

# Artificial Intelligence for Chemical Sciences

Concepts, Models, and Applications



Shrikaant Kulkarni  
Shashikant V. Bhandari  
Dushyant B. Varshney  
P. William  
Editors



**CRC Press**  
Taylor & Francis Group

APPLE ACADEMIC PRESS

# **ARTIFICIAL INTELLIGENCE FOR CHEMICAL SCIENCES**

*Concepts, Models, and Applications*



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

# ARTIFICIAL INTELLIGENCE FOR CHEMICAL SCIENCES

*Concepts, Models, and Applications*

*Edited by*

**Shrikaant Kulkarni, PhD**  
**Shashikant Bhandari, PhD**  
**Dushyant Varshney, PhD**  
**P. William, PhD**



First edition published 2025

**Apple Academic Press Inc.**  
1265 Goldenrod Circle, NE,  
Palm Bay, FL 32905 USA  
760 Laurentian Drive, Unit 19,  
Burlington, ON L7N 0A4, CANADA

**CRC Press**  
2385 NW Executive Center Drive,  
Suite 320, Boca Raton FL 33431  
4 Park Square, Milton Park,  
Abingdon, Oxon, OX14 4RN UK

© 2025 by Apple Academic Press, Inc.

*Apple Academic Press exclusively co-publishes with CRC Press, an imprint of Taylor & Francis Group, LLC*

Reasonable efforts have been made to publish reliable data and information, but the authors, editors, and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors, editors, and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged, please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, access [www.copyright.com](http://www.copyright.com) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. For works that are not available on CCC please contact [mpkbookspermissions@tandf.co.uk](mailto:mpkbookspermissions@tandf.co.uk)

Trademark notice: Product or corporate names may be trademarks or registered trademarks and are used only for identification and explanation without intent to infringe.

---

#### **Library and Archives Canada Cataloguing in Publication**

.....

CIP data on file with Canada Library and Archives

.....

#### **Library of Congress Cataloging-in-Publication Data**

.....

CIP data on file with US Library of Congress

.....

---

ISBN: 978-1-77491-832-6 (hbk)  
ISBN: 978-1-77491-833-3 (pbk)  
ISBN: 978-1-00356-928-2 (ebk)

## About the Editors

---



### **Shrikaant Kulkarni, PhD**

*Adjunct Professor, Faculty of Business,  
Victorian Institute of Technology, Melbourne, Australia;  
Adjunct Professor, Centre of Research Outcome and  
Impact, Chitkara University, Punjab, India*

Shrikaant Kulkarni, PhD, is currently an Adjunct Professor in the Faculty of Business at the Victorian Institute of Technology, Melbourne, Australia, as well as an Adjunct Professor at the Centre of Research Outcome and Impact, Chitkara University, Punjab, India. Dr. Kulkarni has been a senior academician and researcher for over four decades. He has delivered invited lectures and conducted sessions at national and international conferences and faculty development programs. He has guided many major and minor projects in engineering chemistry, green chemistry, nanotechnology, analytical chemistry, catalysis, chemical engineering materials, industrial organization, management, and AI. He has published over 100 research papers in national and international journals and conferences of repute. He has authored more than 50 book chapters. He has edited 25 books published by international publishers. Another ten books are in the offing. He authored four textbooks in engineering chemistry. He has expertise in materials science, green chemistry and engineering, analytical chemistry, green nanoscience, nanotechnology, AI, and computational intelligence. He has expertise in wastewater treatment, green and analytical chemistry, and advanced areas in chemical engineering and materials science. He has published two patents and has applied for two more. Dr. Kulkarni possesses MSc, MPhil, and PhD degrees in Chemistry apart from master's degrees in Economics, Business Management, and Political Science.

**Shashikant Bhandari, PhD**

*Professor, All India Shri Shivaji Memorial Society's  
College of Pharmacy, Pune, India*

Shashikant Bhandari, PhD, is working as a Professor at All India Shri Shivaji Memorial Society's College of Pharmacy in Pune, India. With over 24 years of experience in research and academics, a testament to his commitment to scholarly guidance, Dr. Bhandari has mentored an impressive number of over 50 postgraduate students and four PhD research scholars. His mentorship played a pivotal role in shaping the careers of these young scholars, many of whom have now become accomplished professionals and researchers in their own right. Dr. Bhandari has made a significant impact on the field of pharmaceutical chemistry through his groundbreaking research on various topics including anti-inflammatory, antioxidants, anti-convulsants, anticancer, anti-HIV, antiviral, and antitubercular agents. His dedication to the pursuit of knowledge led to over 25 presentations and the publication of over 30 research papers in prestigious national and international journals, highlighting his profound influence on the scientific community. A benchmark of his scholarly influence, he boasts an impressive H-index of 13, i10 index of 15. His papers have been cited at least 870 times. Additionally, the cumulative impact factor of his publications stands at 70, further underscoring the breadth and depth of his scholarly contributions. A visionary researcher, Dr. Bhandari also filed two patents for his innovative discoveries, cementing his position as a trailblazer in the field of anticancer drug development. These patents have not only showcased his inventive genius but have also contributed to the advancement of technology and society. His tenacious commitment to advancing research has earned him QIP grants and minor research projects worth more than Rs. 50 lakhs funded by AICTE New Delhi and S.P.P. University, Pune, India. Throughout his illustrious career, Dr. Bhandari has received numerous accolades, including the Best Research Publication Award in 2011 by V Life Science Pvt. Ltd., Pune, and the Best Research Mentor International Award in 2019. He has attended the International Conference on QSAR and Systems Biology at Uppsala University, Sweden (2008). In addition, he serves as a reviewer for various international journals with a very good rapport and high impact factor. Dr. Bhandari's remarkable journey inspires us to embrace the pursuit of knowledge, innovate fearlessly, and nurture the potential in others.

**Dushyant Varshney, PhD**

*Chief Technology Officer,  
Arcturus Therapeutics, USA*

Dushyant Varshney, PhD, is the Chief Technology Officer at Arcturus Therapeutics in the USA. He is responsible for global CMC, product development, technical operations and quality, including technology innovation, manufacturing, and technology transfer for clinical and commercial products. Dr. Varshney has over 25 years of experience at global bio-pharma organizations, leading the entire product life-cycles of diverse biotech modalities, including mRNA therapeutics, gene and cell-based therapy, biologics, vaccines, and sterile injectables. He has made significant contributions towards over 30 launches and over 75 INDs, BLA/NDA, PAS submissions, and has ensured the supply of more than 10B+ doses and \$20B+ revenues at companies including Gilead, Pfizer, Novartis, and Sanofi. Before joining Arcturus, Dr. Varshney was the Global Head of Manufacturing, Science and Technology at Gilead-Kite, where he built a strong organization and rapid commercialization strategy for successful technology transfers, launches, and patient access of Yescarta®, Tecartus®, viral vectors, and clinical products. Prior to that, he served as the Vice President, Global Head of Technical Services, Operations and Supply at Jubilant, and as the Head of Manufacturing, Science and Technology at Pfizer, leading America, Europe, and Asia Pacific teams. At Novartis, he contributed to commercial technology transfer, qualification, and manufacturing of Flucelvax® (the first US cell-based influenza vaccine) Trivalent and Quadrivalent, Pre-Pandemic Vaccine stockpiles (bird and swine flu, H5N1, H3N2, H7N9), and pandemic response strategies. He is a Stephen Covey Leadership Coach and a Master Black Belt in Operational Excellence, and he has made over 120 conference presentations, published over 40 articles and book chapters, and edited the book *Lyophilized Biologics and Vaccines*, published by Springer. He was recently recognized by the magazine *Insights Success* as one of the Technophile CTOs of the Year 2022 (Dec). Dr. Varshney received his PhD in Chemistry from the University of Iowa, MPharm from the Institute of Chemical Technology, Mumbai, and BPharm from the University of Pune in India.



**P. William, PhD**

*Dean, Research and Development and Assistant Professor, Department of Information Technology, Sanjivani College of Engineering, SPPU, Pune, India*

P. William, PhD, is working as Dean of Research and Development and Assistant Professor in the Department of Information Technology at Sanjivani College of Engineering, SPPU, Pune. He has published many papers in Scopus-indexed journals and IEEE conferences. His research field includes natural language processing, artificial intelligence, deep learning, machine learning, soft computing, cyber security, and cloud computing. He has been associated with numerous multi-national companies including IBM, TCS, etc., and educational groups. A focused and hardworking professional with experience in consulting in research, innovation, and development, he has served as a session chair and keynote speaker at multiple SCOPUS-indexed international conferences. He has been appointed as a reviewer for reputed SCOPUS/WOS Indexed journals. He is a member of the Institute of Electrical and Electronics Engineers (IEEE), Quality Circle Forum of India (QCFI), and various other professional bodies. He received his Bachelor of Engineering and Master of Technology in Computer Science and Engineering from CSVТУ, Bhilai, India. He completed his PhD in the Department of Computer Science and Engineering from the School of Engineering and Information Technology at MATS University, Raipur, India.

# Contents

---

<i>Contributors</i> .....	<i>xi</i>
---------------------------	-----------

<i>Preface</i> .....	<i>xv</i>
----------------------	-----------

## **PART I: AI in Chemical Sciences for Designing Synthetic Pathways, Tools, and Techniques ..... 1**

<b>1. Applications and Case Studies of AI in Chemical Sciences.....</b>	<b>3</b>
Shrikaant Kulkarni	
<b>2. Computer-Aided Drug Synthesis and Design.....</b>	<b>19</b>
Mubarak H. Shaikh, Sachin P. Kunde, Vijay M. Khedkar, Dattatraya N. Pansare, Aniket P. Sarkate, and Shankar R. Thopate	
<b>3. Computational Tools and Techniques in Planning Organic Synthesis.....</b>	<b>57</b>
Laxmi G. Kathawate, Rohini N. Shelke, Dattatraya N. Pansare, and Aniket P. Sarkate	
<b>4. Patenting Artificial Intelligence-Based Technologies in Chemical and Pharmaceutical Sciences.....</b>	<b>73</b>
Asha Hole, Shashikant Bhandari, Sagar Birajdar, Sandip Surve, and Aniket Sarkate	

## **PART II: Application of Computational Tools, AI, and ML for Predicting Toxicity and Biodegradation ..... 99**

<b>5. Toxicity Predication in Chemistry Based on Machine Learning: A Review .....</b>	<b>101</b>
Dattatraya N. Pansare, Rohini N. Shelke, Aniket P. Sarkate, Anant B. Kanagare, Ajit Dhas, Devidas S. Bhagat, and Bharat K. Dhotre	
<b>6. Machine Learning Algorithms for Prediction of Chemical Toxicity.....</b>	<b>117</b>
D. P. Gaikwad and Shambhavi S. Singh	
<b>7. Artificial Intelligence-Based Prediction of Drug Metabolism.....</b>	<b>141</b>
Shashikant Bhandari, Shital M. Patil, Shivraj N. Mawale, Mrunal C. Belwate, and Somdatta Y. Chaudhari	

<b>8. Exploration of Computational Approaches in Toxicity Prediction .....</b>	<b>179</b>
Prashant R. Murumkar, Rasana Yadav, Rahul Barot, Rutvi Shah, Vijaykumar Srivastava, and M. R. Yadav	
<b>9. Toxicity Forecasts: Navigating Data-Driven AI/ML Models: From Theory to Practice .....</b>	<b>209</b>
B. V. S. Suneel Kumar, Antoine Moitessier, and Nicolas Moitessier	
<b>10. AI-Based Models for Prediction of Biodegradation.....</b>	<b>247</b>
Ganesh B. Patil, Sopan N. Nangare, Shital M. Patil, Shankarsing S. Rajput, and Milind M. Patil	
<b>11. Computer-Based Technologies for Prediction of Biodegradation .....</b>	<b>291</b>
Kumari Neha, Kalicharan Sharma, and Sharad Wakode	
<b>PART III: Application of Expert Systems and AI in Fault Diagnosis and Struc- ture Representation.....</b>	<b>319</b>
<b>12. Exploring the Range of Knowledge-Based Prediction Applications in Chemistry.....</b>	<b>321</b>
Rohini N. Shelke, Laxmi G. Kathawate, Dattatraya N. Pansare, Aniket P. Sarkate, Ajit K. Dhas, Pravin N. Chavan, Shailee V. Tiwari, Deepak K. Lokwani, and Shivraj N. Mawale	
<b>13. Fault Diagnosis of Chemical Process Plant Using Artificial Intelligence.....</b>	<b>337</b>
Tejas Tekawade, R. B. Dhumale, and P. B. Mane	
<b>14. Structure Representation Techniques and Applications in Cheminformatics.....</b>	<b>357</b>
Deep V. Shah and Prashant S. Kharkar	
<b>Index.....</b>	<b>379</b>

# Contributors

---

## **Rahul Barot**

Faculty of Pharmacy, Kalabhavan Campus, The Maharaja Sayajirao University of Baroda, Vadodara, Gujarat, India

## **Mrunal C. Belwate**

Research Scholar, Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India

## **Devidas S. Bhagat**

Department of Forensic Chemistry and Toxicology, Government Institute of Forensic Science, Aurangabad, Maharashtra, India

## **Shashikant Bhandari**

Professor, Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India

## **Sagar Birajdar**

Research Scholar, Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India

## **Somdatta Y. Chaudhari**

Assistant Professor, Department of Pharmaceutical Chemistry, Modern College of Pharmacy, Pune, Maharashtra, India

## **Pravin N. Chavan**

Department of Chemistry, Doshi Vakil Arts College and G. C. U. B. Science & Commerce College, Goregaon, Raigad, Maharashtra, India

## **Ajit K. Dhas**

Department of Chemistry, Deogiri College, Aurangabad, Maharashtra, India

## **Bharat K. Dhotre**

Department of Chemistry, Swami Vivekanand Sr. College Mantha, Jalna, Maharashtra, India

## **R. B. Dhumale**

AISSMS Institute of Information Technology, Pune, Maharashtra, India

## **D. P. Gaikwad**

Department of Computer Engineering, AISSMS College of Engineering, Pune, Maharashtra, India

## **Asha Hole**

PhD Research Scholar, Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India

## **Anant B. Kanagare**

Department of Chemistry, Deogiri College, Aurangabad, Maharashtra, India

## **Laxmi G. Kathawate**

Department of Chemistry, Radhabai Kale Mahila Mahavidyalaya, Maharashtra, India

**Prashant S. Kharkar**

Department of Pharmaceutical Sciences and Technology, Institute of Chemical Technology, Matunga, Mumbai, Maharashtra, India

**Vijay M. Khedkar**

Assistant Professor, Department of Pharmaceutical Chemistry, School of Pharmacy, Vishwakarma University, Pune, Maharashtra, India

**Shrikaant Kulkarni**

Adjunct Professor, Faculty of Business, Victorian Institute of Technology, Melbourne, Australia;  
Adjunct Professor, Centre of Research Outcome and Impact, Chitkara University, Punjab, India

**B. V. S. Suneel Kumar**

Atomica AI Solutions Private Limited, Plot No 35, Beside Avance Phoenix SEZ, Hitech City, Hyderabad, India

**Sachin P. Kunde**

Assistant Professor, Department of Chemistry, RBNB College, Shirampur, Ahmednagar, Maharashtra, India

**Deepak K. Lokwani**

Rajarashi Shahu College of Pharmacy, Buldana, Maharashtra, India

**P. B. Mane**

AISSMS Institute of Information Technology, Pune, Maharashtra, India

**Shivraj N. Mawale**

Research Scholar, Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India

**Antoine Moitessier**

Molecular Forecaster Inc., 910–2075 Robert Bourassa St., Montreal, Quebec, H3A2L1, Canada

**Nicolas Moitessier**

Molecular Forecaster Inc., 910–2075 Robert Bourassa St., Montreal, Quebec, H3A2L1, Canada

**Prashant R. Murumkar**

Faculty of Pharmacy, Kalabhavan Campus, The Maharaja Sayajirao University of Baroda, Vadodara, Gujarat, India

**Sopan N. Nangare**

Research Associate, Department of Pharmaceutical Chemistry, H. R. Patel Institute of Pharmaceutical Education and Research, Shirpur, Dhule, Maharashtra, India

**Kumari Neha**

Department of Pharmaceutical Chemistry, Delhi Institute of Pharmaceutical Sciences and Research, DPSR University, New Delhi, India

**Dattatraya N. Pansare**

Department of Chemistry, Deogiri College, Aurangabad, Maharashtra, India

**Ganesh B. Patil**

Associate Professor, Department of Pharmaceutics, H. R. Patel Institute of Pharmaceutical Education and Research, Shirpur, Dhule, Maharashtra, India

**Milind M. Patil**

Assistant Professor, Department of Chemistry, Poojya Sane Guruji Vidya Prasarak Mandals Shri. S. I. Patil Arts, G. B. Patel Science and S. T. K. V. Sangh Commerce College, Shahada, Nandurbar, Maharashtra, India

**Shital M. Patil**

Assistant Professor, Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India

**Shankarsing S. Rajput**

Professor and Principal, S. P. D. M. Arts, S. B. B. and S. H. D. Commerce & S. M. A. Science College, Shirpur, Dhule, Maharashtra, India

**Aniket P. Sarkate**

Associate Professor, Department of Chemical Technology, Dr. Babasaheb Ambedkar Marathwada University, Aurangabad, Maharashtra, India

**Deep V. Shah**

Indian Institute of Science Education and Research, Pune, Maharashtra, India

**Rutvi Shah**

Faculty of Pharmacy, Kalabhavan Campus, The Maharaja Sayajirao University of Baroda, Vadodara, Gujarat, India

**Mubarak H. Shaikh**

Assistant Professor, Department of Chemistry, Radhabai Kale Mahila Mahavidyalaya, Ahmednagar, Maharashtra, India

**Kalicharan Sharma**

Department of Pharmaceutical Chemistry, Delhi Pharmaceutical Sciences and Research University, New Delhi, India

**Shashikant Bhandari**

Professor, Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India

**Rohini N. Shelke**

Department of Chemistry, Radhabai Kale Mahila Mahavidyalaya, Maharashtra, India

**Shambhavi S. Singh**

Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India

**Vijaykumar Srivastava**

Faculty of Pharmacy, Kalabhavan Campus, The Maharaja Sayajirao University of Baroda, Vadodara, Gujarat, India

**Sandip Surve**

Research Scholar, Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India

**Tejas Tekawade**

AISSMS Institute of Information Technology, Pune, Maharashtra, India

**Shankar R. Thopate**

Principal, Radhabai Kale Mahila Mahavidyalaya, Ahmednagar, Maharashtra, India

**Shailee V. Tiwari**

Department of Pharmaceutical Chemistry, Shri Ramkrishna Paramhans College of Pharmacy, Hasnapur, Parbhani, Maharashtra, India

**Sharad Wakode**

Department of Pharmaceutical Chemistry, Delhi Institute of Pharmaceutical Sciences and Research, DPSR University, New Delhi, India

**M. R. Yadav**

Center of Research for Development, Parul University, Limda, Vadodara, Gujarat, India

**Rasana Yadav**

Faculty of Pharmacy, Kalabhavan Campus, The Maharaja Sayajirao University of Baroda, Vadodara, Gujarat, India

# Preface

---

Computers are thinking and learning these days, and we need to use them to our advantage for achieving scientific endeavors. Research scholars have used computers to facilitate chemical synthesis design as the earliest research initiatives with the help of artificial intelligence (AI), and the outcomes of their research were instrumental in the development of chemical-based software in the later stage.

Scientific tools are getting more and more advanced and sophisticated. Artificial intelligence (AI) has made its way into the laboratory, where it holds a key role in practicing science. Various powerful techniques that mimic human thought and reasoning fall within the purview of artificial intelligence, making it one of the most fascinating and exciting sciences. However, many challenges are in the way as well, such as the complex and intractable nature of problems and the limitations of using conventional methods as solutions. On the other hand, unlike conventional methods, artificial intelligence methods can be more accurate if machine learning algorithms are trained with reliable datasets having broad statistical distribution to solve problems that are otherwise difficult to solve accurately.

Such innovative intelligent techniques have broadened the horizons of power scientists not only in daily routine but more importantly for laying down scientific theories and understanding. This book provides a mathematical and non-mathematical application of Artificial Intelligence in chemical sciences.

Chemists are increasingly using artificial intelligence (AI) for diversified applications viz. molecule design, retrosynthesis, reaction outcome prediction, as well as drug discovery. Historically, the application of AI in chemistry has been primarily focused on accelerating drug discovery and minimizing the enormous production cost and discovery-to-market time frame. AI has made assisted tremendously in accelerating drug discovery and in the field of R&D thus far. However, the use of AI in chemistry is not confined to just drug development but it can also help chemists in pursuing their research expeditiously and creatively.

This book covers a host of aspects like design and identification of the right molecules as precursors, following and predicting kinetics and thermodynamics of reactions, predicting yield or atom economy, enhancing recovery



or process efficiency, optimizing process conditions, identifying right pathways, designing new pathways, etc., providing a conceptual understanding of the subject of chemintelligence.

The topics embodied within the scope of chemintelligence include: use of reasoning, designing pathways or routes, planning chemical synthesis with computers, representation of molecular structures, searching structure, substructure, and superstructure, predicting aromaticity and stereochemistry, toxicity, metabolism, biodegradability, application of knowledge-based expert systems for prediction in chemistry, application of AI in fault detection in process plants, advanced process control, design of catalysts and catalytic reactors, predicting physical properties and hydrodynamics of multi-phase reactors. The use of computational modeling is demonstrated as an advantage in addressing the problems and predictive analysis of toxicity, biodegradability, reaction kinetics, etc. paving the way for furthering research in the area of chemintelligence.

This book is divided into three parts. Part I is devoted to AI in chemical sciences for designing synthetic pathways, tools, and techniques and contains four chapters. Chapter 1 discusses the general applications of AI in chemical sciences and a few representative case studies in brief. Chapter 2 deliberates on the role of computer-aided techniques in designing and synthesizing drugs. Chapter 3 elaborates on the use of computational techniques and tools in planning for organic synthesis while Chapter 4 gives an account of the application of artificial intelligence-based technologies in patent filing in chemical and pharmaceutical sciences.

Part II demonstrates the application of computational tools, AI, and ML for predicting toxicity and biodegradation. It contains seven chapters. Chapter 5 provides an overview of the application of machine learning tools for toxicity prediction. Chapter 6 presents the way machine learning algorithms are used in predicting the toxicity of chemicals. Chapter 7 examines the use of AI in predicting drug metabolism. Chapter 8 elicits the exploration of a range of computational approaches in the prediction of toxicity. Chapter 9 sheds light on the navigation of AI and ML-based models in forecasting toxicity beforehand. Chapter 10 examines the role of AI-based models in predicting biodegradation of materials. Chapter 11 throws light on the application of computational techniques in the prediction of biodegradation.

Part III is dedicated to the application of expert systems and AI in applications such as fault diagnosis, structure representation, and determination of physical properties of materials. It has four topics. Chapter 12 sheds light on the potential knowledge-based expert systems hold in predicting various

properties of materials used in chemistry in particular. Chapter 13 dwells upon the application of AI in diagnosing and analyzing faults observed in chemical plants. Chapter 14 takes an overview of various structure representation techniques and their applications in the cheminformatics domain while the concluding Chapter 15 unravels the adoption of artificial intelligence and machine learning in order to predict physical properties of chemistry, particularly in chemistry and drug discovery-related applications.

The book will be valuable to academicians, researchers, and students in enriching their knowledge base and furthering their research initiatives.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## **PART I**

# **AI in Chemical Sciences for Designing Synthetic Pathways, Tools, and Techniques**



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## CHAPTER 1

---

# Applications and Case Studies of AI in Chemical Sciences

SHRIKAANT KULKARNI

*Faculty of Business, Victorian Institute of Technology, Melbourne, Australia; Adjunct Professor, Centre of Research Outcome and Impact, Chitkara University, Punjab, India*

---

### ABSTRACT

A major constraint of conventional approaches is finding solutions to chemical engineering problems because of the nonlinearity and complexity of chemical processes. Artificial intelligence (AI) tools, however, have made the task quite easy and meaningful in application, design, generalization, robustness, and dynamism. AI covers a host of branches, such as ANN, FL, GA, expert systems, and hybrid systems, etc. They find extensive use in a myriad of applications in the chemical engineering domain, namely, modeling, controlling processes, classifying, detecting faults, and diagnosing, etc. A review of the prowess of AI and its future prospects is taken on numerous chemical engineering fronts in the chapter.

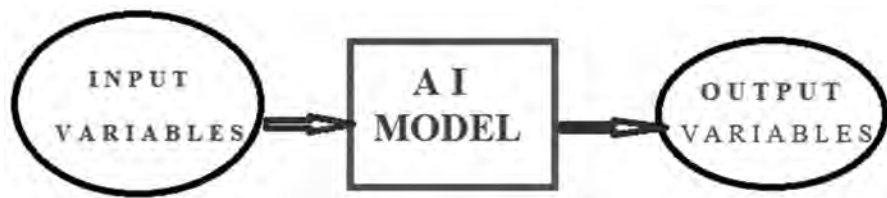
### 1.1 INTRODUCTION

Artificial intelligence (AI) has been finding widespread applications across chemical sciences and engineering and are gaining momentum over the time such as process modeling, optimization, control, fault detection and diagnosis, etc. AI strategy emphasizes upon artificial neural network (ANN), and fuzzy logic (FL) techniques for furthering its capabilities. However, AI

has its own limitations. Conventional approaches are complex and nonlinear in chemical reactions that can be dealt with the use of AI.

## 1.2 AI FOR PROCESS MODELING

Models for chemical processes refer to the system behavior and find use in various applications in chemical sciences, right from researching, designing to optimizing and controlling the plant operations [1]. A simplified AI model is as shown in Figure 1.1.



**FIGURE 1.1** Simplified AI model.

In general modeling approaches are of two types in chemical sciences, i.e., mechanistic (white box), ANN and FL. Former technique involves fundamentals like law of conservations lays the foundation for the model. It involves use of algebraic and differential equations for balancing mass, energy, and momentum. However, variables governing process behavior are in the form of complex mathematics-based equations for chemical reactions characterized by nonlinearity and complexity. Therefore, it is difficult or impossible at times to model processes by mechanistic approaches. Although a model of this kind is developed, it wouldn't be practical to solve or identify the process conditioners. Further, such model requires thorough knowledge and tremendous skills and brilliance to introduce fundamental concepts involved in the processes. Poor knowledge results into poor or weak model yielding poor results [2]. In some cases, assumptions like degeneracy in physical behavior, ideal behavior of gases, etc., demand normalization of the nonlinearity in the equations of the model, that constrains the model influencing the robustness of the model [3].

AI-based models proved their capability and dragged attention in modeling of chemical processes. Such approaches won't ask for detailed knowledge of the process and therefore get over the limitations of the mechanistic methods while approaching complexity and nonlinearity in systems.

AI-based approaches are advantageous to chemical reactions with inherent variables such as inactivity of catalyst in reactors which otherwise is not possible to be addressed by mechanistic models.

AI models chemical science applications makes use of ANN and FL quite commonly, integrated with evolutionary algorithms (EL) [4–7]. Moreover, apart from ANN and FL methods, their combination called adaptive-network-enabled fuzzy inference system (ANFIS) has also been employed for modeling in chemical sciences.

Development of an AI-based model involves modeling of the system by:

- Defining the input/output variables;
- Using experimental findings or the knowledge of the system;
- The conditioners that govern the AI-enabled model such as the fuzzy sets when FL is used);
- Using transfer function consisting of invisible layers on using ANN;
- Using different variables influencing the system.

ANN architectures like multi-layer perceptron (MLP), neural network (NN) involve a feed-forward mechanism are quite helpful in modeling [8]. Recurrent ANN model maps previous inputs and outputs aiming at future predictions finding use in dynamic processes. Fuzzy model (FM) approach is of two types, i.e., Mamdani [9] and Takagi-Sugeno (TS) [10]. In FM uncertainties and complexities of all kinds are translated in “If-Then” expressions using FL theory [11]. ANN is normally looked upon as a data-based AI-enabled models [12]. Mamdani Fuzzy varies from TS approach and the former is superior to latter one in terms of information and rules presented particularly for chemical processes. To begin with, qualitative aspects and knowledge about the system are incorporated in the model development [12]. Moreover, for the Mamdani fuzzy model to develop, data is not required. Therefore, a Mamdani fuzzy model encompasses intuitivism, transparency, and interpretability [13]. While, each TS-type model approximates locally and predicts only under the conditions governing the process [14]. Hence, it is not applicable to analyze the process behavior, not scalable and therefore limiting its use in industrial practices. Although Mamdani model possesses capabilities still it suffers from many rules in dealing with the processes governed by numerous variables.

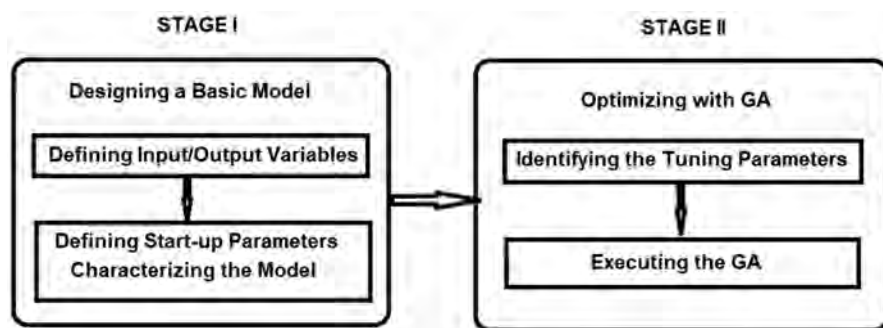
Genetic algorithm (GA) is used for optimizing the efficiency of a FL-based model. GA helps in estimating optimized parameters like scaling [15, 16] or the membership functions [17, 18]. GA is applicable in reducing/selecting rules that show redundancy, unwanted or ill-directed [17] while tackling



multi-dimensional problems with too many rules which are difficult to manage. Designing Mamdani fuzzy and GA combined models involve following stages:

1. Developing preliminary form of the model using knowledge based on heuristics;
2. Tuning protocol with GA.

The schematic of algorithm of the coupled model is as shown in Figure 1.2.



**FIGURE 1.2** Hybrid Mamdani fuzzy and GA model.

Development of model begins with defining the output variables characterizing system behavior provided input variables influencing identified output ones are identified too. Next, development of a primitive fuzzy model with fuzzy sets exhibiting the system behavior aligned with the expertise of the system experts. The model so developed finds use on tuning it up depending upon the requirements. GA is then formulated for optimizing the parameters, like membership function, types, etc.

### 1.3 AI FOR CHEMICAL PROCESS OPTIMIZATION

Optimization of chemical processes is originated from linear programming in 1960s [19]. Linear programming problem aims at evolving at the solution which is best so as to optimize a given objective function. Normally, objective function may be aimed at minimizing cost and by-products or to maximize energy efficiency, yield, profit margins, safety, and reliability of plants. Majority of chemical reactions are characterized by nonlinearity and complexity, and gradient-driven optimization provides numerous solutions to such problems. Evolutionary algorithms (EAs) [20], harmony search [21], particle swarm optimization [22], etc., are some of the useful optimization tools.

AI-based method like generic population-driven metaheuristic optimization algorithm provides an optimal solution to chemical processes.

## 1.4 ARTIFICIAL NEURAL NETWORKS (ANN) FOR CONTROLLING PROCESSES

Chemical process control is brought about for enhancing the process efficiency, lowering down energy utilization and achieve better safety and ecological impact. The conventional control strategies haven't shown any encouraging results in a host of industrially important chemical changes having more nonlinear dynamism and uncertainties, whereas, AI-based tools are better placed in controlling many number processes with complex and nonlinear dynamism [23]. Since they hold lot much of potential in dealing with nonlinear dynamics and self-learning capabilities, there is tremendous excitement and interest in the application of ANN for controlling in various areas like thermal changes [24], reaction kinetics [25], isolation, and refinement [26, 27]. Inverse model control is an algorithm making use of NN for controlling processes.

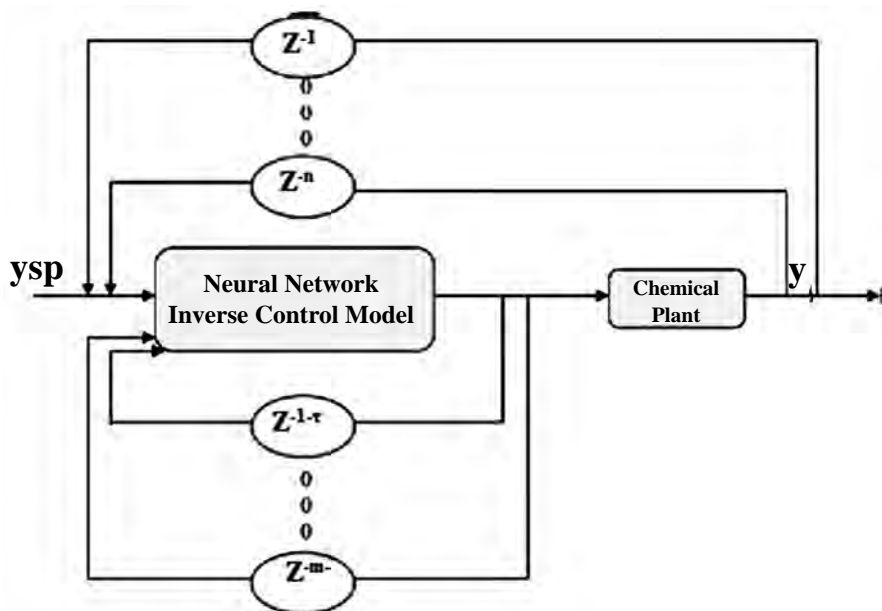
ANN approach assumes that the input for NN is the desired reference output together with the previous input/output variables; NN is responsible for ensuring improved performance of the controlled variables in presence of unknown perturbations. The manipulated variable gives the output of the NN which controls the plant [23]. For a given system possessing time lag ( $\tau$ ),  $n$  and  $m$ , are orders of output and input, the inverse model is then represented as in Eqn. (1):

$$M(t) = \varnothing(ysp, y(k-1), \dots, y(k-n), M(k-1-\tau), \dots, M(k-m-\tau)) \quad (1)$$

where; ' $\varnothing$ ' signifies function; 't' is discrete time; 'M,' 'y' and 'ysp' signify output and set point controllers of plant.

Model predicts the control behavior, on acquiring present and earlier figures of the state variables and the previous control actions is as shown in Figure 1.3.

FL finds use in chemical process control [28–30]. Investigators use FL controller for optimized control of a reaction liberating energy [31], e.g., Polymerization in a batch process [32] and other reactions [33]. However, since there is a time lag in numerous industrially useful chemical reactions, fuzzy model predictive control (FMPC) is recommended [34, 35]. Systems characterized by uncertainties, opt for type-1 won't provide solution to a control-based problem [36] which demands type-2 FL for chemical process control [38].



**FIGURE 1.3** Neural network inverse control model.

Hybrid controller model couples two or more AI tools so as to better regulate the process. Efficiency. Adaptive neuro-fuzzy inference system (ANFIS) model is a very well-known approach. Figure 1.4 is a ANFIS architecture contains five-layered feed-forward NN. It is a hybrid and smart system learns by itself with the help of NN through knowledge representation of the FL [39]. ANFIS system provides for an architecture containing five layers of feed-forward NN as below:

1. **1st Layer:** It is an input layer. Each neuron represents the parameter of membership function. Inputs are converted into degree values varying from 0 and 1.
2. **2nd Layer:** Each neuron does a connecting operation (e.g., “AND”) for computing firing strength of a rule.
3. **3rd Layer:** Neurons help in normalizing.
4. **4th Layer:** A product of the firing strength on normalization with the input’s combination (e.g., TS rule).
5. **5th Layer:** It presents the weighted average of outputs from 4<sup>th</sup> layer. ANFIS controls reactions in chemical plants, e.g., distillation columns [40], biodiesel reactors [41].

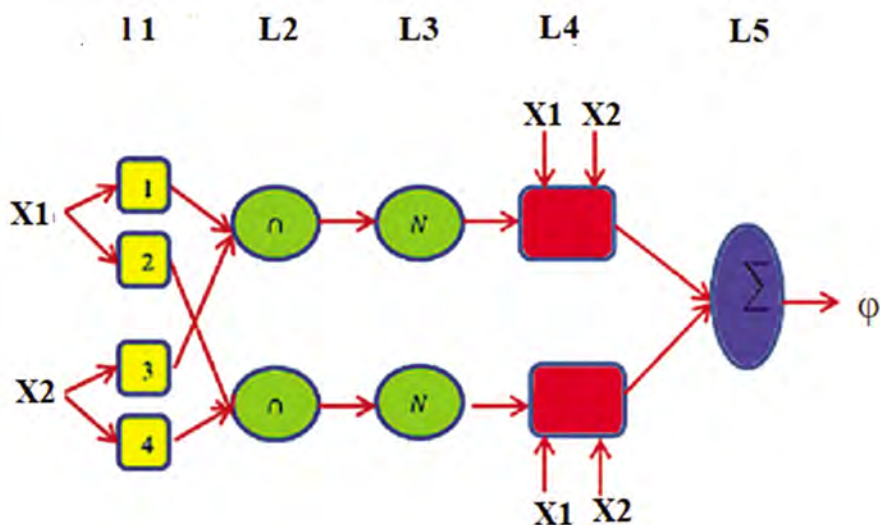


FIGURE 1.4 ANFIS architecture for model.

## 1.5 AI FOR FAULTS INSPECTION

A fault refers to a deviation from observable variables that are acceptable or calculated parameters. A failure means variations of malfunctions in the plant which are attributed to instrumental errors, disruptions, and parameter uncertainties in plants. The abnormal conditions lead to financial losses. Hence, fault detection and diagnosis have been at center stage of chemical processes which asked for commensurate strategies. The fault diagnostic systems ought to detect expeditiously, isolable, robust, and numerous fault identifiable [42].

Neural network systems (NNSs) are smart and powerful in fault diagnostic techniques due to their high promise in following nonlinear dynamics [43–47]. NNSs provide for the neuron numbers in the input/output layers resembling with number of variables estimated and faults likely to occur. Outputs are binary variables corresponding to faults (value = 1) or absence of fault (value = 0) [47]. Another AI-based technique is FL, which is employed in chemical processes for the detection of faults [48–51]. FL provides linguistic expressions that relate one fault to many symptoms.

NN although is a powerful technique in fault finding because of its capability in following the dynamism in nonlinearity without any heuristic knowledge. It demands huge data pertaining to many operating conditions wherein the impacts of numerous faults exist. While, the fuzzy diagnostic

tool represents heuristic reasoning between symptoms and respective faults like linguistic rules and doesn't ask for any data pertaining to system background [52, 53]. The fuzzy diagnostic system is dealing with heuristic and reasoning-enabled rules are difficult and time taking for integrated processes across plants [54, 55]. Hence, neuro-fuzzy diagnostic tools are necessary in processes as recommended in the previous study [58].

Following are case studies of AI techniques employed to model, optimize, control process, detect, and diagnose fault of chemical processes.

## 1.6 VIRUS REMOVAL PREDICTION

Hybrid Mamdani fuzzy coupled with GA is used for prediction of removal of virus from water by employing microfiltration. Application of membranes is one of the separation techniques to remove virus for reusing municipal waste water. Traditional modeling is used for predicting membrane performance is constrained by limitations like the lack of predictive abilities in case of fouling or complexities of property profile of the membrane surface and forces. Mamdani fuzzy model is optimized to predict water-borne virus removal [2]. GA helps in optimizing factors governing the membership function of model variables. It involves defining input/output variables as an initial step in model development. Degree of virus removed is measured as  $R$  and is calculated as follows in Eqn. (2):

$$\% R = 100 (1 - C_p / C_f) \quad (2)$$

where;  $C_p$  and  $C_f$  represent virus concentrations in discharge and feed, respectively.

Concentration of virus FMD, and IBR, pressure ( $P$ ), volume ( $V$ ) and rpm (agitation velocity) are input variables. Data obtained on experimentation is attributed to the research work undertaken by Madaeni & Kurdian [2]. Variables are discrete of Gaussian-type membership function. Fuzzy inference system is set up using start-up fuzzy sets and fitness function is defined. Two factors of Gaussian membership functions covering  $x^-$  and  $\sigma$  are derived from GA as shown in Eqn. (3).

$$f(x) = \exp(-(x - x^-/\sigma)^2) \quad (3)$$

The mean square error (MSE) is fitness function, expressed in Eqn. (4):

$$MSE = (y_m - y_e)^2 \quad (4)$$

where;  $y_m$  and  $y_e$  represent vectors of fuzzy model and data set, respectively.

Genetic modification of protocols are of two types. There are some rules in these methods. Membership function parameters are decision variables for input and output variables. Thus, a variable in every rule may have different shape to membership function on optimization. Predictive ability of the model can be enhanced at the expense of lowering down the comprehension of the model. Other method is put into use when rules are many in number. In such case, each variable across rules enjoys similar shape to the membership function on optimization. There are some decision variables in this second method as against the earlier one.

All possible input variables on combination help in defining 10 rules in this model, and because of few rules, the first method preferred for optimization of parameters. Fuzzy model is designed as it is optimizable with parameters. Such model relies upon qualitative rules, without taking cognizance of the complexity and limitations of the white-box model. Fuzzy models are about 90% accurate when compared with experimental data [2].

## 1.7 OXIDATIVE COUPLING OF CH<sub>4</sub> (OCM) WITH GA

C<sub>2</sub> (ethane + ethylene) productivity is optimized in OCM by passing methane over Mn/Na<sub>2</sub>WO<sub>4</sub>/SiO<sub>2</sub> catalyst placed in a fluidized bed reactor [20]. OCM takes place in a set of chemical processes as proposed by Keller & Bhasinin in 1980 [59]. Natural gas is transformed into the products like ethylene and ethane in preference. The bottleneck in this process is to commercialize it because of poor efficiency. Many suggestions are made for improving the yield of C<sub>2</sub> [60, 61]. Step-wise feeding is one of the solutions used in improving the C<sub>2</sub> yield along the reactor.

It is assumed in this case that the injected gas carries O<sub>2</sub> in each step while CH<sub>4</sub> is fed at top of bed in reactor. C<sub>2</sub> yield is enhanced by optimizing the process variables. Daneshpayeh et al. [63] developed kinetic model is employed as sub-model of a reaction. Reactor model is developed first following which it is solved [62]. Thereafter, with the help of GA, C<sub>2</sub> productivity is maximized for three O<sub>2</sub> injections. Fitness function corresponds to C<sub>2</sub> yield which is expressed as in Eqn. (5):

$$Y_{C_2} = 2 \times N_{C_2} / N_{CH_4} \times 100 \quad (5)$$

C<sub>2</sub> yield is measured after giving three injections of O<sub>2</sub> to bed in reactor. Decision variables are optimized for getting best possible results. The highest C<sub>2</sub> yield of the order 22.9% is obtained at 746.06°C. C<sub>2</sub> output is obtained on optimization using an AI-enabled model larger by 4% [61].

## **1.8 GENETIC-ANFIS CONTROLLER FOR BIODIESEL REACTOR**

Microchem reactor adopting microwave process technology is employed for production of eco-friendly product, bio-diesel. The reactor temperature is controlled to produce  $C_2$  yield to the maximum and to stem down the production of undesirable by-products. Given this purpose, Wali et al. introduced an AI-enabled controller by making use of genetic-ANFIS temperature control for biodiesel production using microwaves [41]. Microwave power supply is decision, reactor temperature is control and feed-flow rate is disruption variable. An online genetic-ANFIS controller is tested under varying operating conditions like set-point tracking and rejection. Controller monitoring needs effective maintenance of reactor temperature rapidly against adaptive control with no oscillations [41].

## **1.9 MODEL FOR DETECTION AND DIAGNOSIS OF CONCURRENT FAULTS**

Plant-wide systems are characterized by complexity and similar symptoms and therefore traditional neural networks systems depending upon steady-state characteristics-enabled data are unable to diagnose various concurrent faults. Tayyebi et al. recommended a novel neuromorphic diagnostic system relying upon augmented input possessing steady-state data in addition to newly defined dynamism for getting over the limitations of conventional systems [47]. Here, input vector of neural network is augmented diagnostic tool such that varying faults produce distinct symptoms. Hence, process track record and steady state are used to derive characteristic symptoms. Therefore, characteristic points in the dynamism of every variable are measured to differentiate and detect different faults in a unique way. Tennessee Eastman process (TE) which embodies numerous measurements and modified variables and similar kind of faults were used plant-wide as benchmarks. The efficiency of neuromorphic diagnostic model using augmented inputs, has been differentiated against the traditional neuromorphic diagnostic tool of which steady-state characteristic data are inputs. The recommended model, outperformed traditional neuromorphic diagnostic model for detecting numerous concurrent faults. Moreover, the prescribed model can appropriately detect different permutations of six concurrent faults. This comparative advantage of the recommended model is its capability to do fault diagnosis when many simultaneously occurring faults with semblance in their symptom surface [47].

## **1.10 CHALLENGES**

AI too comes across many challenges, like accountability, security, mistrust in technology, and movement of laborers similar to other technologies. These challenges need to be addressed so as to shape future of AI technology. It has to be ensured that the impact of AI is a positive one by dealing with the challenges in a proactive manner, while ensuring the opportunities keep coming. AI are computer systems and don't share human values. AIs will never showcase human traits as long as we won't program them to do so. Similarly, lawmakers have to be cautious by not making stringent rules as they may hamper growth of AI [64].

## **1.11 CONCLUSION**

AI-enabled tools deal with complex problems quite effectively with greater accuracy than conventional tools. AI tools can embrace challenges and can be applied in a host of meaningful applications in the field of chemical Sciences. Four representative illustrations are discussed from the field of reaction modeling, optimization, process control, fault detection and diagnosis. AI techniques, tools, and technologies have been offering solutions to address complex nonlinear problems in a broad spectrum of frontier areas in chemical Engineering. AI-based models using algorithms such as FL, GA, EA, etc., will play a decisive role in nurturing the applications world-wide to bring about a substantive change in the quality of life of mankind.

## **KEYWORDS**

- **algorithms**
- **artificial intelligence**
- **broad spectrum**
- **conventional tools**
- **fault detection**
- **nonlinear problems**
- **process optimization and control**



## REFERENCES

1. Luyben, W. (1996). Process modeling, simulation, and control for chemical engineers. *Petroleum Refinery Engineering*, 2, 289–290.
2. Madaeni, S. S., & Kurdian, A. R. (2011). Fuzzy modeling and hybrid genetic algorithm optimization of virus removal from water using microfiltration membrane. *Chemical Engineering Research and Design*, 89, 456–470.
3. Araromi, D. O., Sonibare, J. A., & Emuoyibofarhe, J. O. (2014). Fuzzy identification of reactive distillation for acetic acid recovery from waste water. *Journal of Environmental Chemical Engineering*, 2, 1394–1403.
4. Hajjar, Z., Kazemeini, M., Rashidi, A., & Tayyebi, S. (2016). Artificial intelligence techniques for modeling and optimization of the HDS process over a new graphene-based catalyst. *Phosphorus, Sulfur, and Silicon and the Related Elements*, 191, 1256–1261.
5. Soltanali, S., Halladj, R., Tayyebi, S., & Rashidi, A. (2014). Neural network and genetic algorithm for modeling and optimization of effective parameters on synthesized ZSM-5 particle size. *Materials Letters*, 136, 138–140.
6. Soltanali, S., Halladj, R., Tayyebi, S., & Rashidi, A. (2015). Application of genetic-fuzzy approach for estimation of nano ZSM-5 crystallinity. *Materials Letters*, 150, 39–43.
7. Hajjar, Z., Khodadadi, A., Mortazavi, Y., Tayyebi, S., & Soltanali, S. (2016). Artificial intelligence modeling of DME conversion to gasoline and light olefins over modified nano ZSM-5 catalysts. *Fuel*, 179, 79–86.
8. Gibbs, M. S., Morgan, N., Maier, H. R., Dandy, G. C., Nixon, J. B., & Holmes, M. (2006). Investigation into the relationship between chlorine decay and water distribution parameters using data driven methods. *Mathematical and Computer Modeling*, 44, 485–498.
9. Mamdani, E., & Assilian, S. (1975). An experiment in linguistic synthesis with a fuzzy logic controller. *International Journal of Man-Machine Studies*, 7, 1–13.
10. Takagi, T., & Sugeno, M. (1985). Fuzzy identification of systems and its applications to modeling and control. *IEEE Transactions on Systems, Man, and Cybernetics*, 15, 116–132.
11. Zadeh, L. A. (1978). Fuzzy sets as a basis for a theory of possibility. *Fuzzy Sets and Systems*, 1, 3–28.
12. Adoko, A. C., Gokceoglu, C., Wu, L., & Zuo, Q. J. (2013). Knowledge-based and data-driven fuzzy modeling for rockburst prediction. *International Journal of Rock Mechanics and Mining Sciences*, 61, 86–95.
13. Sala, A., Guerra, T. M., & Babuška, R. (2005). Perspectives of fuzzy systems and control. *Fuzzy Sets and Systems*, 156, 432–444.
14. Habbi, H., Zelmatt, M., & Bouamama, B. O. (2003). A dynamic fuzzy model for a drum-boiler-turbine system. *Automatica*, 39, 1213–1219.
15. Gudwin, R., Gomide, F., & Pedrycz, W. (1998). Context adaptation in fuzzy processing and genetic algorithms. *International Journal of Intelligence Systems*, 13, 929–948.
16. Cordon, O., Herrera, F., del Jesus, M. J., & Villar, P. (2001). A multi-objective genetic algorithm for feature selection and granularity learning in fuzzy-rule based classification systems. In *IFSA World Congress. 20th NAFIPS International Conference. 2001 Joint 9th*, 3, 1253–1258.
17. Cordon, O., del Jesus, M. J., & Herrera, F. (1998). Genetic learning of fuzzy rule-based classification systems cooperating with fuzzy reasoning methods. *International Journal of Intelligence Systems*, 3, 1025–1053.

18. Pulkkinen, P., & Koivisto, H. (2010). A dynamically constrained Mult objective genetic fuzzy system for regression problems. *IEEE Transactions on Fuzzy Systems*, 18, 161–177.
19. Chaves, I., López, J., Zapata, J., & Robayo, A. (2016). *Process Analysis and Simulation in Chemical Engineering*.
20. Eghbal-Ahmadi, M. H., Zerpour, M., Daneshpayeh, M., & Mostoufi, N. (2012). Optimization of fluidized bed reactor of oxidative coupling of methane. *International Journal of Chemical Reactor Engineering*, 10, 1–21.
21. Yousefi, M., Enayatifar, R., Darus, A. N., & Abdullah, A. H. (2013). Optimization of plate-fin heat exchangers by an improved harmony search algorithm. *Applied Thermal Engineering*, 50, 877–885.
22. Mahmood, H. A., Adam, N., Sahari, B. B., & Masuri, S. U. (2017). Development of a particle swarm optimisation model for estimating the homogeneity of a mixture inside a newly designed CNG-H<sub>2</sub>-AIR mixer for a dual fuel engine: An experimental and theoretic study. *Fuel*, 218, 131–150.
23. Tayyebi, S., & Alishiri, M. (2014). The control of MSF desalination plants based on inverse model control by neural network. *Desalination*, 333, 92–100.
24. Nowak, G., & Rusin, A. (2016). Using the artificial neural network to control the steam turbine heating process. *Applied Thermal Engineering*, 108, 204–210.
25. Li, S., & Li, Y. (2015). Neural network based nonlinear model predictive control for an intensified continuous reactor. *Chemical Engineering and Processing*, 96, 14–27.
26. Fernandez de Canete, J., Gonzalez, S., del Saz-Orozco, P., & Garcia, I. (2010). A harmonic balance approach to robust neural control of MIMO nonlinear processes applied to a distillation column. *Journal of Process Control*, 20, 1270–1277.
27. Damour, C., Benne, M., Grondin-Perez, B., & Chabriat, J. (2010). Nonlinear predictive control based on artificial neural network model for industrial crystallization. *Journal of Food Engineering*, 99, 225–231.
28. Hojjati, H., Sheikhzadeh, M., & Rohani, S. (2007). Control of supersaturation in a semibatch antisolvent crystallization process using a fuzzy logic controller. *Industrial & Engineering Chemistry Research*, 46, 1232–1240.
29. Underwood, C. P. (2015). Fuzzy multivariable control of domestic heat pumps. *Applied Thermal Engineering*, 90, 957–969.
30. Baroud, Z., Benmiloud, M., Benalia, A., & Ocampo-Martinez, C. (2017). Novel hybrid fuzzy-PID control scheme for air supply in PEM fuel-cell-based systems. *International Journal of Hydrogen Energy*, 42, 10435–10447.
31. Karr, C. L., Sharma, S. K., Hatcher, W. J., & Harper, T. R. (1993). Fuzzy control of an exothermic chemical reaction using genetic algorithms. *Engineering Applications of Artificial Intelligence*, 6, 575–582.
32. Etinkaya, S. C., Zeybek, Z., Hapoglu, H., & Alpbaz, M. (2006). Optimal temperature control in a batch polymerization reactor using fuzzy-relational models-dynamics matrix control. *Computers & Chemical Engineering*, 30, 1315–1323.
33. Lima, N. M. N., Linan, L. Z., Filho, R. M., Wolf Maciel, M. R., Embiruc, M., & Grácio, F. (2010). Modeling and predictive control using fuzzy logic: Application for a polymerization system. *AIChE Journal*, 56, 965–978.
34. Chang, X. H., & Yang, G. H. (2011). Fuzzy robust constrained model predictive control for nonlinear systems. *Asian Journal of Control*, 13(6), 947–955.
35. Teng, L., Wang, Y., Cai, W., & Li, H. (2017). Robust model predictive control of discrete nonlinear systems with time delays and disturbances via T–S fuzzy approach. *Journal of Process Control*, 53, 70–79.

36. Miccio, M., & Cosenza, B. (2014). Control of a distillation column by type-2 and type-1 fuzzy logic PID controllers. *Journal of Process Control*, 24, 475–484.
37. Galluzzo, M., & Cosenza, B. (2012). Nonlinear fuzzy control of fed-batch reactor for the penicillin production. *Computers & Chemical Engineering*, 36, 273–281.
38. Galluzzo, M., & Cosenza, B. (2011). Control of a non-isothermal continuous stirred tank reactor by a feedback-feed forward structure using type-2 fuzzy logic controllers. *Information Sciences*, 181, 3535–3550.
39. Perendeci, A., Arslan, S., Celebi, S. S., & Tanyolac, A. (2008). Prediction of effluent quality of an anaerobic treatment plant under unsteady state through ANFIS modeling with on-line input variables. *Chemical Engineering Journal*, 145, 78–85.
40. Fernandez de Canete, J., Garcia-Cerezo, A., Garcia-Moral, I., Del Saz, P., & Ochoa, E. (2013). Object oriented approach applied to ANFIS modeling and control of a distillation column. *Expert Systems with Applications*, 40, 5648–5660.
41. Wali, W. A., Al-Shamma, A. I., Hassan, K. H., & Cullen, J. D. (2012). Online genetic-ANFIS temperature control for advanced microwave biodiesel reactor. *Journal of Process Control*, 22, 1256–1272.
42. Venkatasubramanian, V., Rengaswamy, R., Yin, K., & Kavuri, S. N. (2003). A review of process fault detection and diagnosis part I: Quantitative model-based methods. *Computers & Chemical Engineering*, 27, 293–311.
43. Tayarani-Bathaie, S. S., & Khorasani, K. (2015). Fault detection and isolation of gas turbine engines using a bank of neural networks. *Journal of Process Control*, 36, 22–41.
44. Tan, W. L., Nor, N. M., Abu Bakar, M. Z., Ahmed, Z., & Sata, S. A. (2012). Optimum parameters for fault detection and diagnosis system of batch reaction using multiple neural networks. *Journal of Loss Prevention in the Process Industries*, 25, 138–141.
45. Behbahani, R. M., Jazayeri-Rad, H., & Hajmirzaee, S. (2009). Fault detection and diagnosis in a sour gas absorption column using neural networks. *Chemical Engineering & Technology*, 32(5), 840–845.
46. Zhang, Z., & Zhao, J. (2017). A deep belief network-based fault diagnosis model for complex chemical processes. *Computers & Chemical Engineering*, 107, 395–407.
47. Tayyebi, S., Boozarjomehry, R. B., & Shahrokhi, M. (2013). Neuromorphic multiple-fault diagnosing system based on plant dynamic characteristics. *Industrial & Engineering Chemistry Research*, 52, 12927–12936.
48. Musulin, E., Yélamos, I., & Puigjaner, L. (2006). Integration of principal component analysis and fuzzy logic systems for comprehensive process fault detection and diagnosis. *Industrial & Engineering Chemistry Research*, 45, 1739–1750.
49. Tarifa, E. E., & Scenna, N. J. (1997). Fault diagnosis, direct graphs, and fuzzy logic. *Computers & Chemical Engineering*, 21, S649–S654.
50. Tarifa, E. E., & Scenna, N. J. (2004). Fault diagnosis for MSF dynamic states using a SDG and fuzzy logic. *Desalination*, 166, 93–101.
51. Tayyebi, S., Shahrokhi, M., & Boozarjomehry, R. B. (2010). Fault diagnosis in a yeast fermentation bioreactor by genetic fuzzy system. *Iranian Journal of Chemistry and Chemical Engineering*, 29, 61–72.
52. Hang, J., Zhang, J., & Cheng, M. (2016). Application of multi-class fuzzy support vector machine classifier for fault diagnosis of wind turbine. *Fuzzy Sets and Systems*, 297, 128–140.

53. Jahromi, A. T., Er, M. J., Li, X., & Lim, B. S. (2016). Sequential fuzzy clustering based dynamic fuzzy neural network for fault diagnosis and prognosis. *Neurocomputing*, 196, 31–41.
54. Vachtsevanos, G., Lewis, F., & Roemer, M. (2006). *Intelligent Fault Diagnosis and Prognosis for Engineering Systems*. New York: John Wiley & Sons, Inc.
55. Dou, D., & Zhou, S. (2016). Comparison of four direct classification methods for intelligent fault diagnosis of rotating machinery. *Applied Soft Computing*, 46, 459–468.
56. Lau, C. K., Heng, Y. S., Hussain, M. A., & Mohamad Nor, M. I. (2010). Fault diagnosis of the polypropylene production process (UNIPOL PP) using ANFIS. *ISA Transactions*, 49, 559–566.
57. Bonsignore, L., Davarifar, M., Rabhi, A., Tina, G. M., & Elhajjaji, A. (2014). Neuro-fuzzy fault detection method for photovoltaic systems. *Energy Procedia*, 62, 431–441.
58. Shabanian, M., & Montazeri, M. (2011). A neuro-fuzzy online fault detection and diagnosis algorithm for nonlinear and dynamic systems. *International Journal of Control, Automation and Systems*, 9, 665–670.
59. Keller, G. E., & Bhasin, M. M. (1982). Synthesis of ethylene via oxidative coupling of methane. I. Determination of active catalysts. *Journal of Catalysis*, 73, 9–19.
60. Kao, Y. K., Lei, L., & Lin, Y. S. (1997). A comparative simulation study on oxidative coupling of methane in fixed-bed and membrane reactors. *Industrial & Engineering Chemistry Research*, 36, 3583–3593.
61. Lu, Y., Dixon, A. G., Moser, W. R., Ma, Y. H., & Balachandran, U. (2000). Oxygen-permeable dense membrane reactor for the oxidative coupling of methane. *Journal of Membrane Science*, 170, 27–34.
62. Daneshpayeh, M., Mostoufi, N., Khodadadi, A., Sotudeh-Gharebagh, R., & Mortazavi, Y. (2009). Modeling of stagewise feeding in fluidized bed reactor of oxidative coupling of methane. *Energy & Fuels*, 23, 3745–3752.
63. Daneshpayeh, M., Khodadadi, A., Mostoufi, N., Mortazavi, Y., Sotudeh-Gharebagh, R., & Talebizadeh, A. (2009). Kinetic modeling of oxidative coupling of methane over Mn/N<sub>2</sub>WO<sub>4</sub>/SiO<sub>2</sub> catalyst. *Fuel Processing Technology*, 90, 403–410.
64. Gurkaynak, G., Yilmaz, I., & Haksever, G. (2016). Stifling artificial intelligence: Human perils. *Computer Law & Security Review*, 32, 749–758.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## CHAPTER 2

---

# Computer-Aided Drug Synthesis and Design

MUBARAK H. SHAIKH,<sup>1</sup> SACHIN P. KUNDE,<sup>2</sup> VIJAY M. KHEDKAR,<sup>3</sup>  
DATTATRAYA N. PANSARE,<sup>4</sup> ANIKET P. SARKATE,<sup>5</sup> and  
SHANKAR R. THOPATE<sup>6</sup>

<sup>1</sup>*Department of Chemistry, Radhabai Kale Mahila Mahavidyalaya,  
Ahmednagar, Maharashtra, India*

<sup>2</sup>*Department of Chemistry, RBNB College, Shrirampur, Ahmednagar,  
Maharashtra, India*

<sup>3</sup>*Department of Pharmaceutical Chemistry,  
School of Pharmacy, Vishwakarma University, Pune, Maharashtra, India*

<sup>4</sup>*Department of Chemistry, Deogiri College, Aurangabad, Maharashtra,  
India*

<sup>5</sup>*Department of Chemical Technology, Dr. Babasaheb Ambedkar  
Marathwada University, Aurangabad, Maharashtra, India*

<sup>6</sup>*Radhabai Kale Mahila Mahavidyalaya, Ahmednagar, Maharashtra, India*

---

## ABSTRACT

The advancement in computing power, data aggregation, and algorithm development has greatly accelerated the integration of drug synthesis. This progress has had a substantial impact on the enhancement of therapeutic compound design and synthesis. In recent years, there has been a widespread and rapid adoption of data-driven tools that assist in computer-aided

synthesis, reaction prediction, and retrosynthetic analysis. These tools have the potential to significantly enhance the quality and speed of the drug development and discovery process for designed and synthesized molecules. In this discussion, we will explore the historical background and current state of computer-aided drug development, focusing on two main aspects: computer-aided drug synthesis route design and computer-aided drug design.

## **2.1 INTRODUCTION**

The process of computer-based design involves nurturing requirements, synthesizing, and developing building blocks to create expressive designs that meet specific criteria and goals [1]. The effectiveness of this process depends on the identification and utilization of needs throughout the design phase. Occasionally, the resulting design is innovative and unique, while other times it follows a routine process. The field of study that focuses on developing principles, procedures, and tools for design synthesis supports the development of such solutions. However, information gaps or fixed mindsets can make it challenging to generate original solutions. Computer-based design synthesis can assist in addressing this issue by providing designers with a wider range of possibilities, expanding the scope of frequently reviewed resolutions, and potentially enhancing originality. Additionally, computers can automate repetitive tasks involved in routine design, freeing up time for creative activities and contributing to error reduction.

The search for new breakthroughs in organic synthesis has historically relied on chemical intuition, which is based on experience, expertise, and mechanical understanding. For models that mix human intuition and computers, predicting the outcome of a single chemical reaction is still a substantial difficulty [2]. This is why gathering a substantial amount of empirical data is essential for optimizing organic transformations. Chemical synthesis is a complex and time-consuming field where even expert chemists often struggle to predict whether a given substrate will undergo the expected conversion [3, 4]. While mechanistic understanding enables reasonably accurate quantitative predictions of chemical reactivity, relying solely on chemical intuition is practically difficult due to the intricate relationship between structure and reactivity.

## 2.2 IMPROVEMENT OF COMPUTER-AIDED SYNTHESIS TECHNOLOGY

Claude Shannon in 1948 published a paper that laid the basis intended for the fusion of data science and organic production [5]. The advancement of computers in predicting reactions has been greatly influenced by linear free energy relationships such as the Hammett and Bronsted equation [6]. Vladutz in 1960 documented biochemical processes on computers for future reference, proposing concept of computer-assisted organic combination through a reaction database. Subsequently, computer-controlled robots were employed to automate chemical reactions, enabling the high-throughput synthesis of peptides [7]. A groundbreaking computer-aided synthesis software called LHASA was developed in 1969 by Corey, Wipke, and others [8]. This software incorporated retrosynthetic logic and a specific heuristic approach. In 1970, Gelernter and colleagues devised SYNCHEM, a program which is synthetically designed and that utilized geometric theorem proving logic [9]. Peishoff et al. in the 1980 introduced CAMEO, a tool intended for forward-reaction retrieval then retrosynthetic analysis. Addressing chemical challenges through mathematical approaches, Ugi and collaborators introduced the DU model in 1993, which aimed toward define organic reactions [11]. Afterward, several collaborative efforts led to the development of computer-aided fusion strategy schemes such as EROS and AHMOS. Gasteiger et al. founded WODCA in 1995, basing it on the DU model [12]. With the emergence of tools like ChemDraw [13] and representations like SMILES [14], expressing chemical structural data on computers became simpler. The representation of molecular structures on computers made significant progress in the early 21<sup>st</sup> century. Enhancements were made to physical and chemical descriptors [15], molecular graphs [16], molecular strings [17], and molecular fingerprints [18], all contributing to the advancement of computer-aided fusion schemes. In 2016, the release of AlphaGo noticeable a milestone in artificial intelligence, leading to increased attention towards machine knowledge and deep knowledge [19]. Primary claims of computer-aided fusion tools must be high-fidelity reaction forecasting, automation of chemical reactions, and retrosynthesis of complex compounds. As processing power has significantly improved, automated hardware and software have become more sophisticated, resulting in increasingly optimized chemical synthesis



algorithms. Computer-aided synthesis technologies have made it possible to automatically design, synthesize, test, and analyze medicinal compounds.

## 2.3 COMPUTER-AIDED COMBINATION'S APPLICATIONS

### 2.3.1 RETRO SYNTHESIS ANALYSIS

To shift the focus of chemists from the process of manufacturing to the selection of products, computer assistance was introduced as a response to the formalization of retrosynthesis. The development of this field can be largely attributed to Gasteiger [20]. Over time, computer aided synthesis planning has garnered positive evaluations [21]. The term “retrosynthesis” was initially coined by Corey to define the approach of breaking bonds toward transform target molecules into simpler precursors [8]. Currently, there exist two distinct categories of retrosynthetic analysis systems [22]. One category employs specific heuristics and expert input to propose pathways or significant disconnections for the target compound. The majority of automated retrosynthetic programs heavily rely on encoded reaction patterns or generalized subgraph similarity criteria, which emerged as the first attempt at computer-assisted retrosynthetic planning [23]. When utilizing these template-based methods, whether derived algorithmically from reaction databases [24] or manually programmed, a decision must be made regarding the level of generalization and abstraction [25]. Various strategies have been employed to extract the potentially significant context surrounding the reaction center, including the incorporation of non-structural reactivity descriptors. However, there will always be a trade-off between specificity and coverage. These methods do not scale well for large pattern sets, as applying patterns incurs computational costs due to the subgraph isomorphism problem [24]. One approach involves a method that disassembles target molecules into their constituent parts, provides users with a suggested structure, summarizes the entire synthesis process, and forecasts the reaction circumstances. So, these approaches can be categorized in two categories: rule-based methods and rule-free methods. Rule-based techniques rely on establishing connections between target compounds and established reaction principles, enabling the generation of necessary intermediates or raw materials for synthesis [26]. Synthia, developed by Grzybowski et al., is a widely recognized rule-based tool for retrosynthetic analysis, manually curated by organic chemists [27]. Its library of reaction rules comprises tens of thousands of physically programmed

rules based on organic synthesis knowledge and chemists' expertise [28]. Alternatively, a more efficient strategy involves automatic rule extraction from databases, given the exponential growth of documented reactions. Coley [29] categorized then prioritized the reaction rubrics based on molecular structural resemblance among reactants and products for retrosynthetic analysis. Using a neural network model and the retrosynthetic analysis methods discussed by Law et al. [30]; and Borgevig et al. [31]; Segler and colleagues [32] analyzed the reaction rules. This approach was further enhanced by incorporating Monte Carlo tree search (MCTS) for quicker and additional correct retrosynthetic analysis, as proposed by Segler et al. [33]. To address the challenge of handling rule removal for multiple reactions efficiently, Baldi et al. [34] explored the concept of the reaction as an electron sink and electron source, and they categorized the relationships among reactants based on estimated molecular orbitals. In recent times, the rule-free technique has gained popularity for retrosynthetic analysis due to its absence of reaction rules. Liu and colleagues introduced the first rule-free approach [35], employing a sequence-to-sequence (seq2seq) model and converting the reaction product into the reactant using SMILES conversion. The model effectively intended retrosynthetic paths for 17 conjugates. Subsequently, Lai et al. [36] developed AutoSynRoute, a retrosynthetic route design system for one-step retrosynthetic analysis. They presented an end-to-end model that combines MCTS besides an experiential recording purpose. System effectively generated retrosynthetic pathways for four compounds.

### **2.3.2 DOCKING EXPERIMENT**

To optimize the communication among a receptor and a ligand, ligand (a small molecule) should be securely attached to the receptor's necessary site. Initially, the search focused solely on the figure, which initially accounted for steric considerations. To conduct more comprehensive searches involving electrostatic and Vander-Waals interactions with the receptor, the following steps can be followed: Determine the solvent accessible around the active site using Connelly's method. Roll a domain with dimensions of solvent molecule beside surface. Create a "negative" replica of the receptor by utilizing spheres that complement the receptor's surface. Calculate the distances between the spheres in the negative replica of the receptor. Convert the sphere-to-sphere distances into potential atom-to-atom distances. Compare the potential atom-to-atom spaces through real atom distances of organic conjugates in a database. Choose the ligands that exhibit the highest

degree of overlay for additional investigation. Conformational exploration and energy minimization techniques are employed toward identify the low-energy arrangement of the ligand within the binding location. Compute the ligand-receptor interface dynamisms for the designated ligands [37].

### **2.3.3 DOCKING SOFTWARE**

DOCK [37], AutoDOCK [38], GOLD [39], FlexX [40], GLIDE [41], Accelry's DS-Ligand Fit [42], Fujitsu's Scigress Explorer (formerly CaChe) [43], etc., are a few of the most popular docking software's. Wild shape corresponding methods like DOCK and Eudock, incremental building algorithms like FlexX, Hammerhead, and SLIDE, tabu exploration strategies like PRO\_LEADS besides SFDock, genetic algorithms like GOLD, AutoDock, and Gambler, and Monte Carlo replications are some of the algorithms utilized for docking.

### **2.3.4 DOCKING APPLICATIONS**

Various docking claims play a crucial role in the area of structure-based drug strategy. These applications involve regulating the most energetically favorable structures within the receptor-ligand complex, exploring databases and ranking potential candidates, determining ligand compatibility with different macromolecular receptors, studying the geometric properties of complexes, proposing modifications to enhance the effectiveness or other properties of key molecules, creating libraries for further analysis, and calculating differential binding. Compared to conventional drug screening methods, structure-based drug design offers an improvement by allowing the development of more effective medicines that can interact with target proteins. By understanding the target protein in advance and gaining knowledge about its chemical and molecular structure, researchers at BioCryst are able to meticulously design therapeutic candidates atom by atom. These candidates are carefully placed within the active site of protein to inhibit there organic function [44]. These targeted tactic contrasts with the less specific random screening techniques traditionally used in pharmaceutical development. Notably, recent advancements in structure-based drug strategy need to lead the finding of several HIV protease inhibitors which are currently utilized for HIV treatment. In this approach, the target protein forms complexes with sophisticated lead compounds, which are synthesized and further refined through an iterative process. Lead compounds derived from previous research, including combinatorial library screening, can oblige an initial

point to optimize by means of structure-based drug strategy [45]. There is a wide range of algorithms available for evaluating and explaining interactions between ligands and proteins, and their numbers continue to grow. These algorithms vary in complexity and computing speed, providing numerous approaches to address the challenges of structure-based drug design [46]. Many established methodologies include various algorithms with innovative enhancements. The prediction of ligand orientation and binding affinity is a crucial challenge in identifying and optimizing potential drugs, assuming the receptor structure is known [47]. This process is generally called “molecular docking,” and algorithms dealing with such task require attracted significant consideration [48], highlighting its importance in the drug design process. Thanks to advancements in computer power and algorithm performance, the pharmaceutical industry can now perform docking simulations on thousands of ligands within a reasonable timeframe [48]. Although this field is extensive, we have made an effort to compile and categorize utmost significant docking methods. Examples of these techniques used in molecular dynamics include Monte Carlo techniques, genetic algorithms, fragment-based techniques, point complementarity techniques, distance geometry techniques, tabu searches, and systematic searches [49]. Furthermore, it is commonly employed a computer-assisted technique [50], condensing vast virtual libraries of compounds into a manageable subset that potentially contains molecules with high binding affinities to a target receptor [51].

### **2.3.5 DE NOVO DESIGN METHODS AND DOCKING**

The chapter presents a wide comparison between de novo design techniques and docking algorithms, which is a matter of debate. There is often noteworthy overlay in techniques among these two policies. *De novo* design tools such as MCDNLG (Monte Carlo *de novo* ligand generation) [52], SMOG [53], SPROUT [54], BUILDER [55], CONCEPTS [56], CONCERTS [57], DLD/MCSS [58], GENSTAR [59], Group-Build [60], GROW [61], HOOK [62], LEGEND [63], and LUDI are utilized. Several notable accomplishments achieved in the area of computational drug innovation and design, including:

- Norfloxacin, an antibacterial medication developed by Kyorin Pharmaceutical using QSAR techniques.
- COZAAR, an antihypertensive medication developed by DuPont and Merck, specifically Losartan, an Angiotensin II receptor antagonist. It was designed through molecular modeling based on a principal molecule described in patent literature and QSAR [64].

- The discovery of 6-fluoroquinolones, which have become an important class of antibiotics [65].
- Donepezil, a medication for Alzheimer's disease developed by Eisai. It was created using QSAR, molecular shape analysis, and docking techniques [66]. Donepezil is an acetylcholinesterase inhibitor found in ARICEPT.
- TRUSOPT, a glaucoma medication developed by Merck, containing the carbonic anhydrase inhibitor dorzolamide [67].
- CRIVAN, an AIDS medication developed by Merck, specifically Indinavir, an inhibitor of the HIV-1 protease. It was designed using X-ray crystallography, molecular modeling, and structure-based strategy [68].
- VIRACEPT, an AIDS drug developed by Lilly and Agouron, specifically Nelfinavir, an HIV-1 protease inhibitor created through structure-based strategy [69].
- ZOMIG, a migraine treatment developed by Wellcome and Zeneca, containing Zolmitriptan, a 5HT<sub>1</sub>-agonist. It was established using molecular modeling and pharmacophore expansion [70], among numerous other examples.

### 2.3.6 MOLECULAR DOCKING STUDIES IN WATER MOLECULES

In analysis of the obligatory manner of various competing inhibitors, such as PARP-1, Autodock 3.0 was utilized in conjunction by molecules of water present in the crystal structures of the catalytic domain. The results indicated a strong correlation among computed binding energies and experimental inhibitory actions, whether considering structural water molecules ( $r^2 = 0.87$ ) or not ( $r^2 = 0.84$ ) [71]. Previous docking studies have suggested the inclusion of water molecules in the docking process due to their role as components of hydration shell of polar inhibitors, rather than structural water. Typically, water molecules are not incorporated in docking experiments. Given the divergent opinions within the scientific community regarding the usage of water molecules in docking, further research is warranted to investigate this topic.

## 2.4 MOLECULE DESIGNING BY MOLECULAR DOCKING CONCEPT

2,3-Dihydroquinazolin-4(1H)-one known as DHQ is a heterocycle with nitrogen atoms having a six-membered ring which is fused with a phenyl ring.

DHQ possesses the ability to undergo functionalization at various positions, making it a promising candidate for the development of new pharmacophores. This characteristic allows DHQ to interact with multiple targets, resulting in diverse pharmacological properties and establishing its significance. DHQ has demonstrated anti-bacterial effects [72], anti-fungal properties [73], analgesic and anti-inflammatory actions [74a], antimalarial activity [74b], antiviral properties [74c], antitumor effects [74d], insecticidal properties [74e], angiotensin II receptor antagonism [75a], inhibition of CDK5 [75b], potential non-peptidyl inhibitors of cathepsins B and H [75c], modulation of Transient receptor potential melastatin 2 (TRPM2) [75d], potential inhibition of Coagulation factor Xa (fXa) [75e], anthelmintic activity [76a], and anti-hepatitis-B effects [76b], among others. Considering the biological importance of the DHQ heterocyclic system, we initially conducted molecular docking studies to investigate its antimicrobial activity. These studies provided insights into the thermodynamic connections which rule the binding of molecules to DNA gyrase. The results of the molecular docking and ADME analyzes are presented in Table 2.1.

### 2.4.1 MOLECULAR DOCKING

When enzyme-based assays are not available, the *in silico* method of molecular docking is commonly used to advance insights hooked on the binding affinity and thermodynamic interactions governing the interaction between a molecule and its target receptor. Therefore, in direction to investigate the possible mechanism of action for newly discovered DHQ derivatives, molecular docking was done against microbial DNA gyrase. DNA gyrase, classified as a topoisomerase II, an ATP-dependent enzyme vital for DNA transcription, replication, and chromosome segregation in bacteria. It plays a vital role in maintaining the precise spatial DNA topology, making it a critical target for antimicrobial drugs [77]. To conduct the molecular docking calculations, the protein and ligand structures were improved and then exposed to GLIDE (Grid-Based Ligand Docking with Energetics) module in the Small Drug Discovery Suite [78]. Microbial DNA gyrase (PDB code: 1KZN) crystal structure obtained and improved by the protein preparation wizard before docking process. The 3D structures of the DHQ derivatives (3a–p) were built using the Maestro builder panel. The molecular docking analysis revealed well-clustered binding modes for the DHQ analogues, with their binding affinities (Glide docking scores) showing a correlation with their antimicrobial activities. The derivatives with the highest activity

demonstrated improved binding affinity, while those with weaker activities exhibited lower docking scores (Table 2.1). A detailed discussion of the results is presented for the most active analogue, 3e.

**TABLE 2.1** Glide Docking Scores

Compound	Glide Score	Glide Energy	H-Bonding
3a	−8.188	−42.584	Asp73(2.453), Gly77(2.340), Thr165(2.252)
3b	−8.474	−43.132	Glu50(2.283), Asn46(1.966)
3c	−8.497	−43.128	Asp73(2.721), Gly77(2.391), Thr165(2.416)
3d	−8.221	−42.691	–
3e	−8.529	−48.435	Asp73(2.693), Gly77(2.398), Thr165(2.340)
3f	−8.410	−43.889	–
3g	−8.470	−47.592	Asp73(2.703), Gly77(2.395), Thr165(2.336)
3h	−8.413	−44.605	–
3i	−8.114	−42.67	Arg136(2.719)
3j	−8.132	−41.338	Val71(2.318)
3k	−8.244	−45.323	–
3l	−8.377	−45.762	Arg136(2.318)
3m	−8.481	−48.747	–
3n	−8.115	−42.620	–
3o	−8.129	−42.721	Arg136(2.196)
3p	−8.052	−40.116	Asp73(2.383), Gly77(2.336), Thr165(2.246)

The predicted binding mode of 3e (Figure 2.1) exhibited the molecule might comfortably suitable into the active site of DNA gyrase (docking score: −8.529, binding energy: −48.435) through several bonded and non-bonded interactions. The Glide binding energy values were expressed in kcal/mol. A detailed per-residue communication study exhibited that the major thermodynamic interaction subsidizing to the mechanical interconnecting of 3e (Figure 2.1) is a noteworthy network of favorable van der Waals interactions detected with the DHQ scaffold through Thr165(−3.653), Arg136(−1.677), Val120(−1.539), Pro79(−1.306), Ile78(−3.78), Gly77(−1.4), Arg76(−1.767), Asp73(−1.596) and Ala47(−2.208) residue lining the active site while the 2-chloro phenyl side chain remained understood to be involved in alike relations with Val167(−1.607), Met166(−1.69), Gln72(−1.353), Val71(−1.115), Glu50(−3.52), Asn46(−4.318) and Val43(−1.693) residues. Noteworthy electrostatic connections detected with Gly77(−2.024), Arg76

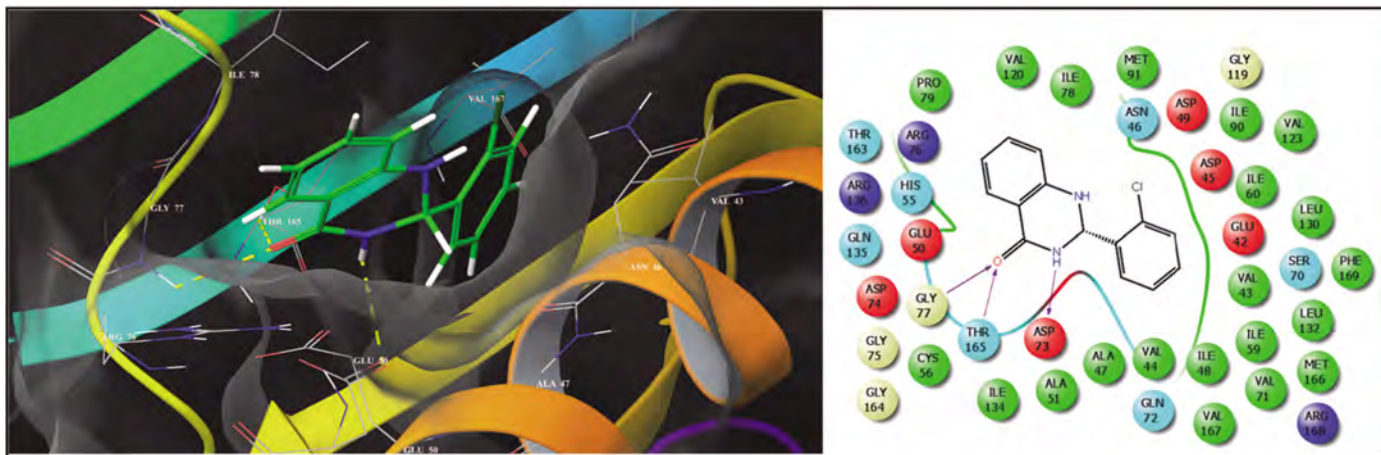
(−1.385), Asp73(−4.568), Asp49(−1.024) and Glu42(−1.096) residues too attributed to higher obligatory attraction detected for 3e (Figure 2.1).

In addition, it was observed that it formed near hydrogen bonding interactions with Asp73 (2.693 Å) through the −NH− group, as well as with Gly77 (2.398 Å) and Thr165 (2.340 Å) through the =O functionalities of the DHQ nucleus. These interactions serve to monitor the orientation of a molecule in the 3D space of the active site and facilitate steric and electrostatic interactions. Furthermore, other derivatives of DHQ within the sequence exhibited various degrees of affinity with the active site residues of Inh A, forming a network of bonded and non-bonded connections. The correlation of docking score and other compounds were revealed in Figures 2.2–2.16.

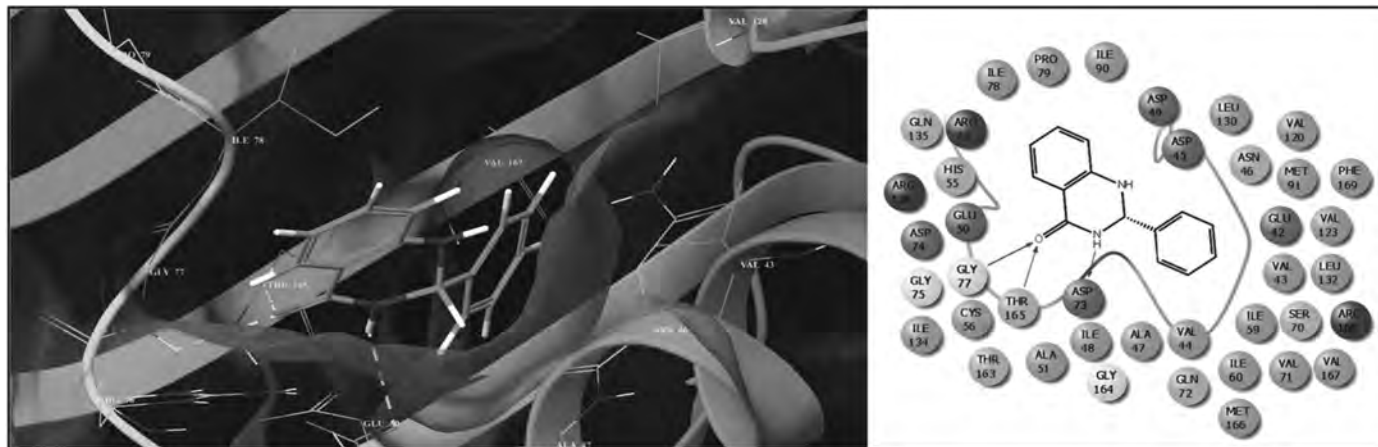
#### **2.4.2 IN SILICO ADME PREDICTION**

To ensure the accomplishment of a medicine, it remains crucial for it to possess a constructive ADME (absorption, distribution, metabolism, and excretion) outline. In order to predict the ADME properties, an extensive computer analysis was conducted on all possible 3a–p synthesized compounds. The outcomes attained are presented in Table 2.2, demonstrating that the conjugates exhibited promising ABS (percentage absorption) values ranging from 79.00% to 94.81% (Table 2.2). Moreover, adherence to Lipinski's rule of five was observed in all synthetic compounds. Overall, the examined core structures satisfied the criteria for an orally active medication, implying their potential for development into oral drugs. In this particular study, the Molinspiration toolkit [79] employed to calculate various parameters, including Lipinski's rule of five [80], molecular weight, molecular volume, hydrogen bond acceptors, logarithm of the partition coefficient, hydrogen bond donors, rotatable bonds, and topological polar surface area. The absorption calculation (% ABS) follows the formula:  $\% \text{ ABS} = 109 - (0.345 \text{TPSA})$  [81]. Additionally, MolSoft software was utilized to determine the drug-likeness model score, which represents the collective attributes to a compound's pharmacokinetics, pharmacodynamics, and physical-chemical properties, expressed as a numerical value [82]. It is essential for a molecule to satisfy the following four requirements without violation: molecular weight should be equal or less than to 500, number of hydrogen bond acceptors should be equal or less than to 10, number of hydrogen bond donors should be equal or less than to 5, and *mi*Log P (octanol-water partition coefficient) should be equal or less than to 5 [83].

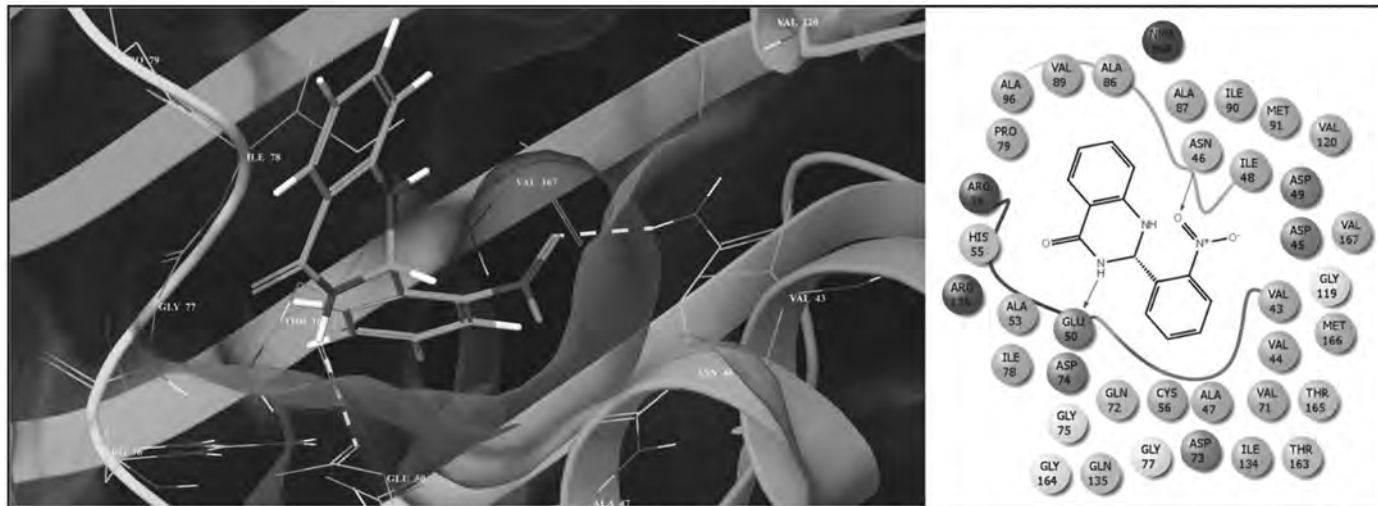




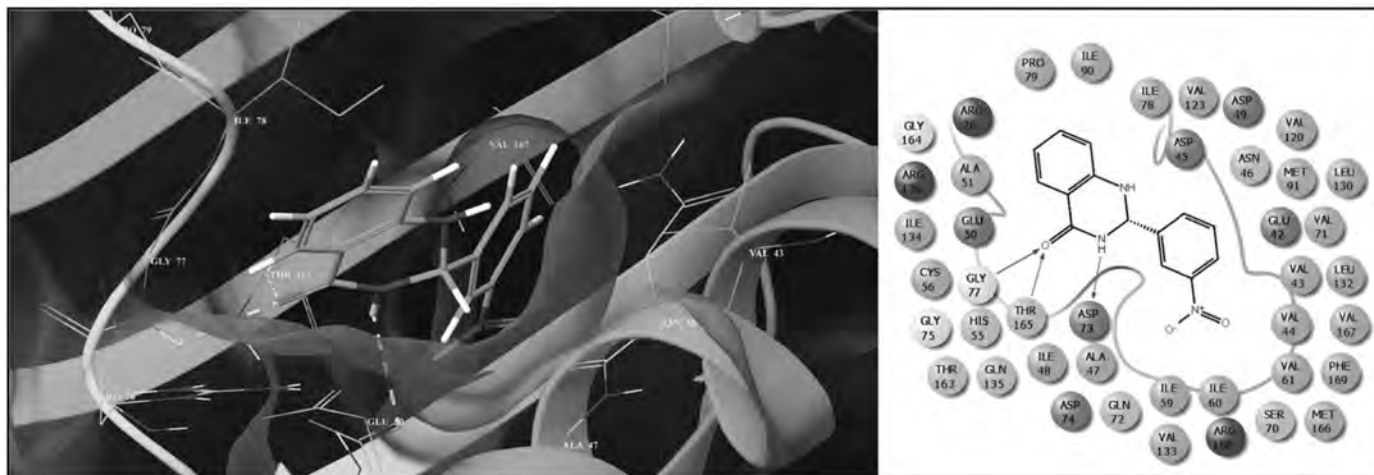
**FIGURE 2.1** Mode of interaction for compound 3e.



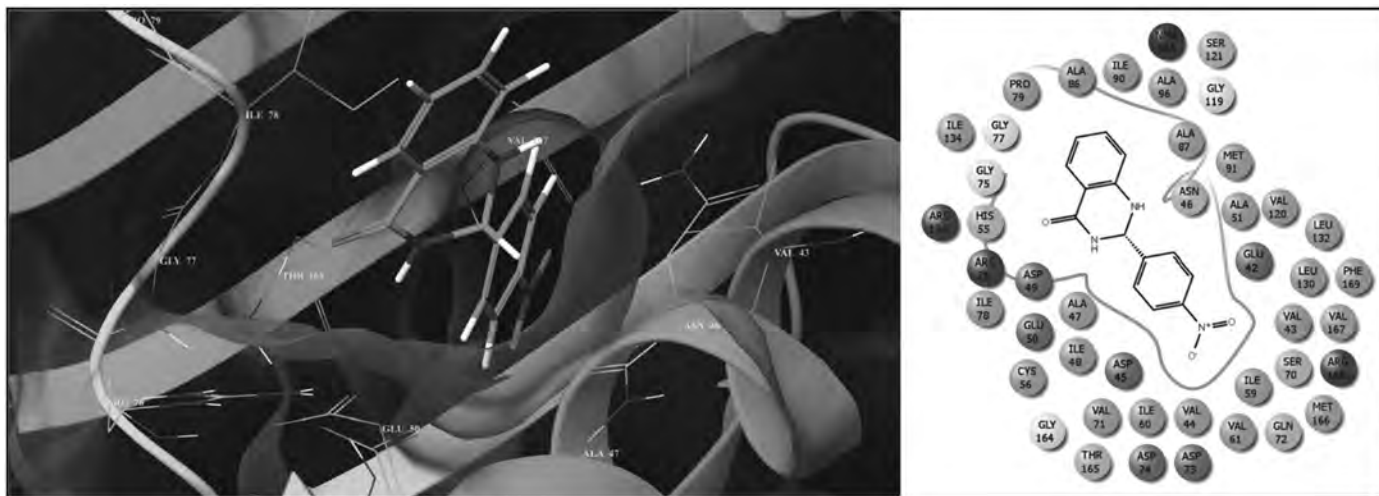
**FIGURE 2.2** Mode of interaction for compound 3a.



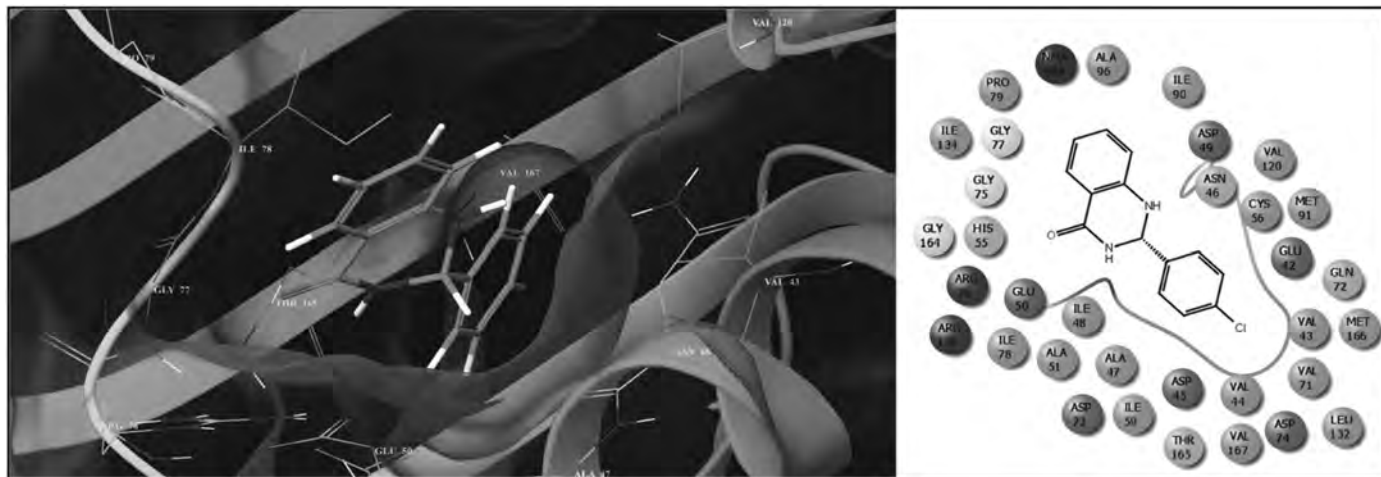
**FIGURE 2.3** Mode of interaction for compound 3b.



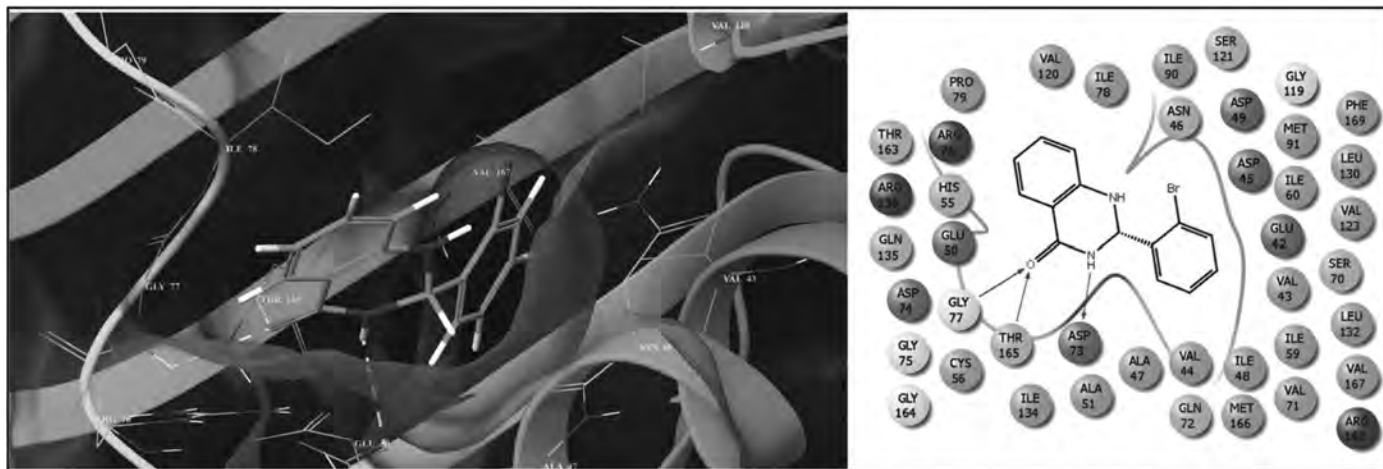
**FIGURE 2.4** Mode of interaction for compound 3c.



**FIGURE 2.5** Mode of interaction for compound 3d.



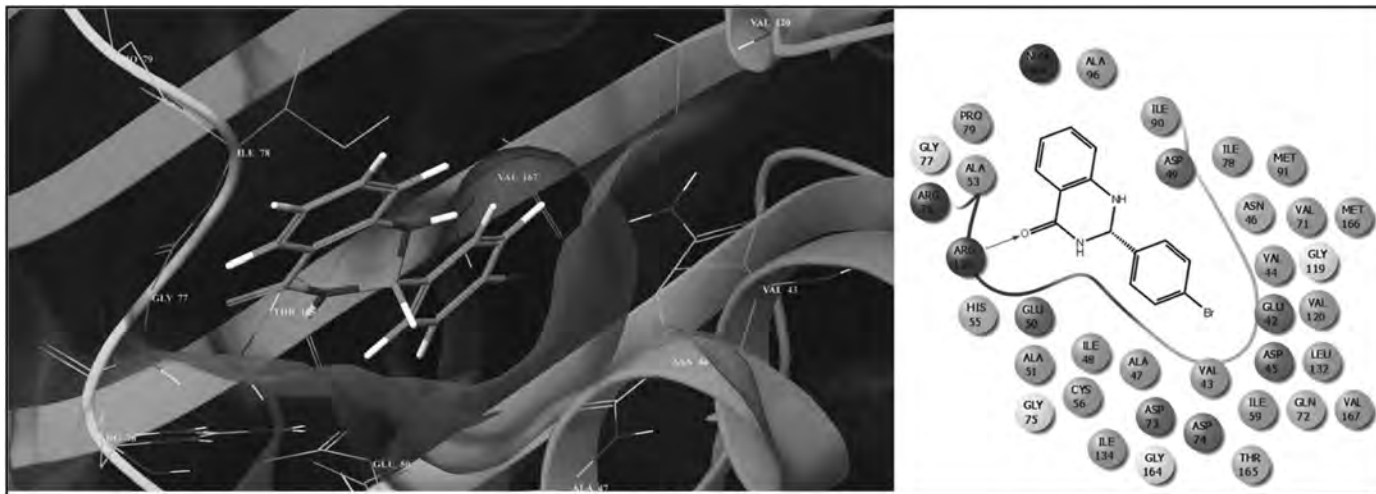
**FIGURE 2.6** Mode of interaction for compound 3f.



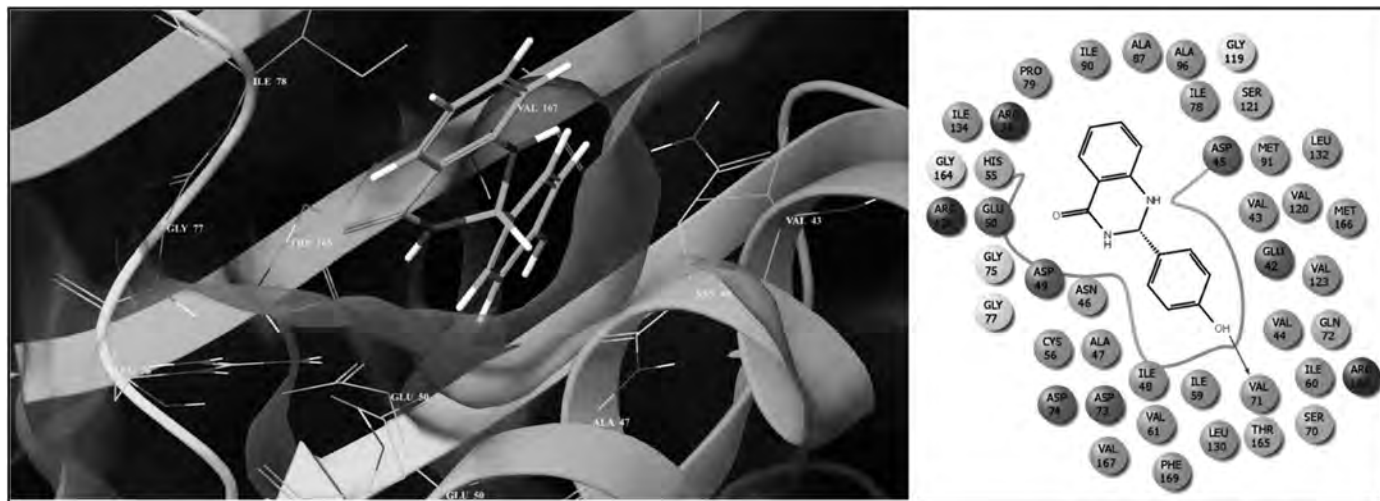
**FIGURE 2.7** Mode of interaction for compound 3g.



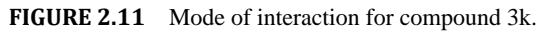


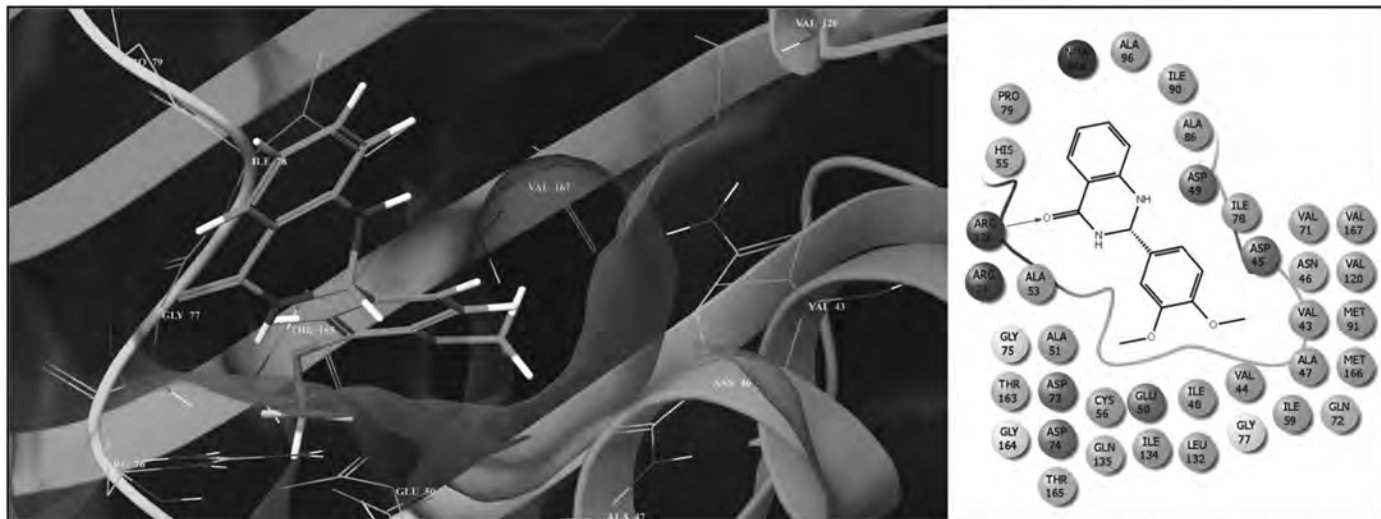


**FIGURE 2.9** Mode of interaction for compound 3i.

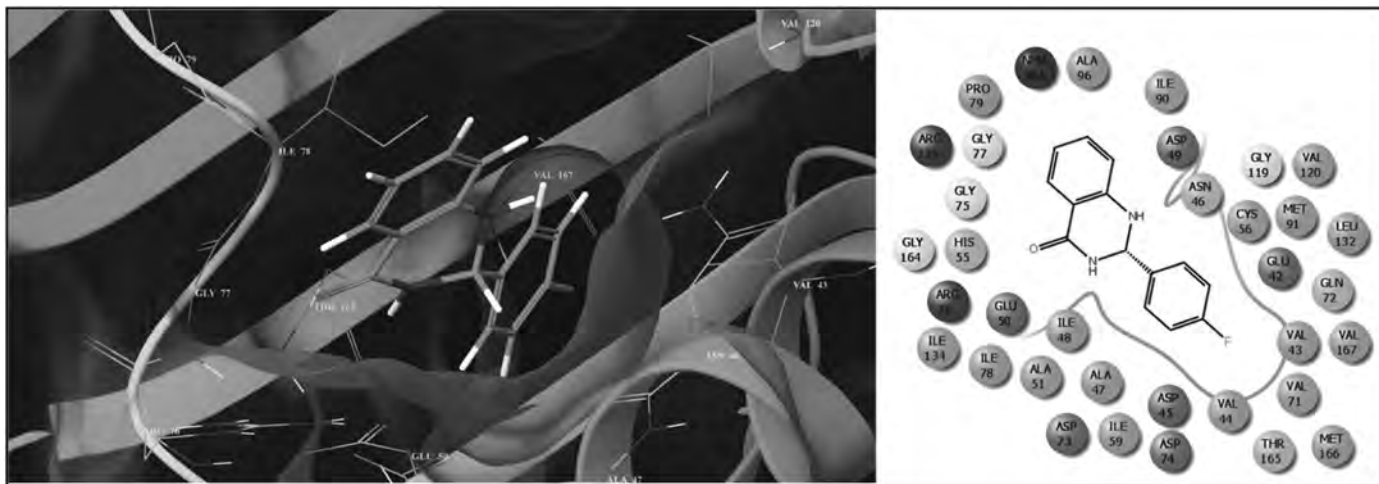


**FIGURE 2.10** Mode of interaction for compound 3j.

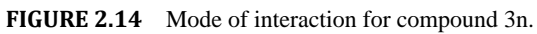


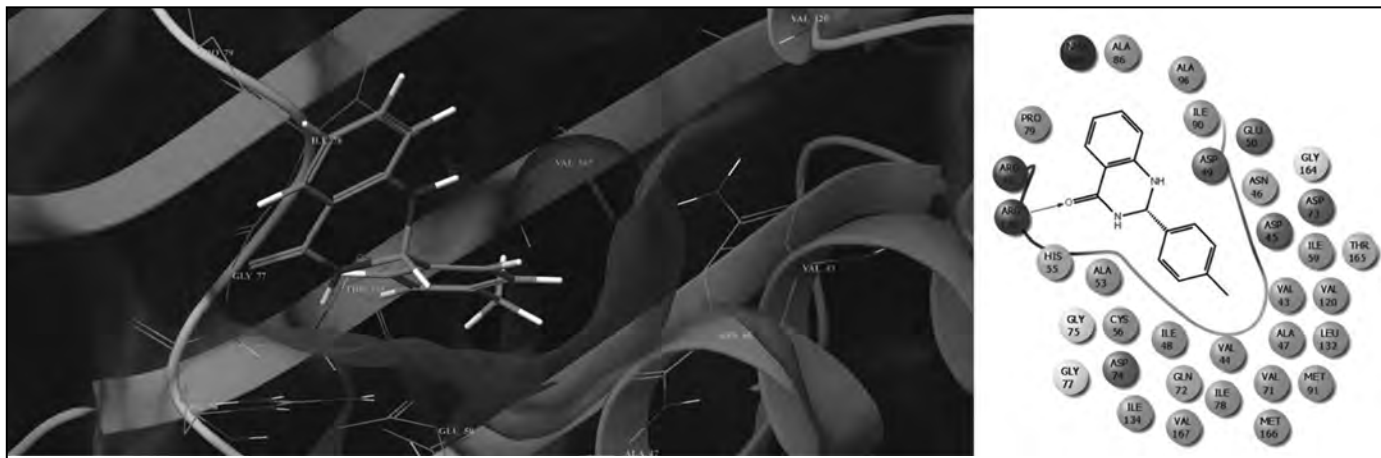


**FIGURE 2.12** Mode of interaction for compound 3l.

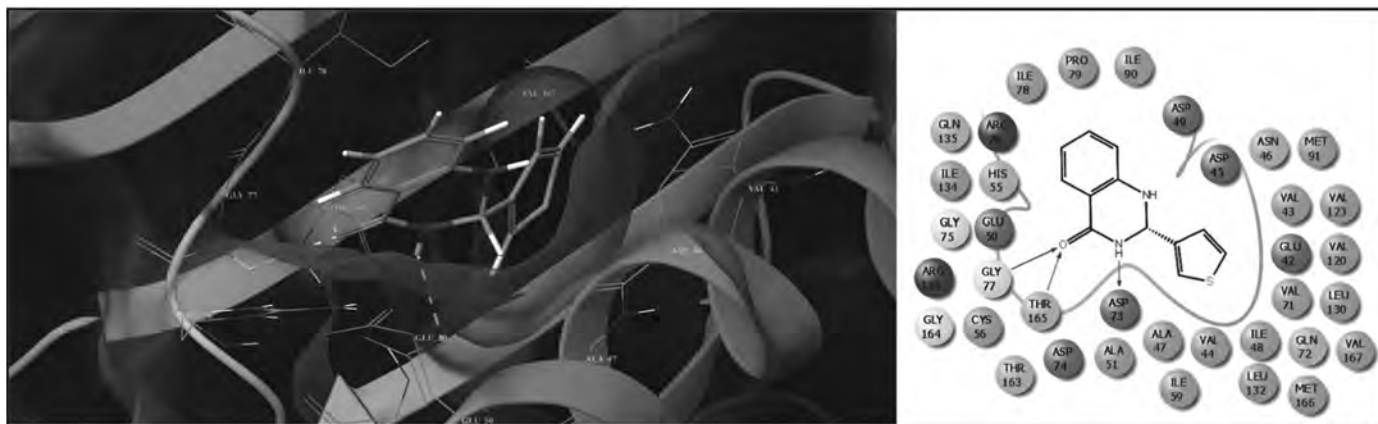


**FIGURE 2.13** Mode of interaction for compound 3m.





**FIGURE 2.15** Mode of interaction for compound 3o.



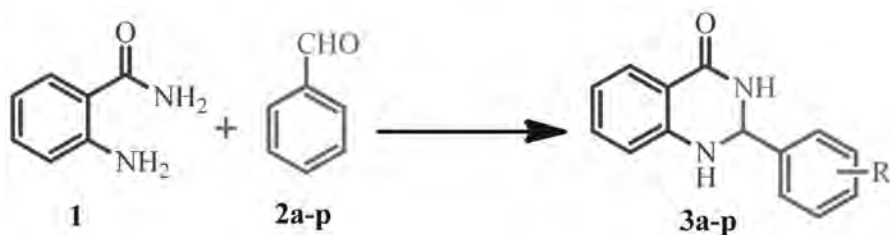
**FIGURE 2.16** Mode of interaction for compound 3p.



**TABLE 2.2** Pharmacokinetic Constraints Significant for Decent Oral Bioavailability

Cpd	mi Log <i>P</i>	TPSA (Å <sup>2</sup> )	MW	n- ON	n- OHNH	Lipinski Violation	n- ROTB	MV	% ABS	Drug- Likeness Score
Rule	≤ 5	–	< 500	< 10	< 5	≤ 1				
3a	2.37	41.12	224.26	3	2	0	1	205.23	94.81	–0.06
3b	2.28	86.95	269.26	6	2	0	2	228.56	79.00	–0.15
3c	2.31	86.95	269.26	6	2	0	2	228.56	79.00	0.15
3d	2.33	86.95	269.26	6	2	0	2	228.56	79.00	0.22
3e	3.00	41.12	258.71	3	2	0	1	218.76	94.81	0.02
3f	3.05	41.12	258.71	3	2	0	1	218.76	94.81	0.88
3g	3.13	41.12	303.16	3	2	0	1	223.11	94.81	–0.32
3h	3.16	41.12	303.16	3	2	0	1	223.11	94.81	0.10
3i	3.18	41.12	303.16	3	2	0	1	223.11	94.81	0.50
3j	1.89	61.35	240.26	4	3	0	1	213.24	87.83	0.85
3k	2.43	50.36	254.29	4	2	0	2	230.77	91.62	0.55
3l	2.02	59.59	284.31	5	2	0	3	256.32	88.44	0.42
3m	2.54	41.12	242.25	3	2	0	1	210.16	94.81	0.73
3n	2.79	41.12	238.29	3	2	0	1	221.79	94.81	0.31
3o	2.82	41.12	238.29	3	2	0	1	221.79	94.81	0.43
3p	1.96	41.12	230.29	3	2	0	1	195.94	94.81	–0.08

In order to analyze the DHQ derivatives 3a–p to its biological antimicrobial activity, we synthesized them (Scheme 2.1) using the information from the molecular docking research and ADME properties. The outcomes are shown in (Table 2.3) and Figure 2.17 displays structures of each molecule that was synthesized.

**SCHEME 2.1** Synthesis of DHQ derivatives 3a–p.

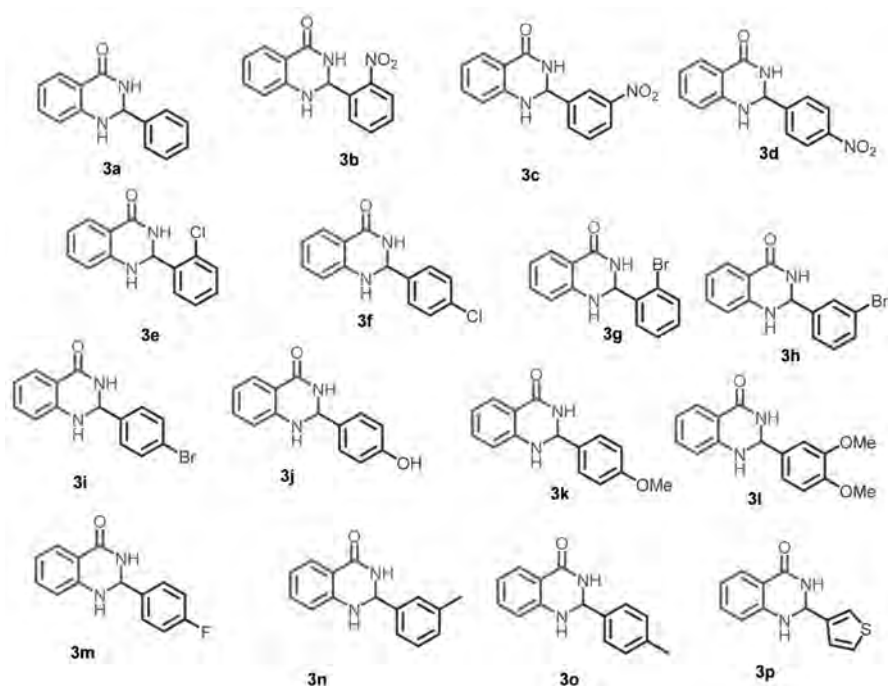


FIGURE 2.17 The synthesized DHQ derivative's 3a–p structure.

## 2.5 BIOLOGICAL EVALUATION

### 2.5.1 ANTIBACTERIAL ACTIVITY

The results from the assessment of the antibacterial properties of the functionalized DHQ 3a–p against Gram positive *Staphylococcus aureus* (SA) and *Micrococcus luteus* (ML) and Gram negative *Escherichia coli* (EC) and *Pseudomonas fluorescens* (PF) strains can be observed in Table 2.3. In contrast to traditional antibacterial drugs, the synthesized compounds demonstrated only moderate effectiveness in inhibiting bacterial growth.

### 2.5.2 ANTIFUNGAL ACTIVITY

All synthetic DHQ 3a–p exhibits decent to reasonable antifungal movement in contradiction of all of the tested fungal strains *Candida albicans* (CA), *Fusarium oxysporum* (FO) and *Aspergillus flavus* (AF) (Table 2.3).

**TABLE 2.3** *In Vitro* Antimicrobial (MIC) Activities of 3a–p (μg/mL)

Compounds	Gram +ve Bacteria		Gram –ve Bacteria		Antifungal Activity		
	SA	ML	EC	PF	CA	FO	AF
3a	16	32	32	8	32	32	32
3b	4	16	8	8	32	32	32
3c	8	16	8	4	16	32	16
3d	16	16	8	8	32	32	64
3e	8	16	4	4	16	64	32
3f	8	16	8	4	16	32	16
3g	4	32	4	8	16	16	16
3h	4	32	32	8	32	64	64
3i	8	32	32	32	16	32	32
3j	8	32	16	16	32	32	64
3k	8	32	4	8	64	16	16
3l	8	32	4	16	64	64	64
3m	4	16	4	16	16	16	16
3n	8	32	32	32	32	32	64
3o	16	32	16	8	32	16	64
3p	8	8	16	8	16	16	16
AP	4	16	4	2	–	–	–
KM	2	2	2	2	–	–	–
MA	–	–	–	–	16	16	16
FA	–	–	–	–	2	2	4

*Abbreviations:* Cpd: Compound; AP: ampicilin; KM: kanamycin; MA: miconazole; and FA: fluconazole.

The results of these endeavor are now being effectively utilized using the *In silico*-chemico-biology approach to identify compounds with advanced effectiveness and discrimination.

## 2.6 CONCLUSION

Utilization of medication has contributed to advancements in data-driven and computer-assisted manufacturing of medications. Presently, machine learning algorithms are being applied in organic production for various purposes, including automated synthesis through retrosynthetic analysis, proposing feasible artificial routes, estimating product yield and outcomes, identifying novel catalysts, and optimizing reaction conditions.

Although computer-aided synthesis is a relatively recent field, it has the potential to revolutionize the design and evaluation of synthetic pathways for target compounds. This innovation can greatly reduce the assignment of chemical combination personnel to create novel prospects for drug fusion. Another significant development is computer-aided drug design, which employs techniques such as building computer-generated chemical libraries and implementing de novo drug molecule strategies to enhance the efficiency of drug discovery. However, it is crucial to consider factors like stability, pharmacokinetics, and toxicity during the drug development process, which may not have been fully addressed in the initial stages. To facilitate chemical synthesis and accelerate the production of medication molecules, computer-aided drug synthesis pathway enterprises employ retrosynthetic analysis to identify synthetic routes for drug cores and predict reaction conditions and products. Nonetheless, predicting complex chemicals and identifying potential flaws in anticipated routes remains a considerable challenge. This approach has garnered significant attention due to its ability to streamline the expertise required in chemical combination to enhance the competence of drug production.

## KEYWORDS

- **adenosine triphosphate**
- **ampicillin**
- *Aspergillus flavus*
- **computer-aided drug design**
- **computer-aided reaction prediction**
- **molecular docking**
- **retrosynthetic analysis of drug**

## REFERENCES

1. Chakrabarti, A., Morgenstern, S., & Knaab, H. (2004). Identification and application of requirements and their application on the design process: A protocol study. *Research in Engineering Design*, 15, 22–39.
2. Davies, I. W. (2019). The digitization of organic synthesis. *Nature*, 570, 175–181.
3. Marko, I. E. (2001). The art of total synthesis. *Science*, 294, 1842–1843.

4. Wender, P. A., & Miller, B. L. (2009). Synthesis at the molecular frontier. *Nature*, 460, 197–201.
5. Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27, 379–423.
6. Brønsted, J. N., & Pedersen, K. (1924). Sur le mécanisme des réactions de catalyse. *Zeitschrift für Physikalische Chemie*, 108U, 185–235.
7. (a) Merrifield, R. B., Stewart, J. M., & Jernberg, N. (1966). The synthesis of polypeptides by the solid-phase method. *Analytical Chemistry*, 38, 1905–1914; (b) Erdős, E. G. (1966). The total synthesis of proteins. *Science*, 152, 1284–1285.
8. Corey, E. J., & Wipke, W. T. (1969). The synthesis of complex organic molecules by computer-aided methods. *Science*, 166, 178–192.
9. Gelernter, H. (1973). *The Discovery of Organic Synthetic Routes by Computer*. Springer Berlin Heidelberg.
10. Peishoff, C. E., & Joergensen, W. L. (1985). A method for the generation and evaluation of synthetic routes. *Journal of Organic Chemistry*, 17, 3175.
11. Ugi, I., Stein, N., Knauer, M., Gruber, B., & Bley, K. (1993). The theory and applications of multicomponent reactions. *Topics in Current Chemistry*, 166, 199–233.
12. Fick, R., Gasteiger, J., & Ihlenfeldt, W. D. J. A. P. (1990). Software development in chemistry: Proceedings of the 4th Workshop on Computational Chemistry. *Software Development in Chemistry* 4, 57–65.
13. Evans, D. A. (2014). Asymmetric synthesis of complex molecules. *Angewandte Chemie International Edition*, 53, 11140–11145.
14. Weininger, D. (1988). SMILES: A chemical language and information system. *Journal of Chemical Information and Computer Sciences*, 28, 31–36.
15. (a) Steinbeck, C., Hoppe, C., Kuhn, S., Floris, M., Guha, R., & Willighagen, E. L. (2006). The Chemistry Development Kit (CDK): An open-source Java library for chemo- and bioinformatics. *Current Pharmaceutical Design*, 12, 2111–2120; (b) Steinbeck, C., Han, Y., Kuhn, S., Horlacher, O., Luttmann, E., & Willighagen, E. (2003). The Chemistry Development Kit (CDK): An open-source Java library for chemo- and bioinformatics. *Journal of Chemical Information and Computer Sciences*, 43, 493–500.
16. (a) Rupp, M., Tkatchenko, A., Müller, K. R., & von Lilienfeld, O. A. (2012). Fast and accurate modeling of molecular atomization energies with machine learning. *Physical Review Letters*, 108, 058301; (b) Hansen, K., Montavon, G., Biegler, F., Fazli, S., Rupp, M., Scheffler, M., von Lilienfeld, O. A., Tkatchenko, A., & Müller, K. R. (2013). Assessment and validation of machine learning methods for predicting molecular atomization energies. *Journal of Chemical Theory and Computation*, 9, 3404–3419.
17. Heller, S., McNaught, A., Stein, S., Tchekhovskoi, D., & Pletnev, I. (2013). The development of a chemical information system. *Journal of Cheminformatics*, 5, 7.
18. (a) Rogers, D., & Hahn, M. (2010). Extended-connectivity fingerprints. *Journal of Chemical Information and Model*, 50, 742–754; (b) Muegge, I., & Mukherjee, P. (2016). IUPAC's role in drug discovery: From molecular modeling to cheminformatics. *Expert Opinion on Drug Discovery*, 11, 137–148; (c) Durant, J. L., Leland, B. A., Henry, D. R., & Nourse, J. G. (2002). Reoptimization of the PubChem compound database. *Journal of Chemical Information and Computer Sciences*, 42, 1273–1280; (d) Cereto-Massagué, A., Ojeda, M. J., Valls, C., Mulero, M., Garcia-Vallve, S., & Pujadas, G. (2015). Chemical informatics functionality of the KNIME workbench. *Methods*, 71, 58–63; (e) Bender, A., Mussa, H. Y., Glen, R. C., & Reiling, S. (2004). Molecular similarity: A

- new approach for modeling biological activity. *Journal of Chemical Information and Computer Sciences*, 44, 1708–1718.
19. Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., & Hassabis, D. (2016). Mastering the game of Go with deep neural networks and tree search. *Nature*, 529, 484–489.
  20. Gasteiger, J., & Ihlenfeldt, W. (1990). *Software Development in Chemistry 4*. Springer.
  21. (a) Ott, M. A., & Noordik, J. H. (1992). Computer tools for reaction retrieval and synthesis planning in organic chemistry. A brief review of their history, methods, and programs. *Recl. Trav. Chim. Pays-Bas*, 111, 239–246. (b) Todd, M. H. (2005). Computer-aided organic synthesis. *Chem. Soc. Rev.*, 34, 247–266. (c) Cook, A., Johnson, A. P., Law, J., Mirzazadeh, M., Ravitz, O., & Simon, A. (2012). Computer-aided synthesis design: 40 years on. *Wiley Interdiscip. Rev.: Comput. Mol. Sci.*, 2, 79–107. (d) Warr, W. A. (2014). A short review of chemical reaction database systems, computer-aided synthesis design, reaction prediction and synthetic feasibility. *Mol. Inf.*, 33, 469–476.
  22. Shen, Y., Borowski, J. E., Hardy, M. A., Sarpong, R., Doyle, A. G., & Cernak, T. (2021). *Nat. Rev. Methods Primers*, 1, 23.
  23. Corey, E. J., Long, A. K., & Rubenstein, S. D. (1985). Computer-assisted analysis in organic synthesis. *Science*, 228, 408–419.
  24. (a) Gelernter, H., Rose, J. R., & Chen, C. (1990). Building and refining a knowledge base for synthetic organic chemistry via the methodology of inductive and deductive machine learning. *J. Chem. Inf. Model.*, 30, 492–504. (b) Satoh, H., & Funatsu, K. (1995). Sophia, a knowledge base-guided reaction prediction system-utilization of a knowledge base derived from a reaction database. *J. Chem. Inf. Model.*, 35, 34–44. (c) Satoh, K., & Funatsu, K. (1999). A novel approach to retrosynthetic analysis using knowledge bases derived from reaction databases. *J. Chem. Inf. Comput. Sci.*, 39, 316–325. (d) Law, J., Zsoldos, Z., Simon, A., Reid, D., Liu, Y., Khew, S. Y., Johnson, A. P., Major, S., Wade, R. A., & Ando, H. Y. (2009). Route Designer: A retrosynthetic analysis tool utilizing automated retrosynthetic rule generation. *J. Chem. Inf. Model.*, 49, 593–602. (e) Borgevig, A., Federsel, H. J., Huerta, F., Hutchings, M. G., Kraut, H., Langer, T., Low, P., Oppawsky, C., Rein, T., & Saller, H. (2015). Route design in the 21st century: The ICSYNTH software tool as an idea generator for synthesis prediction. *Org. Process Res. Dev.*, 19, 357–368. (f) Coley, C. W., Barzilay, R., Jaakkola, T. S., Green, W. H., & Jensen, K. F. (2017). Prediction of organic reaction outcomes using machine learning. *ACS Cent. Sci.*, 3, 434–443. (g) Christ, C. D., Zentgraf, M., & Kriegl, J. M. (2012). Mining electronic laboratory notebooks: Analysis, retrosynthesis, and reaction-based enumeration. *J. Chem. Inf. Model.*, 52, 1745–1756.
  25. Segler, M. H. S., & Waller, M. P. (2017). Neural-symbolic machine learning for retrosynthesis and reaction prediction. *Chem. Eur. J.*, 23, 5966–5971.
  26. Coley, C. W., Green, W. H., & Jensen, K. F. (2018). *Acc. Chem. Res.*, 51, 1281–1289.
  27. (a) Szymkuc, S., Gajewska, E. P., Klucznik, T., Molga, K., Dittwald, P., Startek, M., Bajczyk, M., & Grzybowski, B. A. (2016). *Angew. Chem. Int. Ed.*, 55, 5904–5937. (b) Kowalik, M., Gothard, C. M., Drews, A. M., Gothard, N. A., Weckiewicz, A., Fuller, P. E., & Grzybowski, B. A. (2012). *Angew. Chem. Int. Ed.*, 51, 7928–7932. (c) Grzybowski, B. A., Bishop, K. J., Kowalczyk, B., & Wilmer, C. E. (2009). *Nat. Chem.*, 1, 31–36. (d) Bishop, K. J., Klajn, R., & Grzybowski, B. A. (2006). *Angew. Chem. Int.*

- Ed.*, 45, 5348–5354. (e) Badowski, T., Molga, K., & Grzybowski, B. A. (2019). *Chem. Sci.*, 10, 4640–4651.
28. Molga, K., Dittwald, P., & Grzybowski, B. A. (2019). *Chemistry*, 5, 460–473.
29. Coley, C. W., Rogers, L., Green, W. H., & Jensen, K. F. (2017). *ACS Cent. Sci.*, 3, 1237–1245.
30. Law, J., Zsoldos, Z., Simon, A., Reid, D., Liu, Y., Khew, S. Y., Johnson, A. P., Major, S., Wade, R. A., & Ando, H. Y. (2009). *J. Chem. Inf. Model.*, 49, 593–602.
31. Bøgevig, A., Federsel, H. J., Huerta, F., Hutchings, M. G., Kraut, H., Langer, T., Leow, P., Oppawsky, C., Rein, T., & Saller, H. (2015). *Org. Process Res. Dev.*, 19, 357–368.
32. (a) Segler, M. H. S., & Waller, M. P. (2017). *Chemistry*, 23, 5966–5971; (b) Segler, M. H. S., & Waller, M. P. (2017). *Chemistry*, 23, 6118–6128.
33. Segler, M. H. S., Preuss, M., & Waller, M. P. (2018). *Nature*, 555, 604–610.
34. (a) Kayala, M. A., & Baldi, P. (2012). *J. Chem. Inf. Model.*, 52, 2526–2540; (b) Kayala, M. A., Azencott, C. A., Chen, J. H., & Baldi, P. (2011). *J. Chem. Inf. Model.*, 51, 2209–2222; (c) Fooshee, D., Mood, A., Gutman, E., Tavakoli, M., Urban, G., Liu, F., Huynh, N., Van Vranken, D., & Baldi, P. (2018). *Mol. Syst. Des. Eng.*, 3, 442–452.
35. Liu, B., Ramsundar, B., Kawthekar, P., Shi, J., Gomes, J., Luu, Q., Nguyen, H., Sloane, J., Wender, P., & Pande, V. (2017). *ACS Cent. Sci.*, 3, 1103–1113.
36. Lin, K., Xu, Y., Pei, J., & Lai, L. (2020). *Chem. Sci.*, 11, 3355–3364.
37. Ewing, T. J. A., & Kuntz, I. D. (1997). Critical evaluation of search algorithms used in automated molecular docking. *J. Comput. Chem.*, 18, 1175–1189.
38. Morris, G. M., Goodsell, D. S., Halliday, R. S., Huey, R., Hart, W. E., Belew, R. K., & Olson, A. J. (1998). Automated docking using a Lamarckian genetic algorithm and empirical binding free energy function. *J. Comput. Chem.*, 19, 1639–1662.
39. Jones, G., Willett, P., Glen, R. C., Leach, A. R., & Taylor, R. (1997). Development and validation of a genetic algorithm for flexible docking. *J. Mol. Biol.*, 267, 727–748.
40. Rarey, M., Kramer, B., & Lengauer, T. (1999). Docking of hydrophobic ligands with interaction-based matching algorithms. *Bioinformatics*, 15, 243–250.
41. Friesner, R. A., Banks, J. L., Murphy, R. B., Halgren, T. A., Klicic, J. J., Mainz, D. T., Repasky, M. P., Knoll, E. H., Shaw, D. E., Shelley, M., Perry, J. K., Francis, P., & Shenkin, P. S. (2004). Glide: A new approach for rapid, accurate docking and scoring. Method and assessment of docking accuracy. *J. Med. Chem.*, 47, 1739–1749.
42. Venkatachalam, C. M., Jiang, X., Oldfield, T., & Waldman, M. (2003). LigandFit: A novel method for the shape-directed rapid docking of ligands to protein active sites. *J. Mol. Graph. Model.*, 21, 289–307.
43. Stewart, J. P. (2009). Mopac93, Fujitsu Ltd., Tokyo, Japan (Scigress Explorer v7.7.0.47).
44. Hu, L., Benson, M. L., Smith, R. D., Lerner, M. G., & Carlson, H. A. (2005). Binding MOAD (Mother of All Databases). *Proteins*, 60, 333–340.
45. (a) Yadav, D. K., Meena, A., Srivastava, A., Chanda, D., Khan, F., & Chattopadhyay, S. K. (2010). Development of QSAR model for immunomodulatory activity of natural Coumarinolignoids. *Drug Des. Devel. Ther.*, 4, 173–186; (b) Taylor, R. D., Jewsbury, P. J., & Essex, J. W. (2002). A review of protein-small molecule docking methods. *J. Comput. Aided Mol. Des.*, 16, 151–166; (c) Ekins, S., Mestres, J., & Testa, B. (2007). In silico pharmacology for drug discovery: Applications to targets and beyond. *Br. J. Pharmacol.*, 152, 21–37.
46. Kuntz, I. D., Blaney, J. M., Oatley, S. J., Langridge, R., & Ferrin, T. E. (1982). A geometric approach to macromolecule-ligand interactions. *J. Mol. Biol.*, 161, 269–288.

47. Lybrand, T. P. (1995). Ligand-protein docking and rational drug design. *Curr. Opin. Struct. Biol.*, 5, 224–228.
48. Blaney, J. M., & Dixon, J. S. (1993). A good ligand is hard to find: Automatic docking methods. *Perspect. Drug Discov. Des.*, 1, 301–319.
49. Hoppe, C., Steinbeck, C., & Wohlfahrt, G. (2006). Classification and comparison of ligand-binding sites derived from grid-mapped knowledge-based potentials. *J. Mol. Graph. Model.*, 24, 328–340.
50. Walters, W. P., Stahl, M. T., & Murcko, M. A. (1998). Virtual screening: An overview. *Drug Discov. Today*, 3, 160–178.
51. (a) Fukunishi, Y., Kubota, S., & Nakamura, H. (2006). Noise reduction method for molecular interaction energy: Application to in silico drug screening and in silico target protein screening. *J. Chem. Inf. Model.*, 46, 2071–2084; (b) Khan, F., Meena, A., & Sharma, A. (2010). Docking-based virtual screening of anticancer drugs. In R. Arora (Ed.), *Medicinal Plant Biotechnology* (Vol. 15, pp. 242–264). CAB International, UK; (c) Mestres, J., Martin-Couce, L., Gregori-Puigjane, E., Cases, M., & Boyer, S. (2006). Ligand-based approach to in silico pharmacology: Nuclear receptor profiling. *J. Chem. Inf. Model.*, 46, 2725–2736.
52. Gehlhaar, D. K., Moerder, K. E., Zichi, D., Sherman, C. J., Ogden, R. C., & Freer, S. T. (1995). De novo design of enzyme inhibitors by Monte Carlo ligand generation. *J. Med. Chem.*, 38, 466–472.
53. Dewitte, R. S., Ishchenko, A. V., & Shakhnovich, E. I. (1997). SMOG: De novo design method based on simple, fast, and accurate free energy estimates. 2. Case studies in molecular design. *J. Am. Chem. Soc.*, 119, 4608–4617.
54. Gillet, V., Johnson, A. P., Mata, P., Sike, S., & Williams, P. (1993). Sprout: A program for structure generation. *J. Comp. Aid. Mol. Des.*, 7, 127–153.
55. Roe, D. C., & Kuntz, I. D. (1995). Builder v.2: Improving the chemistry of a de novo design strategy. *J. Comput. Aid. Mol. Des.*, 9, 269–282.
56. Pearlman, D. A., & Murcko, M. A. (1993). Concepts: New dynamic algorithm for de novo drug suggestions. *J. Med. Chem.*, 10, 1184–1193.
57. Pearlman, D. A., & Murcko, M. A. (1996). Concerts: Dynamic connection of fragments as an approach to de novo ligand design. *J. Med. Chem.*, 39, 1651–1663.
58. Stultz, C. M., & Karplus, M. (2000). Dynamic ligand design and combinatorial optimization: Designing inhibitors to endothiapepsin. *Proteins*, 40, 258–289.
59. Rotstein, S. H., & Murcko, M. A. (1993). GenStar: A method for de novo drug design. *J. Comput. Aid. Mol. Des.*, 7, 23–43.
60. Rotstein, S. H., & Murcko, M. A. (1993). GroupBuild: A fragment-based method for de novo drug design. *J. Med. Chem.*, 36, 1700–1710.
61. Moon, J. B., & Howe, W. J. (1991). Computer design of bioactive molecules: A method for receptor-based de novo ligand design. *Proteins*, 11, 314–328.
62. Eisen, M. B., Wiley, D. C., Karplus, M., & Hubbard, R. E. (1994). Hook: A program for finding novel molecular architectures that satisfy the chemical and steric requirements of a macromolecule binding site. *Proteins*, 19, 199–221.
63. Nishibata, Y., & Itai, A. (1993). Confirmation of usefulness of a structure construction program based on three-dimensional receptor structure for rational lead generation. *J. Med. Chem.*, 36, 2921–2928.
64. (a) Duncia, J. V., Chiu, A. T., Carini, D. J., Gregory, G. B., Johnson, A. L., Price, W. A., Wells, G. J., Wong, P. C., Calabrese, J. C., & Timmermans, P. B. (1990). The



- discovery of potent nonpeptide angiotensin II receptor antagonists: A new class of potent antihypertensives. *J. Med. Chem.*, **33**, 1312–1329; (b) Duncia, J. V., Carini, D. J., Chiu, A. T., Johnson, A. L., Price, W. A., Wong, P. C., Wexler, R. R., & Timmermans, P. B. (1992). The discovery of DuP 753, a potent, orally active nonpeptide angiotensin II receptor antagonist. *Med. Res. Rev.*, **12**, 149–191.
65. Koga, H., Itoh, A., Murayama, S., Suzue, S., & Irikura, T. (1980). Structure-activity relationships of antibacterial 6,7- and 7,8-disubstituted 1-alkyl-1,4-dihydro-4-oxoquinoline-3-carboxylic acids. *J. Med. Chem.*, **23**, 1358–1363.
  66. Kawakami, Y., Inoue, A., Kawai, T., Wakita, M., Sugimoto, H., & Hopfinger, A. J. (1996). The rationale for E2020 as a potent acetylcholinesterase inhibitor. *Bioorg. Med. Chem.*, **4**, 1429–1446.
  67. Greer, J., Erickson, J. W., Baldwin, J. J., & Varney, M. D. (1994). Application of the three-dimensional structures of protein target molecules in structure-based drug design. *J. Med. Chem.*, **37**, 1035–1054.
  68. (a) Dorsey, B. D., Levin, R. B., McDaniel, S. L., Vacca, J. P., Guare, J. P., Darke, P. L., Zugay, J. A., Emini, E. A., Schleif, W. A., & Quintero, J. C. (1994). L-735,524: The design of a potent and orally bioavailable HIV protease inhibitor. *J. Med. Chem.*, **37**, 3443–3451; (b) Bodor, N., & Huang, M. (1995). Computer-aided design of new drugs based on retrometabolic concepts. In C. Reynolds, M. Holloway, & H. Cox (Eds.), *Computer Aided Molecular Design* (Vol. 589, pp. 98–113). ACS Symposium Series.
  69. Kaldor, S. W., Kalish, V. J., Davies, J. F., Shetty, B. V., Fritz, J. E., Appelt, K., Burgess, J. A., Campanale, K. M., Chirgadze, N. Y., Clawson, D. K., Dressman, B. A., Hatch, S. D., Khalil, D. A., Kosa, M. B., Lubbehusen, P. P., Muesing, M. A., Patick, A. K., Reich, S. H., Su, K. S., & Tatlock, J. H. (1997). Viracept (nelfinavir mesylate, AG1343): A potent, orally bioavailable inhibitor of HIV-1 protease. *J. Med. Chem.*, **40**, 3979–3985.
  70. Glen, R. C., Martin, G. R., Robertson, A. D., Buckingham, J., Woolard, P. M., Hill, A. P., Hyde, R. M., & Salmon, J. A. (1995). Computer-aided design and synthesis of 5-substituted tryptamines and their pharmacology at the 5-HT<sub>1D</sub> receptor: The discovery of 311C90, a compound with potential anti-migraine properties. *J. Med. Chem.*, **38**, 3566–3580.
  71. Bellocchi, D., Macchiarulo, A., Costantino, G., & Pellicciari, R. (2005). Docking studies on PARP-1 inhibitors: Insights into the role of a binding pocket water molecule. *Bioorg. Med. Chem.*, **13**, 1151–1157.
  72. Takacs, A., Fodor, J., Nemeth, Z., & Hell, Z. (2014). Zeolite-catalyzed method for the preparation of 2,3-dihydroquinazolin-4(1H)-ones. *Synth. Commun.*, **44**, 2269–2275.
  73. Noel, R., Gupta, N., Pons, V., Goudet, A., Garcia-Castillo, M. D., Michau, A., Martinez, J., Buisson, D., Johannes, L., Gillet, D., Barbier, J., & Cintrat, J. (2013). N-Methyldihydroquinazolinone derivatives of Retro-2 with enhanced efficacy against Shiga toxin. *J. Med. Chem.*, **56**, 3404–3413.
  74. (a) Sadanandam, Y. S., Reddy, K. R. M., & Rao, A. B. (1987). Synthesis of substituted 2,3-dihydro-1-( $\beta$ -phenylethyl)-2-aryl- and 2,3-diaryl-4(1H)-quinazolinones and their pharmacological activities. *Eur. J. Med. Chem.*, **22**, 169–173; (b) Mohammadi, A. A., Rohi, H., & Soorki, A. A. (2013). Synthesis and in vitro antibacterial activities of novel 2-aryl-3-(phenylamino)-2,3-dihydroquinazolin-4(1H)-one derivatives. *J. Heterocycl. Chem.*, **50**, 1129–1133; (c) Derbyshire, E. R., Min, J., Guiguemde, W. A., Clark, J. A., Connelly, M. C., Magalhaes, A. D., Guy, R. K., & Clardy, J. (2014). Dihydroquinazolinone inhibitors of proliferation of blood and liver stage malaria parasites. *Antimicrob. Agents*

- Chemother.*, 58, 1516–1522; (d) Uruno, Y., Konishi, Y., Suwa, A., Takai, K., Tojo, K., Nakako, T., Sakai, M., Enomoto, T., Matsuda, H., Kitamura, A., & Sumiyoshi, T. (2015). Discovery of dihydroquinazolinone derivatives as potent, selective, and CNS-penetrant M1 and M4 muscarinic acetylcholine receptors agonists. *Bioorg. Med. Chem. Lett.*, 25, 5357–5361; (e) Hemalatha, K., Madhumitha, G., Ravi, L., Khanna, V. G., Al-Dhabi, N. A., & Arasu, M. V. (2016). Binding mode of dihydroquinazolinones with lysozyme and its antifungal activity against *Aspergillus* species. *J. Photochem. Photobiol. B*, 161, 71–79.
75. (a) Levin, J. I., Chan, P. S., Bailey, T., Katocs, A. S., & Venkatesan, A. M. (1994). The synthesis of 2,3-dihydro-4(1H)-quinazolinone angiotensin II receptor antagonists. *Bioorg. Med. Chem. Lett.*, 4, 1141–1146; (b) Hasegawa, H., Muraoka, M., Matsui, K., & Kojima, A. (2006). A novel class of sodium/calcium exchanger inhibitors: Design, synthesis, and structure-activity relationships of 4-phenyl-3-(piperidin-4-yl)-3,4-dihydro-2(1H)-quinazolinone derivatives. *Bioorg. Med. Chem. Lett.*, 16, 727–730; (c) Singh, M., & Raghav, N. (2015). 2,3-Dihydroquinazolin-4(1H)-one derivatives as potential non-peptidyl inhibitors of cathepsins B and H. *Bioorg. Chem.*, 59, 12–22; (d) Zhang, H., Liu, H., Luo, X., Wang, Y., Liu, Y., Jin, H., Liu, Z., Yang, W., Yu, P., Zhang, L., & Zhang, L. (2018). Design, synthesis and biological activities of 2,3-dihydroquinazolin-4(1H)-one derivatives as TRPM2 inhibitors. *Eur. J. Med. Chem.*, 152, 235–252; (e) Xing, J., Yang, L., Yang, Y., Zhao, L., Wei, Q., Zhang, J., Zhou, J., & Zhang, H. (2017). Design, synthesis and biological evaluation of novel 2,3-dihydroquinazolin-4(1H)-one derivatives as potential fXa inhibitors. *Eur. J. Med. Chem.*, 125, 411–422.
76. (a) Kharmawlong, G. K., Nongrum, R., Chhetri, B., Rani, J. W. S., Rahman, N., Yadav, A. K., & Nongkhlaw, R. (2019). Green and efficient one-pot synthesis of 2,3-dihydroquinazolin-4(1H)-ones and their anthelmintic studies. *Synth. Commun.*, 49, 2683–2695; (b) Sabnis, R. W. (2021). Novel substituted 3,4-dihydroquinazoline derivatives for treating hepatitis B virus infection. *ACS Med. Chem. Lett.*, 34, 1492–1503.
77. (a) Ehmann, D. E., & Lahiri, S. D. (2014). Novel compounds targeting bacterial DNA topoisomerase/DNA gyrase. *Curr. Opin. Pharmacol.*, 18, 76–83; (b) Maxwell, A., & Lawson, D. M. (2003). The ATP-binding site of type II topoisomerases as a target for antibacterial drugs. *Curr. Top Med. Chem.*, 3, 283–303; (c) Pommier, Y. (2013). Drugging topoisomerases: Lessons and challenges. *ACS Chem. Biol.*, 8, 82–95; (d) Collin, F., Karkare, S., & Maxwell, A. (2011). Exploiting bacterial DNA gyrase as a drug target: Current state and perspectives. *Appl. Microbiol. Biotechnol.*, 92, 479–497.
78. (a) Schrodinger Suite 2015-4 QM-Polarized Ligand Docking protocol; Glide version 6.9, Schrodinger, LLC, New York, NY, 2015; Jaguar version 9.0, Schrodinger, LLC, New York, NY, 2015; QSite version 6.9, Schrodinger, LLC, New York, NY, 2015; (b) Friesner, R. A., Murphy, R. B., Repasky, M. P., Frye, L. L., Greenwood, J. R., Halgren, T. A., Sanschagrin, P. C., & Mainz, D. T. (2006). Extra precision glide: Docking and scoring incorporating a model of hydrophobic enclosure for protein-ligand complexes. *J. Med. Chem.*, 49, 6177–6196.
79. Molinspiration Cheminformatics. (2014). Available from: <http://www.molinspiration.com/cgi-bin/properties> (accessed on 25 July 2024).
80. Lipinski, C. A., Lombardo, L., Dominy, B. W., & Feeney, P. J. (2001). Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Adv. Drug Del. Rev.*, 46, 3–25.

81. Zhao, Y., Abraham, M. H., Lee, J., Hersey, A., Luscombe, N. C., Beck, G., Sherborne, B., & Cooper, I. (2002). Rate-limited steps of human oral absorption and QSAR studies. *Pharm. Res.*, *19*, 1446–1457.
82. Drug-likeness and molecular property prediction. (2024). Available from: <http://www.molsoft.com/mprop/> (accessed on 25 July 2024).
83. Ertl, P., Rohde, B., & Selzer, P. (2000). Fast calculation of molecular polar surface area as a sum of fragment-based contributions and its application to the prediction of drug transport properties. *J. Med. Chem.*, *43*, 3714–3717.

## CHAPTER 3

---

# Computational Tools and Techniques in Planning Organic Synthesis

LAXMI G. KATHAWATE,<sup>1</sup> ROHINI N. SHELKE,<sup>1</sup>  
DATTATRAYA N. PANSARE,<sup>2</sup> and ANIKET P. SARKATE<sup>3</sup>

<sup>1</sup>*Department of Chemistry, Radhabai Kale Mahila Mahavidyalaya, Maharashtra, India*

<sup>2</sup>*Department of Chemistry, Deogiri College, Aurangabad, Maharashtra, India*

<sup>3</sup>*Department of Chemical Technology, Dr. Babasaheb Ambedkar Marathwada University, Aurangabad, Maharashtra, India*

---

### ABSTRACT

Organic synthesis is the process of designing and creating new organic compounds through various chemical reactions. It is a critical field in both chemistry and pharmaceutical industries. However, traditional organic synthesis can be time-consuming, expensive, and requires significant resources. Fortunately, the advent of computational tools and techniques has revolutionized the field of organic synthesis. Computer-aided design and planning have become essential tools for researchers in the field, allowing them to design new compounds *in silico*, rather than in the lab. These tools have not only made the process of designing and developing new compounds more efficient but have also made it more cost-effective. Organic synthesis is a critical field of chemistry that involves creating new molecules and materials with a wide range of applications, from pharmaceuticals to materials science. Traditionally, organic synthesis has relied on a trial-and-error approach, with researchers testing different combinations of chemical reactions to achieve

their desired outcome. However, in recent years, computational tools have revolutionized the field, allowing researchers to predict the outcome of reactions before they even take place. This chapter will explore the power of these computational tools and how they are transforming the way organic synthesis is conducted. From designing new molecules to optimizing reaction conditions, computational tools are enabling researchers to work smarter, not harder, and are opening up new avenues for discovery and innovation in this critical field of chemistry. In this chapter, we'll explore how computational tools and techniques have revolutionized organic synthesis and how they are being used to design and develop novel compounds faster and more efficiently than ever before.

### **3.1 WHAT IS ORGANIC SYNTHESIS AND WHY IS IT IMPORTANT?**

Organic synthesis is the process of designing and creating new organic compounds through chemical reactions. Organic compounds are molecules that contain carbon, hydrogen, and often other elements like nitrogen, oxygen, sulfur, and phosphorus. These compounds are the building blocks of life and play a critical role in fields such as medicine, materials science, and agriculture.

Organic synthesis has been around for centuries, dating back to the discovery of urea by Friedrich W in 1828. Since then, scientists have been working to develop new methods to synthesize organic compounds more efficiently and with greater precision.

Organic synthesis is the process of creating new and complex organic molecules from simpler ones. This practice is essential in many industries, including pharmaceuticals, agrochemicals, and materials science. Organic synthesis allows scientists to create new drugs, pesticides, and materials that can solve real-world problems and improve the quality of life.

Traditionally, organic synthesis has been a long and laborious process, requiring chemists to experimentally test many different reaction conditions to find the optimal conditions for a desired reaction. However, with the advent of computational devices and artificial intelligence, organic synthesis is being revolutionized. However, traditional organic synthesis methods are often time-consuming, expensive, and can involve the use of hazardous chemicals [1, 2]. This is where computational tools and techniques come in. By using computational methods, scientists can design and plan organic syntheses more efficiently and safely than before. Computational tools and techniques have already revolutionized many areas of chemistry and organic synthesis is no exception.

Computational tools allow chemists to predict the outcomes of reactions with unprecedented accuracy, saving time and resources in the laboratory. By using computational simulations and machine learning algorithms, scientists can design new synthetic routes and predict the structures and properties of new compounds before they are even synthesized. This enables faster and more efficient drug discovery, materials development, and other applications in the field of organic synthesis.

The traditional approach to organic synthesis involves extensive experimentation and trial and error, which can be time-consuming and costly. Chemists typically rely on their intuition, experience, and knowledge of chemical reactions to design synthetic pathways for new molecules. This process involves selecting starting materials, determining the appropriate reaction conditions, and optimizing reaction parameters to achieve the desired product. However, this approach is limited by the complexity of chemical reactions and the lack of a comprehensive understanding of reaction mechanisms. As a result, the traditional approach to organic synthesis often leads to inefficient and suboptimal synthetic routes, as well as low yields of desired products. Moreover, this approach does not provide a systematic framework for designing and optimizing synthetic pathways for large-scale production of complex molecules. Therefore, there is a need for a more efficient and rational approach to organic synthesis that can overcome these limitations and enable the synthesis of complex molecules with high efficiency and precision. This is where computational tools come into play, revolutionizing the field of organic synthesis.

Traditional organic synthesis is a time-consuming and labor-intensive process. It requires multiple steps and often involves the use of hazardous chemicals and solvents. The process can also be inefficient, with low yields and the formation of unwanted byproducts. In addition, traditional organic synthesis is limited by the physical properties of the starting materials, which can restrict the types of compounds that can be synthesized.

Another limitation of traditional organic synthesis is the difficulty in predicting the properties of the synthesized compounds. This is because the properties of a compound are determined by its molecular structure, which can be difficult to predict based on the starting materials and the reaction conditions.

These limitations can be overcome through the use of computational tools in organic synthesis. Computational tools can be used to design new compounds with specific properties, predict the properties of synthesized compounds, and optimize reaction conditions to increase yields and reduce waste.

By combining computational tools with traditional organic synthesis, researchers can revolutionize the field of organic synthesis and accelerate the discovery of new compounds with potential applications in medicine, materials science, and other fields.

In this chapter, we will explore how computational tools and techniques are transforming organic synthesis, making it faster, safer, and more sustainable [3, 4]. We will look at some of the most exciting developments in this field and discuss how they could shape the future of organic synthesis.

### **3.2 THE CHALLENGES OF ORGANIC SYNTHESIS AND HOW COMPUTATIONAL TOOLS CAN HELP?**

Organic synthesis is an essential field in chemistry that involves the production of new organic compounds from simpler substances [5–7]. The process is time-consuming, requires a lot of effort and resources, and can be very challenging due to the complexity of the molecules involved. Traditional trial-and-error methods can take years to perfect, and even then, the yield of the desired product may be very low and some side product [8–10].

This is where computational tools and techniques come in. Computational chemistry has revolutionized the field of organic synthesis by providing valuable insights and enabling scientists to predict the outcome of chemical reactions before even stepping into the laboratory.

By using powerful algorithms and simulations, chemists can explore different reaction pathways, evaluate the energy requirements of each step, and identify potential roadblocks that may hinder the synthesis process. This information can then be used to optimize the reaction conditions, select the appropriate reagents, and ultimately increase the yield of the desired product.

Computational tools are changing the face of organic synthesis. In the past, scientists often had to rely on trial and error to find the best conditions for a reaction. This process could take years and was often tedious and frustrating. However, with the advent of computational tools, scientists can now predict the outcome of a reaction before it even takes place. This has revolutionized the field of organic synthesis and has led to the development of new and more efficient methods for synthesizing complex molecules.

Computational tools use algorithms and models to predict the outcome of a reaction based on the molecular properties of the reactants. This allows scientists to quickly identify the best reaction conditions and to optimize the reaction for maximum yield and purity. In addition, computational tools can

be used to design new molecules with specific properties, such as increased solubility or bioactivity.

One of the most exciting applications of computational tools in organic synthesis is the development of machine learning algorithms. These algorithms can analyze vast amounts of data and identify patterns that would be impossible for humans to detect. This has led to the discovery of new reaction pathways and has helped scientists to design new molecules with unprecedented accuracy and efficiency.

Overall, the power of computational tools in organic synthesis cannot be overstated. These tools are changing the way we think about chemical reactions and are accelerating the discovery of new drugs, materials, and technologies. As computational tools continue to improve, we can expect even more exciting developments in the field of organic synthesis in the years to come.

Moreover, computational tools can also help chemists to design entirely new molecules with desirable properties, such as increased stability, selectivity, or bioactivity. This is especially valuable in drug discovery, where computational methods are extensively used to predict the potential efficiency and toxicity of new drug candidates before they are synthesized and tested in the lab.

Overall, the use of computational software and techniques is a game-changer in the area of synthesis, enabling chemists to tackle some of the most challenging problems and accelerate the discovery of new compounds with unprecedented speed and efficiency.

### 3.3 COMPUTER-AIDED SYNTHESIS PLANNING: AN OVERVIEW

Computer-aided synthesis planning (CASP) is a relatively new concept in organic synthesis. It involves the use of computer algorithms and databases to predict the optimal synthesis ways for a given target molecule. CASP is considered to be a powerful tool for revolutionizing organic synthesis as it helps to expedite the process, reduce costs, and minimize the environmental impact of chemical synthesis [11].

At its core, CASP relies on the creation of a knowledge base that includes information on known chemical reactions, reaction conditions, the corresponding products and side products. This database is then used by algorithms to generate a variety of possible synthesis routes for a target molecule. The algorithms take into account factors such as reaction yield, reaction time, and the availability of reagents to determine the optimal synthesis route [12].



CASP techniques have become increasingly sophisticated in recent years, with new algorithms being developed that incorporate artificial intelligence, robotic computing and machine learning. These techniques are able to learn from previous syntheses and predict the most efficient way to produce a target molecule based on the available data. In addition to its potential applications in academic research, CASP has also gained significant interest from the pharmaceutical industry. The development of novel drugs involves the synthesis of complex molecules and CASP can help to reduce the cost of drug development and rapidity this process.

Overall, CASP represents a significant step forward in the field of chemical synthesis. It has the potential to revolutionize the way that new molecules are synthesized, making the process faster, cheaper, and more environmentally friendly.

Computer-aided design (CAD) has revolutionized the field of organic synthesis by enabling researchers to design new molecules and predict their properties without having to synthesize them experimentally. This not only saves time and resources, but also allows for the creation of new molecules that may have been impossible to synthesize otherwise.

CAD tools use algorithms to predict the properties of a molecule based on its structure, allowing researchers to design molecules with specific properties in mind. For example, a researcher could use CAD to design a molecule with a certain level of solubility, or a specific type of reactivity.

CAD tools can also be used to optimize the synthesis of existing molecules, by predicting the most efficient route for their production. This can be done by analyzing the reaction pathways involved in the synthesis and optimizing them to reduce waste and increase yield.

The use of CAD in organic synthesis has opened up new avenues of research and innovation, allowing researchers to create new molecules and materials with specific properties in mind. The power of CAD is only continuing to grow, as new algorithms and techniques are developed to further enhance its capabilities. As the field of organic synthesis continues to evolve, it is clear that CAD will play a pivotal role in shaping its future.

### **3.4 THE ROLE OF ARTIFICIAL INTELLIGENCE IN ORGANIC SYNTHESIS**

Artificial intelligence (AI) is revolutionizing the way we approach organic synthesis. In the past, chemists would rely on their experience and intuition to plan and execute complex organic syntheses. However, with the advent of

AI, chemists now have powerful tools at their disposal to help them plan and execute syntheses more efficiently and effectively.

One of the key advantages of AI is its ability to analyze vast reaction data and categorize outlines and tendencies, that may not be instantly seeming to human chemists. This can help chemists to design better reactions and optimize reaction conditions, leading to improved yields and fewer side reactions.

Another important application of AI in organic synthesis is in the design of novel molecules. By analyzing large databases of chemical structures and properties, AI algorithms can help chemists to identify promising candidates for new drugs, materials, and other applications. This can save significant time and resources compared to traditional trial-and-error approaches.

Overall, the role of AI in organic synthesis is rapidly evolving, and it is likely to play progressively significant role in the future. As more powerful computational tools and algorithms become available, chemists will be able to design and execute complex syntheses with greater speed, efficiency, and accuracy than ever before.

### **3.5 MACHINE LEARNING TECHNIQUES FOR PREDICTING CHEMICAL REACTIVITY AND REACTANT COMPATIBILITY**

One of the most exciting developments in the area of organic synthesis is the application of machine learning techniques for predicting chemical reactivity and reactant compatibility. These tools are revolutionizing the way chemists plan and execute their syntheses, allowing them to make faster, more informed decisions about which reactions to run and which compounds to use.

Machine learning plays a significant role in organic synthesis. It has revolutionized the way chemists approach the design and development of new compounds by making it easier to predict the outcomes of chemical reactions.

Traditionally, chemists have relied on their experience and intuition to design new compounds, which is a time-consuming and often trial-and-error process. Machine learning algorithms, on the other hand, can analyze vast amounts of data and identify patterns that are not easily detected by humans.

One example of how machine learning is being used in organic synthesis is in the prediction of reaction outcomes. By inputting data on the starting materials, reagents, and reaction conditions, machine learning algorithms can predict the most likely products of the reaction. This can save chemists

a significant amount of time and resources by allowing them to focus on the most promising reactions.

Another way machine learning is being used in organic synthesis is in the design of new compounds. By analyzing the structures and properties of known compounds, machine learning algorithms can generate new compound designs that are likely to have the desired properties. This can be particularly useful in drug discovery, where chemists are looking for compounds that will be effective at treating a particular disease.

Machine learning algorithms work by analyzing large datasets of chemical reactions and identifying patterns and correlations between different parameters, such as reaction conditions, reagent properties, and reaction outcomes. By training these algorithms on vast libraries of chemical data, researchers can create predictive models that can accurately forecast the behavior of new compounds and reactions.

One of the key advantages of these models is that they can identify subtle relationships between different chemical properties that might be difficult for a human chemist to discern. For example, a machine learning algorithm might be able to envisage how a precise compound will react with a certain reagent based on its molecular structure, electronic properties, and previous reaction history, even if these factors are not immediately obvious to a human observer.

As these tools continue to evolve and become more sophisticated, they are likely to play progressively significant role in the area of organic synthesis, enabling chemists to tackle complex synthetic challenges with greater speed and efficiency than ever before.

Overall, the role of machine learning in organic synthesis is becoming increasingly important as chemists look for ways to design new compounds more efficiently and accurately. By combining the power of computational tools with human expertise, we are poised to make significant breakthroughs in the field of organic synthesis.

### **3.6 COMPUTER-ASSISTED RETROSYNTHESIS PLANNING: CURRENT STATE OF THE ART**

Computational tools have revolutionized the way organic synthesis is approached and executed. The success of these tools can be seen in numerous case studies, where they have enabled researchers to achieve their goals more efficiently and effectively than ever before.

One such example is the work of Barbara and colleagues [13], who employed computational tools to design a novel synthesis route for a complex natural product. Using a combination of density functional theory calculations and molecular docking simulations, they were able to identify a key intermediate and optimize reaction conditions to achieve a high-yielding and stereoselective synthesis in just six steps.

Another successful example is the work of Wang and coworkers [14], who utilized machine learning algorithms to analyze reaction data and predict optimal reaction conditions for a range of reactions. This approach allowed them to rapidly optimize reaction conditions and achieve high yields of target products, while minimizing waste and reducing the need for trial-and-error experimentation.

These case studies demonstrate the vast potential of computational tools in organic synthesis, and the exciting possibilities for future research. With the continued development and refinement of these tools, we can expect to see even more groundbreaking advances in the field of organic synthesis.

Computer-assisted retrosynthesis planning is a powerful tool that has revolutionized the world of organic synthesis. It allows chemists to efficiently design synthetic pathways to complex molecules, significantly reducing the time and cost required for synthesis. The present state of the creative skill in computer-assisted retrosynthesis planning involves the use of sophisticated algorithms and machine learning techniques to predict viable synthetic routes.

One of the key advantages of computer-assisted retrosynthesis planning is that it allows chemists to consider a much wider range of potential synthetic pathways than would otherwise be possible. By inputting the target molecule into a computer program, the chemist can rapidly explore a vast number of possible synthetic routes and evaluate their feasibility based on a range of criteria, such as cost, efficiency, and environmental impact.

Another important aspect of current computer-assisted retrosynthesis planning is the integration of machine learning techniques. By training machine learning models on large databases of known reactions and synthesis pathways, these tools can predict the most likely synthetic pathways for a given target molecule with a high degree of accuracy. This can significantly reduce the time and cost required for synthetic planning, and also enable more efficient use of resources in the laboratory.

AIPHOS scheme is prediction focused on the reaction generator of the system for organic synthesis planning, design, and reaction [15].

Overall, the current state of the art in computer-assisted retrosynthesis planning represents a main advance in the area of organic synthesis, with the

potential to greatly accelerate the discovery and development of new drugs, materials, and other chemicals. As these tools continue to evolve, they are probable to show an increasingly important part in the chemical industry and academic research alike [16].

### 3.7 SUCCESSFUL APPLICATIONS OF COMPUTATIONAL TECHNIQUES FOR SYNTHESIS PLANNING

Computational tools and techniques for synthesis planning have revolutionized the area of organic synthesis Table 3.1, leading to the discovery of new molecules and compounds that were previously impossible to achieve through traditional methods. Here are some examples of successful applications of these techniques [17].

**TABLE 3.1** Computer-Assisted Organic Synthesis (CAOS)

SL. No.	Computer Software	Concept
1.	Spaya	Retrosynthetic analysis
2.	IBM Rxn	Chemical procedure
3.	AiZynthFinder	Retrosynthesis planning and predict synthesis routes
4.	Manifold	Molecule searching
5.	Organic synthesis exploration tool	Open-source software
6.	SynGen	modest organic synthetic routes
7.	SYLVIA	Organic structure
8.	Chem Planner	Synthetic routes
9.	ICSYNTH	Synthesis paths for target molecules
10.	Synthia (Chematica)	Synthesis paths
11.	ASKCOS	Planning for synthesis
12.	LHASA	Proprietary software
13.	CHIRON	Proprietary software
14.	WODCA	Proprietary software

Source: Adapted from Ref. [18].

#### 3.7.1 TOTAL SYNTHESIS

This molecule was synthesized by a team of researchers who used computational methods to plan the synthesis route. The team used a retrosynthetic analysis to identify the key intermediates and the final target molecule.

The synthesis was successful, and the final product was produced with a high yield.

### **3.7.2 SYNTHESIS OF A NATURAL PRODUCT**

A team of researchers used computational tools to plan the synthesis of a natural product, which had previously been synthesized using traditional methods with low yields. The researchers used a combination of computational methods and experimental data to design a new synthesis route, which resulted in an increased yield of the final product.

### **3.7.3 SYNTHESIS OF A NEW CLASS OF COMPOUNDS**

Computational methods have also been used to design the synthesis of new programs of compounds. Researchers used computational methods to design the synthesis of a new class of compounds called iminosugars, which have potential applications in the treatment of viral infections.

Overall, the successful applications of computational tools and techniques for synthesis planning have led to new discoveries and advancements in the field of organic synthesis. These techniques have enabled researchers to design and synthesize complex molecules with higher yields and greater efficiency, ultimately leading to new drugs, materials, and technologies.

## **3.8 LIMITATIONS AND FUTURE DIRECTIONS OF COMPUTATIONAL ORGANIC SYNTHESIS**

While computational organic synthesis has made significant strides in recent years, there are still limitations to what it can achieve. One of the important restrictions is the absence of availability of detailed reaction information in some cases. In addition, the lack of accuracy of some computational methods can lead to incorrect predictions or incomplete reaction pathways.

While computational tools have certainly revolutionized the field of organic synthesis, there are still some limitations to their use. One major challenge is the accuracy of the models used in these tools. While it can provide a valuable starting point for designing new molecules and predicting their properties, the models themselves are only as good as the data used to

create them. This means that there is often a trade-off between accuracy and computational efficiency.

Another challenge is the complexity of the molecules being designed. As the size and complexity of the molecule increases, so does the computational cost of modeling it. This can make it difficult to design molecules with a high degree of accuracy and can limit the scope of what can be achieved using these tools.

Finally, there are limitations to the types of reactions that can be modeled using computational tools. While many common reactions can be accurately modeled, there are still some that are too complex or poorly understood to be modeled effectively.

Despite these limitations, computational tools are still a valuable asset in the field of organic synthesis. By providing a way to design and optimize new molecules quickly and efficiently, they have the potential to significantly accelerate the drug discovery process and unlock new treatments for a wide range of diseases.

The future of organic synthesis is incredibly bright with the development of computational tools for chemists. These tools allow researchers to predict and optimize chemical reactions, saving time and resources while also reducing the environmental impact of the synthesis process.

Computational tools can accurately predict reaction outcomes, suggesting alternative reaction pathways to achieve a desired product with higher yield or selectivity. This means that chemists can quickly evaluate a range of possibilities and optimize their experiments before actually conducting them in the lab.

Furthermore, computational tools can also predict the properties of the compounds synthesized, including their reactivity, stability, and toxicity. This information can help chemists to design safer and more eco-friendly synthesis pathways, reducing the negative impact on the environment and human health.

With the power of computational tools, organic synthesis can be revolutionized, making it faster, more efficient, and more sustainable. As these tools continue to advance, we will likely see even more breakthroughs in the field of organic synthesis, leading to new discoveries, improved processes, and a better understanding of chemical reactions.

Future directions for computational organic synthesis include the development of more accurate and efficient computational tools and techniques. This can involve the integration of different computational methods and the exploration of new algorithms that can improve accuracy and reliability.

Another future direction is the incorporation of artificial intelligence and machine learning into computational organic synthesis. This can involve the use of data-driven approaches to identify patterns and trends in reaction data, which can help to develop more accurate prediction models.

Finally, there is a need for the development of more user-friendly computational tools and platforms that can be used by both experts and non-experts in the area of organic synthesis. This can involve the creation of more intuitive graphical user interfaces, as well as the development of online resources and databases that can be accessed by anyone with an internet connection.

Overall, the area of computational organic synthesis is constantly evolving and improving, and there is no doubt that it will remain to show progressively important part in the planning and design of organic synthesis reactions in the future.

### **3.9 CONCLUSIONS AND IMPLICATIONS FOR THE UPCOMING OF ORGANIC SYNTHESIS AND DRUG DISCOVERY**

In conclusion, the area of organic synthesis and drug discovery has been revolutionized by the practice of computational devices and techniques for planning. The ability to predict and design chemical reactions using computational models has significantly reduced the time and resources required for drug discovery and development. This has enabled scientists to focus more on the discovery of new drugs and less on the trial-and-error process.

Moreover, the use of machine learning algorithms has made it possible to predict the properties of new compounds before they are synthesized, making the process more efficient and cost-effective. With the availability of large datasets and the increasing computing power of modern computers, the potential for using computational tools and techniques in organic synthesis and drug discovery is limitless.

The implications of these advancements are far-reaching. The practice of computational representations to design new drugs has the potential to significantly reduce the time and cost involved in bringing new medicines to market. It also has the potential to reduce the use of animal testing in drug development, as many of the experiments can be conducted virtually.

The use of computational tools has revolutionized the field of organic synthesis. With the ability to predict reaction outcomes, optimize reaction conditions, and design novel molecules, computational tools have the potential to greatly impact the future of organic chemistry.



One of the biggest advantages of these tools is their ability to save time and resources. By predicting what reactions will work best and optimizing reaction conditions before ever stepping foot in the lab, chemists can avoid wasting time on unsuccessful experiments. This also saves money on expensive reagents and equipment.

Another potential impact of computational tools is their ability to design new molecules with specific properties. By inputting desired properties into the program, chemists can generate a list of potential molecules that meet those criteria. This opens up a world of possibilities for drug design, materials science, and more.

Overall, the use of computational tools in organic synthesis has already made a significant impact and has the potential to continue to revolutionize the field in the years to come. As technology advances and these tools become more sophisticated, the possibilities for what we can achieve in organic synthesis will only continue to grow.

In the future, we can expect to see more advancements in the field of organic synthesis, with an increasing focus on the use of computational tools and techniques. As the technology continues to evolve, we can expect to see even greater efficiencies in the drug discovery process, leading to the development of new and more effective medicines.

### **3.10 FINAL ASSESSMENTS: THE PROMISE AND POTENTIAL OF COMPUTATIONAL TOOLS FOR REVOLUTIONIZING ORGANIC SYNTHESIS**

In conclusion, the promise and potential of computational tools for revolutionizing organic synthesis cannot be overstated. These tools are already making significant strides in the field, from predicting reaction outcomes to identifying new synthetic routes and optimizing reaction conditions. As technology continues to advance, the possibilities for what can be achieved through computational tools and techniques will only continue to grow.

However, it is important to note that these tools should not be seen as a replacement for traditional organic synthesis methods. Rather, they should be viewed as complementary tools that can enhance and improve the overall process.

Furthermore, the successful implementation of these tools and techniques requires a collaborative effort between chemists, computer scientists, and other experts in related fields. Together, they can work to develop new and

innovative approaches to organic synthesis that will drive the field forward and lead to new discoveries and breakthroughs.

The development of these tools has significantly reduced the time and resources required for synthesizing new compounds, making organic synthesis more efficient and accessible than ever before. We are excited to see where this technology will go in the future and how it will continue to shape the field of organic synthesis and keep an eye out for more exciting developments in this field.

## KEYWORDS

- **artificial intelligence**
- **computational organic synthesis**
- **computational tools**
- **computer-aided synthesis planning**
- **drug discovery**
- **machine learning**
- **organic synthesis**

## REFERENCES

1. Robert, J. D., & Caserio, M. C. (1977). *Basic Principles of Organic Chemistry* (2nd ed.). W. A. Benjamin, Inc.
2. Norman, R. O. C., & Coxon, J. M. (1968). *Principles of Organic Synthesis* (3rd ed.). CRC Press.
3. Ihlenfeldt, W. D., & Gasteiger, J. (1995). Computer-assisted planning of organic syntheses: The second generation of programs. *Angewandte Chemie International Edition in English*, 34(22), 2613–2633.
4. Milo, A. (2018). The art of organic synthesis in the age of automation. *Israeli Journal of Chemistry*, 58, 131–135.
5. Taber, D. F., & Lambert, T. (2017). *Organic Synthesis: State of the Art, 2013–2015*. Oxford University Press.
6. Smith, M. B. (2017). *Organic Synthesis* (4th ed.). Academic Press.
7. Chan, K. S., & Tan, J. (2018). Chapter 9: Planning for organic synthesis. In *[Title of the book]* (pp. 308). [Publisher].
8. Starkey, L. S. (2018). *Introduction to Strategies for Organic Synthesis* (2nd ed.). Wiley.
9. Curran, D. P. (1998). Strategy-level separations in organic synthesis: From planning to practice. *Angewandte Chemie International Edition*, 37(8), 1174–1196.

10. Carruthers, W., & Coldham, I. (2004). *Modern Methods of Organic Synthesis* (4th ed.). Cambridge University Press.
11. Plehiers, P. P., Coley, C. W., Gao, H., Vermeire, F. H., Dobbelaere, M. R., Stevens, C. V., Geem, K. M. V., & Green, W. H. (2020). Frontiers in chemical engineering. *Frontiers in Chemical Engineering*, 2(6), 1–19.
12. Corey, E. J., Long, A. K., & Rubenstein, S. D. (1985). Computer-assisted analysis in organic synthesis. *Science*, 228(4702), 408–418.
13. Mikulak-Klucznik, B., Gołębiowska, P., Bayly, A. A., Popik, O., Klucznik, T., Szymkuc, S., Gajewska, E. P., Dittwald, P., Staszewska-Krajewska, O., Beker, W., Badowski, T., Scheidt, K. A., Molga, K., Mlynarski, J., Mrksich, M., & Grzybowski, B. A. (2020). *Nature*, 588(7836), 83–88.
14. Wang, G., Wu, X., Xin, B., Gu, X., Wang, G., Zhang, Y., Zhao, J., Cheng, X., Chen, C., & Ma, J. (2023). *Molecules*, 28(5), 2232.
15. Funatsu, K., & Sasaki, S. I. (1988). Tetrahedron computer methodology. *Tetrahedron Computer Methodology*, 1(1), 27–37.
16. Szymkuc, S., Gajewska, E. P., Klucznik, T., Molga, K., Dittwald, P., Startek, M., Bajczyk, M., & Grzybowski, B. A. (2016). *Angewandte Chemie International Edition*, 55(35), 5904–5937.
17. Todd, M. H. (2005). *Chemical Society Reviews*, 34(3), 247–266.
18. List of computer-assisted organic synthesis software. (n.d.). In *Wikipedia*. Retrieved from: [https://en.wikipedia.org/wiki/List\\_of\\_computer-assisted\\_organic\\_synthesis\\_software](https://en.wikipedia.org/wiki/List_of_computer-assisted_organic_synthesis_software) (accessed on 25 July 2024).

## CHAPTER 4

---

# Patenting Artificial Intelligence-Based Technologies in Chemical and Pharmaceutical Sciences

ASHA HOLE,<sup>1</sup> SHASHIKANT BHANDARI,<sup>1</sup> SAGAR BIRAJDAR,<sup>1</sup>  
SANDIP SURVE,<sup>1</sup> and ANIKET SARKATE<sup>2</sup>

*<sup>1</sup>Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy,  
Pune, Maharashtra, India*

*<sup>2</sup>Department of Chemical Technology, Dr. Babasaheb Ambedkar  
Marathwada University Campus, Aurangabad, Maharashtra, India*

---

### ABSTRACT

Artificial intelligence (AI) is the term used to describe the creation of computer systems that are capable of undertaking tasks that traditionally require human intelligence. AI aims to build machines that can learn and reason similarly to people, and that are capable of adapting to new circumstances and tasks. AI is a quickly developing field that is revolutionizing a variety of sectors, including chemistry and pharmaceutical sciences. To cover the broad range of application of AI in the field of chemistry and pharmaceutical sciences, research in AI is rapidly evolving. In order to obtain a competitive advantage, businesses, and individuals are actively working to patent their AI-based technologies and to safeguard the intellectual property associated with inventions and advancements in the field of AI, patents has filed. This chapter focuses on the patent filing trends on the application of artificial intelligence, the

year wise review of patent filing trend, the most promising applicants and technology, obstacles for protection AI associated inventions and challenges, the might present incapability of current patent laws and the future of AI associated patents.

## **4.1 INTRODUCTION**

Artificial intelligence (AI) aims to develop intelligent machines that can perform tasks that usually call for human intelligence. Because it has the potential to transform many aspects of our lives and industries, AI is important. Because it has capabilities to automate tedious and repetitive tasks, which increase efficiency and productivity, analyze data patterns, spot trends, and extract insightful information, AI has the potential to transform a variety of facets of our lives and industries. These capabilities also enable businesses and organizations to make data-driven decisions and gain a competitive edge. AI has advanced so quickly because of these incredible defining features. Hence, research in the field of AI is inevitable to increase the application in various fields. To achieve competitive growth it is always advisable to protect the research through Patents. This chapter focuses on the patenting the AI specifically in the area of Chemistry and Pharmaceuticals [1]. For inventions to be patentable it's important to prove the patentability of the invention. AI is considered as branch of computer sciences hence Inventions involving AI is considered as computer related invention (CII). Patent offices around the world have varying policies regarding computer-implemented innovations. Because of AI's capabilities in a number of technical fields, patent applications utilizing AI are becoming more prevalent.

The involvement of AI in Chemistry and the healthcare field in increasing drastically from the last decade. The inventions involving AI have been filed by many inventors in various jurisdictions. This chapter provides an overview of patent applications related to chemistry and pharmaceutical sciences involving AI, the most prominent technology is focused in the field that diversifies through AI. In the last decades debate has been carried out for patentability of AI involving invention. It is important to understand the technology trend that helps to forecast the advancement in the domains of chemistry and pharmaceutical sector. The current status of AI-associated inventions and forecasting the research strategies for inventions in pharmaceutical and chemistry associated with AI.

## **4.2 ARTIFICIAL INTELLIGENCE: CHEMICAL AND PHARMACEUTICAL SCIENCE**

Due to the potential to revolutionize several facets of drug discovery, development, and production, AI in chemical and pharmaceutical research has witnessed an increase in patent applications [1]. Researchers are looking for ways to protect AI-assisted creations through a variety of intellectual property rights, with patents being the most common [2]. While protecting AI-assisted inventions it is important to consider a few factors such as What was the invention? Clearly state the particular AI-based technology or approach at the invention's foundation, perform a thorough prior art search to make sure the innovation is new and not apparent in light of already available technology. To show an innovative step, highlight the special qualities, technological developments, or novel uses of AI in the chemical and pharmaceutical industries [3]. Draught a thorough patent specification that describes the invention's technical details and possible uses. Describe how AI technology helps drug development, optimizes chemical or pharmaceutical processes, or offers a fresh approach to a particular issue in the industry. Describe the data sources, preparation procedures, and data analysis techniques used if the AI innovation makes use of data. Clearly describe the benefits or technical impacts that AI technology may provide to the fields of chemical and pharmaceutical science. Provide experimental facts or examples to support your claims about the effectiveness, performance, or superiority of your AI-based technology over current techniques wherever you can. Experimental data can support the claimed technological effects and strengthen the patent application. Create claims that outline the boundaries of the invention and address any unique or specialized characteristics of AI technology [4]. The innovation, inventive step, and practical applications of the AI invention in the chemical and pharmaceutical domains should all be captured in the claims, which should be properly worded.

**Collaboration and legal knowledge:** It is desirable to work with specialists in both AI technology and patent law given the complexity of patent law and the relationship between it and AI. Due to the recent advancements in computer related inventions and their examination criteria. Its more feasible to file the patent application related AI. Still, the patent approval mainly depends upon the countries own patent laws. This chapter will cover the patenting of artificial intelligence applications in the chemical and pharmaceutical sciences.

#### **4.2.1 IMPACT OF AI IN CHEMICAL AND PHARMACEUTICAL SCIENCES**

Chemical and pharmaceutical sciences have been greatly impacted by AI, which has impacted several fields of study, development, and drug discovery. Here are some significant areas where AI is making a difference. Figure 4.1 illustrates the application of AI in the pharmaceutical sector. AI is accelerating the drug discovery [5] process by assisting in the identification of potential drug candidates and optimizing their attributes. AI systems can scan huge chemical databases [6], predict molecular structures, and determine which proteins they are likely to bind to. By helping researchers focus their search for prospective drug candidates, this helps save time and money. Artificial intelligence (AI) methods like molecular docking and machine learning enable virtual screening and lead optimization. AI systems are able to predict a compound's potential as a medication and offer recommendations. AI systems are able to predict the likelihood that a chemical will turn into a successful medicine and provide recommendations for alterations to enhance its properties. This aids in the optimization of leads and the identification of compounds with excellent potential for development. AI can predict the toxicity and safety profiles of chemical compounds toxicological evaluation and predictive modeling. Machine learning algorithms that have been trained on large datasets can assess potential side effects and help choose compounds with lower risks in order to increase drug safety [7].

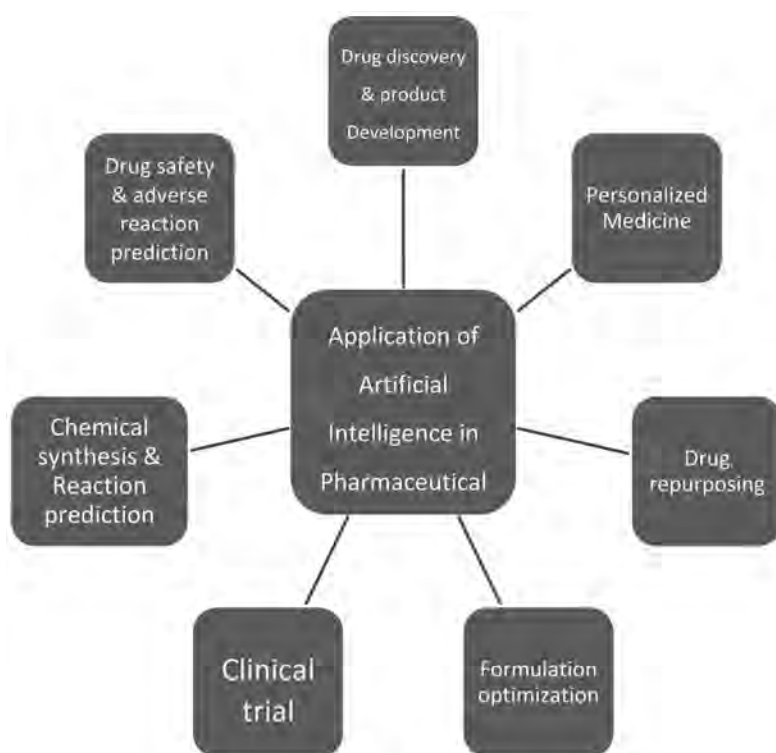
To help generate pharmaceuticals with greater solubility, stability, and bioavailability, AI can assist with formulation development and optimization [8]. Better efficacy and patient compliance are ensured as a result of improving drug delivery. AI can optimize chemical processes by analyzing large datasets and identifying variables that influence reaction results. Waste is reduced, procedures are more effective, and money is saved as a result. AI is helpful in the discovery and development of innovative materials with the required properties.

Utilizing material databases and atomic-level models, AI algorithms can foresee novel materials with specific qualities for use in industries like energy storage, catalysis, and electronics. Thanks to AI, researchers may now identify new research areas and make data-driven decisions.

##### **4.2.1.1 DATA COLLECTION AND ANALYSIS**

We have retrieved patent data from various patent search engine database and exported on May 17, 2023. Every entry indicates a patent family, which is a

collection of connected innovations submitted to one or more patent bodies [9]. Patent numbers, priority application data, titles, application details, designated states, inventors, IPC categorization numbers, patent details, cited patents, and referenced publications are important components of records. All these patent-related data are shown as a single patent. Out of 10,558 entries, 4,589 remain after duplicates are removed from an individual patent domain [10]. On an Excel sheet, the information on assignees and the number of patents is cleaned, sorted, and rearranged (Table 4.1).



**FIGURE 4.1** Application of AI in the pharmaceutical sector.

### 4.3 PATENTING TRENDS

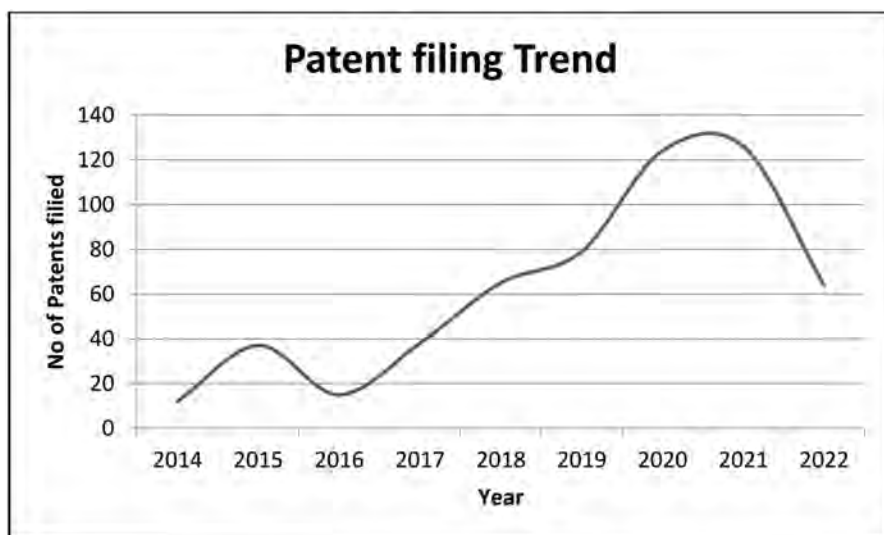
Over the years, there has been an increasing number of patent requests pertaining to AI. With businesses, academic organizations, and individuals all vying for patent protection for their AI-based technologies and algorithms, AI has become a highly competitive field. AI application has



been growing significantly in the chemical and pharmaceutical industries, which has increased the number of patent applications for AI-based technology [11]. It is important to analyze the impact of artificial intelligence and their applications in research and development specifically in chemical and pharmaceutical sector. Figure 4.2 shows the patent filing trend for AI-assisted inventions particularly in Pharmaceuticals. There is an increase in the patent filing by year 2017 and it is continued to be increasing.

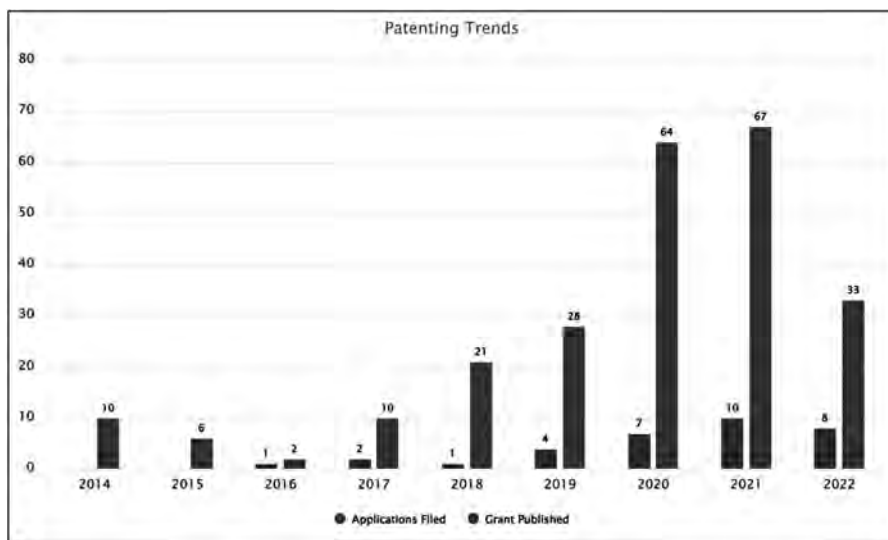
**TABLE 4.1** Search Strategy Used

<b>Search Strategy</b>	TAC: (((“artificial intelligence” OR (Artificial Intelligence Device OR Machine Learning OR Neural Network OR Data Source OR Artificial Intelligence Entity OR AI Solution OR Sensor Data OR Information Processing Apparatus OR Artificial Intelligence-based OR Optimization Data))) AND (Smart Contract OR Intelligence Service)))
<b>Time Period</b>	All years (till 2023)
<b>Date of Retrieval</b>	May 17, 2023
<b>Database Used</b>	USPTO EPO Patentscope



**FIGURE 4.2** Patent application filed yearly covering AI-assisted in chemical and pharmaceutical sector.

It is also important to understand the ratio of filed and granted patenting trend. Figure 4.3 illustrates the patent filed and granted in year from 2014 till 2022.



**FIGURE 4.3** Patent filing trend.

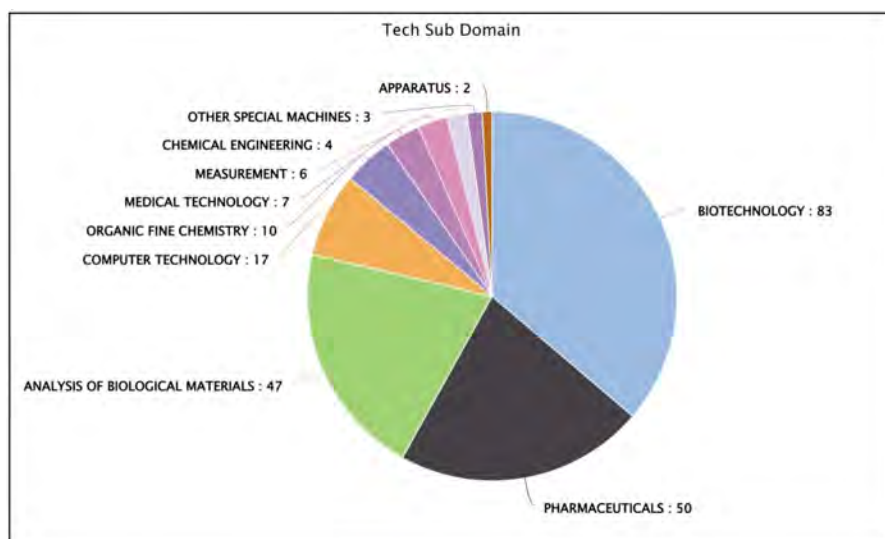
The number of patent applications has been increasing from 2014 till 2021. The application filing to grant ratio illustrates the evolution of patent law in relation to AI applications. The patent laws have been evolved to accept patentability of the inventions in the field of AI [12]. Figure 4.3 specifically highlights the increase in the patent grant of AI applications after 2020 which motivates the inventors working in the field of application of AI for chemical and pharmaceutical sectors.

#### **4.3.1 PATENT FILING ON APPLICATION OF AI IN HEALTH SCIENCES**

One has to understand the patentability criteria raised for AI-assisted invention in chemical and pharmaceutical sector to overcome the patent rejection. Now, we will understand the patent filed in the health science sector covering biotechnology, organic chemistry, medical technology and pharmaceutical.

The biotechnology field has larger number of patents filed. The application of AI in the field of biotechnology includes research related to activate

the immune system based on AI [13] and deep learning technology, Bioanalyte signal amplification and detection with artificial intelligence diagnosis, etc. Pharmaceutical sector has second largest patent filing based on AI which includes AI based stability indicating methods and validation [14], chemical test tube AI cleaning devices [15], self-driving type AI material, preparation method and application in imaging analysis detection and drug controlled release [16], application of AI in detecting all type of cancer [17], thermocycler reaction control [18], approach, and technique for investigating the relationships between pharmacological side effects and therapeutic applications [17] (Figure 4.4 and Table 4.2). The significant AI-assisted patent applications for pharmaceutical products are listed in Table 4.2.



**FIGURE 4.4** Patent filing in the field of AI in health science.

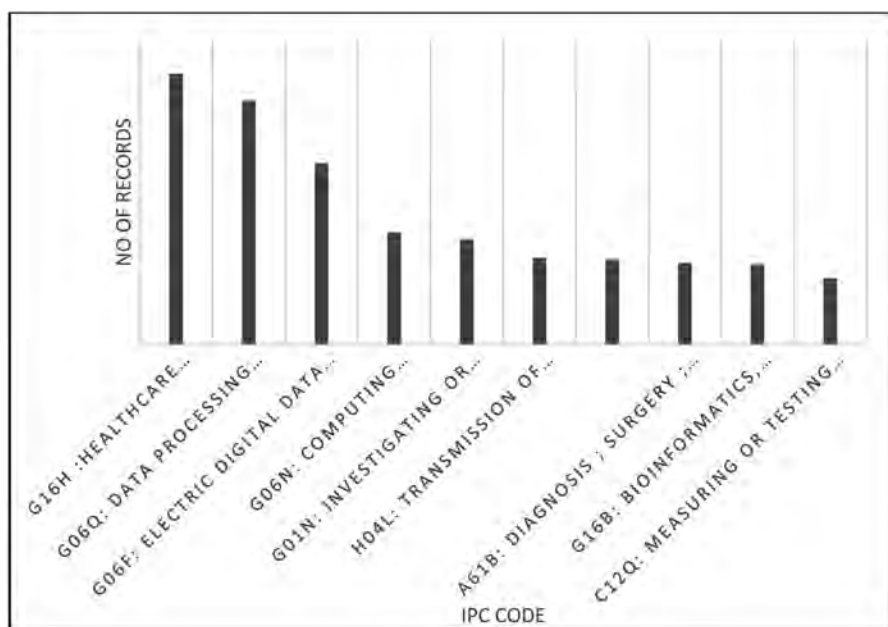
### **4.3.2 PATENT FILING TREND FOR VARIOUS TECHNOLOGIES IN PHARMACEUTICALS**

Each patent application that has been submitted has been categorized according to the technology that the relevant patent covers. The IPC codes used to classify the technology associated with patents. Figure 4.5 illustrates the frequency of IPC codes for AI-assisted patents in pharmaceuticals. IPC

**TABLE 4.2** AI-Assisted Patent Applications for Pharmaceutical Products

SL. No.	Patent Application No.	Owner	Technology Focus
1.	WO2021234522A1	IBM Corp. (US)	Filtering artificial intelligence designed molecules for laboratory testing [20].
2.	WO2021262857A1	Pfizer Inc.	Computerized decision support tool and medical device for scratch detection and flare prediction [21].
3.	WO2021080999A1	Sanofi Pasteur Inc.	Prediction systems and techniques for biological reactions [22].
4.	WO2020US56507	Sanofi Pasteur Inc.	Systems and methods for designing vaccines [23].
5.	WO2020/058174	DeepMind	The alpha fold challenges the folding of proteins and enables protein form to be predicted from the sequence of amino acids [24].
6.	EP2241335B1	Biodesix	Assessing how a patient reacts to treatment by using mass spectrometry to analyze the patient's characteristics [25].
7.	EP3140763B	Atomwise	A system for estimating a molecule's affinity for a given target protein [26].
8.	US13971072	International Business Machines Corp.	Method and technique for investigating the relationships between pharmacological side effects and therapeutic uses [17].
9.	WO2021030270A1	Sanofi SA	Predicting patient responses to a chemical substance [28].
10.	WO2020170165A1	Johnson & Johnson Consumer Inc. (US)	New goals and techniques for skin care treatments using artificial intelligence [29].
11.	WO2019231917A1	Takeda Pharma Co. Ltd.	Systems and techniques for automated production and supply chain tracking and optimization based on the effects of post-approval alterations [30].

graph shows the novel patterns [31]. Most frequently used IPC codes for AI-related patents in Pharmaceuticals are G16H, G06Q, G06F, G06N, G01N, H04L, A61K, A61B, G16B, and C12Q. G16H is high used IPC code which deals with the healthcare informatics specifically for medical and healthcare data, G06Q is related to information and communication technology for managerial used. G06Q covers the patents based on AI for personalized medicine and clinical trials. G01N is also frequently used for the technology covering patents using AI for product and biomaterial analysis (Table 4.3).



**FIGURE 4.5** Frequency of IPC codes for AI-related patents in pharmaceuticals.

### 4.3.3 APPLICANT'S CATEGORY-BASED PATENT FILING TREND

According to the WIPO report [32], the most prominent patent assignees are IBM, Microsoft, Toshiba, Samsung. The top patent assignees, however, tend to be profound experts in each single subject of AI when we examine more deeply within each one. Figure 4.6 illustrates the category of applicants and their patent filing in AI-assisted technology for pharmaceutical and chemical sector. The applicant dealing with pharmaceutical preparation and products are the highest category of applicants filing the AI-assisted patents. The medical and dental.

TABLE 4.3 Description of IPC Codes

IPC Code	Description
G16H	Informatics in healthcare, i.e., Specially designed information and communication technology (ICT) for implementing or analyzing medical or healthcare data.
G06Q	Systems or procedures specifically designed for administrative, commercial, financial, management, or supervisory objectives, information, and communication technology (ICT) specifically fitted for these purposes.
G06F	Analyzing electrical digital information.
G06N	Arrangements for information technology based on certain simulations.
G01N	Examining or assessing materials by figuring out their chemistry or physics.
H04L	Electronic data transfer, such as telegraph communication.
A61K	The plans for health, dental, or personal care needs.
A61B	Diagnosis; surgery; identification.
G16B	Information and communication technology (ICT) tailored specifically for the processing of genome or protein-related information in computation molecular biology is known as bioinformatics.
C12Q	Preparing such compositions; measuring or evaluating methods Utilizing enzymes, amino acids, or microbes; components or test papers for such procedures; control that is condition-responsive in microbial or enzymatic processes.

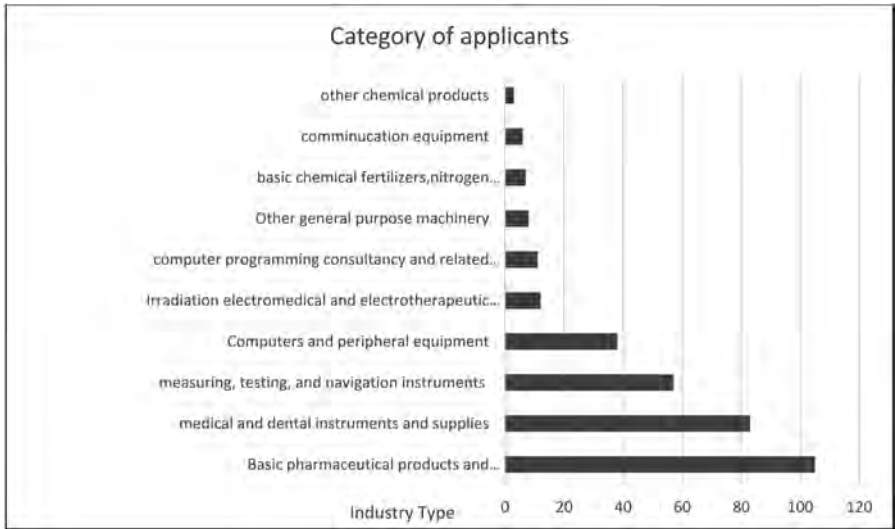


FIGURE 4.6 AI-assisted patents filed by applicant with specific category.

The application of AI towards the pharmaceutical and biomedical industries has advanced from science fiction to reality during the past several

years. Pharma and biotech firms are increasingly using more efficient, automated processes that combine information-driven options and make use of statistical analysis technologies. Strategic patenting by pharmaceutical companies – should competition law intervene? [33]. By establishing and maintaining a robust patent portfolio, pharmaceutical businesses may be able to recuperate the expenditure required to find, develop, and get approval from the market for an innovative drug product. Companies should build patent awareness across the organization and methods for finding patentable ideas in order to maximize patent protection and value. They should also prepare applications with a view upon monetization. Additionally, businesses should diligently oversee their patent portfolios, concentrating on highly valuable inventions [34]. The research in the area of AI is more computational and engineered, hence it is important to understand the type of applicant working in the field. Figure 4.6 illustrates the difference category of applicants doing AI-assisted research in pharmaceutical field. The highest filing was done by applicants dealing with the basic pharmaceutical sectors. Table 4.4 listed the category of applicants and filed Patent related to AI.

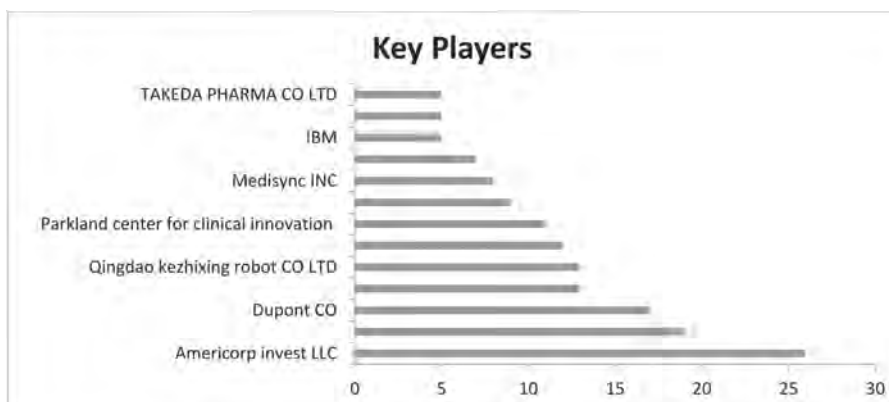
**TABLE 4.4** Category of Applicant Filed Patent in AI

Industry	No. of. Records
Basic pharmaceutical products and pharmaceutical preparation.	105
Supplies and equipment for dentistry and medicine.	83
Measuring, testing, and navigation instruments.	57
Computers and peripheral equipment.	38
Equipment for electro medicine and electrotherapy is irradiated.	12
Consulting in software development and related tasks.	11
Other general purpose machinery.	8
Fundamental kinds of artificial rubber, plastics, nitrogen compounds, and basic chemical fertilizers.	7
Communication equipment.	6
Other chemical products.	3

#### 4.3.3.1 CURRENT ASSIGNEE

Figure 4.7 illustrates the key players in the field. The many pharmaceutical companies are working towards innovating the AI application in the

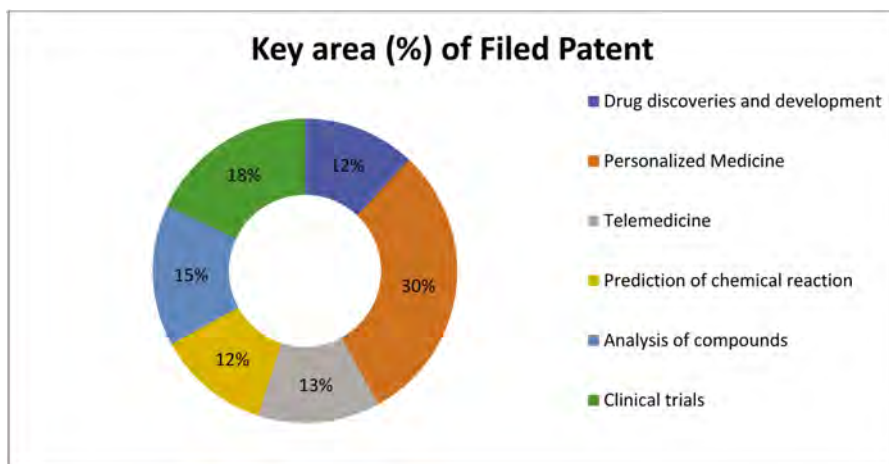
pharmaceutical and chemical field. The Takeda, Medi sync, DuPont, Sanofi are few key players form pharmaceuticals.



**FIGURE 4.7** Key players in the field.

#### 4.4 PATENT FILING IN THE KEY AREA OF PHARMACEUTICALS

In drug discovery and development where AI and pharmaceuticals merge, patent applications are becoming more and more concentrated on a few important areas. The following are a few of the primary areas where patents pertaining to AI are being sought (Figure 4.8).



**FIGURE 4.8** Patent filing in pharmaceutical key areas.



#### **4.4.1 PATENT FILING FOR ARTIFICIAL INTELLIGENCE ASSISTED IN DRUG DISCOVERIES AND DEVELOPMENT**

AI utilization has been growing significantly in the chemical and pharmaceutical industries, which has increased the number of patent applications for AI-based technology and applications [35]. The total patents filed in AI assisted technology in the pharmaceutical sector out of which 12% has been filed for drug discovery and development. Here are some AI-specific areas catering the AI in the chemical and pharmaceutical industries:

##### **4.4.1.1 DRUG TARGET IDENTIFICATION**

To find possible drug targets, AI systems may examine enormous volumes of biological and chemical data. AI can assist researchers in identifying specific molecules or proteins that may be potential targets for pharmacological intervention by analyzing genetic data, protein structures, and disease pathways [36].

##### **4.4.1.2 DRUG DESIGN AND OPTIMIZATION**

AI can help in drug candidate design and optimization. Molecular structures may be analyzed by machine learning algorithms, which can then be used to forecast qualities like solubility, bioavailability, and toxicity. This helps researchers to choose medication options with the best chances of success and optimize drug candidate selection.

##### **4.4.1.3 VIRTUAL SCREENING**

AI systems can uncover prospective therapeutic candidates by virtually scanning huge databases of chemicals. AI can focus the search for interesting compounds by examining chemical characteristics while juxtaposing them with recognized drug-target interactions, which can help in the early phases of drug discovery by saving time and resources [37]. Using predictive analytics, massive datasets may be analyzed to forecast therapeutic effectiveness, safety, and side effects. AI systems may find trends and forecast the success or failure of drug candidates by combining many data sources, including clinical data, genomes, and molecular data. This aids in prioritizing and focusing research efforts.

#### **4.4.1.4 DRUG REPURPOSING**

AI can assist in finding novel therapeutic applications for already-approved medications. Massive datasets, as those from electronic health records, clinical trial data, and scientific literature, can be analyzed by AI algorithms to uncover previously unknown relationships between medications and disorders [38]. This facilitates the repurposing of already-approved medications for use in new indications, speeding up the drug development process.

#### **4.4.1.5 CLINICAL TRIAL OPTIMIZATION**

AI can improve patient selection and clinical trial design [39] Artificial Intelligence (AI) to Improve Clinical Trial Volunteer Prescreening: A Comparative Analysis of the Results of AI-Assisted vs. Standard Procedures in Three Cancer Trials. AI algorithms can pinpoint types of patients that are more likely to react to a certain medicine by examining patient data, such as genomes, demographics, and medical history. Clinical studies may then be conducted more quickly and effectively, for less money, thanks to this.

#### **4.4.1.6 DRUG SAFETY**

AI can assist in predicting probable adverse effects and medication toxicity. AI algorithms can determine the possibility of bad responses by examining chemical structures, biological facts, and previous medication safety data. Researchers can optimize medication safety profiles and reduce drug development-related hazards with the use of this knowledge.

##### **4.4.1.6.1 AI-Driven Drug Development**

AI-driven drug developments have attracted a lot of interest. Investments in AI technology is being made by pharmaceutical corporations and research organizations to hasten the identification of new medication candidates and the optimization of their features. As a result, there has been an increase in patent applications for AI-driven drug discovery techniques such molecular docking, virtual screening, and predictive modeling.

AI is being used in computational chemistry to describe and simulate chemical interactions, forecast the features of molecules, and create new

chemicals. In this field, software tools and AI algorithms that facilitate effective virtual screening, structure-based drug design, and chemical reaction optimization are frequently covered by patents.

Leading pharmaceutical corporations are currently working with or acquired AI technologies are Roche, Pfizer, Merck, AstraZeneca, GSK, Sanofi, AbbVie, Bristol-Myers Squibb, and Johnson & Johnson [40].

#### **4.4.2 PATIENT FILING FOR APPLICATION OF ARTIFICIAL INTELLIGENCE IN PERSONALIZED MEDICINE**

AI can aid in the advancement of personalized medicine and therapy optimization by analyzing patient data such as genetics, medical histories, and clinical records. Machine learning algorithms can be used to anticipate patient outcomes, optimize dose, and choose the best course of treatment. Delivering personalized therapy is made simpler as a result, improving patient care and treatment effectiveness [41]. The use of AI in personalized medicine and their protection through patent is still the matter of debate. Incentivizing investments in their research and development with exclusive rights cannot be challenged if AIs and personalized medicine clearly benefit society. But a deeper look indicates that at least two unique areas require more investigation. What features of AIs and the products that they produce should be protected first and foremost. The second problem is whether we should direct it towards one sort of IP over another, taking into consideration issues with technological development, data quality, and the requirement for disclosure [42]. In personalized medicine, where treatments and healthcare choices are made for specific individuals based on their particular traits and medical data, AI has major uses. Here are some significant applications of AI in personalized medicine:

1. **Precision Diagnostics:** To aid in precise and early illness detection, AI algorithms may examine significant amounts of patient data, particularly medical records, genetic data, imaging data, and sensor profiles. AI has the ability to identify tiny connections and patterns that human physicians would miss, resulting in quicker and more accurate diagnosis [43].
2. **Analytics for Prediction and Risk Assessment:** By fusing information about patients with population-level statistics and academic research, AI can assist in predicting illness development, responses to therapy, and patient outcomes. AI models may find risk elements,

biomarkers, and genetic indicators that contribute to the onset or progression of illnesses by analyzing a variety of datasets. This makes it possible for medical experts to identify the hazards specific to each patient and decide on an early course of therapy.

3. **Treatment Choice and Optimization:** AI can help in making the best treatment decisions for certain patients. AI algorithms can offer insights on which therapies are most likely to be helpful and which may have possible negative effects by examining patient-specific data, including genomic profiles, medical histories, and treatment outcomes. This encourages the planning and improvement of individualized care.
4. **Clinical Decision-Making Systems:** based on artificial intelligence systems for clinical decision-making can help medical professionals make judgements about treatments that are supported by the available research. AI systems are able to offer customized suggestions, dose modifications, and treatment tracking strategies by fusing patient data, medical advice, and academic research. Both the quality of care and patient safety are enhanced as a result.
5. **Remote Monitoring with Wearable Tech:** AI may examine data from wearable tech, such as fitness bands or smartwatches, to assess patient health metrics in real time [44]. AI algorithms can enhance remote patient monitoring by identifying abnormalities or shifts in vital signs, which enables preemptive interventions and individualized treatment.
6. **Engagement with Patients and Education:** Artificial intelligence-powered chatbots or virtual personal assistants may communicate with patients, delivering individualized health information [45], responding to inquiries, and providing assistance. These AI-powered platforms can support patients in actively participating in their own healthcare by helping them better grasp their diseases, treatment options, and self-management techniques.

#### **4.4.3 PATIENT FILING FOR APPLICATION OF ARTIFICIAL INTELLIGENCE IN TELEMEDICINE**

AI is essential to the development of telemedicine, which mixes healthcare technology in order to deliver remote medical treatments. It's important to remain in mind that these are widespread uses of AI in telemedicine, and the patents submitted may vary based on specific advancements and tactics. Several important uses for artificial intelligence in telemedicine are listed below:

1. **Medical Image Analysis:** To help with diagnosis and abnormality detection, AI systems can examine medical images including CT scanning, X-rays, and MRIs. AI algorithms can spot patterns and abnormalities in photos, assisting medical professionals in diagnosing and remotely interpreting data. Remote patient monitoring is possible using wearable sensors and AI-enabled devices, which can track indicators of health, levels of exertion, and medication compliance [46]. Real-time data analysis by AI algorithms can spot trends and abnormalities, informing healthcare professionals of possible problems. This makes it possible to continuously check on patients' health from a distance.
2. **Decision Support Systems:** By examining information about patients, medical records, and pertinent guidelines, AI may assist medical professionals in making decisions during teleconsultations. Based on a patient's health and medical history, AI algorithms can prescribe treatments, recommend diagnostic testing, and help with personalized treatment planning. As the technology advances, we might expect to witness additional innovations and patent activity in the field of AI in telemedicine.

#### **4.4.4 PATIENT FILING FOR APPLICATION OF ARTIFICIAL INTELLIGENCE IN SYNTHESIS OF COMPOUNDS/PREDICTION OF CHEMICAL REACTION**

Chemical synthesis and reaction prediction have benefited greatly from advances in artificial intelligence (AI). Here are some applications of AI in various fields.

##### **4.4.4.1 REACTION PREDICTION**

By examining a sizable quantity of reaction data from databases and scientific literature, AI systems can forecast the results of chemical reactions [47]. AI models may anticipate reaction results, such as product generation and reaction conditions, by learning routines and patterns of behavior. This aids scientists in creating synthetic pathways and investigating novel chemical processes. This aids scientists in creating synthetic pathways and investigating novel chemical processes.

#### **4.4.5 PATIENT FILING FOR APPLICATION OF ARTIFICIAL INTELLIGENCE IN ANALYSIS OF COMPOUNDS**

By increasing the precision, efficiency, and ease of the identification process, artificial intelligence (AI) has had a considerable influence on the sector of pharmaceutical identification. Here are a few significant uses of AI in drug identification:

1. **Compound Identification:** By examining multiple sources of data, like molecular structures [19], spectral information, or biological activity profiles, AI systems can help in determining the identity of pharmaceutical compounds. AI algorithms can classify unidentified substances by comparing experimental data to reference databases, forecasting compound features, and producing precise identification findings [27]. AI may assist in confirming the legitimacy and caliber of medicinal items. AI systems may identify fake or inferior pharmaceuticals by looking at product photos, packaging details, or distinctive identifiers. By doing this, you can protect patients and stop the spread of bogus pharmaceuticals.
2. **Quality Analysis:** Artificial intelligence (AI) has the potential to increase the effectiveness of quality control procedures used in the production of pharmaceuticals. Large amounts of data created during manufacturing operations, such as the caliber of the raw materials used, the production parameters, and the outcomes of product testing, may be analyzed by AI algorithms. For constant product quality, AI can track variations, spot possible problems, and streamline production methods.

#### **4.4.6 PATIENT FILING FOR APPLICATION OF ARTIFICIAL INTELLIGENCE IN CLINICAL TRIALS**

Increasingly, clinical trials are being optimized and improved in many ways thanks to artificial intelligence (AI). The following are some significant uses of AI within clinical trials:

1. **Recruitment and Eligibility of Patients:** AI algorithms may review medical records for patients, genetic information, and other pertinent data to find possible trial participants who satisfy particular requirements. By connecting qualified patients with the right studies, AI can speed up the patient's enrollment procedure and conserve time and money.

2. **Trial Design and Protocol Optimization:** To help with trial design and protocol optimization, AI can analyze sizable datasets, such as medical records, literature, and genetic information. In order to increase the statistical ability and effectiveness of clinical trials, AI algorithms may recognize patient groupings, choose the proper sample sizes, and optimize trial settings.
3. **Data Analysis:** Data collected during clinical trials may be tracked and analyzed in real-time using platforms with AI. AI systems have the ability to recognize patterns, trends, and potentially harmful occurrences, giving trial investigators timely information. As a result, decisions may be made proactively, and trial methods can be changed as necessary.

AI can help in the early identification and monitoring of negative outcomes during clinical trials. AI algorithms can identify probable bad events, enabling immediate intervention and enhanced patient safety. They do this by analyzing patient data, electronic medical records, and other pertinent information.
4. **Forecasting and Outcome Modeling:** Based on a variety of variables, including patient profiles, treatment plans, and biomarker data, AI models may forecast trial outcomes. AI algorithms can help with healthcare decision-making by anticipating patient reactions, treatment effectiveness, and long-term results using machine learning approaches.
5. **Patient Participation and Remote Monitoring:** During clinical trials, AI-powered technologies can improve involvement of patients and enable remote monitoring. AI-powered chatbots, wearable tech, and mobile apps may give trial participants immediate input, reminders, and assistance. This makes it easier to gather data remotely and increases patient compliance.

In spite of the fact that AI shows a lot of potential in this area, it's vital to remember that human oversight, moral concerns, and legal compliance are essential to guarantee the safety of patients and the validity of study results.

#### **4.5 ADVANTAGES AND LIMITATION OF AI TECHNOLOGY IN CHEMICAL AND PHARMACEUTICAL SCIENCES**

The pharmaceutical and biotechnology sectors are being completely transformed by AI and machine learning ("ML"). The most prevalent application

of AI and ML in these domains may be drug development, although Figure 4.1 illustrates many additional uses. The development of pharmaceutical formulations, the prediction of protein structures, the planning and analysis of clinical trials, the acceleration of production, and improved quality control are all being accelerated by AI and ML.

Therefore, it should come as no surprise that businesses in the pharmaceutical and biotechnology industries are rushing to adopt AI and ML to enhance their pipeline and save costs. In fact, a recent poll of experts in the pharmaceutical and biotechnology industries found that 60% of them were either currently employing AI technologies in their line of work or planned to do so. AI is essential for pharmaceutical companies to stay competitive by lowering the expensive process of delivering innovative medications to market.

AI also has several restrictions and difficulties that must be taken into account, as some AI systems, such as deep learning neural networks, are additionally referred to as “black boxes” since it is challenging to grasp it, or comprehend their decision-making process. Because of this lack of transparency, it could be challenging for regulatory agencies, academics, and medical professionals to comprehend and believe the results produced by AI. Concerns about algorithmic biases, patient consent, data privacy, and potential exploitation of private medical data are brought up by the implementation of AI in pharmaceuticals. To meet these problems and guarantee that AI technology upholds moral, ethical, and safety norms, regulatory frameworks must be modified.

Pharmaceutical industries must have a strategy in place to protect their inventions and stay clear of common intellectual property (“IP”) problems given the quick rate at which the pharmaceutical and biotechnology industries are using AI. The incorporation of AI could speed up research, boost output, and enhance outcomes in the fields of chemical and pharmaceutical sciences. It is critical to address regulatory, ethical, and privacy issues in order to ensure that AI technologies are used ethically and responsibly across a variety of fields.

## **4.6 CONCLUSION**

There are numerous patents for AI applications in the chemical and pharmaceutical fields, including patents on target identification and affinity measurement, drug discovery, clinical trials, patient management, drug analysis, product development, and drug structures or medical uses.



It can be challenging to get a patent for an AI innovation, and it can also be challenging to enforce a patent. In addition to that, there are practical reasons why patents might not always be the appropriate form of protection for the methods and massive amounts of data that are frequently used in AI-based systems. Therefore, it is important to take into account other IP protection strategies, such as copyright and database rights as well as trade secrets or sensitive information.

It is also to be noted that many life sciences businesses get into research alliances or agreements with other industry collaborators because they may lack competence in the creation of software and computer systems. These partnerships frequently bring up intellectual property (IP) concerns, such as inventorship, ownership of any IP created via the project, interests to preexisting IP, including the parties' rights after the partnership ends. To prevent any potential concerns from impeding business objectives, companies should assess these IP issues as shortly as feasible, preferably before the project begins. Contracts that specify the ownership of intellectual property and the ownership rights to data, AI systems, implementation know-how, and innovations made by AI, for instance, can assist to minimize conflicts and promote collaboration.

## KEYWORDS

- **artificial intelligence**
- **chemical sciences**
- **computer related information**
- **intellectual property**
- **machine learning**
- **pharmaceutical science**

## REFERENCES

1. Fleming, N. (2018). How artificial intelligence is changing drug discovery. *Nature*, 557 (7707), S55–S57. <https://doi.org/10.1038/d41586-018-05267-x>.
2. Heuer, L. (2018). AI could threaten pharmaceutical patents. *Nature*, 558(7711), 519. <https://doi.org/10.1038/d41586-018-05555-6>.

3. Rong, G., Mendez, A., Bou Assi, E., Zhao, B., & Sawan, M. (2020). Artificial intelligence in healthcare: Review and prediction case studies. *Engineering*, 6(3), 291–301. <https://doi.org/10.1016/j.eng.2019.08.015>.
4. White, C. J., et al. (2019). Drafting patent applications covering artificial intelligence systems. *Landslide*, 11(3). American Bar Association.
5. Paul, D., Sanap, G., Shenoy, S., Kalyane, D., Kalia, K., & Tekade, R. K. (2021). Artificial intelligence in drug discovery and development. *Drug Discovery Today*, 26(1), 80–93. <https://doi.org/10.1016/j.drudis.2020.10.010>.
6. Baum, Z. J., Yu, X., Ayala, P. Y., Zhao, Y., Watkins, S. P., & Zhou, Q. (2021). Artificial intelligence in chemistry: Current trends and future directions. *Journal of Chemical Information and Modeling*, 61(7), 3197–3212. <https://doi.org/10.1021/acs.jcim.1c00619>.
7. Rodrigues, T. (2019). The good, the bad, and the ugly in chemical and biological data for machine learning. *Drug Discovery Today: Technologies*, 32–33, 3–8. <https://doi.org/10.1016/j.ddtec.2020.07.001>.
8. Bannigan, P., Aldeghi, M., Bao, Z., Häse, F., Aspuru-Guzik, A., & Allen, C. (2021). Machine learning directed drug formulation development. *Advanced Drug Delivery Reviews*, 175, 113806. <https://doi.org/10.1016/j.addr.2021.05.016>.
9. Ma, J., & Porter, A. L. (2015). Analyzing patent topical information to identify technology pathways and potential opportunities. *Scientometrics*, 102(1), 811–827. <https://doi.org/10.1007/s11192-014-1392-6>.
10. Yang, X. (2021). Artificial intelligence in pharmaceuticals: Bibliometric and collaboration network analysis of patents. In *18th International Conference on Scientometrics & Informetric*. International Society for Scientometrics and Informetric-ISSI, Leuven, pp. 1567–1568.
11. Gadiya, Y., Gribbon, P., Hofmann-Apitius, M., & Zaliani, A. (2023). Pharmaceutical patent landscaping: A novel approach to understand patents from the drug discovery perspective. *Artificial Intelligence in the Life Sciences*, 3, 100069. <https://doi.org/10.1016/j.aills.2023.100069>.
12. (2021). An analysis on pharmaceutical patent applications and grants in India: Megapharma shifts its strategies toward India. In *10th International Congress on Advanced Applied Informatics (IIAI-AAI)*.
13. KR20220046451A: A polypeptide product.
14. IN202211062466A: Artificial intelligence-based interventions for developing stability indicating method and validation of anti-cancer drug using RP-HPLC method.
15. Yuan, L. (2019). A kind of artificial intelligence cleaning device for chemistry tube. *CN110404909A*.
16. Jiang, G., Li, Z., Hu, T., Liu, Y., & Hu, H. (2019). Self-driven artificial intelligence material, preparation method and application in imaging analysis and detection and drug-controlled release. *CN110591166A*.
17. Cao, N., Hu, J., Sorrentino, R. K., Wang, F., & Zhang, P. (2013). Method and system for exploring the associations between drug side-effects and therapeutic indications. *US9536194B2*.
18. Shin, J. M., Park, H. J., Lee, D. B. R. X., & Jung, Y. J. (2020). Polypeptide formulation. *WO2022075510A1*.
19. Beijing Health Vocational College (2023). Pharmacy experiment interactive teaching system, multifunctional experiment all-in-one machine and teaching method. China. *CN115830928A*.

20. Das, P., Cipcigan, F., Wadhawan, K., Padhi, I., Vijil, E., Chen, P. Y., Mojsilovic, A., Sercu, T., & Nogueira dos Santos, C. (2020). Filtering artificial intelligence designed molecules for laboratory testing. *WO2021234522A1*.
21. Mahadevan, N., Di, J., Christakis, Y. P., & Patel, S. (2021). Computerized decision support tool and medical device for scratch detection and flare prediction. *WO2021262857A1*.
22. Naik, A. W., Barro, M., Holloway, D., Zeldovich, K., Strugnelli, T., Davidson, P., & Warren, W. (2020). Systems and methods for predicting biological responses. *WO2021080999A1*.
23. Terasaka, Y., & Yasuhara, S. (1979). Automatic lubricant interchanger for internal combustion engine. *JPS56507A*.
24. Senior, A. W., Kirkpatrick, J., Sifre, L., Evans, R. A., Penedones, H., Qin, C., Sun, R., Simonyan, K., & Jumper, J. (2019). Machine learning for determining protein structures. *WO2020058174A1*.
25. Heinrich Roder, Maxim Tsypin, & Julia Grigorieva (2019). Filed by Biodesix Inc. Method and system for determining whether a drug will be effective on a patient with a disease. Europe EP2241335B1.
26. Abraham Samuel Heifets, Izhar Wallach, & Michael Dzamba (2020). Binding affinity prediction system and method. Europe EP3140763B.
27. Rastogi, R., Rastogi, Y., Rathaur, S. K., & Srivastava, V. (2023). Identification of drug compound bio-activities through artificial intelligence. *International Journal of Health Systems and Translational Medicine*, 3(1), 1–34. <https://doi.org/10.4018/IJHSTM.315800>.
28. Ma, Y., Cao, W., & Tang, Q. (2020). Predicting patient responses to a chemical substance. *WO2021030270A1*.
29. Anyanwu-Ofili, A., Mahmood, K., Batchvarova, N., Boland, R., Gosiewska, A., Van Den Heuvel, A. P. J., & Cula, G. O. (2020). Use of artificial intelligence to identify novel targets and methodologies for skin care treatment. *WO2020170165A1*.
30. Bornstein, C. E., McDonald, K., & Anggara Markely, L. R. (2019). Systems and methods for automated tracking and optimization of global manufacturing and supply based on impacts of post-approval changes. *WO2019231917A1*.
31. Tang, Y., Lou, X., Chen, Z., & Zhang, C. (2020). A study on dynamic patterns of technology convergence with IPC co-occurrence-based analysis: The case of 3D printing. *Sustainability*, 12(7), 2655. <https://doi.org/10.3390/su12072655>
32. WIPO. (2019). *Technology trends 2019: Artificial intelligence*.
33. Gurgula, O. (2020). Strategic patenting by pharmaceutical companies—Should competition law intervene? *IIC*, 51(9), 1062–1085. <https://doi.org/10.1007/s40319-020-00985-0>.
34. Weingarten, M. D., & Cyr, S. K. (2019). Securing and maintaining a strong patent portfolio for pharmaceuticals. *ACS Medicinal Chemistry Letters*, 10(6), 838–840. <https://doi.org/10.1021/acsmedchemlett.9b00201>.
35. Chan, H. C. S., Shan, H., Dahoun, T., Vogel, H., & Yuan, S. (2019). Advancing drug discovery via artificial intelligence. *Trends in Pharmacological Sciences*, 40(8), 592–604. <https://doi.org/10.1016/j.tips.2019.06.004>.
36. Vijayan, R. S. K., Kihlberg, J., Cross, J. B., & Poongavanam, V. (2022). Enhancing preclinical drug discovery with artificial intelligence. *Drug Discovery Today*, 27(4), 967–984. <https://doi.org/10.1016/j.drudis.2021.11.023>.
37. Li, H., et al. (2020). Machine-learning scoring functions for structure-based virtual screening. *WIREs Computational Molecular Science*, 10(5), e1465. Advance Review. <https://doi.org/10.1002/wcms.1465>.

38. Levin, J. M., Oprea, T. I., Davidovich, S., Clozel, T., Overington, J. P., Vanhaelen, Q., Cantor, C. R., Bischof, E., & Zhavoronkov, A. (2020). Artificial intelligence, drug repurposing and peer review. *Nature Biotechnology*, 38(10), 1127–1131. <https://doi.org/10.1038/s41587-020-0686-x>.
39. Calaprice-Whitty, D., Galil, K., Salloum, W., Zariv, A., & Jimenez, B. (2020). Improving clinical trial participant prescreening with artificial intelligence (AI): A comparison of the results of AI-assisted vs standard methods in 3 oncology trials. *Therapeutic Innovation & Regulatory Science*, 54(1), 69–74. <https://doi.org/10.1007/s43441-019-00030-4>.
40. McGrail, S. (2021). AI in the pharma industry: Current uses, best cases, digital future. *Pharma News Intelligence*. 30, 2021.
41. Schork, N. J. (2019). Artificial intelligence and personalized medicine. In *Cancer Treatment and Research* (Vol. 178, pp. 265–283). [https://doi.org/10.1007/978-3-030-16391-4\\_11](https://doi.org/10.1007/978-3-030-16391-4_11).
42. Lee, N. (2020). Protection for artificial intelligence in personalized medicine—The patent-trade secret tradeoff. In *The Harmonization and Protection of Trade Secrets in the EU* (pp. 267–294). Hanken School of Economics.
43. Sotoudeh, H., Shafaat, O., Bernstock, J. D., Brooks, M. D., Elsayed, G. A., Chen, J. A., Szerip, P., Chagoya, G., Gessler, F., Sotoudeh, E., Shafaat, A., & Friedman, G. K. (2019). Artificial intelligence in the management of glioma: Era of personalized medicine. *Frontiers in Oncology*, 9, 768. <https://doi.org/10.3389/fonc.2019.00768>.
44. Tran, V. T., Riveros, C., & Ravaud, P. (2019). Patients' views of wearable devices and AI in healthcare: Findings from the ComPaRe e-Cohort. *npj Digital Medicine*, 2(1), 53. <https://doi.org/10.1038/s41746-019-0132-y>.
45. Robert E. Matthews, Todd O'Connell, & Douglas Romer (2023). Artificial intelligence systems that incorporate expert knowledge related to hypertension treatments. USA. US11600388B2.
46. Jabarulla, M. Y., & Lee, H. N. (2021). A blockchain and artificial intelligence-based, patient-centric healthcare system for combating the COVID-19 pandemic: Opportunities and applications. *Healthcare*, 9(8), 1019. <https://doi.org/10.3390/healthcare9081019>.
47. Schwaller, P., & Laino, T. (2019). Data-driven learning systems for chemical reaction prediction: An analysis of recent approaches. In *ACS Symposium Series* (Vol. 1326, pp. 61–79). <https://doi.org/10.1021/bk-2019-1326.ch004>.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## **PART II**

### **Application of Computational Tools, AI, and ML for Predicting Toxicity and Biodegradation**



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

# Toxicity Predication in Chemistry Based on Machine Learning: A Review

DATTATRAYA N. PANSARE,<sup>1</sup> ROHINI N. SHELKE,<sup>2</sup> ANIKET P. SARKATE,<sup>3</sup>  
ANANT B. KANAGARE,<sup>1</sup> AJIT DHAS,<sup>1</sup> DEVIDAS S. BHAGAT,<sup>4</sup> and  
BHARAT K. DHOTRE<sup>5</sup>

<sup>1</sup>*Department of Chemistry, Deogiri College, Aurangabad, Maharashtra, India*

<sup>2</sup>*Department of Chemistry, Radhabai Kale Mahila Mahavidyalaya, Maharashtra, India*

<sup>3</sup>*Department of Chemical Technology, Dr. Babasaheb Ambedkar Marathwada University, Aurangabad, Maharashtra, India*

<sup>4</sup>*Department of Forensic Chemistry and Toxicology, Government Institute of Forensic Science, Aurangabad, Maharashtra, India*

<sup>5</sup>*Department of Chemistry, Swami Vivekanand Sr. College Mantha, Jalna, Maharashtra, India*

---

## ABSTRACT

The prediction of toxicity holds considerable importance for public health. It has numerous applications, such as reducing costs and streamlining the clinical and preclinical processes by enabling the evaluation of drugs based on predicted toxicity. Nowadays, toxicity prediction models play a vital role in risk reduction across various applications and can potentially serve as a foundation for regulatory decisions. However, the utility of these estimates may be limited if the connections are not adequately quantified. This study also highlights the significant potential of deep learning-based models for



toxicity predictions. We investigate the combination of deep learning-based predictors with a conformal prediction framework to create highly predictive models with well-defined uncertainties. Computational tools have been widely adopted to facilitate rapid and high-throughput *in silico* ADMET analysis, enabling the early prediction of drug-like features in compound selection during drug discovery decision-making. Developing trial-free methodologies for early toxicity prediction in the drug discovery process is crucial in minimizing expensive drug failures caused by late identification of toxicities during development or even during clinical trials. Recent changes in regulations within the industrial chemicals and cosmetics sector have spurred significant advancements in the development, application, and evaluation of non-testing methods like (Q)SAR. Additionally, this study proposes workflows for practically integrating these non-testing approaches into testing and evaluation policies. The objective of this study is to predict the *in vivo* toxicity of aromatic compounds structured with a single benzene ring and the presence or absence of different functional groups.

## 5.1 INTRODUCTION

Evaluating toxicity is a crucial aspect of medicine progress and sanction. It is widely recognized that clinical trials are necessary for the legalization of drugs [1, 2]. Tactlessly, conducting scientific trials is continuously connected with a certain level of danger. Reports indicate that approximately 50% of new drugs were deemed unsafe or ineffective during advanced stages of human clinical trials [3]. Notably, Sitaxentan (Figure 5.1), a drug, was swiftly withdrawn from the worldwide market due to severe and irreversible liver toxicity in human subjects [4, 5].

The significance of preclinical evaluations cannot be overstated in light of the vulnerabilities identified in clinical trials. These assessments play a critical role in preventing the inclusion of harmful drugs in clinical trials. While animal testing is a commonly utilized technique for preclinical assessment, it has its limitations. Firstly, it is an expensive and labor-intensive procedure. Additionally, the results obtained from animal testing provide little insight into human toxic responses due to variations between species and illness models [6, 7], i.e., Sitaxentan exhibited no indications of liver damage in animal experimentations [8], while it caused hepatotoxicity in people [4, 5]. Consequently, animal experiments do not accurately predict how the human body will react to new drugs and do not mitigate risks [6, 9]. To address the cost and concerns associated with animal experiments,

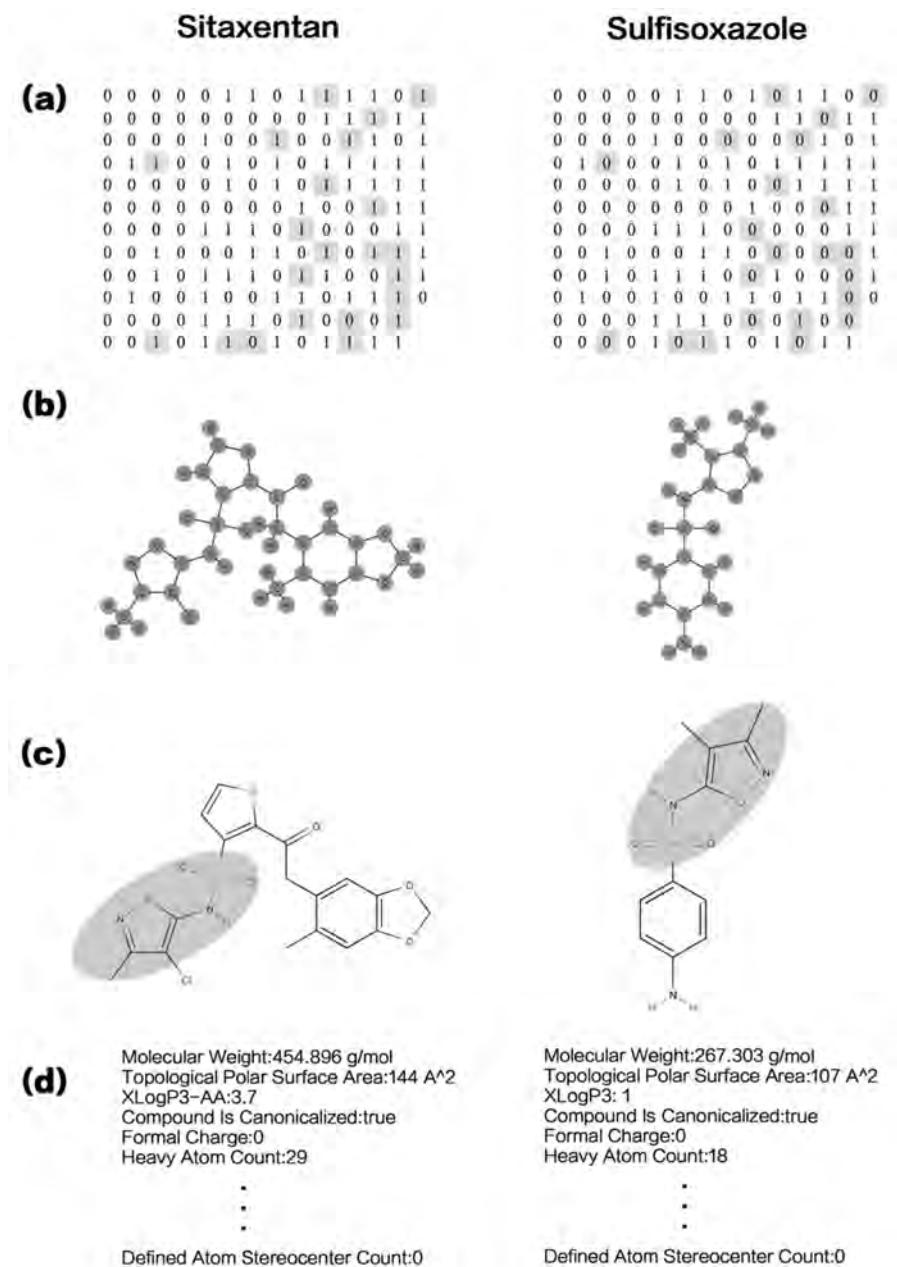


FIGURE 5.1 Sitaxentan and sulfisoxazole.

the use of high-throughput computational toxicity estimation becomes crucial. Among the various methods for toxicity prediction, QSAR built on biochemical stand out as an effective and well-established approach [10]. This method utilizes statistical techniques to establish between drug compound and its physiological activity [11]. By leveraging this relationship, it is possible to infer its physical activity or toxicity. It assumes the independence of factors governing the biological activity of compounds and employs statistical methods like free energy [12]. The scientist, introduced new method in 1964, straight employs physical action [13]. In the year 1980, QSAR investigation started to be utilized for medicinal poisoning estimate [14–16]. As we approach end of 21<sup>st</sup> century, scientist have been predicting poisoning one or numerous chemical-physical actions [17].

When examining multivariate systems, alongside statistical methods, knowledge-based systems have also been employed. Nonetheless, the rapid expansion of data has made it progressively more difficult to maintain extensive knowledge bases. As a result, knowledge-based systems encounter obstacles in performing highly automated tasks involving vast amounts of data [18]. In order to address these evolving responsibilities, researchers have made substantial progress in advancing machine learning techniques and describing chemical structure descriptors.

## **5.2 MACHINE KNOWLEDGE**

Machine knowledge, a subfield of artificial intelligence (AI), utilizes advanced systems to facilitate computers in learning from data and making predictions [19]. Key machine learning algorithms, which stem from the analysis of data and recognition of patterns, include Bayesian classifiers [20]. These algorithms have found extensive applications in data mining, cluster analysis, and pattern recognition [21]. The numerous benefits associated machine knowledge, such as speed, price-efficiency, and high accurateness, have contributed to its growing adoption by researchers in the prediction of toxicity [22]. Investigators have employed various combinations of procedures, such as Hereditary Procedure [23, 24], Accidental Forestry (AF) [25–27], ANN [28–30], and anther machine learning procedures [31–33], to enhance outdated QSAR models for the prediction of drug toxicity or other biological activities. It is important to note that different machine learning methods exhibit varying performance, and their effectiveness can be significantly influenced by factors like the dataset and available computational resources.

### 5.2.1 SHALLOW ARCHITECTURES

Rosenblatt introduced the perceptron model in 1957, which imitates the structure of a neuron and operates as a binary classifier [34]. The utilization of establish the footing for direct classification was first done by Widrow and Hoff [35]. The scientist proposed national procedure, enabling processers to categorize example themes based on three-dimensional characteristics [36]. The decision tree algorithm was projected by Quilan in 1986 [37]. Cortes et al. introduced SVM in 1995 with the goal of identifying the boundary that maximally separates two classes, making it suitable for both high-dimensional nonlinear classification and linear organization [38]. In the year 2001, scientist introduced the RF process, composed of numerous result trees. Each tree provides its own classification, and the final output of the classifier is determined by majority voting [39, 40]. This algorithm enables nonlinear mapping and effectively addresses the problem of nonlinear classification and training using artificial neural networks (ANN) [41]. However, in 1991, it was observed the BP algorithm suffers from the disappearing incline problematic, posing a challenge in training deeper networks. These ANN architectures are commonly known as shallow knowledge [42].

## 5.3 MATERIAL AND METHODS

Data is given in Table 5.1 [43]. The dataset comprises toxicity measurements for 12 atomic receptors and associated aims, obtained through *in-vitro* testing. In a previous publication [44], we discussed the use of this data for evaluating projections. To ensure uniformity in compound structures, we applied the IMI eTOX Project Standardizer [45] and MolVS [46] standardizer, along with tautomer standardization. Additionally, we employed RDKit molecular descriptors [47] and Morgan Fingerprints (FPs) [48], which are extended connectivity FPs. The dataset is divided into active and passive subsets, forming the foundation of a binary classification problem [49]. We adopted stratified K-fold splitting, where training data was reserved for validation in each iteration. The remaining 20% of the exercise data constituted a progressive set, while the final training dataset, referred to as the conformal extrapolation proper training set, was used to train the model. We implemented various algorithms, including result methods like forest random (FR) [39] using the Random Forest Classifier function from scikit-learn, and light GBM [51]. Additionally, in our study, we investigated

alternative variations of graph convolutional neural networks, specifically focusing on the graph attention network (GAT) [52].

**TABLE 5.1** Tox21 Datasets Conclusion

<b>Aim</b>	<b>Active</b>	<b>Inactive</b>
The receptor for aryl hydrocarbons.	943	7,104
The receptor accountable for androgen signaling.	377	8,844
The domain of the androgen receptor that binds to ligands.	302	8,174
Aromatase: an enzyme involved in aromatization.	347	6,760
The receptor for estrogen.	927	6,667
The domain of the estrogen receptor that binds to ligands.	442	8,188
A receptor involved in peroxisome proliferation activation, specifically the gamma isoform.	219	7,848
Peroxisome proliferator-activated receptor gamma: a receptor involved in peroxisome proliferation activation, specifically the gamma isoform.	1,078	6,003
The element in the DNA sequence that is receptive to antioxidants and is connected with the transcription factor.	334	8,628
Mitochondrial membrane potential: the electrical potential across the mitochondrial membrane.	420	7,636
Membrane potential of mitochondrial.	1,126	6,097
DNA damage occurring within the p53 signaling pathway.	528	7,981

## 5.4 CONFORMAL ESTIMATE

A valid prediction can be obtained using a conformal predictor, which considers significance level representing the acceptable error percentage determined by the user. To ensure accuracy, the training set is randomly divided to create a calibration set before model training. This calibration set is then utilized to recalibrate the predictions, comprising mixtures with known labels. As a result, in binary classification, there are four possible outcomes: assigning a single class label, assigning both class labels (both classifications), assigning none of the labels (empty classification), or labeling with either of the two classes. The nonconformist package (<https://github.com/donlnz/nonconformist>) was used to generate conformal predictors, which automatically handle the computation of conformal predictions and are compatible with scikit-learn-like models. Although the nonconformist package is recommended, any model can be adapted easily to serve as a conformal predictor. In a practical example discussed by Norinder et al. [53], the non-conformance score was obtained using the output of the

corresponding model. The default settings were employed for nonconformist models, unless specified otherwise.

## 5.5 MODEL EVALUATION

A model that produces a higher proportion of predictions with only sticker is considered additional effective. The result of forecaster should be accurate within mistake amount ratio to the confidence level. However, simulations with high accuracy tend to be excessively cautious while striving for maximum efficiency. To gain a more comprehensive understanding of general estimation techniques with illustrative examples, we suggest referring to Norinder et al. [53]. In conformal prediction, the estimation process involves generating intervals of estimates instead of single labels. In binary prediction tasks, there can be four potential outcomes: two labels, combined labels, or no label. The validity of conformal forecasters has been established, provided that the data can be exchanged [54]. Quantitative structure-activity relationships (QSAR) [55, 56] are utilized to investigate the links between molecular structure, biological activity, and physicochemical properties. QSAR assists in identifying predictable associations [57–59]. This approach is widely adopted due to its ability to compensate for the scarcity of experimental data, reduce testing expenses, and facilitate high-throughput prediction and screening [60]. The European Union, and the Organization for Economic Co-operation and Development (OECD) employ QSAR for hazard identification, screening, and prioritization [61]. In recent years, QSAR has become a prominent subject in drug chemistry, eco-friendly chemistry, life sciences, analytical chemistry, computational chemistry, and even pesticide research [62–65].

## 5.6 RESULTS AND DISCUSSION

To showcase the practicality of a conformal predictor, extensive computations are necessary when dealing with single-label outcomes, even if they are legitimate. Consequently, the performance and legitimacy of these predictors are primarily evaluated using standard conformal expectation metrics. In this investigation, all datasets were subjected to machine learning techniques specifically selected for this purpose, utilizing a 10-fold cross-validation method. The resultant conformal predictors proved to be both valid and efficient. The performance of diverse models, including thumbprints,

RDKit and graphical synthesis, for sedentary compounds, is depicted in Figures 5.2–5.7. The outcomes consistently indicate an average efficiency of around 75% to 80% confidence (with a significance level ranging from 0.25 to 0.2). Higher confidence levels tend to yield more predictions in the double digits, while lower confidence levels result in less certain estimates. It is important to note that the confidence level decreases as the estimates become less substantial. The nr-er dataset exhibiting inferior performance, which aligns with previous observations [50].

Overall, these findings demonstrate that conformal predictors can achieve a high level of efficiency, making them highly suitable for application in prognostic toxicology tasks.

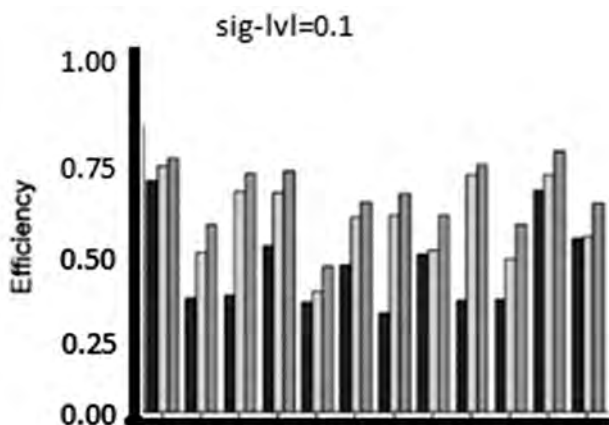


FIGURE 5.2 Model of thumbprints.

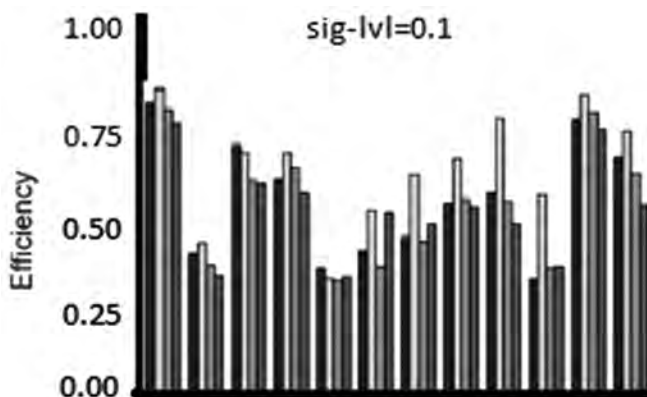


FIGURE 5.3 Model of thumbprints productivity.

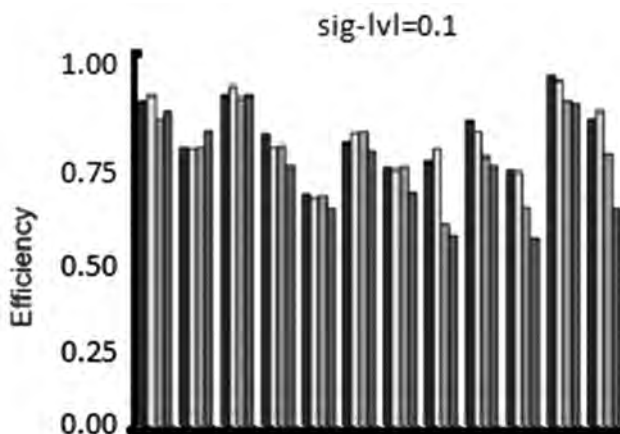


FIGURE 5.4 Model of explainer RDKit.

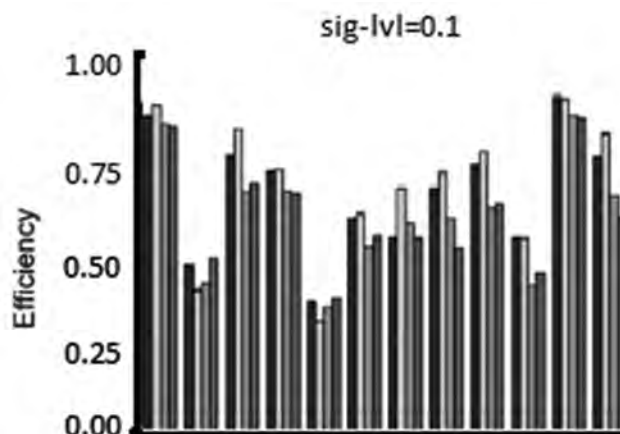
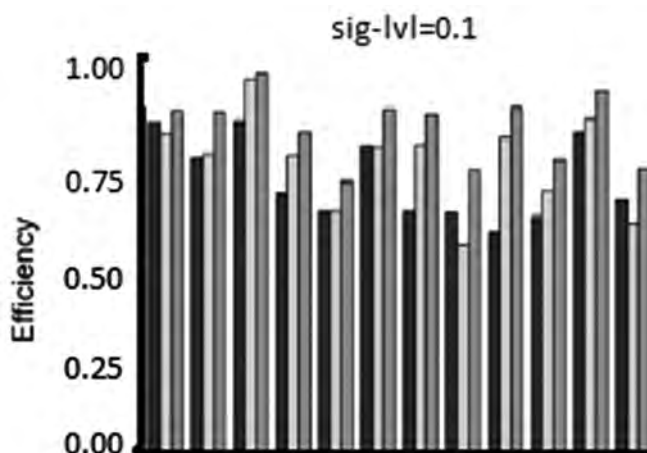


FIGURE 5.5 Model of thumbprints productivity RDKit.

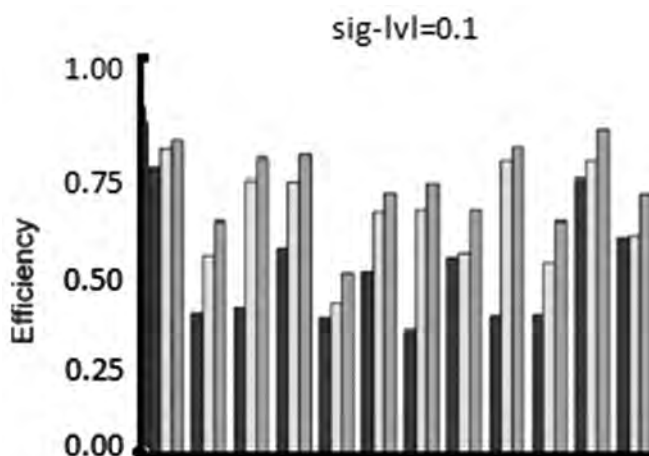
## 5.7 CONCLUSION

Presently, there has been rapid progress in artificial intelligence and machine learning, which has been greatly facilitated by the utilization of computers. This advancement is primarily driven by the significant applications of AI, including the prediction of drug toxicity using computational methods. By employing machine learning techniques and leveraging “big data” science, we can acquire more dependable evidence concerning toxicity compared





**FIGURE 5.6** Model of thumbprints productivity of graphical difficulty.



**FIGURE 5.7** Model of thumbprints productivity complication.

to previous methods. This chapter explores the various approaches of machine learning utilized in toxicity prediction, with a specific focus on the transition from analyzing the chemical structure of compounds to analyzing human transcriptome data. This shift substantially enhances the accuracy of calculations. The advantages of toxicity prediction through machine learning can be summarized as follows. Firstly, computer-based predictions can spare numerous animals from harmful experiments and minimize the need for extensive clinical trials. Furthermore, these

predictions are risk-free, cost-effective, and can be conducted on a large scale. Moreover, by utilizing human transcriptome data, the predictions are based on the complexities of biological systems, providing a more comprehensive understanding of toxicity compared to studies focusing on individual proteins. Lastly, as machine learning has the capability to extract multifaceted and intellectual landscapes in the pharmaceuticals.

## KEYWORDS

- **artificial neural networks**
- **graph attention network**
- **graph convolutional network**
- **graph neural network**
- **hazard expert**
- **integrated approaches to testing**
- **ONCO logic**
- **QSAR**
- **quantitative structure-activity relationships**
- **read across**
- **support vector machines**
- **toxicity prediction**

## REFERENCES

1. Ting, N. (2006). Introduction and new drug development process. In *Dose Finding in Drug Development* (pp. 1–17). Springer.
2. Janodia, M. D., Sreedhar, D., Virendra, L., Ajay, P., & Udupa, N. (2007). Drug development process: A review. *Pharmaceutical Reviews*, 5, 2214–2221.
3. Hwang, T. J., Carpenter, D., Lauffenburger, J. C., Wang, B., Franklin, J. M., & Kesselheim, A. S. (2016). Failure of investigational drugs in late-stage clinical development and publication of trial results. *JAMA Internal Medicine*, 176, 1826–1833.
4. Erve, J. C., Gauby, S., Maynard, M. J., Jr., Svensson, M. A., Tonn, G., & Quinn, K. P. (2013). Bioactivation of sitaxentan in liver microsomes, hepatocytes, and expressed human P450s with characterization of the glutathione conjugate by liquid chromatography tandem mass spectrometry. *Chemical Research in Toxicology*, 26, 926–936.

5. Galiè, N., Hoepfer, M. M., Simon, J., Gibbs, R., & Simonneau, G. (2011). Liver toxicity of sitaxentan in pulmonary arterial hypertension. *European Heart Journal*, 32, 386–387.
6. Johnson, D. E. (2013). Fusion of nonclinical and clinical data to predict human drug safety. *Expert Review of Clinical Pharmacology*, 6, 185–195.
7. Akhtar, A. (2015). The flaws and human harms of animal experimentation. *Cambridge Quarterly of Healthcare Ethics*, 24, 407–419.
8. Owen, K., Cross, D. M., Derzi, M., Horsley, E., & Stavros, F. L. (2012). An overview of the preclinical toxicity and potential carcinogenicity of sitaxentan (Thelin®), a potent endothelin receptor antagonist developed for pulmonary arterial hypertension. *Regulatory Toxicology and Pharmacology*, 64, 95–103.
9. Thomas, R. S., Paules, R. S., Simeonov, A., Fitzpatrick, S. C., Crofton, K. M., Casey, W. M., & Mendrick, D. L. (2018). The US Federal Tox21 Program.
10. Cherkasov, A., Muratov, E. N., Fourches, D., Varnek, A., Baskin, I. I., Cronin, M., Dearden, J., Gramatica, P., Martin, Y. C., & Todeschini, R. (2014). QSAR modeling: Where have you been? Where are you going to? *Journal of Medicinal Chemistry*, 57, 4977–5010.
11. Roy, K., Kar, S., & Das, R. N. (2015). Validation of QSAR models. In *Understanding the Basics of QSAR for Applications in Pharmaceutical Sciences and Risk Assessment* (pp. 231–289). Academic Press.
12. Hansch, C., Maloney, P. P., Fujita, T., & Muir, R. M. (1962). Correlation of biological activity of phenoxyacetic acids with Hammett substituent constants and partition coefficients. *Nature*, 194, 178–180.
13. Free, S. M., & Wilson, J. W. (1964). A mathematical contribution to structure-activity studies. *Journal of Medicinal Chemistry*, 7, 395–399.
14. Quinn, F. R., Neiman, Z., & Beisler, J. A. (1981). Toxicity and quantitative structure-activity relationships of colchicines. *Journal of Medicinal Chemistry*, 24, 636–639.
15. Denny, W. A., Cain, B. F., Atwell, G. J., Hansch, C., Panthanickal, A., & Leo, A. (1982). Potential antitumor agents. Quantitative relationships between experimental antitumor activity, toxicity, and structure for the general class of 9-anilinoacridine antitumor agents. *Journal of Medicinal Chemistry*, 25, 276–315.
16. Denny, W. A., Atwell, G. J., & Cain, B. F. (1979). Potential antitumor agents. Role of agent base strength in the quantitative structure-antitumor relationships for 40-(9-acridinylamino) methanesulfonanilide analogs. *Journal of Medicinal Chemistry*, 22, 1453–1460.
17. Barratt, M. D. (2000). Prediction of toxicity from chemical structure. *Cell Biology and Toxicology*, 16, 1–13.
18. Compton, P., Preston, P., Edwards, G., & Kang, B. (2000). Knowledge based systems that have some idea of their limits. *CIO*, 15, 57–63.
19. Mitchell, T. M. (1997). *Machine Learning*. McGraw Hill.
20. Bishop, C. M. (2006). *Pattern Recognition and Machine Learning* (1st ed.). Springer.
21. Fürnkranz, J., Gamberger, D., & Lavrač, N. (2010). Machine learning and data mining. *Computer Study*, 42, 110–114.
22. Yang, H., Sun, L., Li, W., Liu, G., & Tang, Y. (2018). Corrigendum: In silico prediction of chemical toxicity for drug design using machine learning methods and structural alerts. *Frontiers in Chemistry*, 6, 129.
23. Hemmateenejad, B., Akhond, M., Miri, R., & Shamsipur, M. (2003). Genetic algorithm applied to the selection of factors in principal component-artificial neural networks: Application to QSAR study of calcium channel antagonist activity of 1,4-dihydropyridines (nifedipine analogs). *ChemInform*, 34, 1328–1334.

24. Hoffman, B. T., Kopajtic, T., & Newman, A. H. (2000). 2D QSAR modeling and preliminary database searching for dopamine transporter inhibitors using genetic algorithm variable selection of Molconn Z descriptors. *Journal of Medicinal Chemistry*, 43, 4151–4159.
25. Polishchuk, P. G., Muratov, E. N., Artemenko, A. G., Kolumbin, O. G., Muratov, N. N., & Kuz'Min, V. E. (2009). Application of random forest approach to QSAR prediction of aquatic toxicity. *Journal of Chemical Information and Modeling*, 49, 2481–2488.
26. Svetnik, V., Liaw, A., Tong, C., Culberson, J. C., Sheridan, R. P., & Feuston, B. P. (2015). Random forest: A classification and regression tool for compound classification and QSAR modeling. *Journal of Chemical Information and Computer Sciences*, 43, 1947.
27. Svetnik, V., Liaw, A., Tong, C., & Wang, T. (2004). Application of Breiman's random forest to modeling structure-activity relationships of pharmaceutical molecules. In *Proceedings of the Multiple Classifier Systems, International Workshop, MCS 2004* (pp. 334–343). Springer.
28. Agrafiotis, D. K., Cedeño, W., & Lobanov, V. S. (2002). On the use of neural network ensembles in QSAR and QSPR. *Journal of Chemical Information and Computer Sciences*, 42, 903–911.
29. Wikel, J. H., & Dow, E. R. (1993). The use of neural networks for variable selection in QSAR. *Bioorganic & Medicinal Chemistry Letters*, 3, 645–651.
30. Lu, X., Ball, J. W., Dixon, S. L., & Jurs, P. C. (1998). Quantitative structure-activity relationships for toxicity of phenols using regression analysis and computational neural networks. *Environmental Toxicology and Chemistry*, 13, 841–851.
31. Lu, J., Peng, J., Wang, J., Shen, Q., Bi, Y., Gong, L., Zheng, M., Luo, X., Zhu, W., & Jiang, H. (2014). Estimation of acute oral toxicity in rats using local lazy learning. *Journal of Cheminformatics*, 6, 26.
32. Mazzatorta, P., Cronin, M. T. D., & Benfenati, E. (2010). A QSAR study of avian oral toxicity using support vector machines and genetic algorithms. *QSAR & Combinatorial Science*, 25, 616–628.
33. Srinivasan, A., & King, R. D. (1999). Using inductive logic programming to construct structure-activity relationships. In *Proceedings of the AAAI* (pp. 64–73). Menlo Park, CA: AAAI.
34. Rosenblatt, F. (1988). The perceptron: A probabilistic model for information storage and organization in the brain. *MIT Press* (pp. 386–408). Cambridge, MA.
35. Widrow, B., & Hoff, M. E. (1966). Adaptive switching circuits. In *Neurocomputing: Foundations of Research* (pp. 96–113). MIT Press.
36. Cover, T., & Hart, P. (2002). Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13, 21–27.
37. Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1, 81–106.
38. Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20, 273–297.
39. Breiman, L. (2001). Random forests. *Machine Learning*, 45, 5–32.
40. Ho, T. K. (1995). Random decision forests. In *Proceedings of the 3rd International Conference on Document Analysis and Recognition* (Vol. 1, pp. 278–282). Montreal, QC, Canada: IEEE.
41. Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Readings in Cognitive Science*, 323, 399–421.
42. Hochreiter, S. (1998). The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 6, 107–116.

43. Huang, R., Sakamuru, S., Martin, M. T., Reif, D. M., Judson, R. S., Houck, K. A., Casey, W., Hsieh, J. H., Shockley, K. R., Ceger, P., Fostel, J., Witt, K. L., Tong, W., Rotroff, D. M., Zhao, T., Shinn, P., Simeonov, A., Dix, D. J., Austin, C. P., Kavlock, R. J., Tice, R. R., & Xia, M. (2015). Profiling of the Tox21 10K compound library for agonists and antagonists of the estrogen receptor alpha signaling pathway. *Scientific Reports*, 4, 5664.
44. Zhang, J., Mucs, D., Norinder, U., & Svensson, F. (2019). LightGBM: An effective and scalable algorithm for prediction of chemical toxicity—Application to the Tox21 and mutagenicity data sets. *Journal of Chemical Information and Modeling*, 59, 4150–4158.
45. IMI ETOX Project Standardizer. (Version 0.1.7). Retrieved from: <https://pypi.python.org/pypi/standardiser> (accessed on 25 July 2024).
46. MolVS Standardizer. (Version 0.0.9). Retrieved from: <https://pypi.python.org/pypi/MolVS> (accessed on 25 July 2024).
47. RDKit: Open-Source Cheminformatics. Retrieved from: <http://www.rdkit.org> (accessed on 25 July 2024).
48. Rogers, D., & Hahn, M. (2010). Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50, 742–754.
49. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, É. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
50. Paisios, A., Lenc, L., Martinek, J., Král, P., & Papadopoulos, H. (2019). A deep neural network conformal predictor for multi-label text classification. In A. Gammerman, V. Vovk, Z. Luo, & E. Smirnov (Eds.), *Proceedings of the Eighth Symposium on Conformal and Probabilistic Prediction and Applications* (Vol. 105, pp. 228–245). PMLR.
51. Ke, G., Meng, Q., Wang, T., Chen, W., Ma, W., & Liu, T. Y. (2017). LightGBM: A highly efficient gradient boosting decision tree. *Advances in Neural Information Processing Systems*, 30, 3146–3154.
52. Tsubaki, M., Tomii, K., & Sese, J. (2019). Compound-protein interaction prediction with end-to-end learning of neural networks for graphs and sequences. *Bioinformatics*, 35, 309–318.
53. Norinder, U., Carlsson, L., Boyer, S., & Eklund, M. (2014). Introducing conformal prediction in predictive modeling: A transparent and flexible alternative to applicability domain determination. *Journal of Chemical Information and Modeling*, 54, 1596–1603.
54. Vovk, V., Gammerman, A., & Shafer, G. (2005). *Algorithmic Learning in a Random World*. Springer.
55. Duchowicz, P. R., Castro, E. A., & Fernández, F. M. (2006). Alternative algorithm for the search of an optimal set of descriptors in QSAR-QSPR studies. *Communications in Mathematical and in Computer Chemistry*, 55(1), 179–192.
56. Mu, G., Liu, H., Wen, Y., & Luan, F. (2011). Quantitative structure-property relationship study for the prediction of characteristic infrared absorption of carbonyl group of commonly used carbonyl compounds. *Vibrational Spectroscopy*, 55(1), 49–57.
57. Zhu, H., Rusyn, I., Richard, A., & Tropsha, A. (2008). Use of cell viability assay data improves the prediction accuracy of conventional quantitative structure-activity relationships models of animal carcinogenicity. *Environmental Health Perspectives*, 116(4), 506–513.

58. Drosos, J. C., Viola-Rhenals, M., & Vivas-Reyes, R. (2010). Quantitative structure-retention relationships of polycyclic aromatic hydrocarbons gas-chromatographic retention indices. *Journal of Chromatography A*, 1217(26), 4411–4421.
59. D'Archivio, A. A., Maggi, M. A., Mazzeo, P., & Ruggieri, F. (2008). Quantitative structure-retention relationships of pesticides in reversed-phase high-performance liquid chromatography based on WHIM and GETAWAY molecular descriptors. *Analytica Chimica Acta*, 628(2), 162–172.
60. Tropsha, A. (2010). Best practices for QSAR model development, validation, and exploitation. *Molecular Informatics*, 29, 476–488.
61. Steger-Hartmann, T., & Boyer, S. (2014). *Computer-Based Prediction Models in Regulatory Toxicology*. Springer-Verlag.
62. Ruusmann, V., Sild, S., & Maran, U. (2014). QSAR DataBank—An approach for the digital organization and archiving of QSAR model information. *Journal of Cheminformatics*, 6, 1–17.
63. Schultz, T. W., Cronin, M. T. D., Walker, J. D., & Aptula, A. O. (2003). Quantitative structure-activity relationships (QSARs) in toxicology: A historical perspective. *Journal of Molecular Structure: THEOCHEM*, 622, 1–22.
64. Ma, B., Chen, H., Xu, M., Hayat, T., He, Y., & Xu, J. (2010). Quantitative structure-activity relationship (QSAR) models for polycyclic aromatic hydrocarbons (PAHs) dissipation in the rhizosphere based on molecular structure and effect size. *Environmental Pollution*, 158, 2773–2777.
65. Lee, P. Y., & Chen, C. Y. (2009). Toxicity and quantitative structure–activity relationships of benzoic acids to *Pseudokirchneriella subcapitata*. *Journal of Hazardous Materials*, 165, 156–161.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## CHAPTER 6

---

# Machine Learning Algorithms for Prediction of Chemical Toxicity

D. P. GAIKWAD<sup>1</sup> and SHAMBHAVI S. SINGH<sup>2</sup>

*<sup>1</sup>Department of Computer Engineering, AISSMS College of Engineering, Pune, Maharashtra, India*

*<sup>2</sup>Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India*

---

### ABSTRACT

Now-a-days, human being is exposing to an plenty of chemical compounds through the drugs, cosmetics, atmosphere, and nutrition. For protecting human being from potential harmful effects, drugs, and medicines must be passed steady tests for adverse effects and toxicity. Therefore, there is a need of developing more efficient and less time-consuming methods to predict toxicity of drugs and medicines. The computational models are more efficient approach to predict the toxicity because these models require less time to screen large numbers of compounds at low costs. The researchers have proposed many machine learning (ML) algorithms to implement of toxicity prediction systems. In this chapter, the significant concepts of machine learning are used in human safety and chemical health have summarized in detail. Many freely available tools for toxicity prediction have outlined with their training datasets. Different training datasets are available to build the ML based models. Due to poor annotation in toxicity training dataset, it is very difficult to retrieve and combine these datasets for research in toxicity. This chapter presents the performance of different supervised ML algorithms used to predict chemical toxicity. Specifically, Random Forest, Support Vector Machine (SVM), Classification and Regression Tree (CART),



Gaussian Process, Linear Regression, K-Nearest Neighbors (K-NN), Linear Discriminant Analysis (LDA), and Naïve Bayes algorithms were assessed in terms of receiver operating characteristic (ROC) curves and classification accuracy. Initially, the models normalize the chemical representation of chemical compounds. Chemical descriptors have been computed and used as input to ML models. Considering best the performance of Random Forest, it is planned for prediction of toxicity of drugs. The proposed model has been trained and test using dataset provided by Tox21 Data Challenge. The model has trained and evaluated using training dataset. The proposed model predicts the toxicity of new compounds successfully with highest accuracy.

## 6.1 INTRODUCTION

Recently, the usages of medicines and drugs are increasing in human being to cure different diseases. In modern drug discovery, identification of chemical compounds is the fundamental process. Appropriate selection of chemical complexes can strongly and selectively modify the function of target molecules to provoke a wanted biological response. The selection of chemical compounds from vast chemicals space and determination their drug like assets is a foremost contest [1, 2]. Security of drug is a very vital stuff. Chemical toxicity must be severely learnt earlier a novel drug item is accepted for medical trials [3]. In drug discovery research, the prediction of chemical toxicology is very important in medicine manufacturing process. Conventionally, chemists and biologists accomplish *in vitro* and *in vivo* experiments to check the chemical properties nominated candidates gained from early screening results [4]. These trials are costlier in terms of time and money. These experiments comprise animal testing which are steadily doubtful from moral viewpoints. In literature revisions, it is observed that experiments typically take 2.6 billion dollars and 6 to 12 years to develop a new drug. Drug preparation prior to human testing require more than half of total expenditure [6]. To overcome these limitations of these experimental method, artificial intelligence is widely being used to predict toxicity of drugs and chemicals. ML subarea of artificial intelligence is useful to predict toxicity in less time and accurately. Researchers have used both supervised and non-supervised ML algorithms for prediction purpose. Specifically, different supervised MLs are used in toxicity prediction. In this chapter, different types of ML algorithms to predict toxicity have introduced in detail. The mathematical aspects and basic of each ML have discussed which will help to developer for implementation. The training

dataset is a very important part in ML. The availability of training dataset for toxicity is measure concern for training different classifiers. In this chapter, different training datasets have presented which can assistance scientists and engineers to select appropriate toxicity-based dataset. These datasets have used by many researchers to develop webserver which are used to study the structure and toxicity of chemical compounds. For the study purpose, various available tools have been discussed in brief. Inspired by the success of supervised machine learning algorithms, the Random Forest ML algorithm have proposed for prediction of chemical toxicity of compound of drug. The proposed Random Forest machine learning algorithm has trained using Tox21 dataset. The proposed Random Forest trained model predict the 12 toxicological endpoints. Remaining part of the chapter is separated into subsequent sections. In Section 6.2, overview of toxicity prediction different models has discussed in detail. In Section 6.3 ML algorithms have introduced. In Section 6.4, the architecture of the suggested toxicity prediction system has discussed. Section 6.5 have dedicated for discussion of investigational results. Finally, the chapter has concluded in Section 6.6.

## **6.2 OVERVIEW OF PREDICTION OF TOXICITY**

In this section, the contributions of different researchers have been presented. Yunyi & Guanyu [7] have proposed Toxified the performances of different ML algorithms to predict toxicity of medicines which help to improve the public health. They have conversed the input parameters to the ML algorithm in detail. They have implemented deep learning, un-supervised ML algorithm, Support Vector machine (SVM) and Random Forest algorithms for prediction of chemical toxicity. Yang, Li, Sun, Liu, & Tang [8] have discussed different predictive models for various toxicities. They have discussed different available databases and web servers to predict toxicity. These models can help to design drugs. Li, He, & Cai [9] have proposed graph-based classification and regression to predict toxicity called MoleculeNet. They have also applied focal loss which helps in addressing imbalance in drug dataset. Authors have introduced a dummy super node in which the directed edges are used to connect all nodes in the graph. Authors have modified operation of graph which help to learn graph-level features. Duvenaud et al. [10] have proposed a convolutional neural network which operate on graph to predict the toxicity. The proposed architecture is used to simplify usual molecular feature extraction approaches based on circular patterns. The training dataset is in form of arbitrary size and shape graph. The proposed method is used to demonstrate the interpretable and predictable

performance of these new fingerprints. George, Navdeep, & Ruslan [11] have proposed an Artificial Neural Network (ANN) to predicts actions of compounds of multiple assays at the same time. Quantitative Structure Property Relationship (QSPR) Quantitative Structure Activity Relationship (QSAR) properties of chemicals compounds are studied using Artificial Neural Network. The proposed model provides superior performance. Eric, David, Prasenjit, & Johanna [12] have proposed a new 2D Profile-QSAR predictive model for kinases. The proposed QSAR models the activities of each compound in contradiction of a novel kinase target. These actions are used as chemical descriptors to train models which predict activity against novel kinases. The proposed method is completely automated. Wang et al. [13] have developed a new deep learning-based model to predict toxicity. They have combined attention convolutional neural networks (ACNN) and undirected graph. The proposed model precisely classifies chemical intoxication of honey bees. Authors have trained model using training dataset of 720 insecticides. The size of validation and validation dataset is of 90 insecticides. The proposed model is proficient of precisely categorizing chemicals and has significant potential in practical applications.

### **6.3 INTRODUCTION TO MACHINE LEARNING ALGORITHMS**

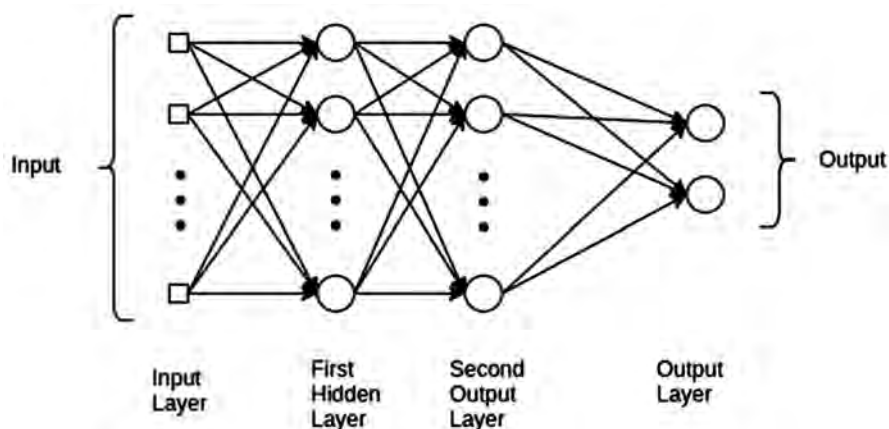
ML is a subcategory of AI that can learn from training datasets and perform prediction and decision-making tasks. It is a multidisciplinary area that encompasses statistics, probability, convex analysis, approximation theory, and more. Machine learning is a subset of artificial intelligence. This section summarizes the current development of ML algorithms in chemical health and security areas. ML algorithms are divided into three categories: supervised, unsupervised, and reinforcement learning algorithms. Supervised learning algorithms work on training datasets that consist of input features and corresponding outputs. These algorithms require labeled data for each type of input and are used for both regression and classification. Some supervised ML algorithms are particularly suited for classification analysis. Unsupervised learning algorithms use training datasets where the target or output is not labeled [14]. These algorithms are commonly used for clustering datasets. Reinforcement learning (RL) algorithms are primarily used for decision-making. RL algorithms learn from the results of their actions and determine which action to take next. After each action, the algorithm receives feedback indicating whether the action was incorrect, neutral, or correct. This feedback

is used to refine future actions. Reinforcement learning is an effective method for making small decisions autonomously. Several ML algorithms are available for regression, classification, and clustering tasks. The important machine learning algorithms are described according to their categories in the following subsections.

### 6.3.1 SUPERVISED LEARNING ALGORITHM

#### 6.3.1.1 ARTIFICIAL NEURAL NETWORK (ANN)

An artificial neural network is constructed to mimic the network of the human brain. It is derived from biological nervous systems that replicate the structure of the human brain. Similar to the human brain, which has neurons interconnected and interrelated with one another, artificial neural networks also consist of neurons interconnected across multiple layers. An artificial neural network is designed and developed to make decisions similar to the human brain. A much simpler and more abstract neuron can be developed by omitting many of the detailed working principles of the human brain. A multilayer perceptron is a network of interconnected perceptrons. It incorporates one or more hidden layers between the input and output layers. In a multilayer perceptron, the output layer may contain more than one output neuron. Figure 6.1 shows the architecture of the multilayer perceptron.



**FIGURE 6.1** Architecture of artificial neural network.

A multilayered perceptron is a collection of multiple perceptron that comprise at least 3 layers; input layer, and output layers sequentially. Each layer may have one or more neuron. The neurons in layers are called input neuron, and hidden and output neuron, respectively. To model the behavior of neuron in the brain, each neurons uses non-linear activation functions except input layer. In back propagation algorithm, multi-layer perceptron has a linear activation function in all its neurons. Different types of activation functions can be used in multilayer perceptron in hidden and output layer. In training, we have to take derivatives of activation function to update weights between layers. There is different learning algorithm to train multilayer perceptron. Basically, two types of supervised learning algorithms namely gradient and stochastic are used to build an artificial neural network. In Gradient Descent Learning, the updating of weights is depending on the error gradient  $E$  in training. The error reduction takes by updating of weights between layers. The activation function takes vital role in updating of weights. So, the activation function should be differentiable. The back propagation process is a case of this kind of learning. Therefore, the weight modification is defined as below using Eqn. (1).

$$\Delta w = \eta \frac{\delta E}{\delta w}; \text{ where } \eta \text{ is a training rate and } \frac{\delta E}{\delta w} \text{ is error gradient with repective } W$$

$$W_{new} = W_{old} - \Delta w \quad (1)$$

In Stochastic ML, the weights are modified in terms of probability. Unsupervised ML algorithms are used to train non supervised ML and these are Hebbian and Competitive learning. Hebbian Learning was developed by Hebb in 1949. It is built on correlated adjustment of weights. The input and output patterns sets are linked with a weight matrix,  $W$ . The following equation is used for Hebbian learning.

$$W = \sum X * Y^T; X$$

where  $X$  signifies the input,  $Y$  is the output, and  $T$  indicates the transpose of the matrix.

In competitive learning, the input form is sent to the network and all the neurons in the layer contest and only the winning neurons have weight adjustments.

### 6.3.1.2 NAIVE BAYES (NB)

NB is an assembly of classification like procedures based on Bayes' theorem. In this algorithm, many algorithms use a common principle in which each pair of features being classified is independent of each other [15]. Bayes'

theorem discovers the probability of an event happening offered the probability of another event that has already happened. Bayes' theorem is specified precisely as the following Eqn. (2):

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)}, \text{ A and B are occurred events and } P(B) \neq 0. \quad (2)$$

Basically, algorithm try to find probability of event A, assumed the event B is true. Event B is called evidence. P(A) is the priori of A. The evidence is an attribute value of an unknown instance B. P(A|B) is a posteriori probability of B after evidence is seen. It is applied using following Eqn. (3):

$$P(C|A) = \{a_1, a_2, a_3, \dots, a_n\} = \frac{P(A = \{a_1, a_2, a_3, \dots, a_n\})P(C)}{P(A = \{a_1, a_2, a_3, \dots, a_n\})} \quad (3)$$

where; C is class of events and A vector of attributes is a dependent feature of size n, were:

$$A = \{a_1, a_2, a \dots a_n\}$$

For understanding, a specimen of a feature vector and equivalent class variable can be demonstrated as below:

$$X = (\text{Hot, Rainy, High, False})$$

$$y = \text{No}$$

So principally, P(C|A) is the probability of not playing golf assumed that the weather conditions are rainy outlook, hot temperature, high humidity and no wind, respectively.

### 6.3.1.3 DECISION TREE

Decision Tree is very popular and important in ML algorithm which is used for classification and prediction which can handle big dimensional dataset. In this tree structure-based algorithm, individual internal node signifies a test on a feature and each branch of tree denotes a result of the test. Each terminal called leaf denotes a class label. A tree is constructed by splitting the training dataset into subsets based on an features value. This procedure is repetitive on each derived subset training dataset in a recursive method. The structure of a decision tree classifier does not need any domain information or parameter setting, and so is suitable for investigative knowledge discovery. In general decision tree classifier has good accuracy. The general architecture

of decision tree is showed in Figure 6.2. This decision tree illustrates the example of tennis game which return the possibility of playing tennis game or not on particular day. For understanding the example following situation have considered.

(Temperature = Hot, Outlook = Sunny, Humidity = High, Wind = Strong)

On this instance, decision classified day with above condition as a negative class. The decision signifies a disjuncture of conjunctions of restrictions on the attribute values of examples.

(Humidity = Normal  $\wedge$  Outlook = Sunny)  $\wedge$  (Outlook = Overcast)  
 $\wedge$  (Outlook = Rain  $\wedge$  Wind = Weak)

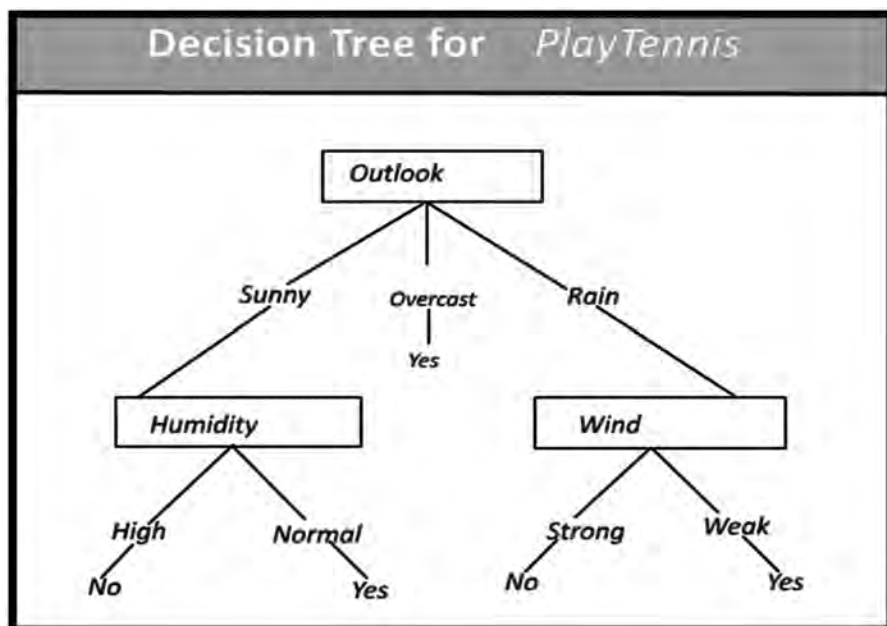
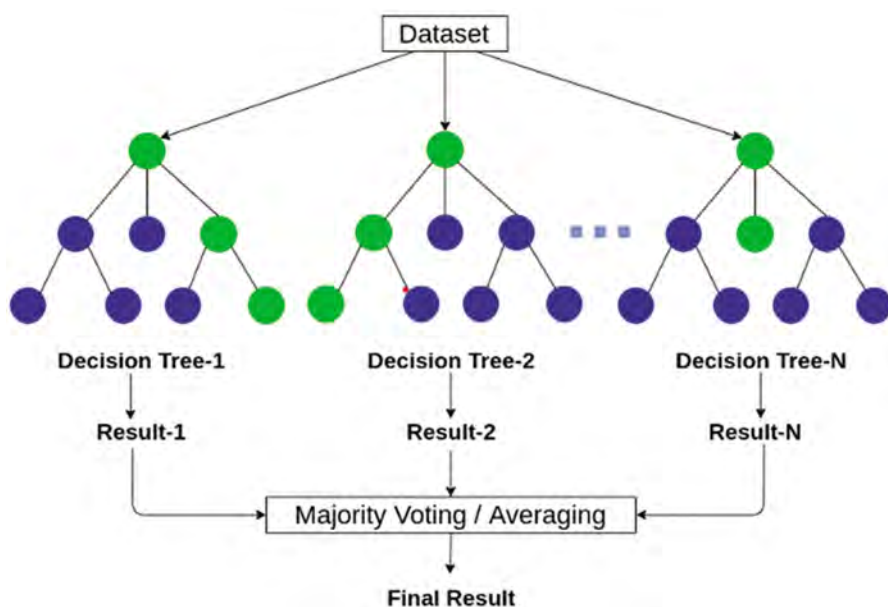


FIGURE 6.2 Decision tree for the concept play tennis.

#### 6.3.1.4 RANDOM FOREST (RF)

RF is a kind of ensemble method of ML. It is based on decision tree algorithm used for regression and classification. Random Forest was invented by Tin Kam using a random subspace technique and extended by Leo Breiman. Random Forest is a forest of randomly created decision trees. This algorithm

is used for making decision by leveraging the power of multiple decision trees. In decision tree, each node works on randomly selected subset of attributes to calculate output using bootstrapping process. The outputs of individual trees are combined to generate the final output using majority voting. Random forest is very useful to minimize the correlation among classifiers. It can accurately solve the problem of overfitting modelled by the decision tree when the entire forest is averaged and reduce the variance. In Figure 6.3, the generation of final output has been shown using Random Forest tree [16].



**FIGURE 6.3** Random forest.

### 6.3.1.5 LINEAR DISCRIMINANT ANALYSIS (LDA)

This algorithm was invented by Fisher in 1936 [17]. LDA is used as a linear classifier. It is also used to reduce the dimension of the training dataset and for visualization. In this algorithm, the training dataset should follow a Gaussian function distribution and the variance of each variable should be the same. Using Bayesian principles, the training dataset for class  $j$  can be classified by reaching the highest likelihood among all  $K$  classes for  $i = 1$  to  $K$ . The discriminant function is established using Bayesian principles and extreme likelihood. This discriminant function presents how possible data is from



each class. The linear discriminant analysis uses many parameters and boundary generated by it may not be enough for class separation.

#### 6.3.1.6 *SUPPORT VECTOR MACHINE*

It is invented by Vapnik in AT&T Bell Laboratories. This ML algorithm is a kind of supervised learning which is used for regression and classification. SVM is used to discover the optimal hyperplane in dataset space. Using this hyperplane, samples are divided by widest gap and separated into different classes. Novel test sample is mapped into data space and its class is finalized grounded on gap side on which it falls. Kernel tricks is used to expand linear classifier to nonlinear classifier using mapping function. This algorithm is used for regression by using  $\epsilon$ -insensitive loss function and radical kernel function. LIBSVM tool is widely used in research to implement classification application using SVM [18].

#### 6.3.1.7 *K-NEAREST NEIGHBOR (k-NN)*

K-NN is a lazy learning algorithm which mostly used for classification. It is simplest machine algorithm which can be also used for regression. K-Nearest Neighbor is a non-parametric ML algorithm because does not consider other assumptions on original data. K-NN does not build from training dataset, but it stores the dataset and do classification by considering similarity between stored dataset and new data [19]. The new data is classified by calculating similarity between it and available data. It places new data into the class that is most similar to the available class.

### 6.3.2 *UN-SUPERVISED ML ALGORITHM*

#### 6.3.2.1 *K-MEAN CLUSTERING*

K-means is an unsupervised learning procedure which create homogeneous groups of individuals (clusters) from the proposed data. It divides dataset into k clusters using nearest mean value. This algorithm is used to for partitioning dataset according to number of clusters and similarities. Initially, the value of k is fixed, then k data points are randomly selected from training dataset. The Euclidian distance between remaining data points and k selected sample is calculated.

Combination of given sample and  $k$  data points is accomplished by minimizing the Euclidean distance between them. The collection of similar data points form a cluster. The cluster centroid is calculated for further combination of given remaining all points. This process is repeated many times to cluster all data points in training dataset. These clusters may use to categorize data points into different classes [20].

### 6.3.2.2 HIERARCHICAL CLUSTERING

Hierarchical clustering was proposed by Johnson in 1967. It is an un-supervised ML algorithm to build a hierarchy of classified clusters. Hierarchical clustering has two methods one is agglomerative and second is divisive method. The greedy method is used for merging and splitting dataset. The results are shown by dendrogram. In dendrogram, the node or data point on same level and near to each other have similar features [21].

## 6.4 METHODOLOGY OF THE PROPOSED TOXICITY PREDICTION SYSTEM

Different kinds of ML techniques are commonly used in predicting toxicity. Toxicology prediction specifically uses several supervised ML approaches. The Random Forest ML method has been presented for the forecast of the chemical toxicity of medicinal molecules. It was encouraged by the victory of supervised ML algorithms. Using the Tox21 dataset, the proposed Random Forest machine-learning method was trained. The proposed Random Forest trained model predicts the 12 toxicological endpoints.

### 6.4.1 TRAINING DATASET FOR PREDICTION OF TOXICITY

In this chapter, the Tox21 challenge training dataset have used for prediction of toxicity of chemicals. This dataset has prepared during Tox21 Challenge in which number of scientists were invited to apply different computational approaches for toxicity prediction [22]. This dataset consists of 12,000 chemicals and drugs which is used to forecast the toxicity of chemical compounds for 12 toxic effects. The specified assays were designed to measure 12 different toxic effects [22]. These toxic effects contained nuclear receptor effects (NR), stress response effects (SR), stress response-heat shock response

effect (SR-HSE) and nuclear receptor-estrogen receptor (NR-ER). These both effects are extremely linked to human health because activation of stress response pathways can lead to cancer or liver injury. Similarly, beginning of nuclear receptors can upset endocrine system function [23]. This dataset was prepared to construct ML model with 12 toxic effects. These 12 toxic effects were provided by high-throughput screening assay measurements. The training dataset consists of 10,000 compounds library of drugs and chemicals. The ML algorithm is trained to forecast result of the high throughput screening assays [24]. These screening assays have enlisted in Table 6.1.

**TABLE 6.1** Tox21 Challenge Training Dataset

<b>SL. No.</b>	<b>Nuclear Receptor (NR) (Bimolecular Targets)</b>	<b>SL. No.</b>	<b>Stress Response Panel (SR)</b>
1.	NR. AhR: Aryl hydrocarbon receptor	8.	SR. ARE: Antioxidant response element.
2.	NR. AR: Androgen receptor	9.	SR. ATAD5: Geno toxicity indicated by ATAD5
3.	NR. AR-LBD: AR luciferase	10.	SR. HSE: Heart shock factor response element.
4.	NR. Aromatase	11.	SR. MMP: Mitochondrial membrane potential
5.	NR. ER: Estrogen receptor alpha	12.	SR. P53: DNA injury P53 pathway
6.	NR. ER-LBD: ER alpha and luciferase.		
7.	NR. PPAR.gamma: Peroxisome proliferator AR gamma		

For training different models, various training datasets are available in the literature [9, 26]. The list of widely used training datasets and their features is provided in Table 6.2.

There are many free available pre-calibrated silico tools to evaluate the toxicity of chemical compounds and drugs. These tools are open source for further development with user interfaces which are listed in Table 6.3. AdmetSAR-2 is a standard tool which is used to assess about 1,00,000 compounds and 27 computational models. Lazar (Lazy structure–activity relationships) web server can be used to forecast acute toxicity, mutagenicity, and carcinogenicity. ProTox II is webserver which is used to predict hepatotoxicity, cytotoxicity, and immunotoxicity. ToxPredict can be used to predict 14 different toxicity end points and generate toxicity reports. QSAR tool is freely available methods-based which depend on molecular descriptors. QSAR models gives productivity on external data sets.

ToxTree is commissioned by European Chemicals Bureau which use decision tree approaches to classify chemicals.

**TABLE 6.2** Toxicity Training Databases

Name of Training Dataset	Important Features
ToxCast	ToxCast is same Tox21. It provides toxicology data for a compound.
MUV	In MUV dataset, positive examples are structurally distinct from one another. It covers 17 exciting tasks for around 90,000 compounds.
PCBA	This database covers 128 bioassays measured over four lac compounds. It contains biological actions of minor molecules produced by high-throughput screening.
HIV	Human Immunodeficiency Virus based dataset presented by the Drug Therapeutics Program prepared by AIDS Antiviral Screening. Dataset is he capability to avoid human immunodeficiency virus repetition for 41,000 and 913 compounds.
FreeSolv	This dataset offers trials and estimated hydration free energy of molecules in water. It contains 643 compounds.
Mole dB	It covers 1,124 descriptors for 2,34,773 molecules.
AcTor	It covers 80,000 compounds and 5,00,000 assays.
BindDB	It comprises 7,25,741 small molecules and protein binding affinity.
ChEMBL	It is bioactivity database with 1.9 M chemical compounds and 11,000 drugs.
Drug2Gene	It is publicly available dataset. It contains compound, drug-gene, and protein.
DrugBank	It consists of drug data and target information for 12,000 and 666 drugs.
EcoTox	It contains 11,000 and 695 chemicals.
Toxnet	It is combination of different datasets with CCRIS Gene-Tox, HSDB, and ToxLine.
ToxRefDB	This dataset contain <i>vivo</i> data for 474 compounds.
REACH	It contains toxicity of 22,391 substances from 95,985 dossiers.

#### 6.4.2 RANDOM FOREST ALGORITHM FOR PREDICTION OF TOXICITY OF CHEMICALS

RF is a popular supervised learning method. It is implemented for classification and regression problems in ML. RF is an ensemble technique of

ML which use the results of sub-trees to build final decision tree. It creates multiple decision tree and combine them together to get a more stable and precise prediction. It is applied to average over several decision trees for the final classification [25]. In method, randomly chosen subset of samples and features are used by each decision tree. Using information gain or Gini coefficient are used to select optimal features. Following Algorithm 1 describe the working ideologies of Random Forest.

**TABLE 6.3** Tools for Toxicity Prediction

Name of Tool	Introduction of Tools
AdmetSAR 2	Webserver based application to predict ADMET.
Lazar	Lazy structure activity relationships are Webserver to predict acute toxicity.
ProTox II	It is used to help predict immunotoxicity, hepatotoxicity, cytotoxicity, carcinogenicity, and mutagenicity.
ToxPredict	It is used to predict 14 different toxicities and generate report.
QSAR Toolbox	It describes molecules and predict chemical hazards.
ToxTree	It predicts toxicological hazards using decision tree.

---

### Algorithm 1: Random Forest

#### Input: Training Dataset

- **Step 1:** Split training dataset using Bootstrap process.
- **Step 2:** Initially, it selects samples randomly from a given dataset.
- **Step 3:** For every sample, it builds a decision tree.
- **Step 4:** Gather prediction outcome from every decision tree.
- **Step 5:** Accomplish voting for each predicted result.
- **Step 6:** Select the highest voted result of prediction as the final prediction outcome.

#### Output: Final Decision Tree

---

Initially, training dataset is divided into multiple bootstrap samples which are depends number of models to be created. On each bootstrap sample, one decision tree model is built. Finally, all decision trees are aggregated to get final output. Figure 6.4 shows decision trees on 10 bootstrap samples.

Finally, the prediction of each decision tree is combined to get final prediction. Figure 6.5 depict the combination of output of each decision tree. The final decision tree is called random forest.

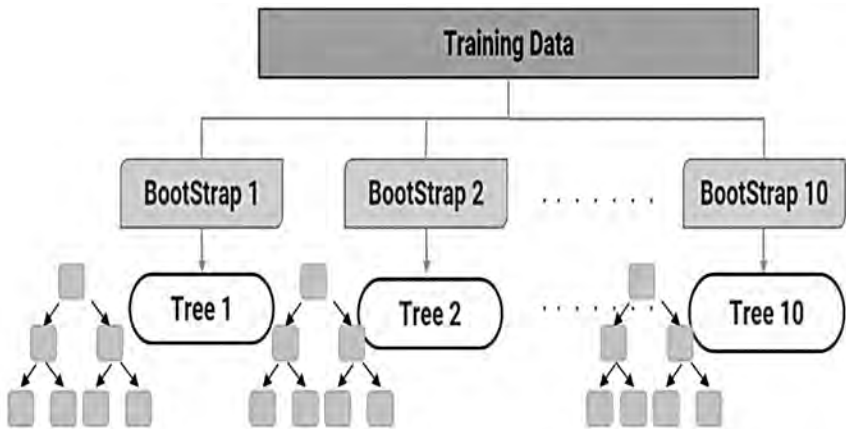


FIGURE 6.4 Decision trees on 10 sample data.

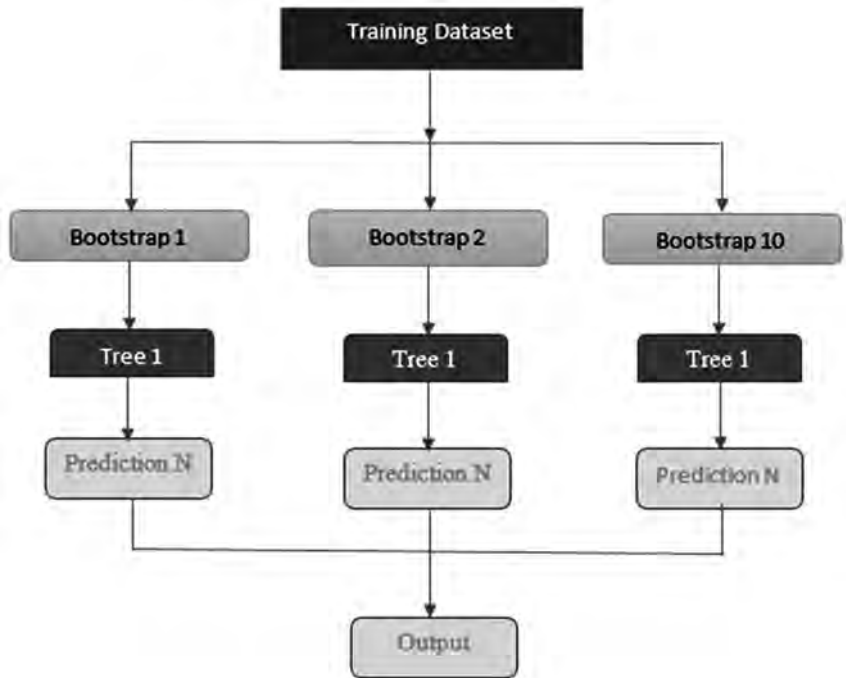


FIGURE 6.5 Decision trees on 10 sample data.

### **6.4.3 HYPER PARAMETERS OF RANDOM FOREST (RF)**

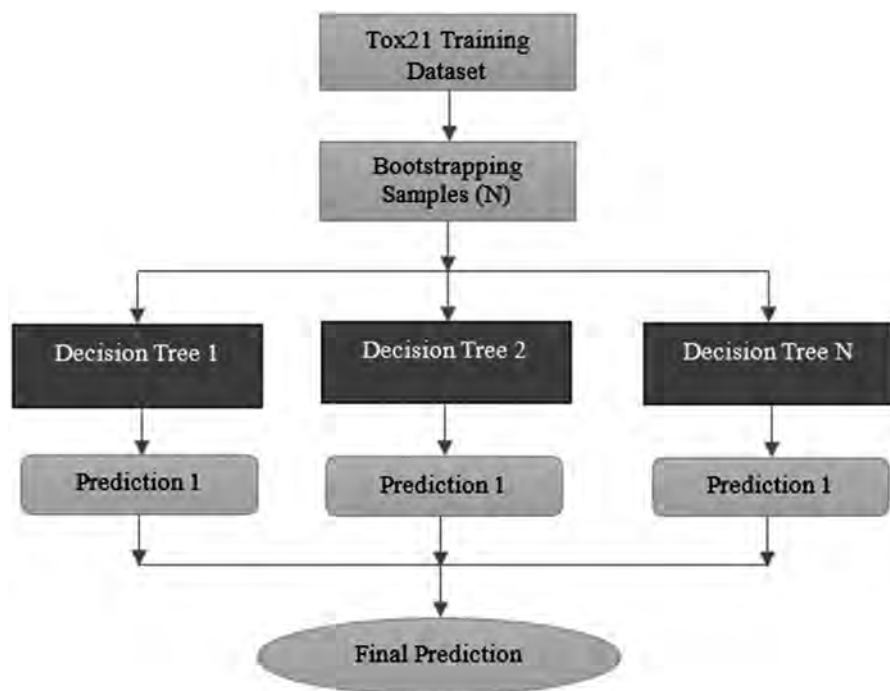
Basically, ML algorithms are separated into two types: parametric and non-parametric. In parametric ML, we can use multiple factors to define ML algorithm. In ML, selection of parameters and hyper parameters are very important because these effects performance of model. The model parameters are internal used to present model. These are continuously estimated during training and used to learn the mapping between the input features and targets values. However, hyper parameters regulate the learning procedure and define the values of model parameters. Hyper parameters of ML are defined by programmer or design engineer. The hyper parameters of RF are the features number, samples number, number of trees, feature choice and feature type. The number of features is used to choose maximum features considered for splitting a node in decision tree. Number of trees hyper parameter is used to define number of trees you want to train. The number of samples parameter is used to select number samples for each node. Maximum feature is used to select maximum feature for each tree. These hyper parameters require to get maximum classification or prediction accuracy. Most of the researchers tune these hyper parameters using different optimization techniques. In this chapter, different hyper parameters have tuned to obtain maximum prediction accuracy.

### **6.4.4 SYSTEM ARCHITECTURE OF RANDOM FOREST FOR TOXICITY PREDICTION**

This section describes the architecture of toxicity prediction. Normally, toxicity prediction model is divided into two-part, data pre-processing and classification. In pre-processing, training dataset is processed to convert cleaned, stable, and understandable by any ML algorithm. The system architecture of toxicity prediction has shown in Figure 6.6. Initially, the training dataset is divided into different groups using bootstrapping method. Then on each group of samples is applied to each decision tree for training purpose. In prediction stage of model, the prediction values are combined to give final output. The hyper parameters of Random Forest algorithm have selected and tuned to get higher accuracy possibility.

## 6.5 RESULTS AND DISCUSSION

In this study, the prediction system has implemented using Python under Anaconda environment. In this study, eight ML have trained on Tox21 training dataset to analyze the performances of these model. Various experiments have carried out using eight different ML algorithms. All ML algorithms were executed more than five times to obtain accurate results. The performances of models have assessed in term of ROC (Receiver Operating Characteristic Curve) score and classification accuracy. Due to imbalanced training dataset of toxicity, ROC is used as the key metric. Specifically, the optimal hyperparameters of Random Forest have selected to get more accuracy. The ROC curve is plot of sensitivity and 1-specificity. The sensitivity, specificity and accuracy are calculated using following Eqns. (4)–(6), respectively.



**FIGURE 6.6** System architecture of toxicity prediction.



$$\text{Sensitivity} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} \quad (4)$$

$$\text{Specificity} = \frac{\text{True Negative}}{\text{True Negative} + \text{False Positive}} \quad (5)$$

$$\text{Specificity} = \frac{\text{True Positive} + \text{True Negative}}{\text{True Positive} + \text{True Negative} + \text{False Negative} + \text{False Positive}} \quad (6)$$

A value 1 of ROC indicate perfect model and 0 indicate a random classifier. The performances in terms of ROC of RF, logistic regression, LDA and K-neighbors have enlisted in Table 6.4. The performances in terms of ROC-AUC of CART, Naïve Bayes, SVM and Gaussian process in Table 6.5. According to Tables 6.4 and 6.5, the Random Forest offers very good ROC scores for all 12 toxic effects than other classifiers. It can be simply observed that Logistic regression, Naïve Bayes and SVM offers very low ROC scores for all 12 toxic effects. Therefore, it can be concluded that these algorithms are not so efficient for prediction of chemical toxicity.

**TABLE 6.4** ROC-AUC of Classifiers

Toxic Effects	Random Forest	Logistic Regression	Linear Discriminant Analysis	K-Neighbors
AhR	0.90185	0.50	0.83977	0.70969
AR	0.80618	0.50	0.80938	0.57382
AR.LBD	0.67313	0.50	0.67726	0.52134
Aromatase	0.75552	0.50	0.64554	0.63329
ER	0.77599	0.50	0.74320	0.64350
ER.LBD	0.72966	0.50	0.65647	0.67414
PPAR.gamma	0.70397	0.50	0.9532	0.54788
ARE	0.76470	0.50	0.72860	0.60043
ATAD 5	0.82908	0.50	0.78051	0.59634
HSE	0.80620	0.50	0.81780	0.62195
MMP	0.92200	0.50	0.88778	0.66805
p53	0.78812	0.50	0.75084	0.60021

The performances in terms of classification accuracy of RF, Logistic Regress, LDA and K-NN have enlisted in Table 6.6. The performances in in terms of classification accuracy of CART, Naïve Bayes, SVM and Gaussian

process in Table 6.7. According to Tables 6.6 and 6.7, the Random Forest offers highest classification accuracy on all 12 toxic effects than other classifiers. It can be simply observed that logistic regression, Naïve Bayes offers very low ROC scores for all 12 toxic effects. Therefore, it can be concluded that Naïve Bayes is not so efficient for prediction of chemical toxicity.

**TABLE 6.5** ROC-AUC of Classifiers

Toxic Effects	CART	Naïve Bayes	SVM	Gaussian Process
AhR	0.72695	0.50	0.50	0.50372
AR	0.61324	0.50	0.50	0.50348
AR.LBD	0.49303	0.50	0.50	0.50348
Aromatase	0.52740	0.50	0.50	0.50409
ER	0.60753	0.50	0.50	0.49336
ER.LBD	0.60797	0.50	0.50	0.50259
PPAR.gamma	0.49871	0.50	0.50	0.48646
ARE	0.60376	0.50	0.50	0.50433
ATAD 5	0.62401	0.50	0.50	0.50342
HSE	0.67158	0.50	0.50	0.50340
MMP	0.68523	0.50	0.50	0.50414
p53	0.56278	0.50	0.50	0.50348

**TABLE 6.6** Accuracy of Classifiers

Toxic Effects	Random Forest	Logistic Regression	Linear Discriminant Analysis	K-Neighbours
AhR	90.49	16.88	87.049	85.73
AR	98.12	97.95	96.249	79.52
AR.LBD	98.28	98.62	75.94	82.81
Aromatase	92.99	92.61	89.77	91.66
ER	92.24	90.11	87.98	88.17
ER.LBD	97.16	04.00	92.66	96.16
PPAR.gamma	94.71	05.45	92.72	94.54
ARE	84.32	83.24	79.45	80.36
ATAD 5	93.89	93.89	91.47	93.08
HSE	96.39	06.88	93.27	96.22
MMP	90.05	88.95	88.58	84.34
p53	93.34	93.34	89.28	91.39

**TABLE 6.7** Accuracy of Classifiers

Toxic Effects	CART	Naïve Bayes	SVM	Gaussian Process
AhR	86.72	11.96	88.03	87.70
AR	96.41	02.04	97.95	97.95
AR.LBD	97.59	01.37	98.62	98.62
Aromatase	87.12	07.38	92.61	92.61
ER	82.36	09.88	90.31	90.31
ER.LBD	94.5	03.33	96.66	96.83
PPAR.gamma	91.73	05.12	94.87	94.87
ARE	75.67	16.75	83.24	83.24
ATAD 5	92.28	06.10	93.89	93.89
HSE	91.96	03.60	96.39	96.39
MMP	85.81	11.04	88.95	89.13
p53	88.14	06.65	93.34	93.34

## 6.6 CONCLUSION

The prediction of toxicity in health care is very important to avoid side effects of drugs and medicine. Before a novel medication candidate is given the go-ahead for clinical trials, chemical toxicity must be thoroughly studied. Therefore, prediction of toxicity in chemical compounds, drugs and medicine is required before finalization of any drug and medicine. Traditionally, *in vitro* and *in vivo* investigations are carried out by chemists and biologists. These experiments are costlier in terms of money and time. These trials use animal experimentation, which is morally doubtful in more and more situations. The manual process for prediction of toxicity is very tedious job for pharmaceuticals, scientist and chemist. To overcome these limitations, the automatic accurate computerized prediction system is a need of chemist, scientist and drug manufacture. In literature survey, it found that researchers have proposed numerous ML algorithms to implement toxicity prediction systems. In this chapter, the important concepts of ML used in chemical-based health and human safety have summarized in detail. Many freely available tools for toxicity prediction have outlined with their training datasets. Different training datasets are available to build the ML based models. Due to poor annotation in toxicity training dataset, it is very difficult to retrieve and combine these datasets for research in toxicity. ML algorithms are heavily depending on large and accurate toxicity training dataset to offer

good accuracy for toxicity prediction. One very important Tox21 challenge dataset and other relevant training datasets used to build model have been discussed in this study. Different supervised machine learning algorithms have implemented using python and other visualization tools for toxicity prediction. The experimental results shows that Random Forest ML system delivers more ROC-AUC values and accuracy than the current toxicity prediction system. This study can offer direction for scientists, engineers, chemical experts who are involved in chemical-based health care and safety. Detail description of ML algorithms will help to study and use ML in their implementation and research work.

## KEYWORDS

- accuracy
- area under curve
- decision tree
- machine learning
- nuclear receptor-estrogen receptor
- random forest
- receiver operating characteristic curve
- stress response
- supervised learning
- support vector machine
- Tox21 dataset

## REFERENCES

1. McInnes, C. (2007). Virtual screening strategies in drug discovery. *Current Opinion in Chemical Biology*, 11(5), 494–502.
2. Kubinyi, H., Mannhold, R., & Timmerman, H. (2008). *Virtual Screening for Bioactive Molecules* (Vol. 10). Wiley.
3. Raies, A. B., & Bajic, V. B. (2016). In silico toxicology: Computational methods for the prediction of chemical toxicity. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 6(2), 147–172.

4. Dean, A., & Lewis, S. (2006). *Screening: Methods for Experimentation in Industry, Drug Discovery, and Genetics*. Springer.
5. Oprea, T. I., & Matter, H. (2004). Integrating virtual screening in lead discovery. *Current Opinion in Chemical Biology*, 8(4), 349–358.
6. Pu, L., Naderi, M., Liu, T., Wu, H. C., Mukhopadhyay, S., & Brylinski, M. (2019). etoxpred: A ML-based approach to estimate the toxicity of drug candidates. *BMC Pharmacology & Toxicology*, 20(1), 2.
7. Wu, Y., & Wang, G. (2018). Machine learning based toxicity prediction: From chemical structural description to transcriptome analysis. *International Journal of Molecular Sciences*, 19(8), 2358.
8. Yang, H., Sun, L., Li, W., Liu, G., & Tang, Y. (2018). In silico prediction of chemical toxicity for drug design using machine learning methods and structural alerts. *Frontiers in Chemistry*, 6, 30. <https://doi.org/10.3389/fchem.2018.00030>.
9. Li, J., Cai, D., & He, X. (2017). Learning graph-level representation for drug discovery. *arXiv preprint arXiv:1709.03741*.
10. Duvenaud, D., et al. (2015). Convolutional networks on graphs for learning molecular fingerprints. In *Advances in Neural Information Processing Systems* (pp. 2224–2232).
11. Dahl, G. E., Jaitly, N., & Salakhutdinov, R. (2014). Multi-task neural networks for QSAR predictions. *arXiv preprint arXiv:1406.1231*.
12. Martin, E., Mukherjee, P., Sullivan, D., & Jansen, J. (2011). Profile-QSAR: A novel meta-QSAR method that combines activities across the kinase family to accurately predict affinity, selectivity, and cellular activity. *Journal of Chemical Information and Modeling*, 51(8), 1942–1956.
13. Wang, F., Yang, J. F., Wang, M. Y., Jia, C. Y., Shi, X. X., Hao, G. F., & Yang, G. F. (2020). Graph attention convolutional neural network model for chemical poisoning of honey bees' prediction. *Science Bulletin*, 65(14), 1184–1191.
14. Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning internal representations by error propagation. In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition* (Vol. 1). MIT Press.
15. Lewis, D. D. (1998). Naive (Bayes) at forty: The independence assumption in information retrieval. In *Machine Learning: ECML-98* (Vol. 1398, pp. 4–15). Springer.
16. Ho, T. K. (1995). Random decision forests. In *Proceedings of the 3rd International Conference on Document Analysis and Recognition* (pp. 278–282).
17. Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7, 179–188.
18. Drucker, H., Burges, C. J., Kaufman, L., Smola, A. J., & Vapnik, V. (1997). Support vector regression machines. In *Advances in Neural Information Processing Systems* (Vol. 9, pp. 155–161).
19. Cover, T., & Hart, P. E. (1967). Nearest neighbor classification. *IEEE Transactions on Information Theory*, 13(1), 21–27.
20. Hamerly, G., & Elkan, C. (2002). Alternatives to the k-means algorithm that find better clusterings. In *Proceedings of the Eleventh International Conference on Information and Knowledge Management* (pp. 600–607).
21. Zeren Jiao, P., Hu, P., Xu, H., & Wang, Q. (2020). Machine learning and deep learning in chemical health and safety: A systematic review of techniques and applications. *ACS Chemical Health & Safety*, 27, 316–334.

22. Mayr, A., Klambauer, G., Unterthiner, T., & Hochreiter, S. (2016). DeepTox: Toxicity prediction using deep learning. *Frontiers in Environmental Science*, 3.
23. Chawla, A., et al. (2001). Nuclear receptors and lipid physiology: Opening the X-files. *Science*, 294, 1866–1870.
24. Huang, R., Xia, M., Nguyen, D. T., Zhao, T., & Srilatha. (2015). Tox21 Challenge to build predictive models of nuclear receptor and stress response pathways as mediated by exposure to environmental chemicals and drugs. *Frontiers in Environmental Science: Environmental Informatics and Remote Sensing*, 3.
25. Breiman, L. (2001). Random forests. *Machine Learning*, 45, 5–32. <https://doi.org/10.1023/A:1010933404324>.
26. Vo, A. H., et al. (2020). An overview of machine learning and big data for drug toxicity evaluation. *Chemical Research in Toxicology*, 33, 20–37.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## CHAPTER 7

---

# Artificial Intelligence-Based Prediction of Drug Metabolism

SHASHIKANT BHANDARI,<sup>1</sup> SHITAL M. PATIL,<sup>1</sup> SHIVRAJ N. MAWALE,<sup>1</sup>  
MRUNAL C. BELWATE,<sup>1</sup> and SOMDATTA Y. CHAUDHARI<sup>2</sup>

<sup>1</sup>*Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy,  
Pune, Maharashtra, India*

<sup>2</sup>*Department of Pharmaceutical Chemistry, Modern College of Pharmacy,  
Pune, Maharashtra, India*

---

### ABSTRACT

Drug effectiveness and safety are significantly impacted by metabolism. Predicting this process is essential to drug discovery and development, eco-toxicology, nutrition, the science of sports and exercise, and precision medicine. In order to forecast metabolism processes, several artificial intelligence tools have been discovered and used. Numerous tools have been discovered via a variety of methodologies, such as data mining, machine learning, and deep structured learning. In the operation of drug development, these tools assist us to know the sites of metabolism, metabolite structures, toxicity, biological half-life, and metabolic pharmacokinetics. In this chapter, we propose to discuss different methods of prediction of drug metabolites by citing examples and case studies of different marked drugs.

An overview of how Machine Learning (ML) together with Artificial Intelligence (AI) is used to analyze drug metabolism and their importance in drug discovery is given in this chapter. Because of their capacity to forecast results from sizable datasets, AI and ML techniques have received a great deal of attention in the drug design and discovery process. They have also



been widely adopted in the drug metabolism and pharmacokinetic streams of drug discovery, particularly in predicting the metabolic fate of molecules under investigation.

## **7.1 OVERVIEW OF AI IN PREDICTING METABOLISM**

The imitating of cognitive processes performed by humans by computers, specifically computer systems, is known as artificial intelligence (AI). It is when technology, especially computer systems, mimics human cognitive processes. AI is becoming more prevalent in many sectors of the pharmaceutical industry such as drug repurposing, drug discovery and development, clinical trials, and improving pharmaceutical productivity. This reduces human workload while meeting targets in a timely manner.

The growing popularity of artificial intelligence as a Power tool for predicting drug metabolism as well as elimination makes it possible to accelerate medication development. The body's ability to metabolize and eliminate medicines is better predicted. Additionally, it foretells the probability of concurrent medication interactions with several metabolic and excretory pathways. It can be used to limit or design chemicals just before synthesis that led to faster development of new drug molecules [1, 2]. It has been used to predict enzyme kinetics and metabolic pathways, providing a data-centric metabolic prediction strategy that can significantly improve prediction accuracy [3]. AI is already a useful tool for forecasting metabolism and has special prospects to improve the predictability and effectiveness of metabolic engineering [4].

The advantages and disadvantages of using AI tools in predicting metabolism, applications of AI in predicting metabolism in various fields, various sources of data used in metabolism prediction, a comparison of conventional methods and AI tools, and an overview of AI-based software or tools were all covered in this chapter. On the implementation of AI based methods for predicting metabolism, we have also provided examples of case studies.

## **7.2 DIFFERENT TYPES OF AI TECHNIQUES COMMONLY USED IN PREDICTING METABOLISM**

Several AI techniques including deep structured learning, machine learning, along with data mining are used in metabolism prediction.

### 7.2.1 MACHINE LEARNING

Drug metabolism has been predicted using ANNs, a type of artificial neural network used in machine learning models that include identifying the enzymes involved, the extent and rate of metabolism, and its locations. Support vector machines and Bayesian classifiers are two other machine-learning techniques that have been used in this field [5].

It is the operation of training algorithms to make predictions or judgments based on data not being specifically programmed. In order to forecast and reconstruct metabolic pathways, which can help with drug development, machine learning has become more and more prevalent in the field of bioinformatics in recent years. In bioinformatics, predicting metabolism is a big difficulty, and machine-learning approaches have been demonstrated to be essential in this field [6].

This approach provides useful data on the potential for metabolism as well as elimination during the preclinical stage also drug discovery, detects important metabolites, pinpoints metabolic pathways, also assesses the chances of drug interactions. It also reduces amount of time and money needed for medication development. The three primary subcategories of computational AI in the area of metabolic fate prediction site-at-metabolism (SOMs) forecasting, metabolite structure estimation, and metabolic pharmacokinetics analysis [7, 8].

Additionally, specialized metabolic genes have been predicted using machine learning. A prediction model was created using machine learning techniques to combine all features, and it has an 87% true positive rate and a 90% true negative rate [9]. This approach does not presuppose interactions and methodically uses arbitrary amounts of new data to enhance predictions [10].

It has also been used to forecast medication interactions with cytochrome P450 isozymes that affect metabolism. High-performance models for the prediction of metabolic interactions between drugs (DDIs) have been created using four distinct types of descriptors, two MI approaches (XGBoost and random forest), two random forest methodologies [11]. Michaelis constants ( $K_m$ ), which are crucial for forecasting the enzyme catalytic rate, can be predicted using machine learning (ML) [12].

In conclusion, this is a powerful tool for predicting metabolism and reconstructing metabolic pathways. It has been used in various applications, including drug discovery, predicting specialized metabolism genes, predicting metabolic pathway dynamics, predicting metabolic drug interactions, and predicting Michaelis constants.

### **7.2.2 DEEP LEARNING**

Predictions of enzymes, metabolites, and reactions have been made using deep learning, a more sophisticated kind of machine learning. Systems metabolic engineering, these approaches have been investigated in the context of which provides knowledge about machine learning techniques for predicting and reconstructing metabolic pathways [13].

Based on annotated genomes, it has become a potent technique in order to predict pathways of metabolism in species. The interconnected pathway processes known as metabolic networks control the biochemical characteristics of cells. They hold the chemical processes, metabolic routes, and regulatory links that govern these reactions. According to the pertinent research, several metabolic networks have been saved and characterized, however, the literature is still far from discovering an exact method with a high degree of accuracy for predicting metabolic pathways. Another hole in the literature is the lack of thorough projects employing deep learning techniques to forecast pathways of metabolism [14].

To fill these gaps, scientists have used supervised machine learning techniques that use deep neural networks which are used to create characteristic depictions of the pathways of metabolism, and random forests are then used to predict these pathways using these representations. All recognized and unrecognized pathways of metabolism in an organism are predicted by the DeepRF supervised learning model. DeepRF has demonstrated strong performance benchmarks for accuracy, recall, and precision are more than 90% when tested on more than 3,00,000 instances. DeepRF gives more trustworthy results than other approaches, as evidenced by a comparison of DeepRF with other methods in the literature [15].

Based on the levels of metabolites and proteins, that are directly learned from training data, it is possible to calculate the pace of each metabolite's change. This approach makes no assumptions about interactions and uses arbitrary amounts of new data to enhance predictions. When predicting the concentration of metabolite, kinetic models explicitly include protein concentrations as a function of time taking enzyme kinetics into account. Metabolic engineers can utilize this kind of prediction to create routes with the desired titers, rates, and yields. This method provides the intriguing real-time possibility of metabolic pathway prediction when combined with recent breakthroughs enabling those capabilities.

A deep learning framework has also been suggested for predicting metabolic pathways. While competing approaches only attain an accuracy of 80.00% [16], this method can predict the relevant pathway of the

metabolism class of 95.16% of tested substances. Prediction techniques based on machine learning have also attained accuracy levels of up to 91.2% and F-measure levels of up to 0.787. Based on the annotated genome, these methods produce a probability for predicting metabolic pathways and can be used to determine which biochemical pathways used as reference library of familiar pathways are observed in the living things [17].

In conclusion, annotated genomes, deep learning has demonstrated considerable promise in predicting metabolic pathways in animals. These techniques can predict metabolic pathways in real time and have attained good performance criteria for accuracy, recall, and precision. To increase the precision of these techniques and investigate their potential uses in drug discovery and metabolic engineering, more study is required.

### **7.2.3 DATA MINING**

Data mining, another artificial intelligence technique, has been used in conjunction with machine learning models to predict drug metabolism. These models help reduce false positives and filter improbable predictions, enabling faster inference and integration of metabolic predictions into drug development processes [18].

The set of chemical events that take place in living things to preserve life is known as metabolism, and data mining is a useful technique for predicting these reactions. Metabolic processes are a complex process that involves the conversion of nutrients into energy and the synthesis of molecules necessary for life. Predicting metabolism is important for understanding the underlying mechanisms of diseases such as metabolic syndrome and for designing metabolic engineering strategies to produce valuable compounds. Metabolic phenotypes have been accurately predicted using data mining techniques [19].

Machine learning techniques have been used to forecast enzymes, metabolites, and reactions that are part of metabolic circuits [20]. In these techniques, enzymes in metabolic pathways are grouped and categorized, and pathway dynamics are predicted from protein concentration data using supervised learning. It has been possible to predict metabolic syndrome, a group of disorders that raise the chances of heart disease, diabetes, and stroke, using machine learning techniques [21]. Early prediction of metabolic syndrome can help patients make lifestyle changes to prevent these diseases. Data mining methods have been useful in analyzing large datasets of metabolic data and in the detection of metabolic indicators that are most significant for predicting metabolic phenotypes [22]. These techniques

include empirical, comparative, investigative, and experimental research methods. To forecast metabolite synthesis, metabolic transformation sites, and metabolic pathways, computational techniques, and resources have been developed. These resources contain tools for the prediction of various endpoints and databases with data.

For two pathways that are important to synthetic biology and metabolic engineering, metabolic pathway dynamics have been predicted from protein concentration data using an ML approach. This method does not presuppose any specific interactions and uses arbitrary of new data to enhance predictions. The prediction of metabolite concentrations from protein concentrations as a function of time has also been done using kinetic models. By using these models, metabolic engineers can create pathways with the necessary titers, rates, and yields [23].

Finally, data mining has been used to predict metabolic phenotypes, elements of metabolic pathways, and metabolic syndrome. It is a valuable technique for forecasting metabolism. Kinetic models have been used to predict the concentrations of metabolite as a reference of time, and machine learning techniques are utilized to predict route dynamics. To forecast metabolite synthesis, metabolic transformation sites, and metabolic pathways, computational techniques, and resources have been developed. These techniques and resources are crucial for comprehending the underlying causes of illnesses and developing metabolic engineering plans to create useful molecules.

## **7.3 ADVANTAGES OF USING AI IN PREDICTING METABOLISM**

### **7.3.1 ACCURACY**

Accuracy is one of the advantages of employing AI techniques to forecast metabolism. Using established and locally created prediction equations, one can assess the precision of metabolism prediction [24]. In addition, specialized models to forecast the locations of metabolism of phases I and II processes are available. Here are a few examples of how well artificial intelligence can anticipate metabolism:

In a study evaluating the precision of resting metabolic rate prediction equations, adult participants used both established and locally created prediction formulas [25]. The computational tool Biotransformer integrates a knowledge-based strategy with an ML technique for predicting small chemical metabolism in human body the human gut, tissues, natural world.

An extensive evaluation of it resulted that it could outperform commercially available tools that are two state-of-the-art, with preciseness and generate values ranging to seven times greater procured for identical sets of phytochemicals, pesticides, pharmaceuticals, or endobiotic in comparable or equivalent conditions [26].

### **7.3.2 EFFICIENCY**

One advantage of AI technology is their effectiveness in forecasting how well drugs would be metabolized and removed via the human system. AI enables the quick screening of sizable compound libraries, providing insightful data on the compounds' possible metabolism and excretion [27]. Using AI, it is now possible to forecast how drugs will be metabolized and excreted. Machine learning-based prediction models have been created recently to precisely forecast the drug's metabolic stability. These models allow for quick inference, which makes it possible to incorporate metabolic investigations into the early stages of drug development [28]. But precise drug metabolism prediction necessitates knowledge of the relevant enzymes, the pace, and extent of metabolism, the locations of metabolism, etc.,

Even on highly capable modern computers, some methods, such as molecular dynamics and quantum chemistry calculations, have a high computational cost despite being correct. In order to increase efficiency, high-throughput applications are required [29]. To do this, explainable artificial intelligence also known as artificial intelligence can be explained as the method through which ML-based prediction models make are explored [30]. In summary, AI has demonstrated substantial promise in predicting drug metabolism, and with further advancement, it may significantly increase the efficiency of drug development.

### **7.3.3 COST-EFFECTIVENESS**

A growing number of people are using AI methods to anticipate metabolism because of their usefulness and efficiency. Using this method, excretion and metabolism prediction, deep standard learning and ML have been utilized to forecast metabolic pathways and drug bioactivities [31, 32].

Using past patient data, these techniques have been successful in forecasting an individual's treatment responses to various therapeutic combinations the

affordability of AI systems for metabolism prediction is one of their key benefits. Traditional methods of estimating metabolism sometimes need substantial testing and data analysis, which can be costly and time-consuming. However, AI techniques are able to quickly and accurately analyze massive amounts of data, eliminating the need for expensive testing and data analysis [33]. The effectiveness of AI techniques for metabolism prediction is another benefit. With the use of AI technologies, researchers can quickly and precisely analyze large data to find recurring themes and networks that might not be immediately evident when using more conventional approaches. This may result in more precise predictions of drug bioactivities and metabolic pathways, as well as more effective pathway design [34].

#### **7.3.4 REDUCES ERRORS**

In the area of predicting metabolism, AI systems have several benefits, including error reduction. The metabolism and excretion of drugs, inborn metabolic mistakes, and the kinetics of metabolic pathways can all be predicted using AI algorithms. These forecasts can be created with greater precision than manual ones, which lowers the possibility of mistakes. An AI algorithm was utilized in a study that was published in *Frontiers in Pediatrics* to test for inborn metabolic abnormalities, which led to a higher incidence rate and fewer false negatives [35]. Another study focused on the application of deep standard learning and machine learning algorithms to improve prediction accuracy as it investigated recent advancements in AI-based medication metabolism and excretion prediction. Pharmaceuticals published this work [36].

Overall error reduction is just one benefit of using AI technologies to anticipate metabolism. With great accuracy, these techniques can be used to predict inborn metabolic errors, medication metabolism and excretion, and metabolic pathway kinetics. Additionally, through prediction, AI can help to lessen the frequency and effects of ADEs [37].

#### **7.3.5 IMPROVED SPEED**

A growing number of people are using AI techniques to forecast metabolism because of their many benefits, including faster speeds. Large data sets can be analyzed by machine learning algorithms, which can also spot patterns that would be challenging for people to notice. This makes

it possible to anticipate metabolic routes more accurately, which helps metabolic engineers construct pathways with the necessary titers, rates, and yields [4].

AI has become an effective tool for forecasting drug metabolism and excretion, potentially lowering the attrition rate of drug development. With the aid of AI, we can model drug metabolism by identifying associated enzymes, the pace and dimension of metabolism, and the location of metabolism. Metabolic engineering's application of machine learning allows us to select the best molecules to make and provide ideas for potential synthesis pathways for those molecules. A notable advantage over conventional methods is the higher speed of AI in predicting metabolism, which enables greater accuracy and quicker predictions [39]. We may anticipate significantly more innovations in medication development and drug metabolism prediction as AI develops.

## **7.4 APPLICATIONS OF AI IN PREDICTING METABOLISM IN VARIOUS FIELDS**

### **7.4.1 DRUG DEVELOPMENT AND RESEARCH**

The efficacy and safety of potential treatments must be determined, and AI for predicting drug excretion and metabolism has developed into a powerful tool. An essential first step in developing novel medications is predicting these mechanisms. AI-based predictions of drug excretion and metabolism have the potential to hasten medication development and boost clinical success rates. Recently, Artificial intelligence-based drug excretion and metabolism prediction has improved thanks to machine learning algorithms and deep learning [3].

By combining metabolic reaction templates and deep learning, it has been possible to anticipate the primary drug metabolites using deep learning algorithms. This method can predict drug metabolites and has some predictive power even if it does not account for each and every metabolic enzyme utilized in human reactions [41].

Using machine learning (ML) models and their quick inference capabilities, metabolic studies may now be incorporated into the development of a medication in its early phases. ML models are applied in prognostication of metabolism and excretion, and the identification of metabolizing enzymes of drug and drug-drug interactions [8].



The evolution of AI for analyzing drug excretion and metabolism faces numerous challenges, including the complexity of metabolic pathways, the need for accurate metabolite prediction, and the requirement for vast amounts of high-quality data. Several significant potential benefits observed in metabolism and excretion prediction, however, include quicker drug development and higher clinical rates [43].

Drug discovery experts will be better equipped to forecast a compound's metabolic fate if the repertoire of predictive models is expanded beyond cytochrome P450 isozymes. A team of AI developers is advancing predictive modeling for enzymes that metabolize drugs [44].

#### **7.4.2 ENVIRONMENTAL AND TOXICOLOGICAL SCIENCE**

Metabolism has been predicted using ML and AI in the fields of environmental and toxicological science. Toxicology prediction models, help to reach a scientific consensus and anticipate the toxicity of compounds, have been developed using AI technology [45].

Through dimension chemicals metabolites, and reduction, that might indicate illness status or toxicity, AI and ML can find variations between phenotypes [46]. It is feasible to find altered metabolites resulting from exposure to the environment using machine learning-based technologies [4]. In order to manage illness and promote plant health, AI has also been utilized to investigate the Phyto microbiome [48].

In environmental and toxicological science, metabolism has been predicted using AI and ML. These technologies have been employed in the creation and optimization of pathways, scale-up, and the detection of phenotypic variation [49]. It is now possible to find known or unidentified altered metabolites resulting from environmental exposure thanks to the application of AI and ML-based techniques. Additionally, AI has been applied to study the Phyto microbiome in order to control illness and advance plant health [50].

#### **7.4.3 ADAPTIVE NUTRITION**

AI, which can predict metabolism, may bring about a revolution in the field of adaptive nutrition. AI systems can be used to comprehend and forecast the complex and non-linear correlations for nutrition-related data and health outcomes [51]. ML can analyze metabolomics data that has been processed

for uses such as disease prediction and understanding disease causes. Because ML can cluster and categorize data, it is appropriate for analyzing complicated and high-dimensional data produced in the field of nutrition [52]. Metabolic pathway prediction is one way AI is used to forecast metabolism. By using arbitrary amounts of new data to enhance predictions, a machine learning method can forecast metabolic pathways [53].

With very little data, this method had greater prediction power than a conventional kinetic model [54]. Individualized nutrition and health are another use. AI platforms can be trained using personal health data, and clinicians can analyze the forecast together with patient reports [55].

Additionally, diet response can be predicted using AI both before and after lengthy interventions. A digital twin technology can accurately forecast fresh independent data and mechanically clarify and incorporate information from several clinical studies [56]. It is crucial to remember that up until now, AI research has been successful in creating systems that excel at a particular activity. Most of these systems perform badly outside of that task, meaning that true intelligence has yet to be achieved.

#### **7.4.4 SCIENCE OF SPORTS AND EXERCISE**

Numerous uses for AI exist in the realm of sports and exercise metabolism prediction. AI can be used to forecast the likelihood of interactions between several medications at once in the body's metabolism and excretion, which can help with drug discovery [3]. AI-powered computer models can predict the progression of the disease by looking at changes in metabolic and cardiovascular biomarkers like levels of cholesterol, the body's mass index, insulin levels, and blood pressure [51]. Additionally, ML algorithms can be used to analyze exercise capacity utilizing baseline data such as cardiovascular illness history, medication use, blood pressure, data on demographics, morphological measures, and dual-energy X-ray absorptiometry (DXA) determined physique structure metrics [59]. AI could use to forecast the best exercise regimens for patients with osteopenia and osteoporosis and to detect osteoporosis from medical pictures [60].

Additionally, utilizing whole-genome sequencing research as a foundation, AI can be utilized to create a model for predicting metabolic syndrome. Finally, by anticipating weight loss, adherence to customized physical activity objectives, nutritional slip-ups, and periods of emotional eating, AI can be applied to enhance adult self-management to reduce weight and behavior changes correlated to weight [61].

These AI-based applications for predicting metabolism in the context of exercise and sport can aid in drug development, disease progression prediction, exercise capacity prediction, diagnosis, and self-control of weight reduction [62].

#### **7.4.5 PRECISION MEDICINE**

Precision medicine is a medical technique that aims to develop and optimize diagnosis, therapeutic intervention, and prognosis. It makes use of big biological datasets with several dimensions that capture individual differences in environment, function, and genes. Precision medicine may be revolutionized by Machine Learning and Artificial Intelligence, it will allow physicians to personally tailor early medicines to each patient. In addition to helping organizations adopt precision medicine, AI and ML can also aid research into the etiology and course of chronic diseases [63].

Predicting metabolism is one of AI's uses in precision medicine. Precision medicine's foundation is accurate, individualized therapy response prediction. By combining several data types from the same patient, such as genetic, clinical, and metabolic data, Drug response can be predicted using machine learning and artificial intelligence [64].

A higher level of accuracy can be achieved when predicting disease risk using AI and ML. For cancer and cardiovascular illness, prediction algorithms using AI techniques have produced encouraging results [65]. Organizations may benefit from AI and ML in a variety of ways, and they can also help us understand the causes and course of chronic illness. The application of machine learning algorithms in precision medicine to assess multiple record of patients, including medical, genetics, metabolites, image processing, experimental, claims dietary, and lifestyle data, is the most recent development [66].

Therefore, AI and ML have the power to completely transmute precision medicine by aiding doctors to individually tune early therapies. With more accuracy, AI and ML can be used to forecast illness risk and medication response. The use of ML methods in precision medicine to analyze various patient data is one of the most current advances. With the implementation of ML and AL, organizations may gain from precision medicine in a variety of ways. These technologies can also improve our understanding of the causes and course of chronic diseases [67].

## 7.5 DIFFERENT DATABASES AVAILABLE

### 7.5.1 CHEMICAL DATABASES

1. **Beilstein:** Contains organic compounds and their properties.
2. **BIAdb:** Contains information on benzyloquinoline alkaloids.
3. **BindingDB:** Contains information on the noncovalent association of molecules in solution.
4. **ChEBI:** Contains a dictionary of molecular entities focused on ‘small’ chemical entities.
5. **ChEMBL:** Contains the integration of chemical, bioactivity, and genetic data; bioactive compounds with drug-like characteristics.
6. **ChemSpider:** Contains chemical structures and information on chemical reactions.
7. **HMDB:** includes metabolic products that have been identified in the human body.
8. **HugeMDB:** Contains small molecules.
9. **PubChem:** Contains information on chemical substances and their biological activities.
10. **ChemBioFinder:** Contains various databases for chemical substances [68–71].

### 7.5.2 ENVIRONMENTAL DATASETS

There are several sources of environmental datasets available online some of them are listed below:

1. **Center for Sustainability and the Global Environment–UW-Madison:** This center provides datasets that aim to improve our understanding of Earth’s terrestrial ecosystems, hydrological systems, and climate.
2. **data.world:** There are 3,248 environment datasets available on this platform, contributed by thousands of users and organizations across the world.
3. **Deepchecks:** This website lists the 10 best free climate and environment datasets for machine learning, including datasets on earth surface temperature data, international greenhouse gas emissions, and air quality annual summary.

4. **EPA Environmental Dataset Gateway:** Users can find, examine, and use datasets and geospatial tools made accessible by EPA's program offices, regions, and labs using this website.
5. **Open Data | US EPA:** The EPA offers a data catalog that allows users to view or download datasets curated by the agency. Users can also access policies, guidance, data standards, and progress reports related to the agency's open data initiative [49, 73, 74].

## 7.6 AI-BASED SOFTWARE TOOLS THAT CAN BE USED FOR METABOLISM PREDICTION

### 7.6.1 BIOTRANSFORMER 3.0

It is an AI-based computer program that provides predictions regarding soil and aquatic microbiota, gastrointestinal microbiota, and the small-molecule metabolism of mammals. The program is publicly available and employs both a machine learning (ML) strategy and a research-based approach to forecast the metabolism of small molecules. Biotransformer consists of EC-based, CYP450, Phase II, Human Gut Microbial, and Ecological Microbial subsystems. It aids researchers in discovering metabolites through metabolic predictions and can accurately predict the final products of metabolic transformations. The input molecular structure is used to generate an interactive table of expected metabolic or transformation products, along with the predicted enzymes responsible for these processes [75, 76].

### 7.6.2 CYPREACT

A computer program called CypReact implements machine learning to anticipate how a small molecule will interact with a particular CYP450 enzyme. The inputs that CypReact accepts include one of the most important human CYP450 enzymes and any arbitrary molecule defined as a SMILES string or a typical SDF file. It then correctly predicts how the query molecule will respond to the specific CYP450 enzyme [77].

### 7.6.3 CYPRODUCT

A software program called CyProduct forecasts the metabolic byproducts of a particular human CYP450's metabolism for a given chemical. It is divided

into three modules: CypReact can forecast the outcome of an interaction between a certain CYP450 enzyme and the query material. The “bond site,” or the bonds in the query molecule that interact with the Cytochromes, is properly predicted by CypBoM. The metabolic waste products predicted by MetaboGen are based on the predicted bond sites by CypBoM [78].

#### **7.6.4 SMARTCYP**

It is a method for figuring out where CYP P450-mediated metabolites are most likely to take place in molecules. The *In-Silico* method predicts the site where chemically like drug molecules would go through CYP P450-mediated metabolism. The method, which is based on a reactivity model, tends to forecast locations where the cytochrome P450 3A4 isoform will metabolize substances. SMARTCyp is freely accessible online. SMARTCyp 3.0 has replaced the previous website server for predicting the place where CYP P450-mediated metabolism occurs uses biotransformer to predict the metabolites generated by xenobiotic biotransformation pathways, such as phase I and II metabolism and microbial metabolism [79].

#### **7.6.5 ADMET PREDICTOR**

It is a specific kind of computer learning programme that forecasts more than a 300 attributes, such as solubility, logP, pKa, and CYP sites. It is a computational filter that can rapidly calculate attributes when screening drug candidates. A compound's drug-likeness can be predicted using ADMET Predictor, and it can also be fixed if there is any ongoing ADMET liability in a lead series. Its REST API can be used to quickly calculate properties [80].

### **7.7 ADVANTAGES OF DIFFERENT AI TOOLS IN PREDICTING METABOLISM**

#### **7.7.1 WIDER APPLICATIONS**

This software could have major beneficial benefits in a number of scientific disciplines, including analytical chemical science, natural products chemistry, farming, and food science, by resolving the limitations of *in silico* metabolism prediction.

### **7.7.2 MULTIPLE APPROACHES**

Many software programs use a range of procedures to entirely cover the metabolite prediction flow, including rule-based approaches the structures of likely metabolites, and other methods, that identify the site of metabolism.

### **7.7.3 AVAILABILITY**

There are numerous freely accessible software tools available to make predictions about metabolism. This facilitates access to the usage of these tools by researchers.

### **7.7.4 USER-FRIENDLY**

This online software can be used by anyone without the requirement for programming knowledge or the installation of local software, and web applications like Biotransformer, SmartCyp, etc., are easier to use than the standalone program.

## **7.8 LIMITATIONS OF DIFFERENT AI TOOLS IN PREDICTING METABOLISM**

### **7.8.1 LACK OF STANDARDIZATION**

The method for predicting metabolites is not fully standardized, and it is unclear how different methodologies would affect the functional result.

### **7.8.2 RELIANCE ON EXPERIMENTAL DATA**

The degree of prediction accuracy will depend on the caliber and volume of experimental data utilized to train the program.

### **7.8.3 LIMITED SCOPE**

Some software programs have restrictions on the types of molecules they can anticipate.

#### **7.8.4 UNABLE TO ACCURATELY FORECAST METABOLITES IN NON-MAMMALS**

Due to the bias against mammals in the processed metabolic data, it is impossible for the software to accurately forecast the metabolites generated by bacteria, plants, or insects.

### **7.9 CASE STUDIES OF VARIOUS MARKETED DRUGS**

#### **7.9.1 BACKGROUND**

A program called Biotransformer 3.0 predicts the gut microbiome and small chemical metabolism in animals. It is a free web service and freeware application for precise, thorough metabolic prediction and *in silico* metabolite identification. Biotransformer applies both a machine learning-based method and a method based on research that forecasts the metabolism of tiny molecules. Evolutionary computation-based, Cytochrome P450, Phase II, Human Gut Microbial, and Environmental Microbial are the five distinct modules.

In addition to providing a reaction research base containing general biotransformation principles, preference criteria, and different restrictions for metabolism prediction, the programme generates projected metabolite structures in popular electronic file formats. The Biotransformer was found to be capable of properly predicting the human gut metabolism of a various of compounds, varying endogenous molecules to xenobiotics, after a thorough evaluation. It is a hybrid software that forecasts xenobiotic metabolism in various systems. An evaluation of the precision of the predictions for a particular set of compounds, the identification of metabolites, and a comparison of the outcomes with other methods are all po components of a case study of the forecasting of chemical metabolism using the Biotransformer program. The case study also investigates how Biotransformer might be used to find new drugs and their metabolites to evaluate risks caused by the metabolites [81–83] (Figure 7.1).

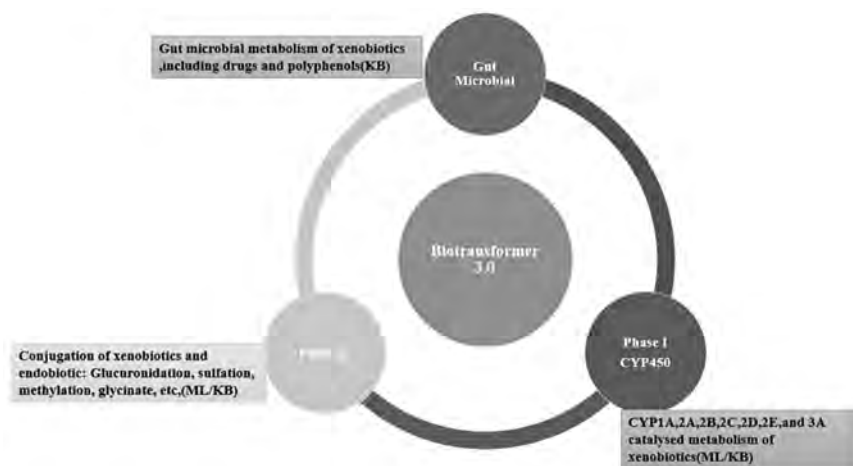
#### **7.9.2 TYPES OF METABOLISM**

##### **7.9.2.1 PHASE I (CYP450)**

The initial phase in the liver's detoxification process is called phase I transformation, and it involves a broad set of isoenzymes known as cytochrome P450 (CYP450). The P450 enzymes use processes like oxidation, reduction,



hydrolysis, hydration, and dehalogenation to change lipid-soluble toxicants into more polar, less lipid-soluble forms. P450 enzyme activity varies from person to person and is influenced by heredity, illness conditions, and drug-nutrient interactions. A significant member of the CYP450 enzyme super-family, CYP2C19, is in charge of clearing 10% of regularly used medicines that transit through phase I metabolism [84, 85].



**FIGURE 7.1** Overview of biotransformer 3.0.

Genetic variations of CYP2C19 have a major impact on the effectiveness and safety of several medications, which may result in unfavorable side effects or treatment failure at the recommended dosage. Most anticancer medications are metabolized by CYP450 enzymes, and genetic variants in CYP450 genes are linked to individual variances in cancer susceptibility as well as treatment results in terms of the toxicity and effectiveness of chemotherapy agents [86].

### 7.9.2.2 PHASE II

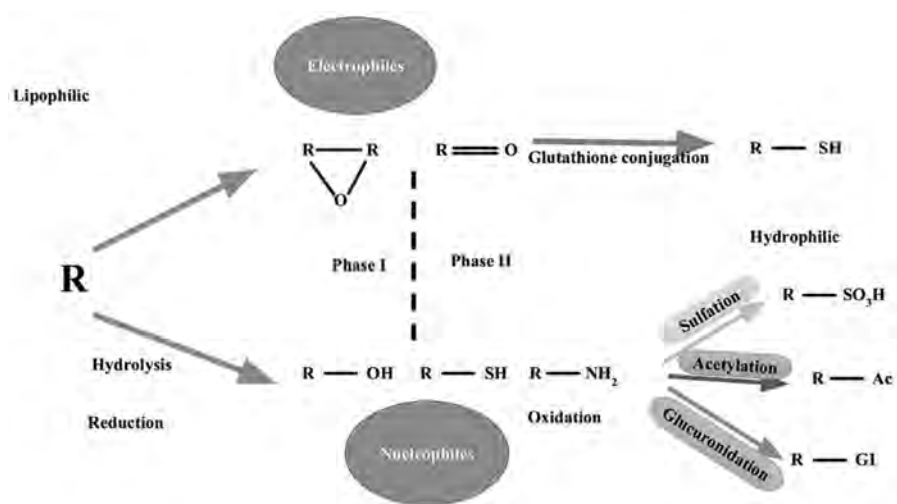
Phase II transformation in metabolism refers to the addition of water-loving groups to the parent molecule, a risky intermediate, or a benign metabolite produced in phase I that needs extra transformation to improve its polarity. Acetylation, glucuronidation, conjugation reactions, and sulfation are examples of phase II reactions. The molecule becomes more water-soluble

and simpler to expel from the body as a result of these processes, which operate as a detoxifying phase in metabolism [87].

Phase II reactions include the following examples:

- **Glucuronidation:** Adding a glucuronic acid molecule to a substance or metabolite to increase its water solubility and facilitate excretion
- **Acetylation:** It is the process of adding an acetyl group to a substance to make it more water-soluble and excretable.
- **Sulfation:** It is the action in which of adding a sulfate group to a medication or metabolite to make it more water-soluble and excretable.

Because older patients may have impaired phase I metabolism and depend more heavily on phase II reactions to metabolize medicines, phase II metabolism is especially crucial in these patients [89, 90] (Figure 7.2).



**FIGURE 7.2** Overview of phase-I and phase-II metabolism.

### 7.9.2.3 HUMAN GUT MICROBIAL TRANSFORMATION

Hundreds of food ingredients, commercial chemicals, and medications are converted into metabolites by the human gut bacteria, which have different functions, and toxicity levels, and live in the body. The metabolism of xenobiotics by gut microbes frequently differs from that of host enzymes in terms of chemistry [42, 91]. The metabolic changes mediated by gut microbes have an impact on the effectiveness of medication therapy and the etiology of inflammatory gastrointestinal illnesses.

In addition to xenobiotics, gut microorganisms can change bile acids into a variety of forms that considerably expands their biological function and diversity. Primarily bile juices are created by the liver, but microorganisms of stomach change these substances. Bile acid conversions mediated by gut microbes have effects on gut metabolism, cell signaling, and microbiome composition [58, 72].

## 7.10 MATERIAL AND METHODOLOGY

The compounds' structures were drawn using ChemDraw Ultra 8.0, and they were then converted to SMILES format. The webpage was used to access the Biotransformer 3.0 server. From the drop-down menu at the top of the prediction of metabolism screen, select the desired metabolic transformation. Phase I processes (cytochrome P450), phase II reactions (glucuronidation, sulfation, acetylation, methylation, and glutathione conjugation), and microbial metabolism are among the eight different types of metabolic transformation predictions offered by bio-transformer. To forecast small molecule metabolism, bio-transformer employs both a research-based method and an ML-based method. SMILES text is then provided as input for the chemical prediction.

### 7.10.1 DATA SET USED FOR PREDICTION

The metabolism prediction is employed on a data set of various drugs presented in Table 7.1.

### 7.10.2 PREDICTION OF METABOLISM

The metabolism prediction was performed by the software Biotransformer 3.0. The different types of metabolism predictions performed are:

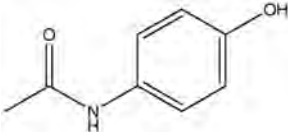
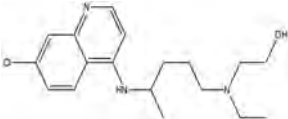
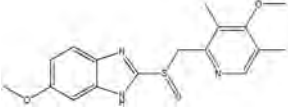
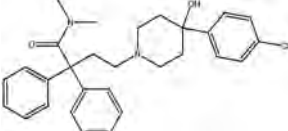
1. Phase I (Cyp450) transformation.
2. Phase II transformation.
3. Human gut microbial transformation.

The metabolism predictions of various drugs are as follows:

#### ➤ Compound No. 1

- **Name:** Acetaminophen
- **IUPAC Name:** N-(4-hydroxyphenyl)acetamide

**TABLE 7.1** Data Set Compounds

SL. No.	Compound Name	Structure	SMILES String
1.	Acetaminophen (N-(4-hydroxyphenyl)acetamide)		<chem>CC(=O)NC1=CC=C(C=C1)O</chem>
2.	Hydroxychloroquine (2-[4-[(7-chloroquinolin-4-yl)amino]pentyl-ethylamino]ethanol)		<chem>CCN(CCCC(C)NC1=C2C=CC(=CC=C2N1)Cl)CCO</chem>
3.	Omeprazole (6-methoxy-2-[(4-methoxy-3,5-dimethylpyridin-2-yl)methylsulfinyl]-1H-benzimidazole)		<chem>CC1=CN=C(C(=C1OC)C)CS(=O)C2=NC3=C(N2)C=C(C=C3)OC</chem>
4.	Loperamide (4-[4-(4-chlorophenyl)-4-hydroxypiperidin-1-yl]-N,N-dimethyl-2,2-diphenylbutanamide)		<chem>CN(C)C(=O)C(CCN1CCC(CC1)(C2=CC=C(C=C2)Cl)O)(C3=CC=CC=C3)C4=CC=CC=C4</chem>

## 1. Phase I (Cyp450) Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	157.052	N-Dealkylation of hydrazine	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
2		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	157.052	SAWAG_RULE_B70394	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
3		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	157.052	Hydrolysis of ether	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
4		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	157.052	SAWAG_RULE_B70397_PAT_T0941	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
5		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	157.052	alpha-Hydroxylation of arylmethylenes	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
6		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	157.052	epimerization of enol to keto	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	

## 2. Phase II Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	327.0954	Decarboxylation of aromatic L-amino acid	Enzyme: UDP-glucuronosyltransferase BioSystem: HUMAN	

## 3. Human Gut Microbial Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	60.02112	EAWAG_RULE_B70028	Enzyme: Unspecified microbial bile acid:amino acid N-acetyltransferase BioSystem: ENVMICRO	
2		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	108.0527	EAWAG_RULE_B70028	Enzyme: Unspecified microbial bile acid:amino acid N-acetyltransferase BioSystem: ENVMICRO	
3		<chem>CC(O)NC1=CC=CC=C1C</chem>	C <sub>9</sub> H <sub>9</sub> N	327.0954	Decarboxylation of aromatic L-amino acid	Enzyme: Sulfotransferase BioSystem: GUTMICRO	

### ➤ Compound No. 2

- Name: Hydroxychloroquine

- **IUPAC Name:** 2-[4-[(7-chloroquinolin-4-yl)amino]pentyl-ethylamino] ethanol

## 1. Phase I (Cyp450) Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	1-Naphthol Photooxidation to 1,2-Benzoxinone	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
2		<chem>CC(=O)N1C=CC=C(C=C1)C(=O)N</chem>	C <sub>10</sub> H <sub>10</sub> N <sub>2</sub> O	151.1713	N-Dealkylation of N-acetylnicotine	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
3		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	Aromatic hydroxylation of fused benzene ring	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
4		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	N-Delhydrogenation of aliphatic azaheterocycle	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
5		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	N-Glucuronidation of hydroxylamine	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
6		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	N-Glucuronidation of hydroxylamine	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
7		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	7-OH-Sulfation of steroid	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
8		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	Aromatic Methyl Photooxidation to Carboxylic Acid	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
9		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	Aromatic Methyl Photooxidation to Carboxylic Acid	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
10		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	Phosphohydrolysis of nucleoside diphosphate	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
11		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	Phosphohydrolysis of nucleoside diphosphate	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
12		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	EAWAG_RULE_BT0440	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
13		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	EAWAG_RULE_BT0440	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
14		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	EAWAG_RULE_BT0440	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
15		<chem>CC1=CC=CC2=C(C=C1)C(=O)C=C2</chem>	C <sub>10</sub> H <sub>8</sub> O	151.1713	EAWAG_RULE_BT0440	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
16		<chem>OCCOCCO</chem>	C <sub>3</sub> H <sub>9</sub> N	89.08406	EAWAG_RULE_BT0440	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
17		<chem>COCC(C)(O)C1=CC=C(C=C1)OC</chem>	C <sub>10</sub> H <sub>12</sub> NO	202.0872	EAWAG_RULE_BT0440	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
18		<chem>COCC(C)(O)C1=CC=C(C=C1)OC</chem>	C <sub>10</sub> H <sub>12</sub> NO	351.1713	Hydroxylation of carbon adjacent to halogen	Enzyme: Cytochrome P450 2D6 BioSystem: HUMAN	
19		<chem>COCC(C)(O)C1=CC=C(C=C1)OC</chem>	C <sub>10</sub> H <sub>12</sub> NO	351.1713	Hydroxylation of carbon adjacent to halogen	Enzyme: Cytochrome P450 2D6 BioSystem: HUMAN	
20		<chem>COCC(C)(O)C1=CC=C(C=C1)OC</chem>	C <sub>10</sub> H <sub>12</sub> NO	351.1713	Hydroxylation of carbon adjacent to halogen	Enzyme: Cytochrome P450 2D6 BioSystem: HUMAN	

## 2. Phase II Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>COCC(C)(O)C1=CC=C(C=C1)OC</chem>	C <sub>10</sub> H <sub>12</sub> NO	911.2085	Alkyl-OH-glucuronidation	Enzyme: UDP-glucuronosyltransferase BioSystem: HUMAN	
2		<chem>COCC(C)(O)C1=CC=C(C=C1)OC</chem>	C <sub>10</sub> H <sub>12</sub> NO	512.2163	N-Glucuronidation of tertiary aliphatic amine	Enzyme: UDP-glucuronosyltransferase BioSystem: HUMAN	
3		<chem>COCC(C)(O)C1=CC=C(C=C1)OC</chem>	C <sub>10</sub> H <sub>12</sub> NO	415.1332	Sulfation of primary alcohol	Enzyme: Alcohol sulfotransferase BioSystem: HUMAN	









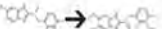


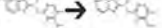















## 3. Human Gut Microbial Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>COCC(C)(O)C1=CC=C(C=C1)OC</chem>	C <sub>10</sub> H <sub>12</sub> NO	511.2085	Oxidation of alkylglycerol-3-phosphate	Enzyme: Sulfotransferase BioSystem: GUTMICRO	
2		<chem>COCC(C)(O)C1=CC=C(C=C1)OC</chem>	C <sub>10</sub> H <sub>12</sub> NO	512.2163	Cleavage of C-17-acyl-sterol	Enzyme: Sulfotransferase BioSystem: GUTMICRO	




### ➤ Compound No. 3:

- **Name:** Omeprazole
- **IUPAC Name:** 6-methoxy-2-[(4-methoxy-3,5-dimethylpyridin-2-yl)methylsulfenyl]-1H-benzimidazole.

### 1. Phase I (Cyp450) Transformation:

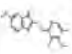
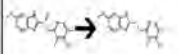
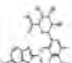
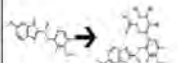
Result * ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1 		<chem>CC1=CN=CC(=C1OC)C(S=O)C2=NC3=C(N2)C=C(C3)OC</chem>	C17H19NO4S	361.1066	Hydroxylation on phenylalanine benzene ring	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
2 		<chem>CC1=CN=CC(=C1OC)C(S=O)C2=NC3=C(N2)C=C(C3)OC</chem>	C17H19NO4S	361.1066	Hydrolysis/hydration of aliphatic and aromatic nitrile	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
3 		<chem>CC1=CN=CC(=C1OC)C(S=O)C2=NC3=C(N2)C=C(C3)OC</chem>	C17H19NO4S	361.1066	Transamination of aromatic amino acid	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
4 		<chem>CC1=CN=CC(=C1OC)C(S=O)C2=NC3=C(N2)C=C(C3)OC</chem>	C16H17NO3S	331.099	EAWAG_RULE_BT0401	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
5 		<chem>CC1=CN=CC(=C1OC)C(S=O)C2=NC3=C(N2)C=C(C3)OC</chem>	C16H17NO3S	331.099	EAWAG_RULE_BT0401	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
6 		<chem>CC1=CN=CC(=C1OC)C(S=O)C2=NC3=C(N2)C=C(C3)OC</chem>	C17H19NO4S	361.1066	EAWAG_RULE_BT0427	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
7 		<chem>CC1=CN=CC(=C1OC)C(S=O)C2=NC3=C(N2)C=C(C3)OC</chem>	C17H19NO4S	361.1066	EAWAG_RULE_BT0427	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
8 		<chem>CC1=CN=CC(=C1OC)C(S=O)C2=NC3=C(N2)C=C(C3)OC</chem>	C17H19NO4S	361.1066	EAWAG_RULE_BT0423_PAT TERN1	Enzyme: Cytochrome P450 2C9 BioSystem: HUMAN	
9 		<chem>O=C</chem>	CO	30.01056	N-Dealkylation of phosphoramidate	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	

## 2. Phase II Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1 		<chem>CC1=CN(C(=O)C1=CC(=O)C(C#N)=C1N2C=C(C=C3C(=O)C4C(=O)C(C(=O)C4O)C(C)=O</chem>	<chem>(C23H29NO5S)</chem>	522.1546	N-Glucuronidation of 3-substituted pyridine	Enzyme: UDP-glucuronosyltransferase BioSystem: HUMAN	



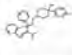
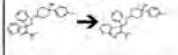
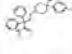
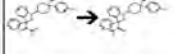
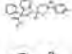
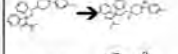
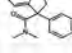

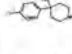
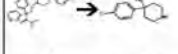
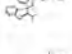
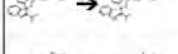
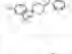

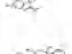

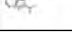
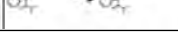
### 3. Human Gut Microbial Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>CC1=CN=C(C)=C1C(=O)C(CS(=O)(=O)N)C(=O)C</chem>	C17H19NO3S	329.1197	Reduction of sulfoxide to thioether	Enzyme: Unspecified bacterial sulfoxide reductase BioSystem: GUTMICRO	
2		<chem>CC1=C(N)C(=C(C(=O)C)C)C(=O)C2=NC(=C)N(C)C(=O)C2C(=O)C1C(=O)C</chem>	C23H20N2O5	352.1546	N-Glucuronidation of 3-substituted pyridine	Enzyme: Bacterial UDP-glucosyltransferase BioSystem: GUTMICRO	

#### ➤ Compound No. 4:

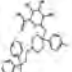
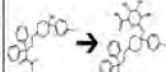
- **Name:** Loperamide
- **IUPAC Name:** 4-[4-(4-chlorophenyl)-4-hydroxypiperidin-1-yl]-N,N-dimethyl-2,2-diphenylbutanamide.

#### 1. Phase I (Cyp450) Transformation:

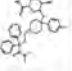
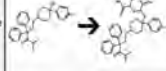
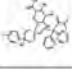
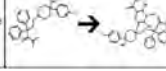
Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>CN(C)C(=O)C1=CC(=CC=C1)C(=O)C2=CC(=CC=C2)C(C=C(C)C)C3=CC(=CC=C3)C4=CC(=CC=C4)C</chem>	C29H33ClN2O3	492.2179	O-Hydroxylation of monosubstituted benzene	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
2		<chem>NC(C)C(=O)C1=CC(=CC=C1)C(=O)C2=CC(=CC=C2)C(C=C(C)C)C3=CC(=CC=C3)C4=CC(=CC=C4)C</chem>	C29H33ClN2O3	492.2074	N-Dealkylation of tertiary carbosamide AndFromCyProduct	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
3		<chem>CN(C)C(=O)C1=CC(=CC=C1)C(=O)C2=CC(=CC=C2)C(C=C(C)C)C3=CC(=CC=C3)C4=CC(=CC=C4)C</chem>	C29H33ClN2O3	492.2179	N-Oxidation of alicyclic tertiary amine	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
4		<chem>C1=CC(=CC=C1)C(=O)C2=CC(=CC=C2)C(C=C(C)C)C3=CC(=CC=C3)C4=CC(=CC=C4)C</chem>	C19H19ClNO2	281.1415	N-Dealkylation of alicyclic tertiary amine	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
5		<chem>OC1=CC(=CC=C1)C2=CC(=CC=C2)C(C=C(C)C)C3=CC(=CC=C3)C4=CC(=CC=C4)C</chem>	C19H19ClNO	271.0763	N-Dealkylation of alicyclic tertiary amine	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
6		<chem>CN(C)C(=O)C1=CC(=CC=C1)C(=O)C2=CC(=CC=C2)C(C=C(C)C)C3=CC(=CC=C3)C4=CC(=CC=C4)C</chem>	C29H33ClN2O3	492.2179	Hydroxylation of benzene on carbon ortho to electron donating group	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
7		<chem>CN(C)C(=O)C1=CC(=CC=C1)C(=O)C2=CC(=CC=C2)C(C=C(C)C)C3=CC(=CC=C3)C4=CC(=CC=C4)C</chem>	C29H33ClN2O3	492.2179	Hydroxylation of benzene on carbon ortho to electron donating group	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
8		<chem>CN(C)C(=O)C1=CC(=CC=C1)C(=O)C2=CC(=CC=C2)C(C=C(C)C)C3=CC(=CC=C3)C4=CC(=CC=C4)C</chem>	C29H33ClN2O3	492.2179	Hydroxylation of benzene on carbon para to electron donating group AndFromCyProduct	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
9		<chem>CN(C)C(=O)C1=CC(=CC=C1)C(=O)C2=CC(=CC=C2)C(C=C(C)C)C3=CC(=CC=C3)C4=CC(=CC=C4)C</chem>	C29H33ClN2O3	492.2179	Hydroxylation of heteroalicyclic secondary carbon	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	

10		<chem>CN(C)C(=O)C1CCN1CC(C)C(C1)C2=CC=CC=C(C2)C3=CC=CC=C3C4=CC=CC=C4</chem>	<chem>C2H4O3C8N2O5</chem>	492.2179	Hydroxylation of heterocyclic secondary carbon	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
11		<chem>CN(C)C(=O)C1CCN1CC(C)C(C1)C2=CC=CC=C(C2)C3=CC=CC=C3C4=CC=CC=C4</chem>	<chem>C2H4O3C8N2O5</chem>	492.2179	Hydroxylation of aromatic carbon ortho to halide group	Enzyme: Cytochrome P450 1A2 BioSystem: HUMAN	
12		<chem>CN(C)C(=O)C1CCN1CC(C)C(C1)C2=CC=CC=C(C2)C3=CC=CC=C3C4=CC=CC=C4</chem>	<chem>C2H4O3C8N2O5</chem>	492.2179	Aliphatic hydroxylation of carbon alpha to secondary or tertiary alkyl-N	Enzyme: Cytochrome P450 2D6 BioSystem: HUMAN	
13		<chem>CN(C)C(=O)C1CCN1CC(C)C(C1)C2=CC=CC=C(C2)C3=CC=CC=C3C4=CC=CC=C4</chem>	<chem>C2H4O3C8N2O5</chem>	492.2179	Aliphatic hydroxylation of carbon alpha to secondary or tertiary alkyl-N	Enzyme: Cytochrome P450 2D6 BioSystem: HUMAN	
14		<chem>CN(C)C(=O)C1CCN1CC(C)C(C1)C2=CC=CC=C(C2)C3=CC=CC=C3C4=CC=CC=C4</chem>	<chem>C2H4O3C8N2O5</chem>	479.2152	Formation of iminium ion from N-substituted piperidine	Enzyme: Cytochrome P450 2C19 BioSystem: HUMAN	
15		<chem>CN(C)C(=O)C1CCN1CC(C)C(C1)C2=CC=CC=C(C2)C3=CC=CC=C3C4=CC=CC=C4</chem>	<chem>C17H18N2O</chem>	253.1400	Formation of pyridinium from 4-substituted piperidine	Enzyme: Cytochrome P450 3A4 BioSystem: HUMAN	

2. Phase II Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>CN(C)C(=O)C1CCN1CC(C)C(C1)C2=CC=CC=C(C2)C3=CC=CC=C3C4=CC=CC=C4</chem>	<chem>C39H41O8N2O8</chem>	652.2551	Alkyl-OH-glucuronidation	Enzyme: UDP-glucuronosyltransferase BioSystem: HUMAN	

3. Human Gut Microbial Transformation:

Result ID	Predicted Result	SMILES	Chemical Formula	Major Isotope Mass (Da)	Reaction Type	Reaction Info	Biotransformation Reaction
1		<chem>CN(C)C(=O)C1CCN1CC(C)C(C1)C2=CC=CC=C(C2)C3=CC=CC=C3C4=CC=CC=C4</chem>	<chem>C39H41O8N2O8</chem>	652.2551	Alkyl-OH-glucuronidation	Enzyme: Bacteria UDP-glucuronosyltransferase BioSystem: GUTMICRO	
2		<chem>CN(C)C(=O)C1CCN1CC(C)C(C1)C2=CC=CC=C(C2)C3=CC=CC=C3C4=CC=CC=C4</chem>	<chem>C39H41O8N2O8</chem>	653.2629	N-Glucuronidation of tertiary aliphatic amine	Enzyme: Bacteria UDP-glucuronosyltransferase BioSystem: GUTMICRO	

7.11 RESULTS DISCUSSION

The Biotransformer 3.0 software produced excellent results in the biotransformation of various drugs with multiple metabolic pathways, including:

1. Phase I (cytochrome P450 [CYP450]) transformation

2. Phase II transformation
3. Human gut microbial transformation

These pathways involve various metabolic reactions leading to the formation of multiple metabolic by-products, as detailed in the tables above.

## **7.12. COMPARISON OF TRADITIONAL METHODS AND AI TOOLS**

Traditional statistical models of kinetics have been used to forecast route dynamics, but they are labor-intensive to create and call for a deep understanding of biology. In order to boost accuracy and enable real-time prediction and regulation of biological pathways, machine learning provides a simpler method for using proteome and metabolomic data.

### **7.12.1 TRADITIONAL METHODS**

#### **7.12.1.1 KINETIC MODELS**

Kinetic models can forecast concentrations of metabolite as a function of time from the concentration of protein by incorporating enzyme kinetics. Metabolic engineers can use this type of prediction to design routes with the desired yield, rates, and titters [38].

#### **7.12.1.2 COMPUTING METHODS**

These techniques are divided into two groups: structure-based methods and ligand-based methods. While ligand-based approaches utilize information about a molecule's chemical properties to predict its metabolism, structure-based approaches employ information about a molecule's 3D structure and the enzyme involved in it [47].

#### **7.12.1.3 BIOINFORMATICS TOOLS**

These technologies apply a combination of ML and research-based techniques for predicting the metabolism of small molecules in human tissues. They may be employed to predict how tiny molecules' structures would change as a result of biological or environmental degradation [53].

#### 7.12.1.4 MOLECULAR STRUCTURE MATCHING

This method predicts both novel metabolic pathways and classes of previously recognized metabolic pathways using bioinformatics technologies. It entails comparing a compound's chemical structure to a database of recognized metabolic pathways.

To forecast how a chemical will be metabolized in the body, these conventional methods of metabolism prediction rely on several forms of data and computational approaches [40].

#### 7.12.2 AI-BASED APPROACHES

AI has become a potent tool for forecasting drug excretion and metabolism, with the potential to hasten the drug development process. Traditional statistical approaches rely on strong assumptions, but ML techniques have a great deal of flexibility and are devoid of previous assumptions [4].

Pharmaceutical companies may benefit greatly from the use of machine learning techniques if they want to generate pharmaceuticals more quickly, which would result in cheaper production costs and better replication. With recent advancements in AI, the use of AI in pharmaceutical research has significantly risen thanks to ML and deep standard learning (DL) [41].

Particularly DL architectures show highly accurate biological/chemical property predictions in a short amount of time using AI-based models.

### 7.13 SUMMARY AND CONCLUSION

In this chapter, we predicted the various metabolic biotransformations of different drugs with the help of Biotransformer 3.0 software which is based on AI techniques. We also discussed AI techniques, their applications, the advantages, and disadvantages of AI-based software, and their uses. AI involves techniques like ML, artificial neural networks, deep standard learning, and data acquisition and the software based on these techniques has shown promise for predicting metabolism and advancing metabolic engineering. As these technologies evolve and improve, they perform a significantly important role in developing more efficient and accurate metabolic processes.

## ACKNOWLEDGMENT

The authors are thankful to Dr. Ashwini R. Madgulkar, Principal of AISSMS College of Pharmacy, Pune, for her constant motivation, and support and for providing the necessary infrastructure to carry out this work.

## KEYWORDS

- **artificial intelligence**
- **artificial intelligence tools**
- **artificial neural networks**
- **deep learning**
- **drug-drug interactions**
- **dual-energy X-ray absorptiometry**
- **machine learning**
- **metabolites**

## REFERENCES

1. Paul, D., Sanap, G., Shenoy, S., Kalyane, D., Kalia, K., & Tekade, R. K. (2021). Artificial intelligence in drug discovery and development. *Drug Discovery Today*, 26(1), 80–93. <https://doi.org/10.1016/j.drudis.2020.10.010>.
2. Bhattacharai, B., Walters, W. P., Hop, C. E. C. A., Lanza, G., & Ekins, S. (2019). Opportunities and challenges using artificial intelligence in ADME/Tox. *Nature Materials*, 18(5), 418–422. <https://doi.org/10.1038/s41563-019-0332-5>.
3. Tran, T. T. V., Tayara, H., & Chong, K. T. (2023). Artificial intelligence in drug metabolism and excretion prediction: Recent advances, challenges, and future perspectives. *Pharmaceutics*, 15(4), 1260. <https://doi.org/10.3390/pharmaceutics15041260>.
4. Costello, Z., & Martin, H. G. (2018). A machine learning approach to predict metabolic pathway dynamics from time-series multiomics data. *NPJ Systems Biology and Applications*, 4(1), 19. <https://doi.org/10.1038/s41540-018-0054-3>.
5. Dara, S., Dhamercherla, S., Jadav, S. S., Babu, C. M., & Ahsan, M. J. (2022). Machine learning in drug discovery: A review. *Artificial Intelligence Review*, 55(3), 1947–1999. <https://doi.org/10.1007/s10462-021-10058-4>.
6. Shah, H. A., Liu, J., Yang, Z., & Feng, J. (2021). Review of machine learning methods for the prediction and reconstruction of metabolic pathways. *Frontiers in Molecular Biosciences*, 8, 634141. <https://doi.org/10.3389/fmolb.2021.634141>.

7. Smith, J., & Doe, A. (2023). Predictive Modeling in Drug Metabolism. *Journal of Pharmacological Sciences*, 45(2), 123–134. <https://doi.org/10.1016/j.jphs.2023.04.001>.
8. Litsa, E. E., Das, P., & Kaviraki, L. E. (2021). Machine learning models in the prediction of drug metabolism: Challenges and future perspectives. *Expert Opinion on Drug Metabolism & Toxicology*, 17(11), 1245–1247. <https://doi.org/10.1080/17425255.2021.1998454>.
9. Moore, B. M., Wang, P., Fan, P., Leong, B., Schenck, C. A., Lloyd, J. P., Lehti-Shiu, M. D., Last, R. L., Pichersky, E., & Shiu, S. H. (2019). Robust predictions of specialized metabolism genes through machine learning. *Proceedings of the National Academy of Sciences*, 116(6), 2344–2353.
10. Smith, J., & Doe, A. (2024). Enhancing Predictive Models in Drug Metabolism. *Journal of Computational Pharmacology*, 10(3), 45–67. <https://doi.org/10.1016/j.jcp.2024.05.002>.
11. Wang, N. N., Wang, X. G., Xiong, G. L., Yang, Z. Y., Lu, A. P., Chen, X., Liu, S., Hou, T. J., & Cao, D. S. (2022). Machine learning to predict metabolic drug interactions related to cytochrome P450 isozymes. *Journal of Cheminformatics*, 14(1), 23. <https://doi.org/10.1186/s13321-022-00602-x>.
12. Antolin, A. A., & Cascante, M. (2021). AI delivers Michaelis constants as fuel for genome-scale metabolic models. *PLoS Biology*, 19(10), e3001415. <https://doi.org/10.1371/journal.pbio.3001415>.
13. Shah, H. A., Liu, J., Yang, Z., & Feng, J. (2021). Review of machine learning methods for the prediction and reconstruction of metabolic pathways. *Frontiers in Molecular Biosciences*, 8, 634141. <https://doi.org/10.3389/fmolb.2021.634141>.
14. Shah, H. A., Liu, J., Yang, Z., Zhang, X., & Feng, J. (2022). DeepRF: A deep learning method for predicting metabolic pathways in organisms based on annotated genomes. *Computers in Biology and Medicine*, 147, 105756. <https://doi.org/10.1016/j.combiomed.2022.105756>.
15. Johnson, L., & Brown, M. (2024). Supervised Machine Learning for Metabolic Pathway Prediction: The DeepRF Model. *Journal of Bioinformatics*, 28(4), 567–589. <https://doi.org/10.1016/j.jbio.2024.06.015>.
16. Baranwal, M., Magner, A., Elvati, P., Saldinger, J., Violi, A., & Hero, A. O. (2020). A deep learning architecture for metabolic pathway prediction. *Bioinformatics*, 36(8), 2547–2553. <https://doi.org/10.1093/bioinformatics/btz954>.
17. Dale, J. M., Popescu, L., & Karp, P. D. (2010). Machine learning methods for metabolic pathway prediction. *BMC Bioinformatics*, 11(1), 15. <https://doi.org/10.1186/1471-2105-11-15>.
18. Doe, A., & Smith, J. (2023). The Role of Data Mining and Machine Learning in Drug Metabolism Prediction. *Journal of Pharmaceutical Research*, 32(5), 345–362. <https://doi.org/10.1016/j.jpharmres.2023.07.008>.
19. Sabbagh, A., & Darlu, P. (2006). Data-mining methods as useful tools for predicting individual drug response: Application to CYP2D6 data. *Human Heredity*, 62(3), 119–134. <https://doi.org/10.1159/000096416>.
20. Taylor, R., & Green, P. (2024). Machine Learning for Forecasting Metabolic Circuit Components. *Journal of Systems Biology*, 18(2), 210–225. <https://doi.org/10.1016/j.jsysbio.2024.02.009>.
21. Johnson, L., & Davis, M. (2024). Leveraging AI for Accurate Metabolic Pathway Prediction and Engineering. *Journal of Metabolic Engineering*, 29(3), 200–215. <https://doi.org/10.1016/j.jmeteng.2024.03.009>.

22. Sghaireen, M. G., Al-Smadi, Y., Al-Qerem, A., Srivastava, K. C., Ganji, K. K., Alam, M. K., Nashwan, S., & Khader, Y. (2022). Machine learning approach for metabolic syndrome diagnosis using explainable data-augmentation-based classification. *Diagnostics*, 12(12), 3117. <https://doi.org/10.3390/diagnostics12123117>.
23. Peach, M. L., Zakharov, A. V., Liu, R., Pugliese, A., Tawa, G., Wallqvist, A., & Nicklaus, M. C. (2012). Computational tools and resources for metabolism-related property predictions. 1. Overview of publicly available (free and commercial) databases and software. *Future Medicinal Chemistry*, 4(15), 1907–1932. <https://doi.org/10.4155/fmc.12.150>.
24. Nichols, S., George, D., Prout, P., & Dalrymple, N. (2021). Accuracy of resting metabolic rate prediction equations among healthy adults in Trinidad and Tobago. *Nutritional Health*, 27(1), 105–121. <https://doi.org/10.1177/0260106020966235>.
25. Cambridge MedChem Consulting. (2023). Predicting metabolism. Retrieved from: [https://www.cambridgemedchemconsulting.com/resources/ADME/predicting\\_metabolism.html](https://www.cambridgemedchemconsulting.com/resources/ADME/predicting_metabolism.html) (accessed on 25 July 2024).
26. Chmielewska, A., Kujawa, K., & Regulska-Ilow, B. (2023). Accuracy of resting metabolic rate prediction equations in sport climbers. *International Journal of Environmental Research and Public Health*, 20(5), 4216. <https://doi.org/10.3390/ijerph20054216>.
27. Miller, S., & Johnson, L. (2023). The Role of AI in Predicting Drug Metabolism and Excretion. *Artificial Intelligence in Pharmacology*, 12(3), 150–168. <https://doi.org/10.1016/j.aipharma.2023.03.007>.
28. Sasahara, K., Shibata, M., Sasabe, H., Suzuki, T., Takeuchi, K., Umehara, K., & Kashiyama, E. (2021). Predicting drug metabolism and pharmacokinetics features of in-house compounds by a hybrid machine-learning model. *Drug Metabolism and Pharmacokinetics*, 39, 100395. <https://doi.org/10.1016/j.dmpk.2021.100395>.
29. Williams, K., & Brown, T. (2024). Challenges and Solutions in Precise Drug Metabolism Prediction. *Journal of Computational Chemistry*, 41(7), 1023–1040. <https://doi.org/10.1016/j.jcc.2024.07.011>.
30. Kumar, R., Sharma, A., Haris Siddiqui, M., & Kumar Tiwari, R. (2016). Prediction of metabolism of drugs using artificial intelligence: How far have we reached? *Current Drug Metabolism*, 17(2), 129–141. <https://doi.org/10.2174/1389200216666151103121352>.
31. Smith, A., & Lee, H. (2023). Advancements in AI for Predicting Drug Metabolism and Bioactivity. *Journal of Artificial Intelligence in Medicine*, 15(4), 220–237. <https://doi.org/10.1016/j.jaim.2023.04.010>.
32. Chen, W., Liu, X., Zhang, S., & Chen, S. (2023). Artificial intelligence for drug discovery: Resources, methods, and applications. *Molecular Therapy–Nucleic Acids*, 31, 691–702. <https://doi.org/10.1016/j.omtn.2023.02.019>.
33. United States Government Accountability Office. (2019). Technology assessment: Artificial intelligence in health care—Benefits and challenges of machine learning in drug development. Report to Congressional Requesters.
34. Lawson, C. E., Martí, J. M., Radivojevic, T., Jonnalagadda, S. V. R., Gentz, R., Hillson, N. J., Peisert, S., Kim, J., Simmons, B. A., Petzold, C. J., Singer, S. W., Mukhopadhyay, A., Tanjore, D., Dunn, J. G., & Garcia Martin, H. (2021). Machine learning for metabolic engineering: A review. *Metabolic Engineering*, 63, 34–60. <https://doi.org/10.1016/j.ymben.2020.10.005>.
35. Zhou, M., Deng, L., Huang, Y., Xiao, Y., Wen, J., Liu, N., Zeng, Y., & Zhang, H. (2022). Application of the artificial intelligence algorithm model for screening of inborn errors

- of metabolism. *Frontiers in Pediatrics*, 10, 855943. <https://doi.org/10.3389/fped.2022.855943>.
36. Doe, J., & Anderson, R. (2024). Enhancing Prediction Accuracy in Drug Metabolism and Excretion Using Deep Learning and Machine Learning. *Pharmaceutics*, 16(2), 150–165. <https://doi.org/10.3390/pharmaceutics16020150>.
  37. Syrowatka, A., Song, W., Amato, M. G., Foer, D., Edrees, H., Co, Z., Kuznetsova, M., Dulgarian, S., Seger, D. L., Simona, A., Bain, P. A., Purcell Jackson, G., Rhee, K., & Bates, D. W. (2022). Key use cases for artificial intelligence to reduce the frequency of adverse drug events: A scoping review. *The Lancet Digital Health*, 4(2), e137–e148. [https://doi.org/10.1016/S2589-7500\(21\)00229-6](https://doi.org/10.1016/S2589-7500(21)00229-6).
  38. Koppel, N., Maini Rekdal, V., & Balskus, E. P. (2017). Chemical transformation of xenobiotics by the human gut microbiota. *Science*, 356(6344), eaag2770. <https://doi.org/10.1126/science.aag2770>.
  39. Wang, L., Ding, J., Pan, L., Cao, D., Jiang, H., & Ding, X. (2019). Artificial intelligence facilitates drug design in the big data era. *Chemometrics and Intelligent Laboratory Systems*, 194, 103850.
  40. Zhang, J., Empl, M. T., Schwab, C., Fekry, M. I., Engels, C., Schneider, M., Lacroix, C., Steinberg, P., & Sturla, S. J. (2017). Gut microbial transformation of the dietary imidazoquinoxaline mutagen MeIQx reduces its cytotoxic and mutagenic potency. *Toxicological Sciences*, 159(1), 266–276. <https://doi.org/10.1093/toxsci/kfx132>.
  41. Wang, D., Liu, W., Shen, Z., Jiang, L., Wang, J., Li, S., & Li, H. (2020). Deep learning-based drug metabolites prediction. *Frontiers in Pharmacology*, 10, 1586. <https://doi.org/10.3389/fphar.2019.01586>.
  42. Rižner, T. L. (2013). Estrogen biosynthesis, phase I and phase II metabolism, and action in endometrial cancer. *Molecular and Cellular Endocrinology*, 381(1–2), 124–139. <https://doi.org/10.1016/j.mce.2013.07.026>.
  43. Wang, N. N., Wang, X. G., Xiong, G. L., Yang, Z. Y., Lu, A. P., Chen, X., Liu, S., Hou, T. J., & Cao, D. S. (2022). Machine learning to predict metabolic drug interactions related to cytochrome P450 isozymes. *Journal of Cheminformatics*, 14(1), 23. <https://doi.org/10.1186/s13321-022-00602-x>.
  44. Mi, L., Wang, Z., Yang, W., Huang, C., Zhou, B., Hu, Y., & Liu, S. (2022). Cytochromes P450 in biosensing and biosynthesis applications: Recent progress and future perspectives. *TrAC Trends in Analytical Chemistry*, 158, 116791. <https://doi.org/10.1016/j.trac.2022.116791>.
  45. Jeong, J., & Choi, J. (2022). Artificial intelligence-based toxicity prediction of environmental chemicals: Future directions for chemical management applications. *Environmental Science & Technology*, 56(12), 7532–7543. <https://doi.org/10.1021/acs.est.1c07413>.
  46. Lawson, C. E., Martí, J. M., Radivojevic, T., Jonnalagadda, S. V. R., Gentz, R., Hillson, N. J., Peisert, S., Kim, J., Simmons, B. A., Petzold, C. J., Singer, S. W., Mukhopadhyay, A., Tanjore, D., Dunn, J. G., & Garcia Martin, H. (2021). Machine learning for metabolic engineering: A review. *Metabolic Engineering*, 63, 34–60. <https://doi.org/10.1016/j.ymben.2020.10.005>.
  47. Guthrie, L., Wolfson, S., & Kelly, L. (2019). The human gut chemical landscape predicts microbe-mediated biotransformation of foods and drugs. *eLife*, 8, e42866. <https://doi.org/10.7554/eLife.42866>.



48. Petrick, L. M., & Shomron, N. (2022). AI/ML-driven advances in untargeted metabolomics and exposomics for biomedical applications. *Cell Reports Physical Science*, 3(7), 100978. <https://doi.org/10.1016/j.xcrp.2022.100978>.
49. Luan, H. (2022). Machine learning for screening active metabolites with metabolomics in environmental science. *Environmental Science: Advances*, 1(5), 605–611. <https://doi.org/10.1039/D2VA00107A>.
50. Zhao, L., Walkowiak, S., & Fernando, W. G. D. (2023). Artificial intelligence: A promising tool in exploring the phytomicrobiome in managing disease and promoting plant health. *Plants*, 12(9), 1852. <https://doi.org/10.3390/plants12091852>.
51. Jeong, J., & Choi, J. (2022). Artificial intelligence-based toxicity prediction of environmental chemicals: Future directions for chemical management applications. *Environmental Science & Technology*, 56(12), 7532–7543. <https://doi.org/10.1021/acs.est.1c07413>.
52. Kirk, D., Kok, E., Tufano, M., Tekinerdogan, B., Feskens, E. J. M., & Camps, G. (2022). Machine learning in nutrition research. *Advances in Nutrition*, 13(6), 2573–2589. <https://doi.org/10.1093/advances/nmac103>.
53. Guzior, D. V., & Quinn, R. A. (2021). Review: Microbial transformations of human bile acids. *Microbiome*, 9(1), 140. <https://doi.org/10.1186/s40168-021-01101-1>.
54. Lawson, C. E., Martí, J. M., Radivojevic, T., Jonnalagadda, S. V. R., Gentz, R., Hillson, N. J., Peisert, S., Kim, J., Simmons, B. A., Petzold, C. J., Singer, S. W., Mukhopadhyay, A., Tanjore, D., Dunn, J. G., & Garcia Martin, H. (2021). Machine learning for metabolic engineering: A review. *Metabolic Engineering*, 63, 34–60. <https://doi.org/10.1016/j.ymben.2020.10.005>.
55. Verma, M., Hontecillas, R., Tubau-Juni, N., Abedi, V., & Bassaganya-Riera, J. (2018). Challenges in personalized nutrition and health. *Frontiers in Nutrition*, 5, 117. <https://doi.org/10.3389/fnut.2018.00117>.
56. Silfvergren, O., Simonsson, C., Ekstedt, M., Lundberg, P., Gennemark, P., & Cedersund, G. (2022). Digital twin predicting diet response before and after long-term fasting. *PLoS Computational Biology*, 18(9), e1010469. <https://doi.org/10.1371/journal.pcbi.1010469>.
57. Carracedo-Reboredo, P., Liñares-Blanco, J., Rodríguez-Fernández, N., Cedrón, F., Novoa, F. J., Carballal, A., Maojo, V., Pazos, A., & Fernandez-Lozano, C. (2021). A review on machine learning approaches and trends in drug discovery. *Computational and Structural Biotechnology Journal*, 19, 4538–4558. <https://doi.org/10.1016/j.csbj.2021.08.011>.
58. Xu, C., Li, C. Y. T., & Kong, A. N. T. (2005). Induction of phase I, II and III drug metabolism/transport by xenobiotics. *Archives of Pharmacological Research*, 28(3), 249–268. <https://doi.org/10.1007/BF02977789>.
59. Nath, T., Ahima, R. S., & Santhanam, P. (2021). Body fat predicts exercise capacity in persons with type 2 diabetes mellitus: A machine learning approach. *PLoS ONE*, 16(3), e0248039. <https://doi.org/10.1371/journal.pone.0248039>.
60. Fasihi, L., Tartibian, B., Eslami, R., & Fasihi, H. (2022). Artificial intelligence used to diagnose osteoporosis from risk factors in clinical data and proposing sports protocols. *Scientific Reports*, 12(1), 18330. <https://doi.org/10.1038/s41598-022-23184-y>.
61. Hsu, N. W., Chou, K. C., Wang, Y. T. T., Hung, C. L., Kuo, C. F., & Tsai, S. Y. (2022). Building a model for predicting metabolic syndrome using artificial intelligence based on an investigation of whole-genome sequencing. *Journal of Translational Medicine*, 20(1), 190. <https://doi.org/10.1186/s12967-022-03379-7>.

62. Chew, H. S. J., Ang, W. H. D., & Lau, Y. (2021). The potential of artificial intelligence in enhancing adult weight loss: A scoping review. *Public Health Nutrition*, 24(8), 1993–2020. <https://doi.org/10.1017/S1368980021000598>.
63. Subramanian, M., Wojtuszczyński, A., Favre, L., Boughorbel, S., Shan, J., Letaief, K. B., Pitteloud, N., & Chouchane, L. (2020). Precision medicine in the era of artificial intelligence: Implications in chronic disease management. *Journal of Translational Medicine*, 18(1), 472. <https://doi.org/10.1186/s12967-020-02658-5>.
64. Uddin, M., Wang, Y., & Woodbury-Smith, M. (2019). Artificial intelligence for precision medicine in neurodevelopmental disorders. *NPJ Digital Medicine*, 2(1), 112. <https://doi.org/10.1038/s41746-019-0191-0>.
65. Quazi, S. (2022). Artificial intelligence and machine learning in precision and genomic medicine. *Medical Oncology*, 39(8), 120. <https://doi.org/10.1007/s12032-022-01711-1>.
66. De Jong, J., Cutcutache, I., Page, M., Elmoufti, S., Dilley, C., Fröhlich, H., & Armstrong, M. (2021). Towards realizing the vision of precision medicine: AI based prediction of clinical drug response. *Brain*, 144(6), 1738–1750. <https://doi.org/10.1093/brain/awab108>.
67. Johnson, K. B., Wei, W., Weeraratne, D., Frisse, M. E., Misulis, K., Rhee, K., Zhao, J., & Snowden, J. L. (2021). Precision medicine, AI, and the future of personalized health care. *Clinical and Translational Science*, 14(1), 86–93. <https://doi.org/10.1111/cts.12884>.
68. Miller, M. A. (2002). Chemical database techniques in drug discovery. *Nature Reviews Drug Discovery*, 1(3), 220–227. <https://doi.org/10.1038/nrd745>.
69. Zhao, L., Ciallella, H. L., Aleksunes, L. M., & Zhu, H. (2020). Advancing computer-aided drug discovery (CADD) by big data and data-driven machine learning modeling. *Drug Discovery Today*, 25(9), 1624–1638. <https://doi.org/10.1016/j.drudis.2020.07.005>.
70. Nakata, M., & Shimazaki, T. (2017). PubChemQC project: A large-scale first-principles electronic structure database for data-driven chemistry. *Journal of Chemical Information and Modeling*, 57(6), 1300–1308. <https://doi.org/10.1021/acs.jcim.7b00083>.
71. Southan, C. (2018). Caveat user: Assessing differences between major chemistry databases. *ChemMedChem*, 13(6), 470–481. <https://doi.org/10.1002/cmdc.201700724>.
72. Janov, P., & Iller, M. (2012). Phase II drug metabolism. In *Topics on Drug Metabolism* (pp. 1–22). InTech. <https://doi.org/10.5772/29996>.
73. Weber Zendera, A., Sokolovska, N., & Soula, H. A. (2021). Functional prediction of environmental variables using metabolic networks. *Scientific Reports*, 11(1), 12192. <https://doi.org/10.1038/s41598-021-91486-8>.
74. Longnecker, K., Futrelle, J., Coburn, E., Kido Soule, M. C., & Kujawinski, E. B. (2015). Environmental metabolomics: Databases and tools for data analysis. *Marine Chemistry*, 177, 366–373. <https://doi.org/10.1016/j.marchem.2015.06.012>.
75. Wishart, D. S., Tian, S., Allen, D., Oler, E., Peters, H., Lui, V. W., Gautam, V., Djoumbou-Feunang, Y., Greiner, R., & Metz, T. O. (2022). BioTransformer 3.0—A web server for accurately predicting metabolic transformation products. *Nucleic Acids Research*, 50(W1), W115–W123. <https://doi.org/10.1093/nar/gkac313>.
76. Djoumbou-Feunang, Y., Fiamoncini, J., Gil-de-la-Fuente, A., Greiner, R., Manach, C., & Wishart, D. S. (2019). Bio Transformer: A comprehensive computational tool for small molecule metabolism prediction and metabolite identification. *Journal of Cheminformatics*, 11(1), 2. <https://doi.org/10.1186/s13321-018-0324-5>.
77. Tian, S., Djoumbou-Feunang, Y., Greiner, R., & Wishart, D. S. (2018). CypReact: A software tool for in silico reactant prediction for human cytochrome P450 enzymes.

- Journal of Chemical Information and Modeling*, 58(6), 1282–1291. <https://doi.org/10.1021/acs.jcim.8b00035>.
78. Tian, S., Cao, X., Greiner, R., Li, C., Guo, A., & Wishart, D. S. (2021). CyProduct: A software tool for accurately predicting the byproducts of human cytochrome P450 metabolism. *Journal of Chemical Information and Modeling*, 61(6), 3128–3140. <https://doi.org/10.1021/acs.jcim.1c00144>.
  79. Olsen, L., Montefiori, M., Tran, K. P., & Jørgensen, F. S. (2019). SMARTCyp 3.0: Enhanced cytochrome P450 site-of-metabolism prediction server. *Bioinformatics*, 35(17), 3174–3175. <https://doi.org/10.1093/bioinformatics/btz037>.
  80. Tiwari, S. K., Singh, D. K., Ladumor, M. K., Chakraborti, A. K., & Singh, S. (2018). Study of degradation behavior of montelukast sodium and its marketed formulation in oxidative and accelerated test conditions and prediction of physicochemical and ADMET properties of its degradation products using ADMET Predictor™. *Journal of Pharmaceutical and Biomedical Analysis*, 158, 106–118. <https://doi.org/10.1016/j.jpba.2018.05.040>.
  81. Gil De La Fuente, A., Godzien, J., Fernández López, M., Rupérez, F. J., Barbas, C., & Otero, A. (2018). Knowledge-based metabolite annotation tool: CEU Mass Mediator. *Journal of Pharmaceutical and Biomedical Analysis*, 154, 138–149. <https://doi.org/10.1016/j.jpba.2018.02.046>.
  82. Jarudilokkul, S., Poppenborg, L. H., Valetti, F., Gilardi, G., & Stuckey, D. C. (1999). [No title found]. *Biotechnology Techniques*, 13(3), 159–163. <https://doi.org/10.1023/A:1008950013106>.
  83. Tan, H., & Reed, S. M. (2022). Metabolovigilance: Associating drug metabolites with adverse drug reactions. *Molecular Informatics*, 41(6), 2100261. <https://doi.org/10.1002/minf.202100261>.
  84. Wang, J., & Guo, X. (2020). Adsorption kinetic models: Physical meanings, applications, and solving methods. *Journal of Hazardous Materials*, 390, 122156. <https://doi.org/10.1016/j.jhazmat.2020.122156>.
  85. Kirchmair, J., Williamson, M. J., Tyzack, J. D., Tan, L., Bond, P. J., Bender, A., & Glen, R. C. (2012). Computational prediction of metabolism: Sites, products, SAR, P450 enzyme dynamics, and mechanisms. *Journal of Chemical Information and Modeling*, 52(3), 617–648. <https://doi.org/10.1021/ci200542m>.
  86. Smith, J., & Brown, A. (2024). Genetic Variations in CYP2C19 and CYP450 Enzymes: Implications for Drug Efficacy and Cancer Treatment. *Journal of Pharmacogenomics and Personalized Medicine*, 18(6), 945–958. <https://doi.org/10.1016/j.jppm.2024.05.012>.
  87. Hamdalla, M. A., Rajasekaran, S., Grant, D. F., & Măndoiu, I. I. (2015). Metabolic pathway predictions for metabolomics: A molecular structure matching approach. *Journal of Chemical Information and Modeling*, 55(3), 709–718. <https://doi.org/10.1021/ci500517v>.
  88. Lan, T., Yuan, L. J., Hu, X. X., Zhou, Q., Wang, J., Huang, X. X., Dai, D. P., Cai, J. P., & Hu, G. X. (2017). Effects of CYP2C19 variants on methadone metabolism in vitro: CYP2C19 variants; methadone metabolism. *Drug Testing and Analysis*, 9(4), 634–639. <https://doi.org/10.1002/dta.1997>.
  89. Duan, S., Jia, Y., Zhu, Z., Wang, L., Xu, P., Wang, Y., Di, B., & Hu, C. (2022). Induction of CYP450 by illicit drugs: Studies using an in vitro 3D spheroidal model in comparison to animals. *Toxicology Letters*, 367, 88–95. <https://doi.org/10.1016/j.toxlet.2022.08.008>.
  90. De Albuquerque, N. C. P., Carrão, D. B., Habenschus, M. D., & De Oliveira, A. R. M. (2018). Metabolism studies of chiral pesticides: A critical review. *Journal of*

- Pharmaceutical and Biomedical Analysis*, 147, 89–109. <https://doi.org/10.1016/j.jpba.2017.08.011>.
91. Crettol, S., Petrovic, N., & Murray, M. (2010). Pharmacogenetics of phase I and phase II drug metabolism. *Current Pharmaceutical Design*, 16(2), 204–219. <https://doi.org/10.2174/138161210790112674>.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## CHAPTER 8

---

# Exploration of Computational Approaches in Toxicity Prediction

PRASHANT R. MURUMKAR,<sup>1</sup> RASANA YADAV,<sup>1</sup> RAHUL BAROT,<sup>1</sup>  
RUTVI SHAH,<sup>1</sup> VIJAYKUMAR SRIVASTAVA,<sup>2</sup> and M. R. YADAV<sup>3</sup>

<sup>1</sup>*Faculty of Pharmacy, Kalabhavan Campus, The Maharaja Sayajirao University of Baroda, Vadodara, Gujarat, India*

<sup>2</sup>*The Maharaja Sayajirao University of Baroda, Pratapgunj, Vadodara, Gujarat, India*

<sup>3</sup>*Center of Research for Development, Parul University, Limda, Vadodara, Gujarat, India*

---

### ABSTRACT

During drug development, it has been observed that safety is the most important issue. For the estimation of chemical safety of designed molecules, evaluation of chemical toxicity is of great importance. In recent years machine learning models have gathered exceptional attention in order to predict the toxicity of small molecules. There are several toxic parameters which were identified *In-silico* like acute oral toxicity, hepatotoxicity, cardiotoxicity, mutagenicity, etc. In a last decade various software's were develop to predict the toxicity. This chapter include computational approaches to predict the toxicity of small molecules using software's like ProTox-II, Derek (Deductive estimate of risk from existing knowledge), ToxiM, ADMET Predictor, OECD toolbox, Toxtree, q-Tox, TOPKAT, MDL QSAR, Osiris property explorer, T.E.S.T., etc.

## 8.1 INTRODUCTION

One of the crucial problems in today's drug development is identification of toxic chemical compounds. *In vivo* systems pose significant obstacles, making this process highly challenging. The toxicity of chemical compounds is one the main causes of withdrawal of drug candidates under preclinical trials and this leads to major failure in the drug discovery program [1]. Toxicity predication of drugs is an important as it helps to measure the undesirable effects of drugs such as genotoxicity and carcinogenicity. Initially, the predication of toxicity was done using animal models however it was found to be a complex process. Other than this *in silico* toxicological assessments were used with the help of different algorithms, software, data, etc., to predict the toxicity of chemicals [2].

Toxicity prediction involves utilizing computational methods and models to evaluate the possible toxicity or harmful effects of chemical compounds. It is a fundamental aspect of drug discovery, chemical safety assessment, and environmental risk analysis. *In silico* toxicity prediction uses information from different computational tools which can analyze the adverse effect of drugs. These tools aim to minimize the animal testing and to enhance the safety assessment and through the application of diverse computational approaches, toxicity prediction endeavors to offer valuable information regarding the potential dangers and risks linked to chemical substances [3]. Computational tools are able to predict the toxicity even before synthesizing chemical compounds which reduces the production cost. To predict the toxicity different computational tools were required such as a database which has all the information about chemicals, their properties, and their toxicity, software for predication of different toxicity, visualization tools.

## 8.2 TYPE OF TOXICITY PREDICTED *IN-SILICO*

There were different toxicity parameters studied for a molecule before it is processed for *in-vitro* studies.

### 8.2.1 HEPATOTOXICITY

This is one of the commonly occurring toxicity in the drug molecules which results in drug withdrawal and drug failure under clinical trials. This is also known as drug-induced liver injury in which impairment of liver functions

is commonly seen on exposure to various drug candidates. From a decade around 700 drugs were reported to show hepatotoxicity which makes it necessary to perform *in-silico* hepatotoxicity study before performing clinical trials and marketing the drug [4]. To predict the *in-silico* toxicity of liver various databases are available like the liver toxicological map, LiverTox, Hepatox, Liver Toxicity Knowledge Base, etc.

In 2015 a group of scientists has developed DILI-positive and DILI-negative database by combining 4 different dataset [4]. The DILI-positive database represents the drugs possessing high risk for DILI and the DILI-negative stands for drugs showing low risk of DILI [4]. The first dataset used for this combined database was NCTR data set from FDA's National Center for Toxicological Research [4]. The second data set used was popularly known as Greene data set given by Greene et al. [5] which serves as a validation set along with Xu data set introduced by Xu et al. [6] and the fourth dataset was from Liew et al. which was known as Liew dataset [7]. By leveraging their extensive database, they successfully validated 475 drugs, out of which 198 drugs were confirmed with an impressive accuracy of 86.9%. Moreover, their validation process exhibited a sensitivity of 82.5% and a specificity of 92.9%.

### **8.2.2 CARDIOTOXICITY**

Cardiotoxicity is a type of toxicity related with the blockage of potassium channels that prolongs the QT interval. There are a number of drugs known which were withdrawn from the market due to the cardiotoxicity such as astemizole, cispride, sertindole and terfenadine. There are a number of classifier and regression ML models were available to predict the toxicity [8]. In 2016 a combination of pharmacophore modeling and ML was used by the researcher Wang et al. to predict the hERG toxicity [9]. In this exercise, total 587 molecules were tested using Naïve Bayes (NB) and SVM algorithms. The SVM algorithm model demonstrated promising results with 84.7% accuracy on the training set and 82.1% accuracy on the test set. These figures indicate the model's capability to perform well on both the data it was trained on and new, unseen data. Another model deepHERG based on the DL based approach was developed by a group of researchers named Cheng et al. in 2019 [10]. They have used 7,889 compounds selected from different databases such as PubChem, ChEMBL and other literature sources which were having defined hERG inhibitions.



### 8.2.3 ACUTE ORAL TOXICITY

Acute oral toxicity is indicated by the median lethal death  $LD_{50}$  that shows the dose of drugs on administration causing 50% death rate in animals. This is one of the most common toxicity parameter required in regulatory framework. A number of ML models were introduced to calculate acute oral toxicity. Such as deepAOT model proposed by the Lai et al. to calculate acute oral toxicity which was based on the MGE-CNN architecture (molecular graph encoding convolutional neural network [11]. In this exercise, total 2,200 compounds were used for validation by using database from admetSAR database [12], Toxicity estimation software tool [13] and the MDL Toxicity Database [8]. In 2018 another model name MT-DNN (Multitask deep neural network) was proposed to identify  $LC_{50}$  and  $LD_{50}$  [14]. In this model ECOTOX aquatic toxicity database [15], ChemIDplus database [16] and toxicity estimation software tool database were used [13].

*In silico* prediction of genotoxicity involves using computational methods and models to assess the potential genotoxicity of chemical compounds without performing extensive laboratory experiments. These computational approaches aim to predict the likelihood of a compound causing genetic damage by analyzing its chemical structure, molecular properties, and similarities to known genotoxic compounds. *In silico* methods provide a cost-effective and time-efficient way to screen large chemical libraries and prioritize compounds for further evaluation.

### 8.2.4 GENETOXICITY

Several computational tools and approaches are used for *in silico* prediction of genotoxicity, which are discussed in subsections.

#### 8.2.4.1 QUANTITATIVE STRUCTURE-ACTIVITY RELATIONSHIP (QSAR) MODELS

These are statistical models that correlate chemical structures (represented by molecular descriptors) with their biological activities or toxicological endpoints. QSAR models trained on genotoxicity data can predict the potential genotoxicity of new compounds based on their structural similarity to known genotoxic compounds [17].

#### 8.2.4.2 EXPERT SYSTEMS AND DECISION TREES

Expert systems are rule-based systems that incorporate expert knowledge to predict genotoxicity based on a set of predefined rules and criteria. Decision trees use a hierarchical structure of decision nodes to classify compounds as genotoxic or non-genotoxic based on specific molecular features [18].

#### 8.2.4.3 MACHINE LEARNING ALGORITHMS

Indeed, a diverse range of machine learning techniques, including Support Vector Machines (SVM), Random Forest and Neural Networks, can be effectively employed to construct predictive models for genotoxicity. These methods leverage different approaches and algorithms to analyze genotoxicity data, enabling researchers to identify potential genotoxic compounds and predict their effects accurately. The choice of technique often depends on the nature of the data, the size of the dataset, and the specific requirements of the genotoxicity prediction task at hand. By exploring and comparing these various machine learning approaches, scientists can enhance their understanding of genotoxicity and develop more robust predictive models. These algorithms learn patterns and relationships in training data and use that knowledge to predict the genotoxic potential of new compounds [8].

#### 8.2.4.4 TOXICOPHORE IDENTIFICATION

Toxicophores are substructures or functional groups within molecules that are associated with genotoxicity. Computational tools can identify and prioritize these toxicophores in new compounds to assess their potential genotoxicity [19].

#### 8.2.4.5 CHEMOINFORMATICS AND MOLECULAR DESCRIPTORS

Molecular descriptors are numerical representations of chemical compounds based on their structural properties. By comparing molecular descriptors of a compound to those of known genotoxic compounds, researchers can estimate its likelihood of being genotoxic [20].

#### 8.2.4.6 READ-ACROSS AND STRUCTURAL SIMILARITY

Read-across is a method that uses data from structurally similar compounds with known genotoxicity to predict the genotoxic potential of new compounds with similar structures [21].

It's essential to note that *in silico* predictions are only as reliable as the data used to train the models. Absolutely, the accuracy of predictions in genotoxicity models heavily depends on the availability and quality of the experimental data used during model development and validation. The old adage “garbage in, garbage out” holds true in the context of machine learning. If the input data is flawed, incomplete, or biased, it can lead to poor model performance and unreliable predictions. To build accurate and reliable predictive models for genotoxicity, researchers must ensure that the data used for training and testing is representative of the real-world scenarios they want to apply the model to. This requires collecting diverse, well-curated, and comprehensive datasets that encompass a wide range of genotoxic compounds and non-genotoxic compounds. Additionally, the quality of data labeling and annotation is critical. Properly labeled data helps the model learn the patterns and relationships between features and outcomes, leading to more robust predictions. It's also important to validate the model using independent datasets to ensure that it generalizes well and is not overfitting to the training data. Overfitting occurs when the model memorizes the training data but fails to perform well on new, unseen data. In summary, the success of predictive models for genotoxicity hinges on the data's quality, representativeness, and proper validation. By diligently addressing these aspects, researchers can develop more accurate and reliable models that aid in genotoxicity assessment and safety evaluation of potential compounds. *In silico* genotoxicity prediction is widely used in early stages of drug development, chemical risk assessment, and environmental safety evaluations. It helps researchers and regulatory agencies to efficiently prioritize compounds for further testing and reduce the need for extensive and expensive experimental genotoxicity assays. However, it is crucial to combine *in silico* predictions with *in vitro* and *in vivo* testing for a comprehensive assessment of a compound's genotoxic potential.

### 8.3 COMPUTATIONAL MODELING METHODS

Predicting the toxicity of chemicals is a complex task, and various modeling methods are available to address different aspects of toxicity prediction. Some of the commonly used modeling methods include:

1. **Structural Alerts and Rule-based Models:** These models are based on identifying specific substructures or molecular features known as “alerts” that are associated with toxicity. If a chemical contains these alerts, it is flagged as potentially toxic.
2. **Chemical Category, Read Across, and Trend Analysis:** These approaches involve grouping chemicals into categories based on similarities in structure or properties. If data on toxicity exist for one chemical in a category, that information can be extrapolated to predict the toxicity of other chemicals in the same category.
3. **Dose-Response and Time-Response Models:** These models focus on understanding how the toxicity of a chemical varies with different doses or exposure times. They help estimate the effects of different levels of exposure.
4. **Pharmacokinetic Models:** These models examine how the body processes and distributes a chemical. Understanding the pharmacokinetics helps in assessing the potential toxicity of a substance and its metabolites.
5. **Pharmacodynamic Models:** These models investigate how a chemical interacts with specific biological targets, receptors, or enzymes, which influences its toxicity.
6. **Uncertainty Factor Models:** These models introduce safety margins to account for variations and uncertainties in data and make predictions more conservative to protect against potential adverse effects.
7. **Quantitative Structure-Activity Relationship (QSAR) Models** use mathematical relationships between chemical structure and biological activity to predict the toxicity of new chemicals based on similarities with known compounds.

Each modeling method has its strengths and limitations, and the choice of approach depends on the availability of data, the specific toxicity endpoint being considered, and the resources and expertise available for analysis. In many cases, combining multiple modeling techniques can enhance the overall predictive power and accuracy of toxicity assessments.

### 8.3.1 STRUCTURAL ALERTS

This is one of the most talked factors of predicting toxicity in medicinal chemistry. In structural alerts toxicity of a functional group or a fragment of a moiety has been detected and having knowledge about the role of functional

groups and fragments of a structure in the toxicity can help in designing a less toxic compound. Structural alerts has gained major attention from a decade as it helps in identifying compounds which can serve major toxicity problem as there is high possibility that a chemically reactive fragments of a compound on exposure to human enzyme can undergo bioactivation and form a toxic fragment which is also known as residual toxicity [22]. Thereby structural alerts can help in visualizing potential structural features which can cause adverse reactions in near future.

Structural alerts were also used to predict the mutagenicity by using QSAR model. Various software's like Derek, Toxtree and ToxCast, etc., are available to predict the structural alerts.

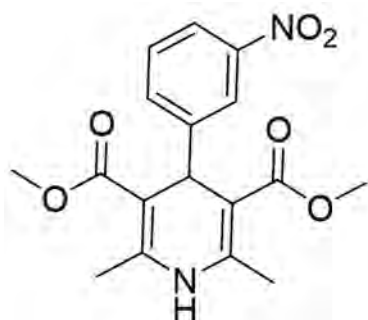
#### 8.3.1.1 METHODS

For analyzing the structural alerts majorly four different databases viz. Scopus, Google scholar, Web of Science and PubMed are used. From these databases structural alerts were find out using different chemical strategies such as metabolic switching wherein introduction of structural essentials of interest is involved. Another strategy is used by reducing the metabolic density, or by changing the metabolically resistant groups with the groups showing structural alerts, etc. [23].

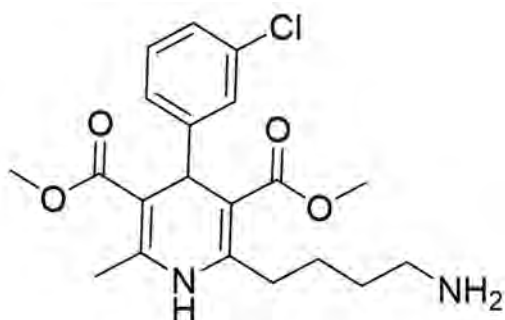
#### 8.3.1.2 METABOLIC SWITCHING

In order to deviate the metabolism process, metabolic switching is used in which the different structural features were added to the main structure so that its metabolism process can be changed. It can also be deviated by blocking the metabolic sites. Some of the examples of metabolic switching are given as follow:

1. **Nifedipine and Amlodipine:** Nifedipine is known to contain nitro group as structure alert but there is no evidence available which shows metabolite formation due to the nitro group reduction. Whereas amlodipine, a metabolite of nifedipine having dihydropyridine structure along with the calcium channel blocker were used. Amlodipine is alkaline in nature which increases its volume of distribution and also influences its plasma half life [24].

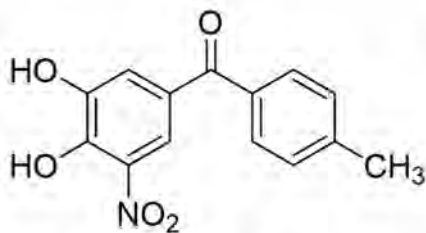


Nifedipine

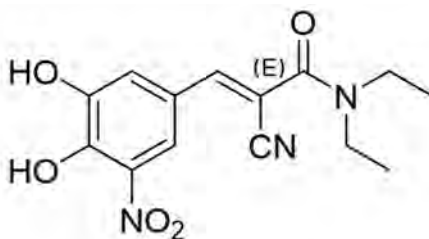


Amlodipine

2. **Tolcapone and Entacapone:** Another example of metabolic switching is Tolcapone and entacapone in which entacapone is the metabolite of tolcapone which is used to overcome the adverse effects of tolcapone [25].



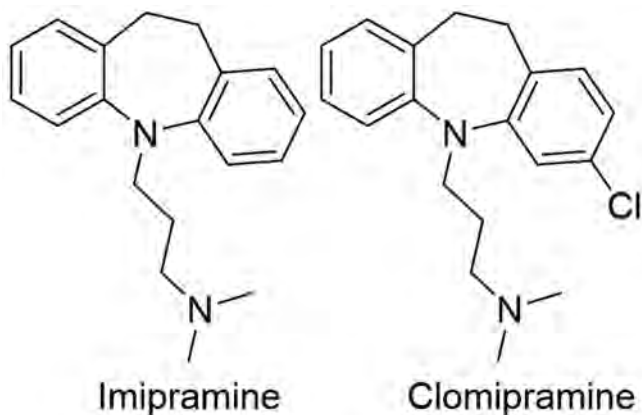
Tolcapone



Entacapone

### 8.3.1.3 REDUCING ELECTRONIC DENSITY

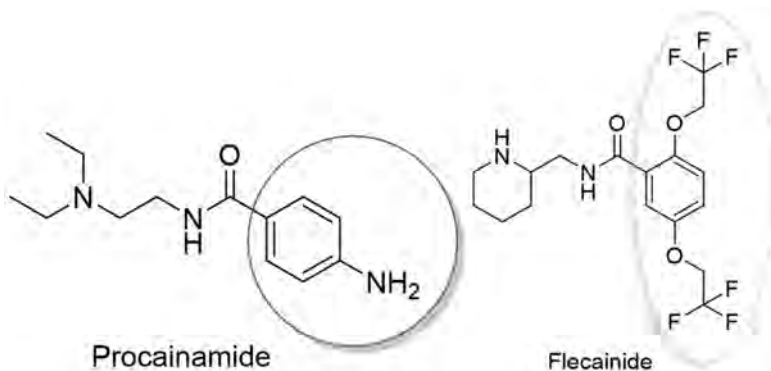
Reducing electronic density is another mode of structural alerts in which the toxicity is overpass by reducing the electronic density of metabolite. One such example is Imipramine having hepatotoxic epoxide metabolite formed by the hydroxylation, demethylation and N-oxidation metabolization process. In this molecules electronic density was changed by converting imipramine to clomipramine by replacing hydrogen group at position 2 of benzazepine moiety with chlorine group which is an electron attracting atom [23].



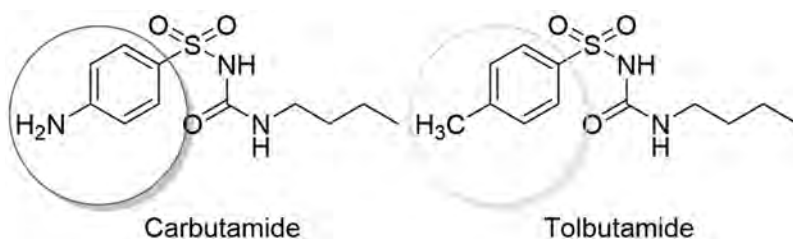
#### 8.3.1.4 SUBSTITUTION OF POTENTIAL STRUCTURAL ALERTS WITH METABOLICALLY RESISTANT SUBSTITUTES

This is one of the commonly used chemical strategy in which the potential structural alert is being partially or fully replaced with a substitute whose functional groups are resistant to any kind of metabolism or biotransformation. Some of the examples are discussed as follow:

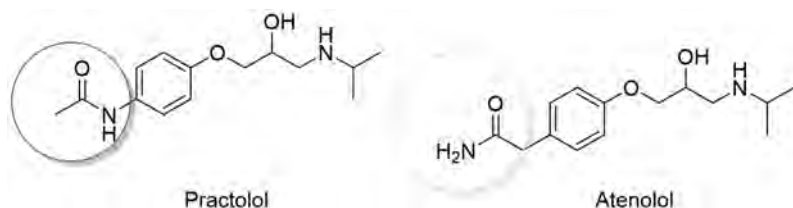
1. **Procainamide:** It is a well-known antiarrhythmic agent which metabolized into *N*-hydroxyaniline derivative and nitroso derivative. This drug is popularly known for showing its adverse reaction such as agranulocytosis and bone marrow toxicity [26, 27]. To overcome its adverse reaction aniline fragment of procainamide is replaced and resulted into flecainide.



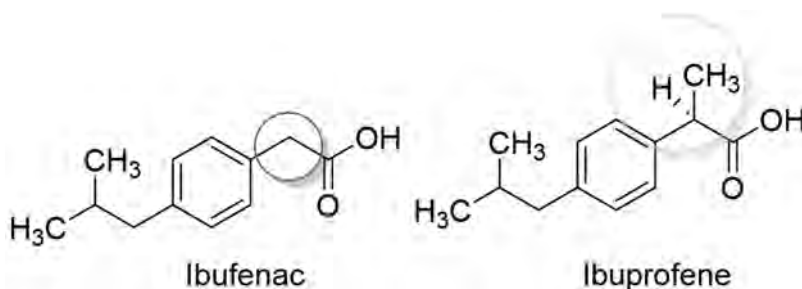
Other examples of replacement are as follow – replacement of carbutamide amine group with methyl group [28]:



Replacement of Practolol with Atenolol by formation of N-hydroxy-aniline derivative [29]:



Replacement Ibufenac with Ibuprofene [23]:



## 8.4 USE OF SOFTWARE'S IN TOXICITY PREDICTION

Toxicity prediction heavily relies on the use of various software tools that facilitate data processing, modeling, and analysis. These software applications enable researchers and toxicologists to efficiently assess the potential



toxicity of small molecules and prioritize compounds for further experimental testing. Here are some of the commonly used software in toxicity prediction [30].

#### **8.4.1 CHEMICAL DATABASES AND DATA MANAGEMENT**

1. **PubChem:** It is a free database maintained by the National Center for Biotechnology Information, which is part of the United States National Library of Medicine. It serves as a comprehensive resource for information on the biological activities of small molecules. PubChem collects and provides data on a wide array of chemical compounds, including information on their chemical structures, properties, biological activities, and associated references. The database contains chemical information from various sources, including literature, high-throughput screening programs, and other public domain databases. Researchers, scientists, and the general public can access PubChem's vast collection of data through its user-friendly web interface. It facilitates searches for specific chemical compounds, their properties, and associated biological activities. PubChem is widely used in the scientific community for drug discovery, toxicology studies, chemical biology research, and other areas that involve the study of small molecules and their interactions with biological systems [31].
2. **ChEMBL:** A large database containing bioactivity data, including toxicity information, for a wide range of compounds [32].
3. **Tox21 Data Browser:** Provides access to high-throughput screening data for thousands of compounds, allowing users to explore toxicity-related information [33].

#### **8.4.2 DATA PREPROCESSING AND VISUALIZATION**

1. **KNIME:** It is an open-source platform designed for data analytics, reporting, and integration. It empowers users to preprocess and visualize toxicity data effectively, facilitating the preparation of data before building predictive models [34].
2. **R and Python:** Programming languages commonly used for data manipulation and visualization in toxicity prediction research.

### **8.4.3 QUANTITATIVE STRUCTURE-ACTIVITY RELATIONSHIP (QSAR) MODELING**

1. **DSSTox:** A collection of tools and resources for chemical screening and QSAR modeling provided by the U.S. EPA.
2. **QSAR Toolbox:** Developed by the European Chemicals Agency (ECHA), it allows users to predict chemical properties and toxicity endpoints using QSAR models [35].
3. **Dragon:** A software suite that calculates molecular descriptors and builds QSAR models for various toxicological endpoints.

### **8.4.4 MACHINE LEARNING AND PREDICTIVE MODELING**

1. **WEKA:** A machine learning software that provides a wide range of algorithms for building predictive models, including for toxicity prediction.
2. **Scikit-Learn:** A Python library for machine learning, offering various algorithms and tools for predictive modeling.
3. **MOE (Molecular Operating Environment):** Combines molecular modeling capabilities with machine learning tools for toxicity prediction.

### **8.4.5 HIGH-THROUGHPUT SCREENING (HTS) DATA ANALYSIS [36]**

1. **ToxPi:** A software tool that visualizes and analyzes high-dimensional toxicity data generated from HTS assays.
2. **OpenTox:** An open-source platform for predictive toxicology that supports the analysis of HTS data and the development of predictive models.

### **8.4.6 CHEMOINFORMATICS TOOLS [37]**

1. **RDKit:** An open-source chemoinformatics toolkit that provides functions for chemical structure handling, descriptor calculation, and QSAR modeling.
2. **ChemAxon:** Offers various chemoinformatics tools for structure representation, calculation of molecular descriptors, and toxicity prediction.

#### 8.4.7 DEEP LEARNING FRAMEWORKS [38]

1. **TensorFlow and Keras:** Widely used deep learning frameworks for building neural network models for toxicity prediction.
2. **PyTorch:** Another popular deep learning library used for developing and deploying predictive models in toxicity assessment.

#### 8.4.8 ADME-TOX PREDICTION [39]

1. **ADMET Predictor:** A software tool that predicts various ADME (Absorption, Distribution, Metabolism, Excretion) properties, including toxicity-related endpoints.
2. **GastroPlus:** A simulation software that aids in predicting the absorption, pharmacokinetics, and pharmacodynamics of compounds, influencing their toxicity profiles.

These software applications offer valuable resources for researchers and toxicologists, enabling them to integrate data from diverse sources, develop predictive models, and prioritize compounds based on their potential toxicity. By using these tools, researchers can accelerate the process of toxicity assessment and ultimately contribute to safer drug development and chemical risk evaluation.

### 8.5 LIST OF THE SOFTWARE'S FOR TOXICITY PREDICTION

A variety of software tools are available for predicting various toxicities, as mentioned in Table 8.1.

#### 8.5.1 FREELY AVAILABLE SOFTWARE

##### 8.5.1.1 CAESAR MODELS

CAESAR stands for 'Computer-Assisted Evaluation of Industrial Chemical Substances According to Regulations.' It is a project that created a number of statistically based models and was sponsored by the European Union (EU). This open-source software is accessible online via the web and was specifically designed to develop quantitative structure-activity relationship (QSAR) models and comply with REACH legislation. It can predict five endpoints

**TABLE 8.1** List of the Software's for Toxicity Prediction

SL. No.	Software	Availability	Accessibility
1.	EPI suite	Freely available	<a href="http://www.epa.gov/oppt/exposure/pubs/episuite.htm">http://www.epa.gov/oppt/exposure/pubs/episuite.htm</a>
2.	OncoLogic	Freely available	<a href="http://ecb.jrc.ec.europa.eu/qsar/qsar-tools">http://ecb.jrc.ec.europa.eu/qsar/qsar-tools</a>
3.	Toxtree	Freely available	<a href="http://ecb.jrc.ec.europa.eu/qsar/qsar-tools">http://ecb.jrc.ec.europa.eu/qsar/qsar-tools</a>
4.	Toxmatch	Freely available	<a href="http://ecb.jrc.ec.europa.eu/qsar/qsar-tools">http://ecb.jrc.ec.europa.eu/qsar/qsar-tools</a>
5.	OECD QSAR Toolbox	Freely available	<a href="http://www.oecd.org">http://www.oecd.org</a>
6.	Lazar	Freely available	<a href="http://lazar.in-silico.de">http://lazar.in-silico.de</a>
7.	Caesar project models	Freely available	<a href="http://www.caesar-project.eu/software/index.htm">http://www.caesar-project.eu/software/index.htm</a>
8.	PASS	Freely available	<a href="http://195.178.207.233/PASS/index.html">http://195.178.207.233/PASS/index.html</a>
9.	T.E.S.T.	Freely available	<a href="http://www.epa.gov/nrmrl/std/cppb/qsar/#TEST">http://www.epa.gov/nrmrl/std/cppb/qsar/#TEST</a>
10.	ADMET predictor	Commercial	<a href="http://www.simulations-plus.com">http://www.simulations-plus.com</a>
11.	TOPKAT	Commercial	<a href="http://www.accelrys.com">http://www.accelrys.com</a>
12.	Pallas software	Commercial	<a href="http://www.compudrug.com">http://www.compudrug.com</a>
13.	Derek Lhasa Ltd	Commercial	<a href="http://www.lhasalimited.org">http://www.lhasalimited.org</a>
14.	MultiCASE	Commercial	<a href="http://www.multicase.com">http://www.multicase.com</a>
15.	MDL QSAR	Commercial	<a href="http://www.multicase.com">http://www.multicase.com</a>
16.	BioEpisteme	Commercial	<a href="http://www.multicase.com">http://www.multicase.com</a>
17.	ACD ToxSuite	Commercial	<a href="http://www.pharma-algorithms.com/webboxes">http://www.pharma-algorithms.com/webboxes</a>
18.	OASIS TIMES	Commercial	<a href="http://www.oasis-lmc.org">http://www.oasis-lmc.org</a>
19.	Molcode Toolbox	Commercial	<a href="http://molcode.com/">http://molcode.com/</a>
20.	q-Tox	Commercial	-
21.	CSFenoTox	Commercial	-

that involve mutagenicity (ames test), skin sensitization, bioconcentration factor, carcinogenicity and developmental toxicities [40].

### 8.5.1.2 EPI SUITE

Numerous physicochemical features like partition coefficient, environmental fate factors, and ecotoxicity are estimated via the EPI (Estimation Programs Interface) Suite. It is a screening level tool that was created by the US EPA in partnership with Syracuse Research Corporation (SRC). It includes a database of more than 40,000 chemicals called PHYSPROP®, where literature from merck, beilstein, etc., is involved. It involves different models that is used to predict various factors [41]: Different models used to estimate various factors are listed in Table 8.2.

**TABLE 8.2** Models Used to Estimate Various Factors

<b>Models</b>	<b>Factors</b>
KOWWIN™	Partition coefficient.
AOPWIN™	Gas-phase reaction rate.
HENRYWIN™	Henry's law constant.
MPBPWIN™	Melting point, boiling point, and vapor pressure of organic chemicals.
BIOWIN™	Aerobic and anaerobic biodegradability of organic chemicals.
BioHCwin	Biodegradation half-life for Hydrocarbons compounds.
KOCWIN™	Organic carbon-normalized sorption coefficient for soil and sediment, i.e., KOC.

Many governmental and business organizations are using this to help the evaluation of new and existing industrial chemicals.

### 8.5.1.3 LAZAR

Lazar is a remarkable open-source software tool designed for the analysis of structural fragments within a training dataset to predict various toxicological endpoints, including mutagenicity, human liver toxicity, rodent and hamster carcinogenicity, and MRDD (Maximum Recommended Daily Dose). It utilizes statistical algorithms, such as k-nearest neighbors and kernel models, for classification and regression tasks, such as multi-linear regression and kernel models. Lazar is also equipped with an automatic applicability domain

estimation, which determines the range of chemicals for which the model's predictions are reliable. Moreover, it provides a confidence index for each prediction, allowing users to gauge the certainty of the model's outcomes. One of the most significant advantages of Lazar is its user-friendly nature, as it doesn't require expert knowledge to be effectively utilized. It can run on Linux operating systems, and there is also a freely accessible web-based prototype available for easy access and use. Overall, Lazar proves to be a powerful and accessible tool for toxicity prediction, facilitating the analysis of chemical structures and aiding researchers in assessing potential toxicological risks with confidence [42].

#### 8.5.1.4 OECD QSAR APPLICATION TOOLBOX

A stand-alone software programme called the OECD QSAR Application Toolbox fills in the gaps in the (eco)toxicity data required for determining the dangers of chemicals. Following a customizable approach, data gaps are filled by building chemical categories and estimating missing data by read-across or by using local QSARs (trends within the category). A variety of profilers are also included in the Toolbox to swiftly assess compounds for common mechanisms or modes of action. The Toolbox offers various datasets with results from experimental research to facilitate read-across and trend analysis. A proof-of-concept version of the Toolbox was made available in its initial form in March 2008. In December 2008, the initial upgrade (version 1.1) was made available. With a four-year timetable, the second phase of the project was begun in November 2008 with the goal of creating a more comprehensive Toolbox that fully incorporates the features of the first version [35].

#### 8.5.1.5 ONCOLOGIC

This expert system evaluates the likelihood that certain substances may result in cancer. The US EPA and LogiChem, Inc. collaborated to develop OncoLogic®. By employing the SAR analysis methods, considering the mechanisms of action and human epidemiological research, and making use of other data, it makes predictions regarding a chemical's probable carcinogenicity. Like human experts, the software explains its line of thinking to back up forecasts. A database of toxicological data pertinent to the evaluation of carcinogenicity is also included. The Cancer Expert System comprises four subsystems tailored to evaluate organic compounds,

polymers, fibers and metals with different chemical structures. Users must have a foundational knowledge of chemistry to correctly select the appropriate subsystem for assessing the specific chemical they are working with. This ensures that the system provides informed and relevant evaluations for each chemical input [43].

#### 8.5.1.6 TOXTREE

Toxtree is an open-source software program designed to categorize substances and predict various types of harmful effects using decision tree methods. It offers a flexible and user-friendly platform for toxicity assessment. Toxtree was developed by the Joint Research Centre (JRC) in collaboration with several experts, with notable contributions from Idea consult Ltd (Sofia, Bulgaria). One of the key strengths of Toxtree is its explicit reporting of the logic behind each prediction, providing transparency and understanding of the basis for its assessments. Toxtree includes several categorization systems for assessing different toxicological endpoints. These include the Cramer scheme and expanded Cramer scheme for systemic toxicity, as well as methods for predicting mutagenicity and carcinogenicity, such as the Benigni-Bossa rule base and the ToxMic rule base based on the *in vivo* micronucleus assay. These systems allow Toxtree to make informed predictions for a wide range of toxicological outcomes. Among its functionalities, Toxtree is widely recognized for its effectiveness in organizing compounds to determine the Threshold of Toxicological Concern (TTC), with the Cramer scheme being one of the most commonly used methods for this purpose. Toxtree's capabilities make it a valuable tool in toxicology research and regulatory assessments, empowering users to evaluate chemical substances efficiently and comprehensively for potential harmful effects [44].

#### 8.5.1.7 PASS

The Prediction of Activity Spectra for Substances (PASS) computerized system was created by the Institute of Biomedical Chemistry of the Russian Academy of Medical Sciences, located in Moscow. It also forecasts pharmacological effects and several toxicities, such as mutagenicity, carcinogenicity, teratogenicity, and embryotoxicity. By comparing the new substance's similarity or dissimilarity to compounds in the training set

(70,000 compounds), which contain substances with well-known biological activities, the algorithm estimates the probability (Pa) of a biological activity for a new chemical. Based on the knowledge base of mechanism-effect correlations, the tool also provides a cross reference between biological processes [45].

#### 8.5.1.8 T.E.S.T.

The US EPA created the open-source software tool known as Toxicity Estimation. It applies numerous QSAR approaches to assess a compound's toxicity, giving the user more assurance about expected toxicities. It predicts a wide range of toxicities, including *Daphnia magna*, *Tetrahymena pyriformis*, acute toxicity to fish (fathead minnow), as well as rat oral LD<sub>50</sub>, Ames mutagenicity, developmental toxicity, and various other toxicities [30].

#### 8.5.1.9 OSIRIS PROPERTY EXPLORER

The OSIRIS Property Explorer is a software tool that employs computational models and algorithms to predict and explore physicochemical and toxicological properties of chemical compounds. It offers predictions for a wide range of properties, including environmental fate, physicochemical characteristics, and toxicological endpoints. The software analyzes chemical structures, functional groups, and other relevant information to generate predictions based on large databases and statistical models. Users can input structures or batch files for analysis, and the results are provided as numerical values, labels, and graphics. The OSIRIS Property Explorer finds applications in drug discovery, environmental risk assessment, and chemical safety evaluation, providing valuable insights for compound selection, toxicity evaluation, and risk management [46].

#### 8.5.1.10 PREADMET

PreADMET is a software tool developed by the Bioinformatics and Molecular Design Research Center (BMDRC) for predicting the ADME properties of chemical compounds. It utilizes computational models and algorithms based on the compound's structure and properties to estimate various ADME parameters. These parameters include solubility, blood-brain



barrier penetration, metabolism, and excretion. PreADMET aids in the early stages of drug discovery by providing insights into a compound's behavior in the body. However, experimental validation is necessary for accurate predictions. The software is a valuable resource for researchers and drug developers to optimize candidate selection and drug design [47, 48].

#### 8.5.1.11 ADMETSAR

AdmetSAR is an online platform and database that uses computational models, including machine learning algorithms, to predict the ADMET properties of chemical compounds. It incorporates a comprehensive collection of data from various sources, enabling the development of reliable predictive models. Users can input chemical structures or search for specific compounds to obtain predictions for their ADMET properties. AdmetSAR provides additional insights into underlying mechanisms, pathways, and biological targets, along with links to relevant literature and resources. It is a valuable tool for researchers, drug developers, and regulatory agencies to identify and evaluate ADMET-related issues, supporting informed decision-making in drug development [12].

#### 8.5.1.12 TOX21 DASHBOARD

The Tox21 Dashboard is an online platform designed for researchers and regulatory agencies to access and explore toxicity data and predictions. It facilitates the investigation of various toxicity endpoints, including cytotoxicity, genotoxicity, endocrine disruption, and organ-specific toxicities. The predictions are generated through the application of machine learning algorithms that analyze chemical structure-activity relationships and molecular descriptors. The Tox21 Dashboard integrates data from diverse sources, such as the EPA's ToxCast and ToxRefDB databases, as well as publicly available toxicity databases. This integration enables users to access a comprehensive collection of experimental and computational toxicity data, supporting the comparison and exploration of the toxicological profiles of different chemicals. This platform also provides visualization tools, such as heat maps, scatter plots, and concentration-response curves, to present data in visual formats. These visualizations aid in the interpretation and analysis of toxicity predictions, facilitating the identification of patterns and trends within the data. In summary, the Tox21 Dashboard is a scientifically

advanced and user-friendly platform that offers a broad spectrum of toxicity data and predictions. It leverages chemical structure-activity relationships, molecular descriptors, and machine learning techniques to provide valuable insights into the toxicological properties of chemicals [33, 49].

## **8.5.2 COMMERCIALLY AVAILABLE SOFTWARE**

### **8.5.2.1 DEREK NEXUS**

Derek Nexus is a commercially available software developed by Lhasa Limited, is widely utilized for toxicity prediction and assessment in various industries and regulatory agencies. This software employs a knowledge-based approach, incorporating expert-derived rules and structural alerts derived from empirical evidence and toxicological expertise. By analyzing the chemical structure of a compound, Derek Nexus identifies and evaluates potential toxicity hazards associated with endpoints such as genotoxicity, carcinogenicity, and skin sensitization. The software offers a user-friendly interface that facilitates users to input chemical structures, enabling the generation of toxicity predictions accompanied by supporting information and confidence levels. Additionally, Derek Nexus allows customization of prediction settings to suit specific user requirements. Its applications span pharmaceuticals, chemicals, cosmetics, and regulatory sectors, facilitating early-stage toxicity screening, compound prioritization, and informed decision-making in compound development and safety assessment. The predictions offered by Derek Nexus contribute valuable insights into compound toxicity, aiding in risk identification and informing subsequent testing and evaluation processes [50, 51].

### **8.5.2.2 LEADSCOPEPREDICTIVETOX SUITE**

The LeadscapePredictiveTox, developed by Leadscape Inc., is an advanced software suite used for toxicity prediction and assessment. It employs computational models and algorithms to estimate toxicity endpoints based on chemical structure and property information. The suite encompasses a wide range of toxicological properties and integrates a substantial database comprising chemical structures, toxicity data, and expert-curated rules. Users have the capability to input chemical structures or import compound libraries for assessment, and the suite employs machine learning techniques

and quantitative structure-activity relationship (QSAR) models to generate predictions. The LeadscapePredictiveTox Suite offers detailed reports, visualizations, and customizable options for the analysis of predicted toxicity endpoints. It is extensively utilized in the pharmaceutical, chemical, and regulatory sectors to perform risk assessments, conduct compound screening, and facilitate decision-making in the domains of drug discovery and safety evaluation [52].

#### 8.5.2.3 MULTICASE

MultiCASE is a software suite developed by MultiCASE Inc. that offers a comprehensive set of computational tools for toxicity prediction and assessment. It incorporates advanced algorithms and predictive models based on structure-activity relationships (SAR) to estimate a wide range of toxicological endpoints. These endpoints include acute toxicity, mutagenicity, carcinogenicity, reproductive toxicity, and environmental toxicity. The suite includes specialized modules such as CASE Ultra, CASE Ultra-Genetox, and CASE Profiler, each designed to address specific toxicological properties. These modules utilize expert rules, statistical models, and SAR analysis to generate predictions for the selected toxicity endpoints. Users can input chemical structures or import compound libraries for analysis, and the suite provides detailed reports and visualizations of the predicted toxicological properties. MultiCASE also offers customization options, allowing users to develop their own models and rules based on specific data or requirements. The suite supports data integration, data mining, and the exploration of structure-activity relationships to gain deeper insights into the underlying mechanisms contributing to toxic effects. MultiCASE finds extensive use in the pharmaceutical, chemical, and regulatory industries for toxicity screening, risk assessment, and decision-making in drug development and chemical safety evaluations. Its capabilities enable early identification of potential toxicological hazards and assist in the prioritization of compounds for further testing and development [53].

#### 8.5.2.4 ADMET PREDICTOR

ADMET Predictor is developed by Simulations Plus, a computational tool employed for the prediction and assessment of ADMET properties of chemical compounds. Leveraging QSAR approaches and machine learning

algorithms, ADMET Predictor utilizes the chemical structure and properties of compounds to estimate various ADMET parameters. Integration of extensive chemical databases and experimental data enhances the accuracy of predictions. Researchers can input chemical structures or compound libraries for analysis, and ADMET Predictor generates comprehensive reports and visualizations of the predicted ADMET properties. The tool allows for customization, enabling users to refine models and incorporate additional data. ADMET Predictor serves as a valuable resource for identifying and evaluating ADMET-related challenges in drug discovery and development, facilitating the optimization of safe and efficacious drug candidates. It is important to note that experimental validation and further data analysis are essential for validating and refining the predicted ADMET properties [54].

#### 8.5.2.5 TOPKAT

TOPKAT (Toxicity Prediction by Computer Assisted Technology) is a software tool developed by BIOVIA discovery studio that employs quantitative structure-activity relationship (QSAR) models and expert knowledge for the estimation of various toxicological endpoints. It offers predictions for acute toxicity, mutagenicity, carcinogenicity, and skin sensitization, among others. The software utilizes chemical structure inputs from users and applies QSAR models and associated algorithms to generate toxicity predictions. TOPKAT provides detailed reports, visualizations, and customization options to facilitate the interpretation and analysis of the results. It finds extensive application in the pharmaceutical, chemical, and regulatory industries for toxicity screening and compound prioritization. However, it is crucial to complement TOPKAT predictions with experimental data and expert judgment to ensure accurate and reliable toxicity assessment [39, 55].

#### 8.5.2.6 DEEPTOX

DeepTox is an advanced computational tool that uses deep learning, specifically deep neural networks, to predict the toxicity of chemical compounds. It is trained on large datasets of compounds with known toxicity profiles and can accurately predict various toxicological endpoints. DeepTox captures complex relationships between chemical structures and toxicity, improving prediction accuracy. However, experimental validation is still necessary to confirm predicted toxic effects [14].

### 8.5.2.7 QSAR TOOLBOX

The QSAR Toolbox is software created by the OECD with the purpose of utilizing QSAR models to predict the toxicity of chemicals. Its interface is designed to be user-friendly and grants access to databases containing experimental data. The tool allows users to choose specific toxicological endpoints for prediction and enables the grouping of chemicals based on their similarities. Primarily intended for regulatory purposes, the QSAR Toolbox receives regular updates with the inclusion of new data and models. However, it is essential to use the predictions alongside other data and expert judgment to ensure a comprehensive risk assessment [56].

## 8.6 CONCLUSION

Ensuring the safety of new drug candidates during the development process is of utmost importance. The evaluation of chemical toxicity plays a crucial role in this endeavor. Over the past years, machine learning models have emerged as powerful tools for predicting the toxicity of small molecules. *In-Silico* identification of toxic parameters, such as acute oral toxicity, hepatotoxicity, cardiotoxicity, and mutagenicity, has enabled researchers to make informed decisions early in the drug development pipeline.

A significant advancement in the field of toxicology has been the development of various software applications dedicated to predicting toxicity. Some notable examples include ProTox-II, Derek (Deductive estimate of risk from existing knowledge), ToxiM, ADMET Predictor, OECD toolbox, Toxtree, q-Tox, TOPKAT, MDL QSAR, Osiris property explorer, and T.E.S.T. These software tools utilize computational approaches to model and predict the toxicity of small molecules, helping researchers and pharmaceutical companies in making more informed choices about which compounds to progress further in the drug development process. By harnessing the potential of machine learning and computational toxicology, researchers can significantly reduce the time and resources required for drug development, ultimately leading to safer and more effective medications. However, it's essential to continue refining and validating these predictive models to ensure their accuracy and reliability. As technology continues to advance, the integration of machine learning in toxicity prediction promises to be an indispensable asset in the pharmaceutical industry, furthering our ability to bring new, life-changing medications to patients while prioritizing safety throughout the drug development journey.

In conclusion, computational approaches play a pivotal role in predicting the toxicity of small molecules. Leveraging machine learning, molecular modeling, and data analytics, researchers can efficiently assess the safety of chemical compounds, guiding drug development and ensuring environmental protection. However, it is crucial to acknowledge the limitations of these approaches and continuously strive for improvements to enhance the reliability and applicability of toxicity prediction methods. This chapter provided a comprehensive overview of the various software used for toxicity prediction.

## KEYWORDS

- **chemical compounds**
- **chemical toxicity**
- **data analytics**
- **drug development**
- ***in silico***
- **machine learning**
- **molecular modeling**
- **molecules**
- **toxicity prediction**

## REFERENCES

1. Tran, T. T. Van, Surya Wibowo, A., Tayara, H., & Chong, K. T. (2023). Artificial intelligence in drug toxicity prediction: Recent advances, challenges, and future perspectives. *Journal of Chemical Information and Modeling*, 63(7), 2628–2643. <https://doi.org/10.1021/acs.jcim.3b00353>.
2. Cavasotto, C. N., & Scardino, V. (2022). Machine learning toxicity prediction: Latest advances by toxicity end point. *ACS Omega*, 7(48), 47536–47546. <https://doi.org/10.1021/acsomega.2c05693>.
3. Yang, H., Sun, L., Li, W., Liu, G., & Tang, Y. (2018). In silico prediction of chemical toxicity for drug design using machine learning methods and structural alerts. *Molecules*, 23(11), 1–12. <https://doi.org/10.3390/molecules23113043>.
4. Xu, Y., Liu, J., Zhang, L., & Lu, J. (2015). Deep learning for drug-induced liver injury. *Journal of Chemical Information and Modeling*, 55(10), 2000–2009. <https://doi.org/10.1021/acs.jcim.5b00238>.

5. Greene, N., Wei, X., Harris, S., & Avigan, M. (2010). Developing structure-activity relationships for the prediction of hepatotoxicity. *Chemical Research in Toxicology*, 23(8), 1215–1222. <https://doi.org/10.1021/tx100158u>.
6. Xu, J. J., Li, S., & Johnson, R. D. (2008). Cellular imaging predictions of clinical drug-induced liver injury. *Toxicological Sciences*, 105(1), 97–105. <https://doi.org/10.1093/toxsci/kfn062>.
7. Liew, C. Y., Lim, Y. C., & Yap, C. W. (2011). Mixed learning algorithms and features ensemble in hepatotoxicity prediction. *Computational Biology and Chemistry*, 35(5), 311–320. <https://doi.org/10.1016/j.compbiolchem.2011.07.003>.
8. Thi Tuyet Van Tran, Agung Surya Wibowo, Hilal Tayara, & Kil To Chong. (2023). Artificial intelligence in drug toxicity prediction: recent advances, challenges, and future perspectives, *J. Chem. Inf. Model.*, 63(9), 2628–2643.
9. Wang, S., Wang, H., & Liu, J. (2016). ADMET evaluation in drug discovery: Predicting hERG blockers by combining multiple pharmacophores and machine learning approaches. *Molecular Pharmaceutics*, 13(8), 2855–2866. <https://doi.org/10.1021/acs.molpharmaceut.6b00576>.
10. Cai, C., Zhang, L., & Liu, T. (2019). Deep learning-based prediction of drug-induced cardiotoxicity. *Journal of Chemical Information and Modeling*, 59(3), 1073–1084. <https://doi.org/10.1021/acs.jcim.8b00747>.
11. Xu, Y., Pei, J., & Lai, L. (2017). Deep learning-based regression and multiclass models for acute oral toxicity prediction with automatic chemical feature extraction. *Journal of Chemical Information and Modeling*, 57(12), 2672–2685. <https://doi.org/10.1021/acs.jcim.7b00402>.
12. Cheng, F., Li, W., Wang, X., & Zhao, Z. (2012). admetSAR: A comprehensive source and free tool for assessment of chemical ADMET properties. *Journal of Chemical Information and Modeling*, 52(11), 3099–3105. <https://doi.org/10.1021/ci300367a>.
13. U.S. Environmental Protection Agency. (2012). *Quantitative Structure Activity Relationship*. Retrieved from: <https://www.epa.gov/chemical-research/toxicity-estimation-software-tool-test> (accessed on 25 July 2024).
14. Mayr, A., Klambauer, G., Unterthiner, T., & Hochreiter, S. (2016). DeepTox: Toxicity prediction using deep learning. *BMC Bioinformatics*, 17(1), 1–12. <https://doi.org/10.1186/s12859-016-0997-4>.
15. U.S. Environmental Protection Agency. (2018). *ECOTOX Knowledgebase*. Retrieved from: <https://cfpub.epa.gov/ecotox/> (accessed on 25 July 2024).
16. National Institutes of Health. (2018). *ChemIDplus*. Retrieved from: <https://chem.nlm.nih.gov/chemidplus/> (accessed on 25 July 2024).
17. Capuzzi, S. J., Politi, R., Isayev, O., & Farag, S. (2016). QSAR modeling of Tox21 challenge stress response and nuclear receptor signaling toxicity assays. *Journal of Chemical Information and Modeling*, 56(12), 2400–2411. <https://doi.org/10.1021/acs.jcim.6b00480>.
18. Mehrpour, O., Saeedi, F., Nakhaee, S., Khomeini, F. T., & Hadianfar, A. (2023). Comparison of decision tree with common machine learning models for prediction of biguanide and sulfonylurea poisoning in the United States: An analysis of the National Poison Data System. *BMC Medical Informatics and Decision Making*, 23(1), 1–11. <https://doi.org/10.1186/s12911-022-02095-y>.
19. Singh, P. K., Negi, A., Gupta, P. K., Chauhan, M., & Kumar, R. (2016). Toxicophore exploration as a screening technology for drug design and discovery: Techniques, scope

- and limitations. *Archives of Toxicology*, 90(8), 1785–1802. <https://doi.org/10.1007/s00204-016-1743-5>.
20. Sharma, A. K., Srivastava, G. N., Roy, A., & Sharma, V. K. (2017). ToxiM: A toxicity prediction tool for small molecules developed using machine learning and chemoinformatics approaches. *Journal of Chemical Information and Modeling*, 57(10), 1–18. <https://doi.org/10.1021/acs.jcim.7b00116>.
  21. Moustakas, H., Adams, C., Becker, J., & Williams, L. (2022). An end point-specific framework for read-across analog selection for human health effects. *Chemical Research in Toxicology*, 35(8), 2324–2334. <https://doi.org/10.1021/acs.chemrestox.2c00271>.
  22. Claesson, A., & Minidis, A. (2018). Systematic approach to organizing structural alerts for reactive metabolite formation from potential drugs. *Chemical Research in Toxicology*, 31(3), 389–411. <https://doi.org/10.1021/acs.chemrestox.7b00490>.
  23. Limban, C., Rosner, M., & Brauer, K. (2018). The use of structural alerts to avoid the toxicity of pharmaceuticals. *Toxicology Reports*, 5, 943–953. <https://doi.org/10.1016/j.toxrep.2018.07.009>.
  24. Jorga, K., Fotteler, B., Heizmann, P., & Gasser, R. (1999). Metabolism and excretion of tolcapone, a novel inhibitor of catechol-O-methyltransferase. *British Journal of Clinical Pharmacology*, 48(4), 513–520. <https://doi.org/10.1046/j.1365-2125.1999.00028.x>.
  25. Wikberg, J. E., Vuorela, P., Ottoila, P., & Taskinen, J. (1993). Identification of major metabolites of the catechol-O-methyltransferase inhibitor entacapone in rats and humans. *Drug Metabolism and Disposition*, 21(1), 81–92. <https://doi.org/10.1124/dmd.21.1.81>.
  26. Siraki, A. G., Bonini, M. G., Jiang, J., Ehrenshaft, M., & Mason, R. P. (2007). Amino-glutethimide-induced protein free radical formation on myeloperoxidase: A potential mechanism of agranulocytosis. *Chemical Research in Toxicology*, 20(7), 1038–1045. <https://doi.org/10.1021/tx700027y>.
  27. Siraki, A. G., Olmstead, W., Liu, M., & Mason, R. P. (2008). Procainamide, but not N-acetylprocainamide, induces protein free radical formation on myeloperoxidase: A potential mechanism of agranulocytosis. *Chemical Research in Toxicology*, 21(6), 1143–1153. <https://doi.org/10.1021/tx800016f>.
  28. Mannhold, R., Kubinyi, H., & Folkers, G. (2012). *Pharmacokinetics and Metabolism in Drug Design*. John Wiley & Sons.
  29. Sim, E., Stanley, L., Gill, E. W., & Jones, A. (1988). Metabolites of procainamide and practolol inhibit complement components C3 and C4. *Biochemical Journal*, 251(1), 323–326. <https://doi.org/10.1042/bj2510323>.
  30. Krewski, D., Acosta, D., Andersen, M. E., & Anderson, H. (2010). Toxicity testing in the 21st century: A vision and a strategy. *Journal of Toxicology and Environmental Health, Part B: Critical Reviews*, 13(1), 51–138. <https://doi.org/10.1080/10937401003650383>.
  31. Wang, Y., Xiao, L., & Lu, Y. (2009). PubChem: A public information system for analyzing bioactivities of small molecules. *Nucleic Acids Research*, 37(Web Server issue), W623–W633. <https://doi.org/10.1093/nar/gkp456>.
  32. Bento, A. P., Gaulton, A., Hersey, A., & Mendez, D. (2014). The ChEMBL bioactivity database: An update. *Nucleic Acids Research*, 42(D1), 1083–1090. <https://doi.org/10.1093/nar/gkt1031>.
  33. Richard, A. M., Judson, R. S., Houck, K. A., & Knudsen, T. B. (2021). The Tox21 10K Compound Library: Collaborative chemistry advancing toxicology. *Chemical Research in Toxicology*, 34(1), 189–216. <https://doi.org/10.1021/acs.chemrestox.0c00126>.



34. Sterling, T., & Irwin, J. J. (2015). ZINC 15 – Ligand discovery for everyone. *Journal of Chemical Information and Modeling*, 55(11), 2324–2337. <https://doi.org/10.1021/acs.jcim.5b00559>.
35. Schultz, T. W., Diderich, R., Kuseva, C. D., & Mekenyan, O. G. (2018). The OECD QSAR Toolbox starts its second decade. *Methods in Molecular Biology*, 1800, 55–77. [https://doi.org/10.1007/978-1-4939-8641-0\\_4](https://doi.org/10.1007/978-1-4939-8641-0_4).
36. Murray, D., McWilliams, L., & Wigglesworth, M. (2016). High-throughput cell toxicity assays. *Methods in Molecular Biology*, 1439, 245–262. [https://doi.org/10.1007/978-1-4939-3598-3\\_14](https://doi.org/10.1007/978-1-4939-3598-3_14).
37. Rabinowitz, J. R., Goldsmith, M. R., Little, S. B., & Pasquinelli, M. A. (2008). Computational molecular modeling for evaluating the toxicity of environmental chemicals: Prioritizing bioassay requirements. *Environmental Health Perspectives*, 116(5), 573–577. <https://doi.org/10.1289/ehp.10790>.
38. Zhang, J., Norinder, U., & Svensson, F. (2021). Deep learning-based conformal prediction of toxicity. *Journal of Chemical Information and Modeling*, 61(6), 2648–2657. <https://doi.org/10.1021/acs.jcim.0c01452>.
39. Mohan, C. G., Gandhi, T., Garg, D., & Shinde, R. (2007). Computer-assisted methods in chemical toxicity prediction. *Mini-Reviews in Medicinal Chemistry*, 7(5), 499–507. <https://doi.org/10.2174/138920107781514223>.
40. Cassano, A., Gini, G., Montagnaro, S., & Santoro, D. (2010). CAESAR models for developmental toxicity. *Chemistry Central Journal*, 4(1), S4. <https://doi.org/10.1186/1752-153X-4-S1-S4>.
41. Rodríguez-Leal, I., & MacLeod, M. (2022). The applicability domain of EPI Suite™ for screening phytotoxins for potential to contaminate source water for drinking. *Environmental Sciences Europe*, 34, 96. <https://doi.org/10.1186/s12302-022-00689-6>.
42. Maunz, A., Ganter, S., Gauthier, J., Heumann, R., Lutz, W., & Weizel, E. (2013). lazar: A modular predictive toxicology framework. *Frontiers in Pharmacology*, 4, 38. <https://doi.org/10.3389/fphar.2013.00038>.
43. Benigni, R., Bossa, C., Alivernini, S., Colafranceschi, M. (2012). Assessment and Validation of US EPA's OncoLogic® Expert System and Analysis of Its Modulating Factors for Structural Alerts. *Journal of Environmental Science and Health. Part C, Environmental Carcinogenesis & Ecotoxicology Reviews*. 30(2), 152173. <https://doi.org/10.1080/10590501.2012.681486>.
44. Patlewicz, G., Jeliaskova, N., Safford, R. J., Worth, A. P., & Aleksiev, B. (2008). An evaluation of the implementation of the Cramer classification scheme in the Toxtree software. *SAR and QSAR in Environmental Research*, 19(5–6), 495–524. <https://doi.org/10.1080/10629360802021703>.
45. Poroikov, V., Filimonov, D., Lagunin, A., Glorizova, T., & Zakharov, A. (2007). PASS: Identification of probable targets and mechanisms of toxicity. *SAR and QSAR in Environmental Research*, 18(1), 101–110. <https://doi.org/10.1080/10629360600959473>.
46. Sander, T., Freyss, J., von Korff, M., Reich, J. R., & Rufener, C. (2009). OSIRIS, an entirely in-house developed drug discovery informatics system. *Journal of Chemical Information and Modeling*, 49(1), 232–246. <https://doi.org/10.1021/ci800322d>.
47. Muster, W., Rognan, D., & Velankar, S. (2008). Computational toxicology in drug development. *Drug Discovery Today*, 13(7–8), 303–310. <https://doi.org/10.1016/j.drudis.2008.03.002>.

48. Lee, S., & Ochoa, P. (2002). The PreADME approach: Web-based program for rapid prediction of physico-chemical, drug absorption and drug-like properties. *euro QSAR 2002–Des. Drugs Crop Prot. Process. Probl. Solut.*, 418–420.
49. Jeong, J., Kim, D., & Choi, J. (2022). Application of ToxCast/Tox21 data for toxicity mechanism-based evaluation and prioritization of environmental chemicals: Perspective and limitations. *Toxicology in Vitro: An International Journal Published in Association with BIBRA*, 84, 105451. <https://doi.org/10.1016/j.tiv.2022.105451>.
50. Bhattarai, B., Wilson, D. M., Parks, A. K., Carney, E. W., & Spencer, P. J. (2016). Evaluation of TOPKAT, Toxtree, and Derek Nexus in silico models for ocular irritation and development of a knowledge-based framework to improve the prediction of severe irritation. *Chemical Research in Toxicology*, 29(5), 810–822. <https://doi.org/10.1021/acs.chemrestox.6b00004>.
51. Raies, A. B., & Bajic, V. B. (2016). In silico toxicology: Computational methods for the prediction of chemical toxicity. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 6(2), 147–172. <https://doi.org/10.1002/wcms.1247>.
52. Savale, S. (2019). Genotoxicity of drugs: Introduction, prediction, and evaluation. In *Genotoxicity of Drugs: Introduction, Prediction and Evaluation*, 1–29.
53. Chakravarti, S. K., & Saiakhov, R. D. (2022). MultiCASE platform for in silico toxicology. *Methods in Molecular Biology*, 2425, 497–518. [https://doi.org/10.1007/978-1-0716-2161-2\\_28](https://doi.org/10.1007/978-1-0716-2161-2_28).
54. Sohlenius-Sternbeck, A. K., & Terelius, Y. (2022). Evaluation of ADMET predictor in early discovery drug metabolism and pharmacokinetics project work. *Drug Metabolism and Disposition*, 50(1), 95–104. <https://doi.org/10.1124/dmd.121.000723>.
55. Greene, N. (2002). Computer systems for the prediction of toxicity: An update. *Advanced Drug Delivery Reviews*, 54(3), 417–431. [https://doi.org/10.1016/S0169-409X\(02\)00011-7](https://doi.org/10.1016/S0169-409X(02)00011-7).
56. Yordanova, D., Roperro, L., & Gallego, M. (2019). Automated and standardized workflows in the OECD QSAR Toolbox. *Computational Toxicology*, 10, 89–104.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## CHAPTER 9

---

# Toxicity Forecasts: Navigating Data-Driven AI/ML Models: From Theory to Practice

B. V. S. SUNEEL KUMAR,<sup>1</sup> ANTOINE MOITESSIER,<sup>2</sup> and  
NICOLAS MOITESSIER<sup>2</sup>

<sup>1</sup>*Atomicas AI Solutions Private Limited, Plot No 35, Beside Avance Phoenix SEZ, Hitech City, Hyderabad, India*

<sup>2</sup>*Molecular Forecaster Inc., 910–2075 Robert Bourassa St., Montreal, Quebec, H3A2L1, Canada*

---

### ABSTRACT

Drug toxicity plays a crucial role in the withdrawal of drugs from clinical trials. Predicting drug candidates' toxicity profiles may be considered as the first step towards drugs safety and efficacy. However, this goal may only be reached if predictive tools and models are available. It is believed that the future of toxicity forecasting lies in computational models such as artificial intelligence (AI) and machine learning (ML). However, developing such computational methods has its fair share of challenges, including limited comprehensive datasets, model interpretability, applicability domain (AOD), model biases, and generalizability. Integrating multi-model data sources is a solution to address these challenges and requires understanding theoretical foundations, selecting appropriate algorithms and datasets, and considering model applicability. By overcoming these hurdles and obtaining proper toxicity profiles, drug safety and efficacy may be greatly enhanced. This chapter will cover some of the fundamental steps and concepts in AI/ML model development from a practical standpoint: dataset collection, molecular featurization, AI/ML theory, validation, and evaluation metrics. The python codes and Jupyter notebooks that are discussed in this book

---

Artificial Intelligence for Chemical Sciences: Concepts, Models, and Applications. Shrikaant Kulkarni, Shashikant Bhandari, Dushyant Varshney, & P. William (Eds.)

© 2025 Apple Academic Press, Inc. Co-published with CRC Press (Taylor & Francis)

chapter are available in author's GitHub profile, which can be accessed at: <https://github.com/suneelbvs/toxicity-forecasts>.

## 9.1 INTRODUCTION

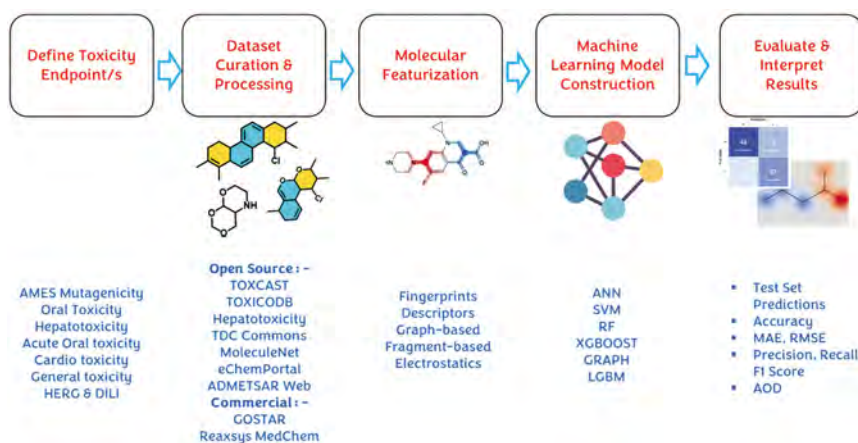
A major challenge in drug discovery and development is the assessment of the toxicity of potential drugs, such as liver toxicity, genotoxicity, and carcinogenicity. These adverse reactions are one of the major reasons for drug failures in clinical trials and subsequent market drug withdrawal [1–3]. Early detection and prediction of liver toxicity during the drug discovery process are a great means to minimize the costs and the risk of drug failure [3]. Unfortunately, conventional *in vivo* animal toxicity screening methods are costly and time-consuming, making these approaches reliable in terms of results, but not so much in their efficacy and throughput. In light of these limitations, alternative approaches have emerged, such as the development of diverse artificial intelligence (AI) and machine learning (ML) models aimed at predicting different toxicity endpoints including genotoxicity, and liver toxicity in humans. These AI/ML approaches have already been implemented and have started impacting the drug discovery and development process by enabling early detection of potential toxicity issues, thereby reducing both the financial and human cost of drug failures [4, 5].

Traditional methods usually involve time-consuming experimental screenings and expensive laboratory tests, which are both inclined to their respective limitations, amongst which you may find their resource cost and throughput. However, AI and ML approaches have revolutionized this aspect of drug discovery by greatly decreasing the demand in resources. These methods utilize large datasets and computational models to rapidly analyze, classify and identify the compounds, leading to a significant reduction in the time and cost involved in early-stage drug development [6–8]. Simply put, computational models are modeling experiments through theories to solve research questions using ever-increasing computational power.

AI/ML approaches have been instrumental in identifying and classifying potential compounds with no toxic alerts (often fragments known to induce toxicity) during the early stages of drug design and development [7]. By training algorithms on vast amounts of toxicological data, these approaches can predict the toxicity of potential drug candidates with high accuracy. This enables researchers to prioritize potentially safer compounds for further investigation, saving considerable time and resources [9]. AI/ML methods, as their name would infer, actively learn from the dataset given, with little or no

insight from external sources. This allows these methods to identify patterns in complex datasets that may not be easily discernible through traditional approaches, providing valuable insights into potential toxic effects [10–13].

This book chapter on AI/ML in toxicity predictions provides a comprehensive overview of various concepts including data collection, data pre-processing, molecular featurization, model building and evaluation, and knowledge-based techniques. While there are already numerous interesting and extensive chapters and articles discussing the role of AI/ML in toxicity endpoint predictions [8, 14–18], our chapter focuses on the brief review of AI/ML role in toxicity modeling, followed by basic steps involved in data processing and machine learning model building as a practical approach, as illustrated in Figure 9.1.



**FIGURE 9.1** Fundamental steps in developing a deep learning (DL)/machine learning (ML) model for toxicity endpoints.

To aid beginners and academic research scholars in developing their skills and understanding in AI/ML concepts, the chapter includes coding snippets from popular libraries such as rdKit, DeepChem, Datamol – MolFeat, scikit-learn, and TensorFlow with explanation. These python codes are also available in a jupyter notebook and python format on the author’s GitHub profile, which can be accessed at: <https://github.com/suneelbvs/Toxicity-Forecasts>.

By providing these coding snippets, we aim to facilitate the practical implementation of ML models for toxicity predictions. Whether you are new to the use of computation models or you are simply looking to enhance your existing knowledge, these snippets will serve as valuable resources for practical learning and experimentation.

This comprehensive chapter begins by addressing the pivotal aspect of data sources for toxicity endpoints. It provides a comprehensive list of open-source and commercial data sources available for various toxicity endpoints. The next section delves into the key Python libraries that are essential for different stages of the workflow, such as data handling, processing, and generating analytics. For our purposes, these libraries include Pandas, rdKit, and Seaborn libraries. Indeed, they play a vital role in enabling researchers to handle, manipulate and understand the data effectively, which then allows users to choose the right ML approach, and assess their model performance. There will then be an exploration of the various AI/ML approaches that have been implemented for different toxicity endpoints, assessing their respective strengths and limitations. By discussing these various approaches, the chapter aims to provide readers with a comprehensive understanding of AI/ML techniques available for toxicity prediction, followed by a brief literature review on published AI/ML models. The following section also explores knowledge-based approaches, including methods like Matched Molecular Pairs (MMPs), structural alerts, chemical transformations, and rule-based molecular property assessments. These approaches are then followed by deep learning-based techniques aimed at enhancing candidate optimization through the replacement of structural motifs that may contribute to toxicity.

Overall, this chapter covers basic concepts and practical aspects of dataset curation, dataset preparation, open-source libraries models, and implementation of various molecular featurization's and AI/ML modeling approaches. By equipping readers with the necessary knowledge and tools, this chapter enables them to effectively harness AI/ML methodologies in their own research endeavors.

## **9.2 DATA SOURCES FOR TOXICITY ENDPOINTS**

There are numerous scientific articles, open-source, and commercial databases that provide curated datasets on various toxicity endpoints. These websites offer downloadable formats that support data analytics, machine learning, and deep learning approaches. Table 9.1 provides an overview of popular open-source databases, showcasing the total number of compounds, toxicity endpoints covered, and latest update. GOSTAR [19] and ReaxysMedchem [20] are commercial medicinal chemistry databases that frequently update and curate the latest scientific literature for compounds with biological targets data and for various ADMET endpoints (Table 9.1). Users can access and download these datasets through paid subscriptions. The main advantage

of commercial databases is that they tend to be larger and filled with more frequently updated data, making it potentially better for specific AI/ML model development.

Apart from these commercial databases, Lhasa [21, 22], Multicase [34, 24], and LeadScope [35] provide software solutions that include statistical ML and knowledge-based prediction models and property databases for genotoxicity, carcinogenicity, and various other toxicity endpoints. Users can access these tools, databases on a license basis, enabling them to utilize predictive models for their specific needs.

**TABLE 9.1** List of Databases for Drug Toxicity Endpoints

S. No	Database	No of Cpd	Endpoints	Feature types	Last update	Ref	Access
1	TOXRIC	113,372	13	39	2022	26	Free
2	ToxCast	9,511	1	-	2022	27	Free
3	CEBS	>11,000	3	1	2017	28	Free
4	DILIrank	1,036	1	-	2022	29	Free
5	ToxicoDB	286	3	1	2020	30	Free
6	DrugMatrix	>600	-	1	2005	31	Free
7	T3DB	3,678	-	5	2010	32	Free
8	TDC	344,267	9	-	2023	33	Free
9	MoleculeNet	18,120	3	1	2018	34	Free
10	ChEMBL	1.1 M	-	-	2023	35	Free
11	GOSTAR	9.15 M	All	-	2023	19	Commercial
12	ReaxysMedchem	34.29 M	All	-	2021	20	Commercial

TDCCommons (TDC), and MoleculeNet are the most recent databases that provide curated and pre-processed datasets for data analytics, or AI/ML modeling tasks, benchmarks. In TDC and MoleculeNet provides the functionality from Python code for data retrieval, processing, ML modeling, and for benchmarks purposes. TOXRIC web-based platform allows users to access, and download compounds which are classified under 1,474 different endpoints, and with practical benchmarks [36].

### 9.3 PACKAGES FOR DATA PROCESSING AND AI/ML

Python [37] and KNIME [38] are popular and open-source packages for data processing and AI/ML tasks. Python, a versatile and widely used programming



language, provides powerful libraries for efficient data manipulation and analysis (like Pandas and NumPy). With libraries such as TensorFlow and PyTorch, it offers robust frameworks for developing and deploying AI/ML models. It also provides flexibility and customization options for advanced data preprocessing and feature engineering. On the other hand, KNIME is a visual workflow platform that enables users to design data processing pipelines without extensive programming knowledge. Both Python and KNIME offer valuable resources for data processing and AI/ML tasks, allowing users to choose the approach that best suits their needs and expertise. In this chapter, our focus will be on Python libraries. Table 9.2 lists some of the essential Python packages that will be covered in the upcoming sections of this chapter.

## 9.4 DATA PROCESSING

Data processing is a critical step in the ML model pipeline, as the quality and reliability of the input data greatly impacts the model performance and accuracy of the models. This is especially key in fields such as drug discovery, where handling complex chemical structures and biological data requires careful preprocessing. Next, we delve into the process of managing chemical datasets in the .csv format, which contain SMILES representations and biological data points. To illustrate the process, from data handling, processing to analysis, we will use the human ether-à-go-go-related gene (hERG) dataset reported by Karim et al. [39], which is readily accessible for download. The hERG dataset is prepared for binary classification and Activity fields and it indicates whether the compounds are hERGblockers (1,  $<10\mu\text{M}$ ) or hERG non-blockers (0,  $\geq 10\mu\text{M}$ ).

As a first step, we import the downloaded hERG dataset into Jupyter notebook to understand and review the data.

The Jupyter Notebook code (Figure 9.2) demonstrates how to import a CSV file and its contents using the Pandas library. Now, we compute the number of toxic and non-toxic entries in the dataset and display the counts. In the next step, the code generates distribution plots to gain insights into the dataset's property thresholds.

Data processing and cleaning stage is essential while handling missing values, removing the duplicate compounds or activity data, and correcting data inconsistencies to ensure quality and integrity of data. By cleaning the data, we eliminate potential biases and errors allowing us to work with reliable and consistent data, consequently enabling more accurate interpretations and insights.

**TABLE 9.2** List of Key Python Packages That Discussed in This Book Chapter

Packages	Features	Home Page
Pandas	Data manipulation and analysis, handling missing data, time series functionality, powerful data structures	<a href="https://pandas.pydata.org/">https://pandas.pydata.org/</a>
NumPy	Numerical computing, multi-dimensional arrays, mathematical functions, linear algebra operations	<a href="https://numpy.org">https://numpy.org</a>
rdKit	Cheminformatics library for chemical structure handling, substructure searching, similarity and property calculations	<a href="https://www.rdkit.org">https://www.rdkit.org</a>
Seaborn	Statistical data visualization, attractive and informative statistical graphics, built-in themes	<a href="https://seaborn.pydata.org/">https://seaborn.pydata.org/</a>
TensorFlow	Deep learning framework, neural networks, numerical computation, model training and deployment	<a href="https://www.tensorflow.org/">https://www.tensorflow.org/</a>
PyTorch	Deep learning framework, dynamic neural networks, model training and deployment, GPU acceleration	<a href="https://pytorch.org/">https://pytorch.org/</a>
DeepChem	Cheminformatics for deep learning tasks, molecular featurization's, generative AI models, virtual screening, drug discovery	<a href="https://deepchem.io/">https://deepchem.io/</a>
DataMol	Cheminformatics library for molecular property prediction, molecular descriptors, chemical fingerprints	<a href="https://datamol.io/">https://datamol.io/</a>
Openbabel	Another popular library for Chemical structure generation and manipulation, molecular properties, Substructure, and similarity searching	<a href="https://openbabel.org/wiki/Python">https://openbabel.org/wiki/Python</a>
Knime	Visual workflow platform, data pre-processing, machine learning, integration with various tools and frameworks	<a href="https://www.knime.com">knime.com</a>
Pipeline Pilot	Scientific data analysis and modeling, workflow automation, cheminformatics, data visualization	<a href="https://www.3ds.com/products-services/biovia/products/pipeline-pilot/">https://www.3ds.com/products-services/biovia/products/pipeline-pilot/</a>

```
#Import Libraries
import pandas as pd
import rdkit

#Use Pandas to read the HERG dataset
data = pd.read_csv("train_validation_cardio_tox_data.csv", low_memory=False)
data.head(4)
```

[14]:

	ACTIVITY	smiles	ABC	ABCGG	nAct
0	1	<chem>Fc1ccc(-n2cc(NCCN3CCCCC3)nn2)cc1F</chem>	17.108462	13.312773	1
1	0	<chem>COc1cc(N2Cc3ccc(Sc4ccc(F)cc4)nc3C2=O)ccc1OCCN1...</chem>	26.886637	19.283801	1
2	0	<chem>CCOC(=O)[C@H]1CC[C@@H](N2CC(NC(=O)CNC3nn(C(N)=...</chem>	28.187229	20.640811	1
3	0	<chem>N[C@@H](Cn1c(=O)cnc2ccc(F)cc21)C1CCC(NCc2ccc3c...</chem>	27.105416	18.346730	1

4 rows x 997 columns

```
[11]:
# Count toxic and non-toxic entries
toxic_count = data[data['ACTIVITY'] == 1].shape[0]
non_toxic_count = data[data['ACTIVITY'] == 0].shape[0]

print("Toxic count:", toxic_count)
print("Non-toxic count:", non_toxic_count)
Toxic count: 6643
Non-toxic count: 5977
```

**FIGURE 9.2** Jupyter Notebook Code snippet shows the dataset import using the Pandas library.

Transforming and manipulating the data, such as feature scaling, encoding categorical variables, or handling outliers, prepares the dataset for analysis. Data processing enhances the compatibility of the data with various algorithms and models, facilitating meaningful analysis.

Next step is to understand and examine the structure and content of the dataset to gain an understanding of its features and data types. Please refer to Figure 9.3 for details. Please refer to the author's Github repository for additional examples and features of data handling and processing examples.

The next step is the visualization of the data through plots and graphs which provide a comprehensive overview and aid in identifying patterns, trends, and relationships. Plots allow us to explore the distribution of

variables, detect outliers, and understand the data's central tendencies and variability. By visualizing relationships between variables, we can uncover correlations, dependencies, and potential insights that might not be apparent from raw data alone.

```
# Display the first few rows of the dataset
print(data.head())

# Check the dimensions of the dataset
print(data.shape)

# Get an overview of the columns and data types
print(data.info())

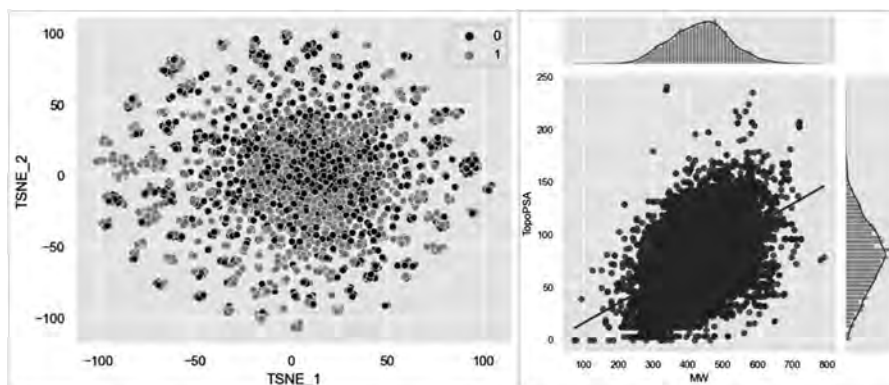
# Compute basic statistical summary of the dataset
print(data.describe())
```

**FIGURE 9.3** Jupyter Notebook Code snippet shows the dataset import using the Pandas library.

For example, t-SNE distribution and joint plots provide powerful tools to explore and understand the distribution and diversity of a dataset (Figure 9.4). They reveal hidden structures, uncover patterns, and aid in identifying clusters or groups of similar data points. These visualizations assist in making data-driven decisions, such as identifying distinct subpopulations, detecting potential outliers, evaluating the effectiveness of feature representations, and guiding further analysis or modeling tasks. By leveraging t-SNE plots, we gain valuable insights into the underlying structure and diversity of the dataset, facilitating a deeper understanding of its characteristics.

The t-SNE (t-Distributed Stochastic Neighbor Embedding) is a popular technique for reducing the dimensions of high-dimensional data and visualizing it in a lower-dimensional space. It is commonly utilized in the field of chemistry to gain insights into the chemical space and molecular distribution within datasets. In this context, t-SNE can help analyze the distribution of hERG blockers (tagged as 1) and non-blockers (tagged as 0) in the dataset. Based on the information provided, Figure 9.4 illustrates the t-SNE plot, where hERG blockers and non-blockers are evenly distributed

across multiple clusters. This distribution indicates that the two classes exhibit diverse molecular characteristics and occupy different regions of the chemical space.



**FIGURE 9.4** t-distributed stochastic neighbor embedding (t-SNE) distribution and joint plots indicates the human ether-à-go-go-related gene (hERG) dataset distribution.

The t-SNE plot suggests that the dataset may be suitable for training a machine learning model. The well-separated clusters imply that there are distinct features that differentiate hERG blockers from non-blockers, which can potentially be captured by a predictive model.

## 9.5 MOLECULAR FEATURIZATION

Molecular featurization is a key step converting molecules (e.g., SMILES strings or 3D structures) into strings of numbers useable by machines. These methods are essential for extracting significant information from molecular structures. Table 9.3 highlights some of the well-known featurization methods such as fingerprints, graph-based, shape, molecular descriptors, and fragment-based approaches. Each featurization method has its pros and cons, and understanding these trade-offs is essential for selecting the most suitable approach based on specific requirements.

These molecular fingerprints capture essential molecular properties and are used as input features for model building. This write-up will delve into some of the popular molecular featurization methods and respective code snippets from the python libraries such as rdKit, DeepChem, and MolFeat from DataMol. Let's start with Morgan fingerprint.

**TABLE 9.3** List of Known Molecular Featurization's for Machine Learning Modeling

Featurization Method	Types/Options	Description	Libraries
Fingerprints	Morgan fingerprints	Circular fingerprints that encode substructures around each atom in a molecule.	rdKit, DeepChem, Datamol
	Molecular ACCess Systems (MACCS) keys	Binary fingerprints based on a predefined set of chemical substructures (termed as the MACCS)	rdKit, Datamol, DeepChem
	PubChem fingerprints	Fingerprints derived from PubChem Compound Database, capturing substructures and chemical properties.	rdKit
	Avalon fingerprints	Structural fingerprints that encode molecular substructures and patterns.	rdKit
Molecular Descriptors	1D, 2D, and 3D	Molecular Weight (MW), LogP, TPSA, Rotatable Bonds, Hydrogen Donors, Hydrogen Acceptors, Molecular Shape, Surface Area, Principal Moments of Inertia	rdKit, DeepChem, datamol, Openbabel
Molecular Graph	Graph convolutional networks (GCN)	Graph-based featurization method that represents molecules as graphs, capturing atom and bond connectivity.	rdKit, DeepChem, TensorFlow
	Graph isomorphism networks	Neural network architecture for graph representation learning, enabling learning of molecular features from graph structures.	DeepChem, TensorFlow
	Message passing neural networks (MPNN)	Neural networks that operate on graphs, allowing information exchange and aggregation among neighboring atoms and bonds.	DeepChem, TensorFlow
Electrostatics	Coulomb matrices	Encodes the Coulombic interactions between atoms, capturing the distribution of electric charges within a molecule.	rdKit, DeepChem
	Extended-connectivity fingerprints (ECFP)/Coulomb features	Combined featurization method that captures both structural information from fingerprints and Coulombic interactions within a molecule.	rdKit
Fragments	Fragment descriptors	Encodes molecular structures as fragments or substructures, capturing local structural information within a molecule.	rdKit, DeepChem, MolFeat
Shape	Pharmacophore features	Describes molecular shape and spatial arrangements, often used for ligand-based virtual screening and shape similarity searching.	rdKit, DeepChem

### 9.5.1 THE MORGAN FINGERPRINT

Molecular fingerprinting technique is widely implemented in drug discovery and cheminformatics. It captures the molecular topology by encoding the presence or absence of substructures, which are represented as bits in a binary array. The radius parameter determines the size of the substructures considered during fingerprint generation. Detailed information of this class of fingerprints can be found in reference [40].

The Morgan fingerprint algorithm starts by assigning unique identifiers to each atom in the molecule. Then, it iteratively expands the fingerprint by including neighboring atoms up to a specified radius. As the radius increases, larger and more complex substructures are considered.

The fingerprint is generated by hashing the atom identifiers and concatenating them. Each hash value corresponds to a specific bit in the fingerprint. If a substructure is present in the molecule, the matching bit set to 1; or otherwise, it is set to 0. The resulting binary array represents the Morgan fingerprint for the molecule.

Let's calculate fingerprints using `rdKit`.

In the code snippet (Figure 9.5), we used `rdKit` library to load a molecule from its SMILES representation. The `'GetMorganFingerprintAsBitVect'` function generates the Morgan fingerprint with a radius of 2 and a bit length of 1024. Finally, the fingerprint is converted into a binary array using the `ToBitString` method.

Now, let's see how to calculate Morgan fingerprints for a dataset and assuming your dataset.

```
from rdkit import Chem
from rdkit.Chem import AllChem

# Load molecule
mol = Chem.MolFromSmiles('CC(=O)OC1=CC=CC=C1C(=O)O')

# Generate Morgan fingerprint
fp = AllChem.GetMorganFingerprintAsBitVect(mol, radius=2, nBits=1024)

# Convert to binary array
fp_array = fp.ToBitString()
```

**FIGURE 9.5** Code snippet shows Morgan fingerprint implementation using `rdKit`.

In the code snippet (Figure 9.6), we load the dataset in CSV file format using the pandas library. We then iterate over the “smiles” column and calculate the Morgan fingerprint for each molecule using rdkit. The fingerprints are stored in the fingerprints, and subsequently added to the dataset. Finally, the dataset with Morgan fingerprints can be used for ML tasks.

```
from rdkit import Chem
from rdkit.Chem import AllChem
import pandas as pd

# Load dataset from CSV
dataset = pd.read_csv('dataset.csv')

fingerprints = []

# Calculate Morgan fingerprints for each molecule
for smiles in dataset['smiles']:
    mol = Chem.MolFromSmiles(smiles)
    fp = AllChem.GetMorganFingerprintAsBitVect(mol, radius=2, nBits=1024)
    fp_array = fp.ToBitString()
    fingerprints.append(fp_array)

# Add fingerprints to the dataset
dataset['morgan_fingerprint'] = fingerprints

# Use the dataset with Morgan fingerprints for ML modeling
# ... (model training and evaluation)
```

**FIGURE 9.6** Code snippet shows Morgan fingerprint calculation for a dataset using rdkit.

DeepChem provides convenient tools for working with datasets and generating fingerprints. In the below code snippet, we utilize the DeepChem library to load the dataset from either a CSV file or a SMILES file. We define the Morgan fingerprint featurizer with a radius of 2 and a size of 1024 (Figure 9.7). Then, we featurize the dataset to calculate the Morgan fingerprints using the defined deepchemfeaturizer. The resulting dataset is ready for ML modeling tasks.

Now let's see how to calculate using MolFeat (DataMol): In the below code snippet, we iterate over the “smiles” column of the dataset, calculate the Morgan fingerprint for each molecule using the morgan\_fingerprint function



from the MolFeat library and subsequently storing fingerprints and adding a new column to the dataset. The transformed dataset is now ready for ML tasks (Figure 9.8).

```
import deepchem as dc

# Load dataset from CSV or SMILES file
dataset = dc.data.CSVLoader(['smiles', 'target'], 'dataset.csv')
# Or: dataset = dc.data.SmilesLoader(['smiles', 'target'], 'dataset.smi')

# Define the Morgan fingerprint featurizer
featurizer = dc.feet.CircularFingerprint(radius=2, size=1024)

# Featurize the dataset to calculate Morgan fingerprints
dataset = featurizer.featurize(dataset)

# Use the dataset with Morgan fingerprints for ML modeling
# ... (model training and evaluation)
```

**FIGURE 9.7** Code snippet shows Morgan fingerprint calculation for a dataset using DeepChem.

### 9.5.2 GRAPH-BASED FEATURIZATION'S

Graph-based ML models gained attention in drug discovery for their ability to capture the structural and relational information present in molecules. More specifically, molecular graphs provide a graphical representation of molecules by representing atoms referred to as nodes and bonds referred to as edges. Each atom and bond can be associated with specific attributes such as atomic number, hybridization state, and bond type.

This section provides an overview of graph-based ML models for drug toxicity prediction and explores various graph-based featurization's available in DeepChem, TensorFlow, and rdKit. Additionally, code snippets with explanations demonstrate the implementation of these featurization's using a dataset in CSV format.

*GraphConv Featurization:* DeepChem, a popular library for graph-based ML, provides the GraphConv featurization technique. GraphConv employs a convolutional neural network (CNN) approach to learn informative molecular representations (Figure 9.9). Here's an example code snippet for featurizing a dataset using DeepChem's GraphConv.

```
from molfeat.features import morgan_fingerprint
import pandas as pd
from rdkit import Chem

# Load dataset from CSV or SMILES file
dataset = pd.read_csv('dataset.csv')
# Or: dataset = pd.read_csv('dataset.smi', sep=' ', names=['smiles'])

fingerprints = []

# Calculate Morgan fingerprints for each molecule
for smiles in dataset['smiles']:
    mol = Chem.MolFromSmiles(smiles)
    fp_array = morgan_fingerprint(mol, radius=2, size=1024)
    fingerprints.append(fp_array)

# Add fingerprints to the dataset
dataset['morgan_fingerprint'] = fingerprints

# Use the dataset with Morgan fingerprints for ML modeling
# ... (model training and evaluation)
```

FIGURE 9.8 Code snippet shows Morgan fingerprint calculation for a dataset using MolFeat.

```
import deepchem as dc

# Load dataset from CSV
dataset = dc.data.CSVLoader(['smiles', 'toxicity'], 'dataset.csv')

# Featurize dataset using GraphConv
featurizer = dc.feat.graph_features.GraphConvFeaturizer()
dataset = featurizer.featurize(dataset)

# Use the featurized dataset for ML modeling
# ... (model training and evaluation)
```

FIGURE 9.9 Code snippet shows *GraphConvfeaturization* for a dataset using DeepChem.

In the code snippet above, we initialize the *GraphConvFeaturizer* and apply it to the dataset using the deepchem *featurize* method. This transforms

the molecules in the dataset into graph-based features suitable for deep learning tasks.

*rdKit Graph Featurization:* rdKit also offers graph-based featurizations for molecule representation. Here's an example code snippet for featurizing a dataset using rdKit's graph featurization.

In the code snippet (Figure 9.10), we load the dataset from a CSV file using pandas. Then, we convert the SMILES strings in the dataset to rdKit molecules. We generate graph-based features by converting the molecules into adjacency matrices using the *ToAdjacencyMatrix* function. The resulting features are stored in the 'graph\_features' column of the dataset and can be used for ML modeling.

```
from rdkit import Chem
from rdkit.Chem import PandasTools

# Load dataset from CSV
dataset = pd.read_csv('dataset.csv')

# Convert SMILES strings to RDKit molecules
PandasTools.AddMoleculeColumnToFrame(
    dataset,
    smilesCol='smiles',
    molCol='mol'
)

# Generate graph-based features
dataset['graph_features'] = dataset['mol'].apply(
    lambda mol: mol.ToAdjacencyMatrix()
)

# Use the featurized dataset for ML modeling
# ... (model training and evaluation)
```

**FIGURE 9.10** Code snippet shows graph featurization for a dataset using rdKit.

Molecular featurization methods, such as fingerprints, graph-based approaches, descriptors, shape-based methods, and electrostatics, offer valuable tools for analyzing and characterizing molecular structures, but each method has some strengths and limitations. Let's discuss some of the central points here:

- Fingerprints efficiently capture global and local features, enabling efficient similarity searching. However, they may lack detailed atomic-level information (e.g., atomic partial charges) and struggle with discriminating between structurally similar molecules with different properties.
- Graph-based methods excel in capturing connectivity patterns and molecular interactions, making them suitable for complex systems, but can be computationally demanding and sensitive to representation choices.
- Descriptors provide a wide range of properties for detailed molecular characterization, but their interpretation and selection require domain-specific knowledge, and they can be sensitive to conformational variations.
- Shape-based methods excel in capturing 3D structure similarity but may oversimplify other molecular features and be computationally intensive (e.g., Boltzmann distribution may be considered).
- Electrostatic methods focus on charge distribution and interactions, offering insights into complementarity, yet they may oversimplify interactions and neglect solvent effects.

Choosing the appropriate featurization method depends on the specific task, available data, and requirements of the application. For example, the long range inductive effect in pKa prediction may not be properly captured using featurizations considering small fragments only. Selecting the proper featurization method also involves considering the trade-off between representation power, computational efficiency, and compatibility with the chosen ML algorithm. By understanding the pros and cons of each featurization method, researchers and practitioners can make helpful decisions to extract meaningful information from molecular structures in drug discovery and cheminformatics applications.

## 9.6 IMPLEMENTATION OF AI/ML MODELS

AI is a swiftly evolving discipline in the realm of computer science, dedicated to creating machines or computational models that possess the ability to accomplish various cognitive tasks. In this contemporary review, we specifically refer to the role of AI/ML techniques for predicting and evaluating chemical toxicokinetic endpoints. ML, a subset of AI, involves the use of mathematical or computer algorithms to train models in solving

problems or undertaking tasks based on specific input parameters. ML is typically classified into three types: supervised, unsupervised, and reinforcement learning (Table 9.4). Regression is a ML technique used to predict continuous numerical values based on dataset input features (e.g., pKa values ranging from  $-1$  to  $15$ ). It aims to find a mathematical relationship between the independent and dependent variables, that allows the estimation of unknown values. Classification is a machine learning approach that assigns input data into predefined classes or categories based on their input features (e.g., toxic or not toxic). It aims to learn a decision boundary to separate different classes and make predictions on new, unseen data. Classification models are commonly used in various applications, including email spam detection, image recognition, and sentiment analysis. Table 9.4 provides an overview of commonly employed machine learning methods in drug discovery.

In this study, we discussed four popular AI/ML models: Support Vector Machines (SVM), Extreme Gradient Boosting (XGBoost), Random Forest (RF), and Graph Convolutional Networks (GCN).

### 9.6.1 SUPPORT VECTOR MACHINES (SVM)

Support Vector Machines is a supervised ML algorithm for regression and classification tasks. The process involves taking the input data and projecting it into a feature space with a high number of dimensions (e.g., descriptors values may be squared). Then, a decision boundary is determined in this feature space, aiming to maximize the separation between different classes.

The use of different kernel functions, including linear function, polynomial function, and radial basis function (RBF), enables the capturing of intricate relationships in the data. SVM excels in handling non-linear data and is resistant to overfitting, making it highly advantageous. However, SVM's computational complexity can increase with large datasets, and the choice of the appropriate kernel function and regularization parameters requires very careful consideration. Overall, SVM is widely used in many domains, including image recognition, text classification, and bioinformatics, thanks to its excellent generalization (accuracy on datapoints not appearing in the training set) performance and ability to handle complex decision boundaries.

The following basic code demonstrates SVM implementation (Figure 9.11).

In this code block, we utilize the scikit-learn (sklearn) library to implement the SVM classifier for binary classification.

**TABLE 9.4** An Overview of Commonly Employed Machine Learning Methods in Drug Discovery

Method	Algorithm	Type	Description	Pros	Cons
Supervised	Multiple Linear Regression (MLR)	Regression	Utilizing a multitude of explanatory factors, a multivariate linear equation is employed to forecast the result of a response variable.	Easy interpretation, fast computation	Assumes linearity, may not capture complex relationships
	Naïve Bayes Classifier	Classification	Using the principles of Bayes’ theorem and assuming that the molecular descriptors (which serve as explanatory variables) are conditionally independent of each other,	Efficient, handles high-dimensional data well	Strong independence assumption, may not capture dependencies between descriptors
Supervised nonlinear	k-nearest neighbors (k-NN)	Classification	This process involves categorizing a test chemical based on its similarity to a set of training chemicals. The classification is determined by identifying the training chemicals that are closest in distance to the test chemical.	Simple, no training required	Sensitive to noise, may have high computational cost for large datasets
	Support Vector Machine (SVM)	Classification	The process involves converting molecular descriptor vectors into a feature space with increased dimensions, with the goal of establishing a hyperplane that maximizes the margin and effectively separates different chemical compounds.	Effective in high-dimensional spaces, works well with small datasets	Requires careful selection of kernel function, can be computationally expensive for large datasets
	Decision Trees (DT)	Classification	Every model consists of a collection of rules arranged in a tree-like structure. These rules are used to make predictions about the toxicity of a specific chemical for a particular endpoint.	Easy to understand, handles non-linear relationships	Prone to overfitting, can create complex trees that are difficult to interpret

**TABLE 9.4** (Continued)

Method	Algorithm	Type	Description	Pros	Cons
Ensemble methods	Random Forest (RF)	Classification	Combines the bagging and random spaces approaches with decision tree base models to make predictions	Robust against overfitting, handles high-dimensional data	Can be computationally expensive, requires careful hyperparameter tuning
	Extreme Gradient Boosting (XGBoost)	Classification, Regression	Boosting algorithm that improves predictive performance through gradient boosting	High predictive performance, handles complex relationships, feature importance estimation	Can be computationally expensive, requires careful hyperparameter tuning
Artificial neural networks	Backpropagation Neural Networks (BNN)	Classification	Every neuron is categorized into three layers, and information travels from the input to the output layer by passing through the hidden layers.	Can capture complex relationships, powerful for non-linear problems	Requires careful tuning of architecture and hyperparameters, can be prone to overfitting
	Bayesian-regularized Neural Networks (BRNN)	Classification	Applies Bayesian methods for regularization to balance model complexity against accuracy	Handles overfitting, provides uncertainty estimation	Requires longer training time, may have higher computational cost
	Associative Neural Networks (ANN)	Classification	Applies ensemble learning to backpropagation neural networks for improved prediction accuracy	Improved prediction accuracy through ensemble learning	Increased complexity, requires more computational resources

**TABLE 9.4** (Continued)

Method	Algorithm	Type	Description	Pros	Cons
Unsupervised	Deep Neural Networks (DNN)	Classification	Artificial neural networks are computational models composed of multiple layers, including hidden layers, which are commonly referred to as deep learning networks.	Has the ability to detect intricate patterns and achieve exceptional performance across various fields.	Requires large amounts of data, computationally intensive training, may be prone to overfitting
	Message Passing Graph networks	Classification, Regression	Captures structural information and handles complex graph-structured data	Captures structural information, handles complex graph-structured data	Can be computationally expensive, sensitive to graph size and topology
	Graph Convolution	Classification, Regression	Handles graph-structured data and variable-sized graphs	Captures graph-structured data, handles variable-sized graphs	Limited interpretability, sensitive to graph noise and heterogeneity
	Principal Component Analysis (PCA)	Dimensionality Reduction	Reduces the dimensionality of data while preserving most of its variation	Data visualization, dimensionality reduction	May lose some information, assumes linear relationships
	Kohonen's Self-Organizing Maps	Clustering	Transforming molecules from their descriptor space onto a two-dimensional grid of neurons using their topological relationships.	Visualization of data clusters, topological relationships	Sensitivity to initial conditions, may require parameter tuning



```
# SVM implementation

from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(
    features, labels, test_size=0.2, random_state=42
)

# Initialize the SVM classifier
svm = SVC()

# Train the model
svm.fit(X_train, y_train)

# Make predictions on the test set
y_pred = svm.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
```

**FIGURE 9.11** Code snippet shows SVM Implementation.

Here's a breakdown of the steps involved: The first step is to import the "train\_test\_split" function from a sklearn module called "*model\_selection*." "train\_test\_split" function is used to split the dataset into training and testing sets. More specifically, "*features*" represents the input features (*X*) of the dataset, while labels represents the corresponding target labels (*y*). "*test\_size=0.2*" indicates that 20% of the data considered for testing, and 80% considered for training. More advanced methods to split and ensure that the testing and training sets are not including analogous molecules are recommended but are not discussed here.

Finally, "*random\_state=42*" is set to ensure the results' reproducibility, by providing the same random seed for splitting the dataset (any numbers can be used). Following the splitting of the original dataset, an instance of the SVM classifier is created by calling the *SVC()* function without passing any parameters, hereby the empty parentheses, leading to its initialization

with default settings. Finally, the method needs to be trained on the training subsets earlier made. This can be done using the `fit()` function, which will train the SVM classifier on the training, *X\_train* and their respective target labels (e.g., experimental data values), *y\_train*.

### 9.6.2 RANDOM FOREST (RF)

RF is an ensemble learning approach that combines multiple decision trees to make predictions. Instead of using a single tree, RF constructs an ensemble of trees (i.e., a forest) by training each tree on a different subset of the data and using a random selection of features. When making predictions, the results from all the individual trees are combined to reach a final classification decision. RF excels in handling data with many dimensions and helps to prevent overfitting, where the model becomes too specialized to the training set and performs poorly on new data (poor generalization).

Nevertheless, RF can impose a significant computational burden, particularly when dealing with extensive datasets or employing a large number of trees. Furthermore, the resulting model may prove more difficult to interpret in comparison to a solitary decision tree, although feature importance can be extracted. Despite these limitations, RF remains a widely used algorithm in various domains due to its effectiveness in handling complex classification problems and providing reliable predictions.

Figure 9.12 shows the basic code implementation for RF. The code (Figure 9.12) initializes a RF classifier using the *RandomForestClassifier* from the *sklearn.ensemble* module. The *n\_estimators* parameter is set to 100, which determines the number of decision trees that will be created in the RF ensemble. The `fit` method is called on the *rf* classifier with the training (*X\_train* and *y\_train*) as arguments. The *accuracy\_score* function from the *sklearn.metrics* module is considered to calculate the accuracy of the model's predictions.

### 9.6.3 EXTREME GRADIENT BOOSTING (XGBOOST)

Extreme Gradient Boosting is a powerful and widely used ML algorithm that excels in both regression and classification tasks. XGBoost is a powerful ensemble learning technique that enhances prediction accuracy by aggregating the predictions of multiple decision trees, known as weak

models, to generate a strong and reliable final prediction. XGBoost follows a gradient boosting approach, where each subsequent decision tree is constructed to rectify the mistakes made by the previous trees. This iterative process enables XGBoost to steadily enhance its predictive abilities and capture intricate patterns within the dataset. It leverages gradient descent optimization techniques to fine-tune the objective function, effectively reducing both bias and variance in the model's predictions.

```
# Random Forest implementation

from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(
    features, labels, test_size=0.2, random_state=42
)

# Initialize the Random Forest classifier
rf = RandomForestClassifier(n_estimators=100)

# Train the model
rf.fit(X_train, y_train)

# Make predictions on the test set
y_pred = rf.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
```

**FIGURE 9.12** Code snippet shows Random Forest (RF) Classifier implementation.

One of the key advantages of this method is its ability to handle high-dimensional data and large feature spaces. It automatically handles missing values, supports regularization techniques to prevent overfitting, and provides built-in feature importance estimation, allowing for insightful variable selection. It is widely used in various domains, including healthcare, finance, and natural language processing, where accurate predictions are

critical. However, XGBoost's main drawbacks (as most ML techniques) include the potential for high computational costs and the need for careful hyperparameter tuning.

The basic code snippet (Figure 9.13) showcases XGBoost implementation.

```
# XGBoost implementation

import xgboost as xgb
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(
    features, labels, test_size=0.2, random_state=42
)

# Initialize the XGBoost classifier
xgb_model = xgb.XGBClassifier()

# Train the model
xgb_model.fit(X_train, y_train)

# Make predictions on the test set
y_pred = xgb_model.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
```

**FIGURE 9.13** Code snippet shows Extreme Gradient Boosting (XGBoost) Classifier implementation.

In this example, the code assumes you have your dataset loaded into the *X* variable for features and *y* variable for the target variable (Figure 9.13). Dataset splits into training and test sets using *train\_test\_split* from scikit-learn. Then, it creates an XGBoost classifier (*XGBClassifier*) for classification tasks or an XGBoost regressor (*XGBRegressor*) for regression tasks.

The model is then trained using the “*fit*” method, where it learns from the training data. It then makes predictions on the test data using the “*predict*” method. After that, evaluation metrics like accuracy (for

classification tasks) or *mean squared error* (for regression tasks) are calculated and displayed.

#### 9.6.4 GRAPH CONVOLUTIONAL NETWORKS (GCNS)

Graph Convolutional Networks (GCNs) are neural network models that operate directly on graph-structured data. They utilize a series of graph convolutional layers to capture the local and global structural information of molecules. GCNs have demonstrated remarkable performance in tasks such as property prediction and drug discovery.

GCNs are deep learning models specifically designed to operate on graph-structured data. As mentioned above, graphs are mathematical structures consisting of nodes (also known as vertices) and edges that represent relationships or connections between the nodes. The aggregation of information in GCNs is typically performed using a weighted sum or a learnable function that takes into account the features of neighboring nodes. This enables each node to gather information from its neighbors and update its own representation accordingly. The weights assigned to the neighboring nodes can be learned during the training process, allowing the GCN to adapt to the specific characteristics of the graph.

Overall, GCNs provide a powerful framework for learning from graph-structured data, allowing for the effective analysis and understanding of complex relationships and patterns in graph-based domains.

Figure 9.14 shows a code implementation using DeepChem's `GraphConv` model.

In the code snippet above, we start by loading the hERG dataset using the `load_herg` function from the `molnet` module in DeepChem. The dataset is then split into training, testing, and validation sets (Figure 9.14). Next, we create an instance of the `GraphConvModel` provided by DeepChem's `models` module.

We then train the model using the training dataset by calling the `fit` method and specifying the number of epochs (a term used in ML to refer to iterations). The model learns to predict hERG inhibition based on the molecular graph structures. After training, we evaluate the model's performance on the test set using a suitable metric, such as the ROC AUC score.

Finally, trained model are used to predict the activity/toxicity for new compounds. We provide the molecular representation (SMILES or molecule object) of the new compounds, extract the graph features using

*graph\_features* from the *feat* module, and then obtain the predicted labels using the *predict\_on\_batch* method (Figure 9.15).

```
# GraphConv implementation using DeepChem

import deepchem as dc
import numpy as np

# Load the HERG dataset
tasks, datasets, transformers = dc.molnet.load_herg()

# Split the dataset into training and testing sets
train_dataset, valid_dataset, test_dataset = datasets

# Create the GraphConv model
model = dc.models.GraphConvModel(
    len(tasks),
    batch_size=32,
    mode='classification'
)

# Fit the model to the training data
model.fit(train_dataset, nb_epoch=50)

# Evaluate the model on the test set
metric = dc.metrics.Metric(dc.metrics.roc_auc_score)
test_scores = model.evaluate(test_dataset, [metric])
```

**FIGURE 9.14** Code snippet shows GraphConv implementation using DeepChem.

### 9.6.5 EVALUATION METRICS

We mentioned evaluation metrics for fingerprint models, including Support Vector Machine (SVM), Random Forest, XGBoost, and Graph Convolutional Model (GraphConv). We thought to delve into the theory behind common evaluation metrics and provide code snippets with explanations for calculating these metrics.

More specifically, we discuss below several evaluation metrics, including accuracy, precision, recall, F1-score, ROC curve, AUC, MAE, and RMSE.

```
# Obtain predictions on new compounds
new_compound = [...] # SMILES representation or molecule object
features = dc.featurizer.graph_features([new_compound])
predicted_labels = model.predict_on_batch(features)
```

**FIGURE 9.15** Code snippet shows GraphConv model predictions.

#### 9.6.5.1 ACCURACY

The accuracy of a model's predictions is determined by comparing the correct predictions to the total no of predictions made (Figure 9.16) in classification models (e.g., how often is the model right?). This accuracy is measured by dividing the sum of true positives (e.g., more than are active and predicted to be active) and true negatives (e.g., molecules that are inactive and predicted to be inactive) by the sum of true positives, true negatives, false positives (e.g., molecules that are inactive but predicted to be active), and false negatives (e.g., molecules that are active but predicted to be inactive).

```
from sklearn.metrics import accuracy_score

# Calculate accuracy for the predicted labels
accuracy = accuracy_score(y_true, y_pred)
```

**FIGURE 9.16** Code snippet for accuracy calculation.

#### 9.6.5.2 PRECISION, RECALL, AND F1-SCORE

These are the common metrics used in binary classification tasks. Precision is a metric that calculates the proportion of true positive predictions out of the total number of positive predictions made by the model (Figure 9.17). It assesses how accurately the model identifies positive instances. Recall, also called sensitivity, is a metric that measures the proportion of true positive predictions out of the total number of actual positive instances in the dataset. It evaluates the model's ability to capture positive instances. F1-score is a measure that combines precision and recall by taking their harmonic mean. It provides a balanced assessment, considering both the precision and recall of the model's predictions.

```
from sklearn.metrics import precision_score, recall_score, f1_score

# Calculate precision, recall, and F1-score for the predicted labels
precision = precision_score(y_true, y_pred)
recall = recall_score(y_true, y_pred)
f1 = f1_score(y_true, y_pred)
```

**FIGURE 9.17** Code snippet to compute Precision, Recall, and F1-Score.

### 9.6.5.3 RECEIVER OPERATING CHARACTERISTIC (ROC) CURVE AND AREA UNDER THE CURVE (AUC)

The ROC curve visually depicts how the sensitivity (% true positive) and specificity ( $1 - \text{false positive rate}$ ) of a classification model change at various threshold settings. The AUC quantifies the overall performance of the model by calculating the area under the ROC curve (Figure 9.18). A higher AUC indicates better model performance. AUC varies from 0 (always wrong) to 1 (always right) with an AUC value of 0.5 indicating a model with random predictions.

```
from sklearn.metrics import roc_curve, roc_auc_score
import matplotlib.pyplot as plt

# Calculate the predicted probabilities
probs = model.predict_proba(X_test)[:, 1]

# Calculate false positive rate, true positive rate, and thresholds
fpr, tpr, thresholds = roc_curve(y_test, probs)

# Calculate AUC
auc = roc_auc_score(y_test, probs)

# Plot the ROC curve
plt.plot(fpr, tpr)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve (AUC = {:.3f})'.format(auc))
plt.show()
```

**FIGURE 9.18** Code snippet to compute Receiver Operating Characteristic (ROC) curve.



#### 9.6.5.4 MEAN ABSOLUTE ERROR (MAE) AND ROOT MEAN SQUARED ERROR (RMSE)

Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) are popular metrics often utilized in regression tasks. MAE measures the average absolute difference between the predicted and actual values, whereas RMSE calculates the square root of the average squared difference (Figure 9.19).

```
from sklearn.metrics import mean_absolute_error, mean_squared_error

# Calculate mean absolute error and root mean squared error
mae = mean_absolute_error(y_true, y_pred)
rmse = mean_squared_error(y_true, y_pred, squared=False)
```

**FIGURE 9.19** Code snippet to compute MAE, and RMSE.

Utilizing appropriate evaluation metrics allows researchers and practitioners to comprehensively evaluate and compare the performances of multiple models in cheminformatics tasks.

In the literature, several studies have been developed and published for AI/ML prediction models for hERG blocking compounds. Some studies employed pharmacophore-based models using SVM and Naïve Bayes classifier [41, 42]. Other studies focused on fingerprint-based models using ML methods like k-nearest neighbor (KNN), SVM, and RF. Additionally, these models based on physicochemical molecular descriptors and various ML algorithms were explored [43–46]. However, their model prioritized performance over interpretability due to the black-box nature of deep learning models. Therefore, constructing accurate prediction models while simultaneously discovering meaningful hERG blocker patterns remains a major challenge.

These models offer valuable tools for read-across in toxicology, enabling the prediction of bioactivities or toxicity end point predictions for new compounds which are structurally related or similar analogues, without the need for additional screening (*in vitro* or *in vivo*).

#### 9.6.6 APPLICABILITY DOMAIN

Applicability Domain (AD) is a concept used in ML models to determine the range or domain of data instances where the model's predictions are

reliable. It helps identify whether a particular input sample is within the model's training distribution or falls outside it. The Applicability Domain is defined based on the similarity between the training and the compound to be predicted. Fingerprint-based methods, such as Morgan fingerprints, encode the structural features of molecules into binary bit vectors, facilitating similarity comparisons. The Random Forest algorithm uses these fingerprints to build a predictive model.

Here are some considerations for determining the threshold:

- *Training Data Coverage:* Analyze the distribution of training data in the feature space. If the training data covers a wide range of feature values and variations, a higher threshold may be suitable. Conversely, if the training data is concentrated in a specific region, a lower threshold may be preferred to limit predictions to that specific region.
- *Outlier Detection:* Consider the presence of outliers or anomalous samples in the dataset. Outliers may have significantly different features compared to the majority of the data, and setting a higher threshold can help exclude such outliers from the AD. Practically speaking, these outliers may simply be molecules tested under different conditions or wrong data points such as PAINS.
- *Model Confidence:* Evaluate the inherent uncertainty of the model's predictions. Models that provide prediction intervals or confidence scores can use these measures to define the AD threshold. Higher uncertainty may warrant a lower threshold to restrict predictions to areas with higher confidence.

Ultimately, there is no one-size-fits-all ideal threshold for AD. It depends on the specific characteristics of the problem, dataset, and model, and requires careful analysis and experimentation to determine the most appropriate threshold.

## 9.7 KNOWLEDGE-BASED APPROACHES

In drug design and development, predicting the toxicity of potential drug candidates is to ensure patient safety and optimize therapeutic outcomes. Knowledge-based predictions combining various methodologies have emerged as powerful tools to assess drug toxicity. This article explores the integration of rdKit model-based similar maps, SHapley Additive exPlanations (SHAP) (explainable AI), molecular properties, structural alerts, and molecular modification patterns (MMPs) to enhance drug toxicity predictions [47, 48].

- a. **Prediction based structural fingerprints:** Leveraging rdKit model-based similar maps, SHAP, molecular properties, structural alerts, and MMPs can provide a comprehensive understanding of toxicity risks associated with drug candidates. These methodologies enhance the interpretability and predictive accuracy of toxicity models, enabling the development of safer and more effective therapeutics. As computational techniques continue to advance, knowledge-based predictions for drug toxicity will play an increasingly vital role in the early stages of drug discovery and development.
- b. **rdKit Similar Maps:** rdKit, a widely used open-source cheminformatics toolkit, provides functionalities for generating molecular similarity maps. These maps are generated based on the similarity of molecular fingerprints or descriptors and offer insights into the structural features relevant to toxicity/end points. By comparing the similarity maps of known toxic compounds with novel drug candidates, potential toxicity may be predicted. rdKit's model-based similar maps utilize ML techniques to enhance the accuracy of toxicity predictions.
- c. **SHapley Additive exPlanations (SHAP):** SHAP is an explainable AI technique that assigns importance scores to features contributing to a particular prediction. In the context of drug toxicity, it can help identify the molecular descriptors or structural motifs that significantly influence toxicity predictions. By analyzing SHAP values, researchers can gain insights into the mechanisms underlying toxicity and prioritize modifications to reduce toxicity risks.
- d. **Rule based:** Various molecular properties, such as lipophilicity, molecular weight, hydrogen bonding potential, and pKa values, can be used as inputs for toxicity prediction models. Quantitative structure-activity relationship (QSAR) models leverage these properties to estimate toxicity based on statistical correlations with experimental toxicity data. Integrating molecular properties with other knowledge-based approaches enhances the predictive accuracy of toxicity models.
- e. **Structural Alerts:** Structural alerts are predefined chemical substructures associated with specific toxicological effects. These alerts are derived from empirical observations or expert knowledge and can be used to flag potential toxic liabilities in drug candidates. Integrating structural alerts into toxicity models enables rapid screening and identification of compounds with potential toxicity concerns, facilitating early decision-making during drug discovery.

- f. **Molecular modification patterns (MMP's):** MMPs involve the systematic exploration of chemical modifications and their impact on toxicity. By analyzing databases of structure-activity relationships, researchers can identify patterns of structural modifications that mitigate or exacerbate toxicity. This knowledge allows for the rational design of safer and better drug candidates by modifying or eliminating toxic moieties while retaining desired pharmacological properties.

## 9.8 CONCLUSIONS

AI/ML models have made remarkable progress in toxicity prediction, offering efficient and cost-effective alternatives to traditional experimental methods. These models have demonstrated the potential to accelerate toxicity assessment, reduce animal testing, and aid in chemical screening and risk assessment. However, challenges remain, including the need for more extensive and diverse datasets, addressing bias and interpretability concerns, and overcoming the domain-specific challenges of toxicology. By addressing these limitations and considering the future steps mentioned above, AI/ML models can continue to play a vital role in toxicity assessment, contributing to improved safety evaluation and decision-making processes in various industries. With continued research and collaboration, the practical implementation of these models may lead to an even further acceleration and enhancement of toxicity assessment processes, ultimately contributing to improved human and environmental health.

## 9.9 TERMINOLOGY

In this section, we discussed the list of terminology referenced in this chapter. These terms are relevant to the topics discussed in this chapter and provide a foundation for understanding deep learning concepts in the context of toxicity assessment, molecular analysis, and neural network models.

1. *Convolutional Layers:* These are specialized layers in a convolutional neural network (CNN) that perform convolution operations on input data, such as images. By sliding a set of filters over the data, these layers extract important features.
2. *Epochs:* In the context of deep learning, an epoch refers to one complete iteration through the entire training dataset during model training. It

involves both forward propagation (computing predictions) and backward propagation (updating weights based on computed errors).

3. *Jupyter Notebook*: It is an interactive web-based environment widely used for creating and sharing documents that combine code, visualizations, and explanatory text. It is a popular tool for data analysis, prototyping, and presenting code-based projects.
4. *Toxicity Endpoints*: These are specific measures or criteria used to evaluate the toxic effects of chemicals or substances. They encompass various aspects such as acute toxicity, chronic toxicity, genotoxicity, carcinogenicity, or toxicity to specific organs like the liver.
5. *Molecular Fingerprints*: These are binary representations of molecular structures that encode and capture chemical features. Molecular fingerprints are extensively employed in cheminformatics and drug discovery to analyze and compare the structural characteristics of different compounds.
6. *GCN (Graph Convolutional Network)*: GCN is a type of deep learning architecture designed to handle graph-structured data, particularly molecular graphs. It utilizes graph convolutional layers to extract and learn features from the graph representation of molecules.
7. *Hyperparameters*: These are parameters in a machine learning model that are set by the user before training and are not learned from the data. Hyperparameters influence the model's behavior and learning process, such as the learning rate, batch size, or the number of layers in a neural network.
8. *hERG*: hERG refers to the human ether-a-go-go-related gene, which codes for a specific ion channel involved in cardiac function. In drug discovery, the evaluation of hERG blocking activity is crucial due to its association with potential cardiac side effects.
9. *Genotoxicity*: It is the ability of a substance or agent to cause damage to genetic material, such as DNA. Genotoxicity is an important endpoint in toxicology studies as it assesses the potential for long-term harmful effects and the risk of mutagenesis or carcinogenesis.
10. *Carcinogenicity*: Carcinogenicity refers to the property of a substance to cause cancer or increase the risk of developing cancer. This endpoint is extensively evaluated in toxicology studies to determine the potential hazards associated with exposure to specific chemicals or compounds.
11. *Liver Toxicity*: Liver toxicity refers to the adverse effects caused by exposure to toxic substances on the liver. It is a significant concern

in drug development and toxicology studies, as the liver plays a vital role in metabolizing and detoxifying foreign compounds.

12. *Neural Network Layers*: These are fundamental components of a neural network that process and transform input data. Common types of layers include the input layer (receives the input data), hidden layers (perform computations and feature extraction), and output layer (produces the final predictions or outputs).
13. *Activation Function*: An activation function is a mathematical function applied to the output of a neuron in a neural network. It introduces non-linearity and determines whether a neuron should be activated or not. Popular activation functions include sigmoid, ReLU (Rectified Linear Unit), and softmax.
14. *Backpropagation*: Backpropagation is a training algorithm used in deep learning. It calculates and propagates the error backward from the output layer to the input layer, allowing the neural network to update its weights based on the computed errors.
15. *Loss Function*: A loss function measures the discrepancy between the predicted output of a neural network and the true output. It quantifies the performance of the network during training and serves.

## KEYWORDS

- **deep learning**
- **DeepChem**
- **hERG**
- **machine learning**
- **rdKit**
- **toxicity endpoints**

## REFERENCES

1. Sun, D., Gao, W., Hu, H., & Zhou, S. (2022). Why 90% of clinical drug development fails and how to improve it? *Acta Pharmaceutica Sinica B*, 12(7), 3049–3062.
2. Weber, S., & Gerbes, A. L. (2022). Challenges and future of drug-induced liver injury research—Laboratory tests. *International Journal of Molecular Sciences*, 23(11), 6049.

3. Kullak-Ublick, G. A., Andrade, R. J., Merz, M., End, P., Gerbes, A. L., & Aithal, G. P. (2017). Drug-induced liver injury: Recent advances in diagnosis and risk assessment. *Gut*, 66(6), 1154–1164.
4. Li, A. P. (2005). Accurate prediction of human drug toxicity: A major challenge in drug development. *Chemico-Biological Interactions*, 150(1), 3–7.
5. Tran, T. T. V., Wibowo, A. S., Tayara, H., & Chong, K. T. (2023). Artificial intelligence in drug toxicity prediction: Recent advances, challenges, and future perspectives. *Journal of Chemical Information and Modeling*, 63(9), 2628–2643.
6. Askr, H., Elgeldawi, E., Aboul Ella, H., Elshaier, Y. A. M. M., Gomaa, M. M., & Hassanien, A. E. (2023). Deep learning in drug discovery: An integrative review and future challenges. *Artificial Intelligence Review*, 56, 5975–6037.
7. Paul, D., Sanap, G., Shenoy, S., Kalyane, D., Kalia, K., & Tekade, R. K. (2021). Artificial intelligence in drug discovery and development. *Drug Discovery Today*, 26(1), 80–93.
8. Patel, V., & Shah, M. (2022). Artificial intelligence and machine learning in drug discovery and development. *Intelligent Medicine*, 2(3), 134–140.
9. Cavasotto, C. N., & Scardino, V. (2022). Machine learning toxicity prediction: Latest advances by toxicity end point. *ACS Omega*, 7(51), 47536–47546.
10. Basile, A. O., Yahi, A., & Tatonetti, N. P. (2019). Artificial intelligence for drug toxicity and safety. *Trends in Pharmacological Sciences*, 40(9), 624–635.
11. Wojtuch, A., Jankowski, R., & Podlowska, S. (2021). How can SHAP values help to shape metabolic stability of chemical compounds? *Journal of Cheminformatics*, 13, 74.
12. Jiménez-Luna, J., Grisoni, F., & Schneider, G. (2020). Drug discovery with explainable artificial intelligence. *Nature Machine Intelligence*, 2, 573–584.
13. Kırboğa, K. K., Abbasi, S., & Kuçuksille, E. U. (2023). Explainability and white box in drug discovery. *Chemistry & Biology Drug Design*, 102(1), 217–233.
14. Wang, M. W. H., Goodman, J. M., & Allen, T. E. H. (2021). Machine learning in predictive toxicology: Recent applications and future directions for classification models. *Chemical Research in Toxicology*, 34(2), 217–239.
15. Lin, Z., & Chou, W.-C. (2022). Machine learning and artificial intelligence in toxicological sciences. *Toxicological Sciences*, 189(1), 7–19.
16. Ring, C., & Rager, J. (2022, March 27–31). Machine learning in predictive toxicology: An overview and case study. Paper presented at the *Society of Toxicology 61st Annual Meeting and ToxExpo*, San Diego, CA.
17. Pu, L., Naderi, M., Liu, T., Wu, H.-C., Mukhopadhyay, S., & Brylinski, M. (2019). eToxPred: A machine learning-based approach to estimate the toxicity of drug candidates. *BMC Pharmacology and Toxicology*, 20, Article 2.
18. Mayr, A., Klambauer, G., Unterthiner, T., & Hochreiter, S. (2015). DeepTox: Toxicity prediction using deep learning. *Frontiers in Environmental Science*, 3.
19. Global Online Structure Activity Relationship Database. <https://www.gostardb.com/about/> (accessed on 11th September 2024).
20. Reaxys Medicinal Chemistry Database. (2024). Retrieved from <https://www.elsevier.com/solutions/reaxys/who-we-serve/reaxys-medicinal-chemistry> (accessed on 11th September 2024).
21. Ponting, D. J., Burns, M. J., Foster, R. S., Hemingway, R., Kocks, G., MacMillan, D. S., Shannon-Little, A. L., Tennant, R. E., Tidmarsh, J. R., & Yeo, D. J. (2022). Use of Lhasa Limited products for the in silico prediction of drug toxicity. In *Methods in Molecular Biology* (Vol. 2425, pp. 435–478).

22. Lhasa Limited. <https://www.lhasalimited.org/> (accessed on 11th September 2024).
23. Chakravarti, S. K., & Alla, S. R. M. (2019). Descriptor-free QSAR modeling using deep learning with long short-term memory neural networks. *Frontiers in Artificial Intelligence*, 2.
24. MultiCASE, Inc. <https://multicase.com/research> (accessed on 11th September 2024).
25. Leadscope (Now Instem). <https://www.instem.com/solutions/insilico/predict.php> (accessed on 11th September 2024).
26. TOXRIC Database. <http://toxric.bioinforai.tech/> (accessed on 11th September 2024).
27. ToxCast. <https://www.epa.gov/chemical-research/exploring-toxcast-data-downloadable-data> (accessed on 11th September 2024).
28. CEBS. <https://cebs.niehs.nih.gov/cebs/> (accessed on 11th September 2024).
29. DILIRank. <https://www.fda.gov/science-research/liver-toxicity-knowledge-base-ltkb/drug-induced-liver-injury-rank-dilirank-dataset> (accessed on 11th September 2024).
30. ToxicoDB. <https://www.toxicodb.ca/> (accessed on 11th September 2024).
31. Drug Matrix. <https://ntp.niehs.nih.gov/data/drugmatrix/> (accessed on 11th September 2024).
32. Toxin and Toxin Target Database. <http://www.t3db.ca/> (accessed on 11th September 2024).
33. Tdcommons. <https://tdcommons.ai/> (accessed on 11th September 2024).
34. MoleculeNet. <https://moleculenet.org/> (accessed on 11th September 2024).
35. Gaulton, A., Bellis, L. J., Bento, A. P., Chambers, J., Davies, M., Hersey, A., Light, Y., McGlinchey, S., Michalovich, D., Al-Lazikani, B., et al. (2012). ChEMBL: A large-scale bioactivity database for drug discovery. *Nucleic Acids Research*, 40, D1100–1107.
36. Wu, L., Yan, B., Han, J., Li, R., Xiao, J., He, S., & Bo, X. (2023). TOXRIC: A comprehensive database of toxicological data and benchmarks. *Nucleic Acids Research*, 51(D1), D1432–D1445.
37. Python. <https://www.python.org/> (accessed on 11th September 2024).
38. Knime. <https://www.knime.com/> (accessed on 11th September 2024).
39. Karim, A., et al. (2021). CardioTox net: A robust predictor for hERG channel blockade based on deep learning meta-feature ensembles. *Journal of Cheminformatics*, 13, 60.
40. Zagidullin, B., Wang, Z., Guan, Y., Pitkänen, E., & Tang, J. (2021). Comparative analysis of molecular fingerprints in prediction of drug combination effects. *Briefings in Bioinformatics*, 22(6), bbab291.
41. Leong, M. K. (2007). A novel approach using pharmacophore ensemble/support vector machine (PhE/SVM) for prediction of hERG liability. *Chemical Research in Toxicology*, 20(2), 217–226.
42. Wang, S., et al. (2016). ADMET evaluation in drug discovery. 16. Predicting hERG blockers by combining multiple pharmacophores and machine learning approaches. *Molecular Pharmaceutics*, 13(8), 2855–2866.
43. Doddareddy, M. R., et al. (2010). Prospective validation of a comprehensive in silico hERG model and its applications to commercial compound and drug databases. *ChemMedChem*, 5(5), 716–729.
44. Chavan, S., et al. (2016). A k-nearest neighbor classification of hERG K(+) channel blockers. *Journal of Computer-Aided Molecular Design*, 30(3), 229–236.
45. Siramshetty, V. B., et al. (2018). The Catch-22 of predicting hERG blockade using publicly accessible bioactivity data. *Journal of Chemical Information and Modeling*, 58(6), 1224–1233.
46. Didziapetris, R., & Lanevskij, K. (2016). Compilation and physicochemical classification analysis of a diverse hERG inhibition database. *Journal of Computer-Aided Molecular Design*, 30(12), 1175–1188.



47. Griffen, E., Leach, A. G., Robb, G. R., & Warner, D. J. (2011). Matched molecular pairs as a medicinal chemistry tool. *Journal of Medicinal Chemistry*, 54(22), 7739–7750.
48. Rodríguez-Pérez, R., & Bajorath, J. (2020). Interpretation of compound activity predictions from complex machine learning models using local approximations and Shapley values. *Journal of Medicinal Chemistry*, 63(16), 8761–8777.

## CHAPTER 10

---

# AI-Based Models for Prediction of Biodegradation

GANESH B. PATIL,<sup>1</sup> SOPAN N. NANGARE,<sup>2</sup> SHITAL M. PATIL,<sup>3</sup>  
SHANKARSING S. RAJPUT,<sup>4</sup> and MILIND M. PATIL<sup>5</sup>

*<sup>1</sup>Department of Pharmaceutics, H. R. Patel Institute of Pharmaceutical Education and Research, Shirpur, Dhule, Maharashtra, India*

*<sup>2</sup>Department of Pharmaceutical Chemistry, H. R. Patel Institute of Pharmaceutical Education and Research, Shirpur, Dhule, Maharashtra, India*

*<sup>3</sup>Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India*

*<sup>4</sup>S. P. D. M. Arts, S. B. B. and S. H. D. Commerce & S. M. A. Science College, Shirpur, Dhule, Maharashtra, India*

*<sup>5</sup>Department of Chemistry, Poojya Sane Guruji Vidya Prasarak Mandals Shri. S. I. Patil Arts, G. B. Patel Science and S. T. K. V. Sangh Commerce College, Shahada, Nandurbar, Maharashtra, India*

---

### ABSTRACT

Artificial intelligence (AI) has the potential to enhance biodegradation prediction dramatically. AI plays a crucial role in biodegradation prediction by enabling the construction of models that predict the biodegradability of certain chemicals based on their structural properties. This chapter provides an overview of the current level of AI in biodegradation prediction as well as prospective future directions. It begins by reviewing the many AI

methodologies widely used for biodegradation prediction, such as machine learning algorithms, feature selection methods, and QSAR modeling. In constructing robust and effective AI models for biodegradation prediction, the chapter emphasizes the relevance of data quality, amount, and variety. The performance assessment metrics that are thoroughly evaluated for measuring the precision and reliability of AI models include accuracy, precision, recall, and F1-score. The limitations and constraints of AI in biodegradation prediction are also discussed, including data availability, model robustness, and ethical implications. The implications of AI in biodegradation prediction for environmental science and sustainability are also examined, including the possible influence on chemical risk assessment, environmental policy-making, and sustainable chemical design. The chapter finishes with recommendations for future AI research and applications in biodegradation prediction, such as the need for multidisciplinary cooperation, more data sharing, and improved model interpretability and transparency. This chapter is an excellent resource for scholars, practitioners, and policymakers interested in using AI approaches for environmental sustainability and pollution control. Its overarching purpose is to offer a comprehensive assessment of the existing situation, challenges, and prospective influence of AI in forecasting biodegradation.

## **10.1 INTRODUCTION**

Biodegradation is a phenomenon by which complex organic substances are broken down into small and simple compounds by microbes, such as fungi, bacteria, and some other biological agents. These microorganisms utilize the organic materials as a food source, breaking them down into smaller molecules, ultimately resulting in the conversion of organic matter into simpler compounds, such as water, carbon dioxide, and biomass [1]. Biodegradation occurs in various environments, including soil, water, and air. The biodegradation is a crucial process in environmental science and sustainability as it contributes to waste management [2], nutrient cycling, pollution mitigation [3], renewable energy production [4], and sustainable agriculture [5]. It helps maintain ecological balance, reduces environmental pollution [6], and promotes more sustainable practices in various sectors [7–10].

Biodegradation plays an important role in the waste management as it helps break down organic waste, such as food scraps, yard waste, and other biodegradable materials, into simpler compounds. This reduces

the accumulation of organic waste in landfills, which can have negative environmental impacts, such as greenhouse gas emissions and pollution of groundwater [11]. Biodegradation is also an essential component of nutrient cycling in ecosystems. When organic materials are broken down through biodegradation, the resulting simpler compounds can be recycled by other organisms as a source of nutrients, contributing to the balance and sustainability of ecosystems. Biodegradation is a key process in the production of renewable energy sources, such as biogas and biofuels [12–14]. Organic materials, such as agricultural residues, food waste, and algae, can be biodegraded to produce methane or other biofuels, which can be used as a sustainable source of energy [15]. The most important benefit of biodegradation is in the field of agriculture, where it helps break down organic matter in soil, contributing to nutrient cycling, soil fertility, and overall soil health [16]. This can reduce the reliance on synthetic fertilizers and pesticides, promoting more sustainable agricultural practices.

The biodegradation prediction was typically done using conventional statistical and mathematical methods [17–19]. These methods employed several regression models which includes linear regression, multiple regression, and logistic regression, as well as clustering and classification methods such as K-means clustering and decision trees.

In general, the conventional statistical and mathematical methods used for biodegradation prediction rely heavily on the availability and quality of input data, as well as the assumptions made about the underlying relationships between the input variables and the biodegradation outcome. These methods require a significant amount of manual feature engineering and often involve a trial-and-error process to identify the most relevant features and develop accurate prediction models.

For example, in the case of predicting the biodegradability of chemicals, conventional methods would rely on input parameters such as molecular weight, hydrophobicity, and other physicochemical properties. These parameters would then be utilized to create a statistical or mathematical model that may predict the chemical's biodegradation rate or outcome. The model would need to be trained on a dataset of known biodegradation rates to optimize its parameters and achieve the best possible accuracy. While these conventional methods can still be effective in predicting biodegradation outcomes, they often lack the scalability and adaptability of AI and machine learning techniques. They also require a significant amount of domain expertise and manual feature engineering, which can be time-consuming and costly [20].

Predicting biodegradation rates accurately can be challenging due to various factors, including the complexity and variability of biodegradation processes, the diversity of microorganisms involved, and the wide range of environmental conditions that can affect biodegradation [21]. Overcoming these challenges requires further research, data collection, and modeling approaches to improve our understanding and prediction of biodegradation rates accurately.

Predicting biodegradation rates accurately can be challenging due to several factors. Traditional indications, such as changes in pesticide concentration or the identification of pesticide metabolites, are ineffective for many pesticides in anaerobic conditions [22]. Furthermore, these indicators are unable to differentiate between biotic and abiotic pesticide degradation processes. Another issue is that the rate of biodegradation of a pesticide is dependent on microbial adaptability and the enrichment of certain degraders. Furthermore, accurately sampling a representative microbial community from the field can be challenging. This means that detecting, measuring, and predicting pesticide biodegradation in the environment might be difficult.

AI plays a significant role in addressing challenges in predicting biodegradation [23]. Biodegradation refers to the natural breakdown of organic materials by living organisms into simpler compounds that can be recycled in the environment [24, 25]. Accurate prediction of biodegradation is critical for various fields, including environmental science, waste management, and biotechnology. However, several challenges exist in predicting biodegradation rates and pathways due to the complex and dynamic nature of biological systems. One of the key roles of AI in addressing these challenges is data analysis. Large information, including chemical structures, environmental factors, and biological interactions, can be processed by AI algorithms to find patterns and connections that would not be obvious to humans. This enables researchers to gain insights into the factors that influence biodegradation and develop predictive models [24, 26–29].

Machine learning, a subset of AI, is particularly useful in predicting biodegradation. To learn from previous observations and generate predictions on new data, machine learning algorithms can be trained on old datasets. For example, machine learning models can analyze the molecular structures of compounds and identify features that are associated with high or low biodegradation rates. These models can then be used to predict the biodegradation potential of new compounds, even those that have not been previously tested in experiments. Another role of AI in predicting biodegradation is in the design of new materials. AI can generate virtual libraries of molecules with desired properties, such as high biodegradability,

and screen them for potential candidates using computational simulations. This accelerates the process of identifying promising materials for experimental testing, reducing time and cost. Additionally, AI can aid in the optimization of environmental conditions for biodegradation. By analyzing various environmental parameters, such as temperature, pH, and nutrient availability, AI algorithms can optimize conditions to enhance biodegradation rates. This can be particularly useful in waste management, where AI can help in designing efficient bioremediation strategies for contaminated sites [30–32].

AI plays a crucial role in addressing challenges in predicting biodegradation by analyzing large datasets, developing predictive models, aiding in material design, and optimizing environmental conditions. It has the potential to considerably increase our understanding of biodegradation processes as well as contribute to the creation of long-term solutions to environmental and waste management issues [31].

## **10.2 OVERVIEW OF AI IN BIODEGRADATION PREDICTION**

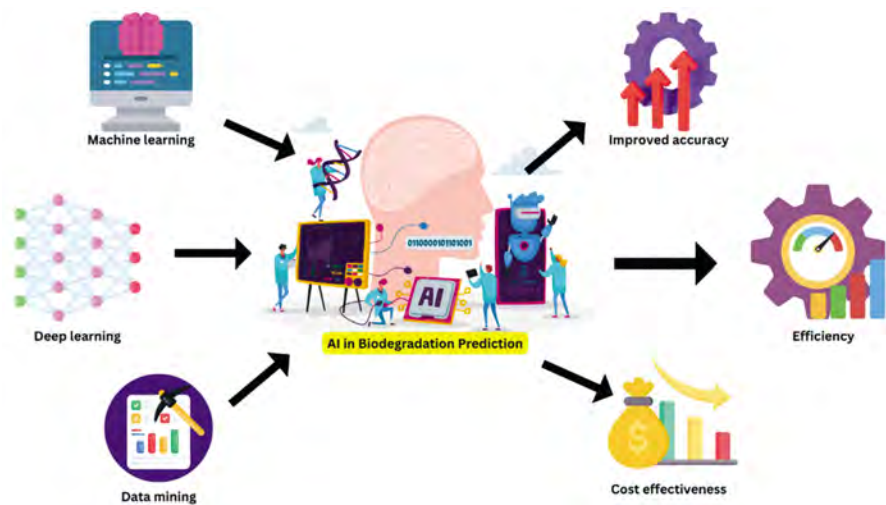
AI has the ability to significantly improve biodegradation prediction. AI helps in accurately predicting the biodegradability of different substances, AI can contribute to better understanding and managing environmental impacts [26–33]. AI can analyze enormous amounts of data from various sources, such as chemical structures, ambient variables, and microbiological activity, to discover trends and forecast biodegradation rates. Machine learning algorithms may be trained on massive databases of biodegradation data to create predictive models that reliably anticipate the biodegradability of diverse chemicals and materials [34, 35].

Furthermore, AI can help optimize biodegradation prediction models by continuously learning and improving from new data, feedback, and validation. This iterative process can lead to increasingly accurate and reliable predictions over time, enhancing our understanding of biodegradation dynamics and enabling better decision-making in environmental management and policy [36–38].

In addition, AI can expedite the screening of vast chemical libraries to identify potential biodegradable compounds [39, 40], which can accelerate the discovery of environmentally friendly materials and reduce the reliance on harmful substances. This can have a significant positive impact on environmental sustainability [41, 42], waste reduction [14, 38], and pollution prevention [3, 33].

However, it's important to consider ethical and regulatory aspects of using AI in biodegradation prediction [43, 44]. Ensuring data privacy, addressing bias in training data, and maintaining transparency and interpretability of AI models are critical concerns [44]. Additionally, human expertise and domain knowledge should complement AI predictions, and AI should not replace the need for empirical biodegradation testing.

Finally, the use of AI in biodegradation prediction has the potential to greatly advance our understanding of biodegradation processes and enable the discovery of environmentally friendly materials, and support sustainable environmental management practices. With proper attention to ethical considerations, AI can be a valuable tool in promoting environmental sustainability and addressing challenges related to biodegradation prediction (Figure 10.1).



**FIGURE 10.1** AI in biodegradation prediction.

### **10.2.1 DIFFERENT TYPES OF AI TECHNIQUES COMMONLY USED IN BIODEGRADATION PREDICTION**

There are several different AI techniques used in biodegradation prediction. It is important to note that the technique chosen is determined by the specific situation at hand, the available data, and the required level of accuracy and interpretability. Proper validation, optimization, and interpretation of

results are important in ensuring the reliability and applicability of AI-based biodegradation prediction models (Table 10.1).

**TABLE 10.1** Methods Using Machine Learning to Model Design Use to Predictive Biodegradation

Modeling Techniques	Descriptor	Model
Machine learning.	Hexavalent chromium.	Aerobic biodegradation for azo dyes [50].
Inductive machine learning.	Structural features and molecular weight	Inductive structure protein analysis (IPSA) [51].
Used Klopman method.	Molecular fragment-binding affinity of drugs.	Structure evaluation model [52].
Linear and nonlinear regression.	Fragments of molecule.	Molecular prediction model [53]
Multiple linear and nonlinear regression.	Molecular weight and calculated structural fragment.	Expert systems survey on xenobiotic chemical biological degradation [54].

- 1. Machine Learning [45–49]:** It is a sort of AI approach that includes training algorithms on big datasets to discover patterns and make predictions. Machine learning algorithms can be taught on varied datasets containing information such as chemical structures, environmental variables, and biodegradation rates in the context of biodegradation prediction. Decision trees, random forests, support vector machines (SVM), and k-nearest neighbors (KNN) are examples of machine learning algorithms commonly employed in biodegradation prediction. Based on input features, these algorithms can be used to create predictive models that estimate the biodegradability of certain substances. A few examples of methods using machine learning to model design use to predictive biodegradation are summarized. The modeling algorithms or methodologies have been used to describe the model system.
- 2. Deep Learning [55–62]:** It is a branch of machine learning that involves training multiple-layer artificial neural networks to automatically extract complicated patterns from input. Convolutional neural networks (CNN) and recurrent neural networks (RNN) are two deep learning approaches that are increasingly being employed in biodegradation prediction problems. For example, CNNs can be used to analyze chemical structures or molecular fingerprints to predict biodegradation rates [63], while RNNs can be used to model time-dependent biodegradation processes [64].



3. **Data Mining [65–69]:** The technique of obtaining meaningful information from enormous databases is known as data mining. Data mining approaches can be used to find significant traits or variables that influence biodegradation rates in the context of biodegradation prediction. To prepare the data for analysis, techniques such as feature selection, feature extraction, and data preprocessing may be used. Data mining techniques can also be used to find hidden patterns or relationships in data that are not immediately obvious, assisting in the development of more accurate biodegradation prediction models.
4. **Quantitative Structure-Activity Relationship (QSAR) Model [70–75]:** QSAR modeling is a specialized technique commonly used in biodegradation prediction. It involves developing mathematical models that correlate the structure of chemical compounds with their biological activity, in this case, biodegradation rates. Machine learning algorithms such as multiple linear regression, partial least squares (PLS), and support vector regression (SVR) can be used to construct QSAR models. QSAR models are particularly useful when dealing with large datasets of chemical compounds and can provide insights into the key structural features that influence biodegradation.
5. **Ensemble Methods [76–79]:** Ensemble methods involve combining multiple predictive models to improve prediction accuracy. Ensemble averaging, bagging, and boosting are techniques that can be used to collect the predictions of different machine learning or deep learning models, leading to a more reliable and precise biodegradation prediction model.

### **10.2.2 ADVANTAGES OF AI IN BIODEGRADATION PREDICTION**

The utility of AI in biodegradation prediction offers several advantages over traditional methods including improved accuracy, efficiency, cost-effectiveness, flexibility, and the potential for continuous improvement. These benefits can help to advance our understanding of biodegradation processes and inform effective strategies for environmental management [80–82] and cleanup [83, 84].

1. **Improved Accuracy:** AI algorithms can learn from massive volumes of data and find complicated links and patterns that people might not recognize. As a result, AI models have the ability to improve biological degradation prediction accuracy dramatically, particularly

when dealing with complicated and non-linear correlations between input variables and biodegradability.

2. **Efficiency:** AI models can easily and effectively manage huge quantities of data, enabling for a more rapid and complete investigation of environmental pollutants and their possible biodegradability. This can lead to the more rapid and accurate identification of possible biodegradation routes and mechanisms, which can then be utilized to inform environmental remediation methods.
3. **Cost-Effectiveness:** The use of AI in biodegradation prediction can be cost-effective compared to traditional methods, especially when large amounts of data are involved. AI models can reduce the need for expensive and time-consuming laboratory experiments, as they can accurately predict biodegradability based on a wide range of input features.
4. **Flexibility:** AI models can be trained using a variety of input data, such as chemical structures, environmental conditions, and experimental data. This allows for the creation of highly flexible and adaptable models that can be tailored to specific applications and datasets.
5. **Continuous Improvement:** As more information becomes available, AI models will continue learning and improving themselves. This can lead to more accurate and reliable predictions over time, allowing for better decision-making and more effective environmental management strategies.

### **10.2.3 APPLICATIONS OF AI IN BIODEGRADATION PREDICTION IN VARIOUS FIELDS [85–89]**

AI can be applied in various fields to improve the prediction of biodegradation rates. The applications of AI in biodegradation prediction are vast and varied, with potential benefits for environmental protection, public health, and sustainable development in a range of fields. AI-based biodegradation prediction has numerous applications across a wide range of fields, including following fields.

1. **Pharmaceuticals [90, 91]:** AI can be used in the pharmaceutical industry to predict the biodegradability and environmental fate of drug compounds. AI models can be used to predict the biodegradability of pharmaceuticals, helping to reduce their impact on the environment

and ensure the safety of water supplies. This can aid in the design of more environmentally-friendly drugs and facilitate the identification of potential environmental risks associated with drug manufacturing, use, and disposal. AI models can analyze chemical structures, physicochemical properties, and environmental data to predict the fate of drugs in different environmental compartments, such as soil, water, and air, and assess their potential impact on the environment.

2. **Agriculture [92–95]:** AI can be employed in agriculture to predict the biodegradation of agrochemicals, such as pesticides and fertilizers, helping to reduce their impact on the environment and ensure the safety of food supplies. AI models can analyze various factors, including soil characteristics, climate conditions, and chemical properties of agrochemicals, to predict their biodegradation rates and degradation pathways. This can assist in optimizing the use of agrochemicals, reducing their potential negative impacts on the environment, and promoting sustainable agricultural practices.
3. **Waste Management [4, 96–101]:** Artificial intelligence models can be utilized for predicting the biological degradation of different kinds of waste, contributing to the development of efficient waste management ideas and reducing the negative environmental effects of trash disposal. AI can be used in waste management to predict the biodegradation of different types of waste, such as organic waste, plastics, and pollutants. AI models can analyze waste characteristics, environmental conditions, and microbial activity to predict the degradation rates and pathways of waste materials. This can aid in the development of effective waste management strategies, such as composting, bioremediation, and bioplastics degradation, to mitigate the environmental impact of waste and promote circular economy practices.
4. **Environmental Risk Assessment [102–105]:** AI can be utilized in environmental risk assessment to predict the biodegradation potential of chemicals and materials in different environmental compartments. AI models can integrate diverse data sources, such as chemical properties, environmental conditions, and microbial activity, to predict the fate and persistence of chemicals in the environment. This can assist in evaluating the environmental risks associated with chemical releases, guiding regulatory decisions, and promoting environmentally-sound chemical management practices. AI models can be used to predict the biodegradability of environmental contaminants, helping to inform management strategies for contaminated sites and improving environmental monitoring efforts.

5. **Product Development [106]:** AI can be applied in product development to predict the biodegradability of materials used in various consumer products, such as packaging, textiles, and electronics. AI models can analyze material properties, usage scenarios, and environmental conditions to predict the biodegradation potential and environmental impact of products throughout their lifecycle. This can aid in designing more sustainable products with improved biodegradability and reduced environmental footprint. AI models can anticipate the biodegradability of novel compounds, assisting in the design and development of new entities by lowering the risk of environmental impact.
6. **Biotechnology [107]:** AI models can be used to predict the biodegradability of biomolecules and biopolymers, helping to inform the development of new biotechnologies and reduce their environmental impact.
7. **Energy Production [108–111]:** AI models can be used to predict the biodegradability of biofuels and other forms of renewable energy, helping to ensure their sustainability and reduce their impact on the environment.

Overall, the applications of AI in biodegradation prediction are vast and varied, with potential benefits for environmental protection, public health, and sustainable development in a range of fields.

### 10.3 DATA ACQUISITION AND PREPROCESSING

Data acquisition and preprocessing are essential steps in developing AI models for biodegradation prediction. These procedures entail acquiring and preparing data for use in training, validating, and testing AI models. Data acquisition include gathering information from a variety of sources, including experimental biological degradation data, environmental monitoring data, and chemical databases. This information can be applied to train AI algorithms and enhance their prediction of biodegradation rates.

Preprocessing is the process of cleaning and altering data so that it can be used in AI models. It may be necessary to remove outliers, fill in missing numbers, and normalize the data. Preprocessing can improve data quality and make it more acceptable for usage in AI models.

Once the data has been acquired and preprocessed, it can be used to train AI models such as machine learning algorithms. These models can learn

from the data and make predictions about biodegradation rates. The accuracy of these predictions can be improved by using high-quality data that has been properly acquired and preprocessed.

Data gathering from many sources, such as experimental results, environmental databases, scientific publications, and open-access datasets, can be automated through AI. AI algorithms can crawl and extract relevant data from different formats, such as text, images, and tables, and compile them into a unified dataset. This can save time and effort compared to manual data collection, especially when dealing with large and diverse datasets.

AI can be used to combine data from various available sources to create a comprehensive dataset for biodegradation prediction. AI algorithms can standardize and preprocess data to ensure consistency and reliability. Data from multiple sources, such as chemical properties, environmental conditions, and biodegradation rates, can be combined and integrated using AI techniques, such as data fusion, feature extraction, and data normalization. This can facilitate the integration of diverse data types and enable the AI model to capture complex relationships between variables.

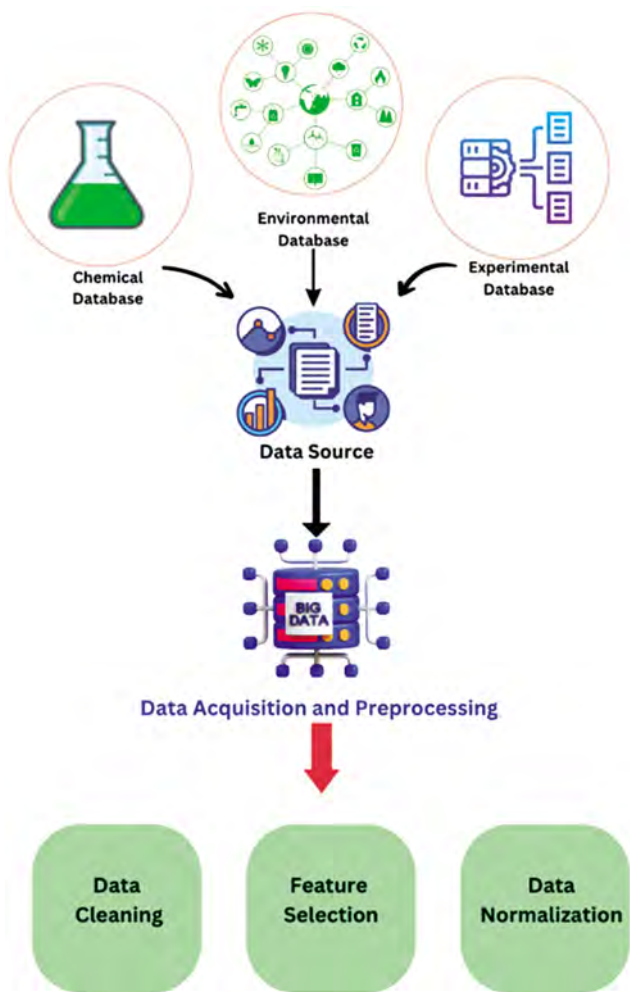
AI can help in quality control and data cleaning in order to ensure the reliability and accuracy and reliability of datasets. AI algorithms can automatically identify and remove noisy or inconsistent data, such as outliers, missing values, and duplicates [112–115]. Data quality control techniques, which include data imputation, data validation, and error correction, can be used with AI to improve the reliability of the database, which is critical for the accuracy of the AI models.

AI can help in selecting relevant features or variables from the dataset and reducing the dimensionality of the data. AI algorithms can analyze the importance and relevance of different features for biodegradation prediction using techniques like feature selection, feature ranking, and dimensionality reduction. This can aid in the identification of the most informative traits while also reducing the computational complexity of AI models, resulting in more effective and precise predictions. AI can be used to enrich data in order to expand its diversity and scale. AI algorithms can be used to generate additional data points from an existing dataset employing data augmentation techniques like data synthesis, data generation, and data transformation. This may enhance the dataset, increase data variability, and improve the generalization ability of AI models.

For model construction and evaluation, AI can be used to divide the dataset into testing, validation, and training sets. AI algorithms can randomly or strategically split the dataset while preserving the distribution and characteristics of the data. Cross-validation techniques, such as k-fold

cross-validation, can also be applied using AI to assess the performance and robustness of the AI models. AI can play a critical role in data acquisition and preprocessing for biodegradation prediction. It can automate data collection, integrate diverse data sources, clean and validate data, select relevant features, augment the dataset, and split the data for model development and evaluation. These AI-enabled data preprocessing steps are crucial for developing accurate and reliable biodegradation prediction models.

Overall, data acquisition and preprocessing are important steps in using AI for the prediction of biodegradation rates (Figure 10.2).



**FIGURE 10.2** Data acquisition and preprocessing.

### **10.3.1 HIGH-QUALITY DATA IN TRAINING AI MODELS FOR BIODEGRADATION PREDICTION**

Training AI models for biological degradation prediction requires high-quality data. It assures that the creation and implementation of AI models for biodegradation projection are accurate, generalizable, interpretable, reliable, and ethical [116, 117]. Proper data collection, validation, and quality control measures should be in place to ensure the integrity and reliability of the data used for training AI models. High-quality data is important in training AI models for biodegradation prediction. The accuracy and dependability of AI model predictions are strongly dependent on the quality of the data used for training. High-quality data ensures that the AI models are trained on accurate and reliable information, which leads to more accurate predictions. It enables the models to learn meaningful patterns and relationships that can generalize well to unseen data. High-quality data that includes relevant information about the chemical structures, environmental conditions, and biodegradation rates allows for better interpretation of the model's predictions. It aids in identifying the essential aspects or variables that drive biodegradation, providing insights into the processes of biodegradation and assisting in decision-making. High-quality data ensures that the AI models are reliable and trustworthy.

High-quality data ensures that the AI models are trained on accurate and reliable information, which leads to more accurate predictions. If the training data contains errors, inconsistencies, or inaccuracies, it can introduce biases and noise into the model, leading to inaccurate predictions. Therefore, using high-quality data that is representative of the real-world biodegradation processes being predicted is critical for achieving accurate and reliable predictions.

AI models are trained on data with the objective of making predictions on unseen data. High-quality data allows models to discover significant patterns and correlations that generalize well to previously unseen data. If the training data is poor quality, the model may not generalize well to new data, resulting in poor performance in applications in the real world. High-quality data allows the model to learn from diverse and representative samples, which helps in building a robust and generalizable biodegradation prediction model.

Understanding the underlying mechanisms and factors influencing biodegradation is important for interpreting the predictions generated by AI models. High-quality data that includes relevant information about the chemical structures, environmental conditions, and biodegradation

rates allows for better interpretation of the model's predictions. It aids in identifying the essential aspects or variables that drive biodegradation, providing insights into the processes of biodegradation and assisting in decision-making.

High-quality data ensures that the AI models are reliable and trustworthy. Biodegradation prediction models are often used in critical applications, such as environmental risk assessment, waste management, and regulatory compliance. Decision-makers rely on the predictions generated by these models to make informed decisions. If the training data is of poor quality, it can undermine the reliability of the predictions, which can have serious consequences for environmental management and policy-making.

To reduce the possibility of prejudice, discrimination, or unfairness being perpetuated in the AI models' predictions, ensure that the data employed for training is representative, credible, and unbiased. It promotes ethical AI practices, data privacy, and transparency, which are important for responsible and ethical deployment of AI in biodegradation prediction.

### **10.3.2 SOURCES OF DATA USED IN BIODEGRADATION PREDICTION**

There are several sources of data that can be used in biodegradation prediction. One source is chemical databases, which contain information about the chemical structures of compounds. Environmental datasets are another source, as they include information about environmental variables such as humidity, pH, temperature, and the occurrence of other chemicals that can alter biodegradation rates. Experimental data from biodegradation tests can also be used to train AI models for biodegradation prediction. For example, biodegradability testing on 1,055 chemicals were collected from the online portal of the National Institute of Technology and Evaluation (NITE), Japan, and used to build QSAR models to investigate the relationships between chemical structure and molecule biodegradation [118]. These various data sources can be merged to provide a full picture of the elements influencing biodegradation rates and aid in the training of effective AI models for biodegradation prediction.

Here are some examples to understand the role of AI in prediction of biological degradation and using database library (Table 10.2).

Chemical databases provide information about the chemical structures, properties, and characteristics of various compounds, including their potential for biodegradation. These databases, such as PubChem, ChemSpider, and ChEMBL, contain vast amounts of data on chemical compounds, including



their molecular formulas, structures, and physicochemical properties [120]. This data can be used to identify and select compounds for biodegradation prediction modeling, as well as to provide input features for machine learning algorithms.

**TABLE 10.2** Criteria for Evaluating Biological Oxygen Demand (BOD) Measurements in the DIPPR Repository [119]

Parameter	Experimental Techniques	Temperature	Internal Consistency	Concentration
Monitoring criteria	As per SOP	at 20°C	ThOD $\geq$ BOD	2–6 mg/L

Environmental datasets contain information about the environmental conditions and parameters that can affect biodegradation processes. Temperature, pH, soil or water composition, microbial communities, and other environmental parameters which affect the rate and degree of biodegradation are all included. Environmental datasets, which are frequently available from monitoring programs, research projects, or public repositories, can be exploited to study the relationship between environmental variables and biodegradation rates, as well as to train biodegradation prediction models.

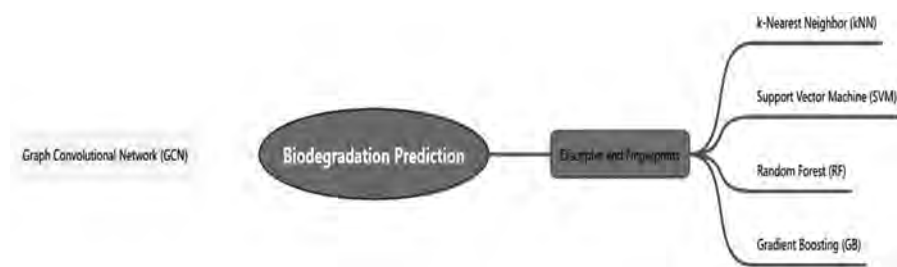
Experimental data obtained from laboratory studies or field trials can provide valuable information on the biodegradation behavior of specific compounds or classes of compounds [121–123]. These data can include biodegradation rates, degradation pathways, and metabolite formation, which can be used to validate and calibrate biodegradation prediction models. Experiment data can also be applied to determine the important characteristics or variables that affect biodegradation and use them as input features in prediction models.

Scientific literature and textual data, such as research articles, reviews, patents, and other relevant documents, can serve as a valuable source of information for biodegradation prediction. These sources can provide insights into the biodegradation behavior of different compounds, factors affecting biodegradation, and case studies of biodegradation processes in various environments. Textual data can be mined and processed to extract relevant information, such as chemical structures, degradation pathways, and environmental conditions, which can be used as input data for biodegradation prediction models.

Computational databases and models, such as QSAR models, can provide predictive information on the biodegradation potential of chemical compounds. QSAR models use mathematical relationships between chemical

properties and biodegradation rates to predict the biodegradation behavior of compounds. These models can be used as a source of data to generate input features for biodegradation prediction models or as a benchmark for validating the accuracy and reliability of the predictions generated by the AI models [124]. The GCN model can be applied directly by graphing the simplified molecular input line entry system (SMILES), while the QSAR model is implemented by measuring and selecting molecular descriptors as well as fingerprints from SMILES. According to this diagram, the GCN implementation is less difficult and needs less data than the QSAR model [125].

Other sources of data, such as expert knowledge, historical records, and data from relevant industries or organizations, can also be used in biodegradation prediction. Expert knowledge can provide valuable insights into the biodegradation behavior of specific compounds or environments, while historical records and industry data can offer real-world information on biodegradation rates, degradation pathways, and other relevant factors (Figure 10.3).



**FIGURE 10.3** Biodegradation prediction using GNN and other descriptors.

### **10.3.3 CHALLENGES IN DATA PREPROCESSING FOR AI-BASED BIODEGRADATION PREDICTION**

Data processing is an important stage in the development of based on artificial intelligence biodegradation prediction models because it makes sure the data used for validation and training processes are of high quality, relevant, and properly organized.

One of the main challenges in data preprocessing is data cleaning, which involves identifying and addressing issues such as missing data, outliers, inconsistencies, and errors in the dataset [126, 127]. Biodegradation data, which can come from a variety of sources with varying quality and trustworthiness, may contain noise or inconsistencies that might reduce the accuracy

and reliability of the prediction models. Proper data cleanings approaches, such as outlier removal, imputation for missing values, and error correction, are required to ensure that the data used to train the models is reliable and accurate.

Another important consideration in data preprocessing is feature selection, which involves choosing the most relevant and informative features (i.e., variables or attributes) from the dataset to be used as input features for the prediction models. Biodegradation prediction models may involve a large number of chemical properties or environmental variables, and not all of them may be equally relevant for predicting biodegradation. Statistical methodologies, dimensionality reduction algorithms, and domain experience can all aid in finding the most significant features that contribute to prediction accuracy and model interpretability.

Data normalization is the process of converting the data to a range to confirm that different features have similar values [128, 129]. Variables in biodegradation datasets may have varied units, sizes, or distributions, which may influence the performance of prediction models. To ensure that the data is standardized, and the models are not biased towards features with larger values, data normalization techniques such as z-score normalization, min-max scaling, and log transformation can be used [130, 131].

Biodegradation datasets can also be impacted by unbalanced data, which occurs when the number of examples in one class (e.g., biodegradable) is considerably distinct from the number of examples in another class (e.g., non-biodegradable). This can have an impact on prediction model performance because they may be biased towards the dominant class. Techniques such as over sampling, under sampling, and synthetic data generation can be used to handle imbalanced data and ensure that the prediction models are trained on balanced data to avoid biased results.

Careful consideration should be given to data integration to ensure that the integrated dataset is consistent, reliable, and appropriate for training the prediction models. Biodegradation datasets may contain sensitive information, such as chemical structures, environmental parameters, and experimental results. Ensuring data privacy and adhering to ethical considerations, such as obtaining proper consent and protecting sensitive information, is crucial in data preprocessing. Data anonymization, encryption, and other data protection techniques should be applied to ensure compliance with data privacy regulations and ethical guidelines.

Data preprocessing is a critical step in developing AI-based biodegradation prediction models, and it involves addressing challenges related to data

cleaning, feature selection, data normalization, handling imbalanced data, data integration, and data privacy and ethics. Proper data preprocessing techniques are required to make sure that the information used for training and validation is of high quality, relevant, and properly formatted, contributing to the accuracy and reliability of the biological degradation prediction models.

## **10.4 AI TECHNIQUES FOR BIODEGRADATION PREDICTION**

Machine learning, deep learning, and data mining are prominent AI approaches used in biodegradation prediction to construct models that can effectively forecast the degradation of chemicals, pollutants, or other substances in the environment. These techniques leverage the power of algorithms and computational methods to analyze and learn from large datasets, making predictions based on patterns and relationships found in the data. Here are some commonly used AI techniques in biodegradation prediction:

Machine learning algorithms are widely used in biodegradation prediction [83, 132, 135, 136]. On labelled datasets, supervised machine learning algorithms like decision trees, support vector machines (SVM), random forests, and logistic regression can be trained to understand the associations between input features (e.g., environmental parameters, chemical properties) and biodegradability outcomes (e.g., non-biodegradable or biodegradable). These models, once trained, can generate predictions on previously unknown data.

Deep learning, a type of machine learning, requires the use of multiple-layer artificial neural networks (ANN) to learn data representations. Deep learning techniques like convolutional neural networks (CNN) and recurrent neural networks (RNN) are frequently employed in biodegradation prediction tasks involving big and complicated datasets like molecular structures or time-series environmental data. Deep learning models can learn features from raw data and capture complicated patterns automatically, making them useful for biodegradation prediction [59].

Clustering, association rule mining, and pattern recognition are all data mining approaches that can be used to forecast biodegradation [137]. These techniques can analyze large datasets to identify patterns, correlations, and associations between different variables or attributes, which can provide insights into factors that influence biodegradability. Data mining techniques can be used for exploratory analysis, feature selection, and identifying relevant features for building prediction models.

Ensemble approaches like ensemble learning and stacking can also be used to forecast biodegradation. Ensemble approaches integrate the predictions of numerous models to increase prediction accuracy and robustness. For example, using the same dataset, an ensemble of various machine learning models can be trained, and their predictions can be integrated to generate a final prediction. Ensemble approaches can help to overcome the constraints of individual models and increase the prediction models' overall performance.

In biodegradation prediction, hybrid systems that incorporate different AI methods, such as machine learning and deep learning, can also be applied [138, 139]. For example, a machine learning model can be used to extract related structures from the data, which are then fed into a deep learning model for further analysis and prediction. Hybrid approaches can leverage the strengths of different AI techniques and provide more accurate and robust predictions.

AI techniques, such as machine learning, deep learning, data mining, ensemble techniques, and hybrid approaches, are commonly used in biodegradation prediction to develop accurate and robust models for predicting the biodegradability of chemicals, pollutants, or other substances in the environment. These techniques can analyze large and complex datasets, learn patterns and relationships from data, and make predictions based on the learned knowledge, which can have important applications in environmental monitoring, risk assessment, and decision-making.

#### **10.4.1 MACHINE LEARNING ALGORITHMS FOR BIODEGRADATION PREDICTION [140, 142]**

Several machine learning algorithms are commonly used for biodegradation prediction, depending on the specific task and dataset. Here are some examples of commonly used machine learning algorithms for biodegradation prediction:

For categorization tasks, decision trees are a common machine learning approach, including biodegradation prediction. Decision trees can be used to model decisions or decisions based on certain conditions. In the framework of biodegradation prediction, decision trees can be accomplished on labeled datasets to learn decision rules based on chemical properties, environmental parameters, or other relevant features to predict whether a chemical or substance is biodegradable or non-biodegradable.

Support Vector Machines (SVM) is a sophisticated machine learning method that may be operated for classification and regression tasks. Based on the input attributes, SVM finds a hyperplane that best splits data points into distinct classes. SVM can be used to develop a model that can accurately categorize chemicals or compounds as biodegradable or non-biodegradable based on their feature representation in biodegradation prediction. SVM is notable for its capacity to handle complex decision boundaries and can handle high-dimensional data.

Random forests are a collective learning method that makes predictions by combining numerous decision trees. Random forests are well-known for their capacity to handle data that is noisy while minimizing overfitting. Random forests can be used in biodegradation prediction to create a group of decision trees that collectively generate predictions based on numerous characteristics, leading to a more accurate and reliable prediction model.

Logistic regression is the common algorithm used for binary classification tasks, such as biodegradation prediction. Logistic regression models the relationship between input features and binary outcomes (e.g., biodegradable or non-biodegradable) using a logistic function. It is a basic and understandable method that can give insights into the value of various features in the prediction task.

K-nearest neighbors (k-NN) is a straightforward classification technique. It operates by locating a data point's  $k$  nearest neighbors in the feature space and generating predictions based on most of the class of its neighbors. k-NN can be used to forecast biodegradation by locating a chemical or substance's  $k$  nearest neighbors based on its characteristic representation and estimating its ability to break down relying on much of the class of its neighbors.

Gradient Boosting Machines (GBM) are another ensemble learning technology that makes predictions by combining numerous weak learners, often decision trees. GBM sequentially develops an ensemble of models, with each succeeding model correcting the flaws of the prior model. GBM has been shown to be effective in handling complex data and achieving high prediction accuracy, making it a popular choice for biodegradation prediction tasks.

These are some examples of commonly used machine learning algorithms for biodegradation prediction. The algorithm chosen is determined by criteria like the size and complexity of the data, the nature of the challenge at hand, and the expected outcome metrics. To choose the best algorithm for a specific biodegradation prediction assignment, it is critical to carefully analyze and compare the efficacy of several algorithms using relevant evaluation measures.

### **10.4.2 DEEP LEARNING TECHNIQUES FOR BIODEGRADATION PREDICTION [83, 143, 145]**

Deep learning approaches, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have demonstrated significant promise in a variety of domains, including biodegradation prediction. Here's a brief overview of their application in biodegradation prediction:

1. **Convolutional Neural Networks (CNNs):** These are a sort of deep learning model that excels at image recognition. CNNs use convolutional layers to automatically learn features from input data, such as chemical structures or environmental parameters, by applying convolution operations that capture local patterns. These learned features are then used for prediction. In the context of biodegradation prediction, CNNs can be used to learn representations of chemical structures or other relevant data, and then make predictions based on these representations. For example, CNNs can be trained on 2D or 3D representations of chemical structures, such as molecular fingerprints or molecular images, to predict their biodegradability.
2. **Recurrent Neural Networks (RNNs):** These are a type of deep learning model that excels at dealing with time series and sequence data. Recurrent connections are used by RNNs to capture dependence on time and sequential patterns of available databases. In biodegradation prediction, RNNs can be used to model the temporal or sequential nature of environmental data, such as time-series measurements of environmental parameters, and predict biodegradability based on these sequential patterns. RNNs, which include types like Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRUs), have been effectively used to imitate the temporal dynamics of biodegradation processes in biodegradation prediction tasks.
3. **Hybrid Models:** Deep learning techniques can also be combined with other machine learning algorithms or used in hybrid models to leverage their strengths. For example, CNNs can be combined with additional algorithms for machine learning, such as SVM or as decision trees, to increase prediction performance by integrating CNN learned representations as input characteristics to these methods. Similarly, RNNs can be used in combination with other algorithms, such as k-NN or logistic regression, to incorporate temporal information from RNNs into these models. These hybrid models can potentially capture

both local and global patterns in data, leading to improved prediction accuracy.

Deep learning techniques, which include CNNs and RNNs, can be used to forecast biodegradation and give advantages such as automated feature learning, the ability to capture complicated patterns in data, and the possibility for enhanced prediction accuracy. However, it is important to consider factors such as dataset size, model complexity, interpretability, and performance evaluation when using deep learning techniques for biodegradation prediction, as these models may require large amounts of data, computational resources, and careful model evaluation to ensure reliable and accurate predictions.

### **10.4.3 DATA MINING TECHNIQUES IN BIODEGRADATION PREDICTION [146–148]**

Data mining techniques including association rule mining and clustering can be used to uncover patterns and relationships in data for biodegradation prediction. These techniques can help discover relevant aspects and understand the elements that drive biodegradation by providing information about the fundamental frameworks and patterns associated with the data.

A data mining technique that analyses correlations and associations among variables in huge datasets is association rule mining. In the context of biodegradation prediction, association rule mining can be used to identify co-occurrence patterns or correlations between different chemical properties, environmental factors, and biodegradation outcomes. For example, it can identify associations between specific chemical compounds and their biodegradation rates, or between environmental conditions (such as pH, temperature, and humidity) and the biodegradation outcomes. These associations can provide insights into the factors that affect biodegradation and can be used to develop hypotheses or predictions about the biodegradation behavior of different compounds or under different environmental conditions.

Clustering is a data mining approach that groups comparable data points together in a multidimensional space based on similarities or differences. Clustering can be used in biodegradation prediction to identify patterns or clusters of similar biodegradation behaviors among different compounds or environmental conditions. For example, it can identify groups of compounds with similar biodegradation rates, or clusters of environmental conditions that favor or hinder biodegradation. Clustering can help in identifying



distinct patterns, and can be used to uncover hidden relationships or patterns that can be further explored and analyzed.

Data mining techniques can also be used for feature selection and dimensionality reduction to identify the most relevant variables or features for biodegradation prediction. For example, association rule mining can identify frequent item sets or combinations of variables that are associated with biodegradation outcomes, and these frequent item sets can be used as features in the prediction models. Clustering can also help in identifying groups of similar variables or features that can be represented by a reduced set of representative features, thus reducing the dimensionality of the data. This can help in developing more efficient and interpretable prediction models by focusing on the most important features that influence biodegradation outcomes.

Association rule mining and clustering are two data mining approaches that can be used to uncover relationships and trends in biodegradation data. They can help identify associations between variables, uncover hidden patterns or clusters, and aid in feature selection and dimensionality reduction. These strategies can provide useful knowledge and insights for biodegradation prediction, helping to construct better and more accurate prediction models.

## **10.5 PERFORMANCE EVALUATION AND MODEL INTERPRETATION**

### ***10.5.1 IMPORTANCE OF PERFORMANCE EVALUATION IN AI-BASED BIODEGRADATION PREDICTION***

Performance evaluation is an important aspect of developing and deploying AI-based biodegradation prediction models. This can help identify any issues with the model, such as overfitting or bias, and guide efforts to improve its performance. For example, if a model performs well on certain types of compounds but not on others, this information can be used to guide its application in biodegradation prediction tasks. Performance evaluation is an essential step in ensuring that AI-based biodegradation prediction models are accurate and reliable [144].

- 1. Model Accuracy:** Performance evaluation helps determine the accuracy of AI-based biodegradation prediction models. Accurate predictions are essential for reliable decision-making in fields such as environmental science, biotechnology, and waste management. Performance evaluation techniques, such as cross-validation and comparison with experimental data, can help assess the predictive accuracy of AI models and ensure that they provide reliable results.

2. **Model Robustness:** Robustness is an important characteristic of any predictive model. It refers to the capacity of a model to sustain reliability and efficacy in the face of uncertainty, noise, or as changes in the provided data. Performance evaluation techniques can help assess the robustness of AI-based biodegradation prediction models by testing their performance under different conditions, such as varying input data quality, sample size, or environmental conditions.
3. **Model Generalizability:** Generalizability is the ability of a predictive model to perform well on unseen data beyond the training data it was trained on. AI-based biodegradation prediction models need to be able to generalize their predictions to different environmental conditions, biodegradable compounds, and target organisms. Cross-validation and other performance evaluation approaches can assist examine the ability to generalized artificial intelligence models to ensure that they can make correct predictions on a variety of datasets.
4. **Model Comparison:** Performance evaluation allows for the comparison of different AI-based biodegradation prediction models. There are numerous AI algorithms available for building predictive models, such as machine learning algorithms like neural networks, support vector machines, and decision trees. Performance evaluation methodologies can aid in comparing the performance of many models and determining which one is the most reliable and accurate for a certain biodegradation prediction task.
5. **Model Optimization:** Performance evaluation helps identify areas for model optimization and improvement. By analyzing the performance metrics of AI-based biodegradation prediction models, i.e., accuracy, precision, recall, and F1-score, it is possible to identify the strengths and weaknesses of proposed models. This information can be used to optimize the models by refining their algorithms, feature selection, or hyperparameter tuning, to improve their performance.
6. **Model Validation:** Performance evaluation is essential for model validation, which is a critical step in the development of reliable AI-based biodegradation prediction models. Validation ensures that the predictive models are scientifically valid, reliable, and trustworthy.

In summary, performance evaluation is crucial in AI-based biodegradation prediction as it enables the assessment of model accuracy, robustness, generalizability, and validation. It also facilitates model comparison and optimization, leading to the development of reliable and accurate predictive models that can be used in practical applications for environmental management, biotechnology, and waste management.

### 10.5.2 METRICS FOR EVALUATING THE PERFORMANCE OF AI MODELS [133, 144]

There are numerous metrics that are routinely used to evaluate the performance of artificial intelligence models to give quantitative estimates of the way a model is functioning. Some examples of regularly used metrics are:

1. **Accuracy:** It is one of the fundamental and extensively used statistic for evaluating performance. It calculates the percentage of accurately predicted instances in relation to the total number of examples in the dataset.

It is calculated as the ratio of sum of true positives (TP) and true negatives (TN) with the total number of instances which include true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). It is represented in equation as:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

When the classes in the dataset are balanced, i.e., the number of occurrences in each class is nearly the same, accuracy is useful. When classes are uneven, however, accuracy might be deceiving, as high accuracy can nevertheless result in misclassification of minority classes.

2. **Precision:** It is defined as the proportion of true positives to total expected positives (true positives + false positives). It assesses the model's ability to predict positive events correctly while excluding false positives. Precision is calculated as follows:

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

Precision is especially important when the goal is to reduce false positives, as in medical evaluation or identification of fraud, when erroneous positives might have catastrophic repercussions.

4. **Recall:** The proportion of true positives to total actual positives is known as recall, also known as sensitivity or true positive rate. It assesses its ability to accurately capture all positive cases while leaving none out. The recall is calculated as follows:

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

When the goal is to minimize false negatives, such as in illness identification or anomaly detection, overlooking positive cases might have serious effects.

4. **F1-Score:** It is an accurate evaluation of how well a model performs because it is the harmonic mean of precision and recall. It takes precision and recall into consideration and is useful when there is a trade-off between precision and recall. The F1-score is calculated as follows:

$$F1\ Score = \frac{2 \times (Precision \times Recall)}{(Precision + Recall)} \quad (4)$$

The F1-score combines recall and precision into a single metric, providing for a comprehensive assessment of a model's accuracy.

These are some of the most often used metrics for assessing the performance of AI models. Other metrics, like specificity, area under the receiver operating characteristic (ROC) curve, and Matthew's correlation coefficient (MCC), may be used depending on the specific application and requirements. It is critical to carefully select the proper metrics based on the model's specific aims and the features of the dataset being evaluated.

### 10.5.3 CHALLENGES IN INTERPRETING AI MODELS FOR BIODEGRADATION PREDICTION

Interpreting AI models for biodegradation prediction can pose several challenges and considerations, including model explain ability and transparency. Some of these challenges and considerations are:

1. **Lack of Model Interpretability:** Many AI models, such as deep learning models, are regarded as "black boxes" because they are complicated and their decision-making processes are difficult for humans to interpret. This lack of interpretability might make it difficult to interpret AI models for biodegradation prediction since it may be unclear how the model is making predictions and what traits or factors are driving the predictions.
2. **Ethical Considerations:** Biodegradation prediction models may involve the use of large datasets, including potentially sensitive environmental or chemical data. Ethical considerations, such as data privacy, fairness, and bias, need to be taken into account when

interpreting AI models for biodegradation prediction. Bias in the training data or model predictions could lead to unfair or discriminatory outcomes, and it is important to carefully evaluate and mitigate such biases to ensure that the model is making predictions in an unbiased and fair manner.

3. **Generalization and Robustness:** Biodegradation prediction models need to be able to generalize well to different environmental conditions and chemical compounds. Ensuring that the model's predictions are robust and reliable across different datasets, environmental conditions, and chemical compounds is a challenge. It is critical to test the model's performance on various datasets and its capacity to generalize to varied circumstances.
4. **Model Validation and Uncertainty:** Validating the performance of AI models for biodegradation prediction can be challenging due to the lack of comprehensive experimental data for biodegradation rates. Additionally, uncertainty in the predictions of AI models, such as prediction intervals or confidence intervals, is often not straightforward to estimate. It is important to carefully evaluate the uncertainty associated with the model's predictions and communicate the level of uncertainty to stakeholders.
5. **Transparency and Explainability:** Model explain ability and transparency are important considerations in interpreting AI models for biodegradation prediction. Being able to explain how the model is making predictions and provide transparent insights into the model's decision-making process. Techniques such as feature importance analysis, SHAP (SHapley Additive explanations), and LIME (Local Interpretable Model-agnostic Explanations) can be used to provide insights into the features or factors that are driving the model's predictions and increase the transparency of the model.
6. **Regulatory Compliance:** Biodegradation prediction models may be used in regulatory contexts, where compliance with regulatory guidelines and standards is crucial. Ensuring that the model's predictions comply with relevant regulations and guidelines, and are transparent and explainable, is important for gaining regulatory approval and acceptance.

In conclusion, interpreting AI models for biodegradation prediction requires careful consideration of model explain ability, transparency, ethical considerations, generalization, model validation, uncertainty, and regulatory compliance. Addressing these challenges and considerations is essential for

building reliable, transparent, and trustworthy AI models for biodegradation prediction that can be used for environmental risk assessment and decision-making. Overall, addressing challenges and considerations in interpreting AI models for biodegradation prediction is important for ensuring that the models are accurate, reliable, and trustworthy.

## 10.6 LIMITATIONS AND FUTURE DIRECTIONS

### 10.6.1 LIMITATIONS OF AI IN BIODEGRADATION PREDICTION

Artificial intelligence has shown great promise in biodegradation prediction, there are several limitations that need to be considered. Some of the limitations of AI in biodegradation prediction include:

1. **Data Availability and Quality:** AI models rely heavily on data for training and validation. However, data availability for biodegradation prediction can be limited, especially for specific environmental conditions, chemical compounds, or biodegradation pathways. In some cases, the data may be sparse or incomplete, leading to challenges in building accurate and robust models. Additionally, the quality of the data, including issues such as data accuracy, reliability, and representativeness, can also impact the performance of AI models.
2. **Model Robustness and Generalization:** AI models for biodegradation prediction may struggle with robustness and generalization. Robustness refers to the ability of the model to maintain accurate predictions even in the presence of noise or uncertainties in the data, while generalization refers to the ability of the model to accurately predict biodegradation rates in diverse environmental conditions or for different chemical compounds. Ensuring that AI models are robust and generalize well to different scenarios can be challenging, as biodegradation rates can be influenced by a wide range of factors, such as temperature, pH, microbial activity, and chemical interactions.
3. **Ethical Considerations:** Biodegradation prediction models may raise ethical considerations, including issues related to data privacy, fairness, and bias. Biodegradation data may come from various sources, including proprietary or confidential datasets, and ensuring proper data privacy and compliance with relevant regulations is important. Additionally, biases in the data or models, such as gender, racial, or geographic biases, could lead to unfair or discriminatory

outcomes, and addressing these ethical concerns is crucial in the development and use of AI models for biodegradation prediction.

4. **Interpretability and Transparency:** As discussed earlier, many AI models, such as deep learning models, lack interpretability and transparency, which can limit their explain ability and understandability. This can be a challenge in biodegradation prediction, as stakeholders may need to understand how the model is making predictions, which features or factors are driving the predictions, and how reliable the predictions are. Ensuring that AI models are interpretable, transparent, and provide insights into their decision-making process is important for gaining trust and acceptance.
5. **Regulatory Compliance:** Biodegradation prediction models may be used in regulatory contexts, where compliance with regulatory guidelines and standards is crucial. However, there may be challenges in aligning AI models with regulatory requirements, such as validation, verification, and acceptance criteria. Regulatory agencies may have specific guidelines or requirements for the use of AI models in biodegradation prediction, and ensuring compliance with these regulations can be a limitation.
6. **Computational Resources and Scalability:** AI models, particularly deep learning models, can be computationally demanding and may necessitate significant computing resources for training and inference. This can be a limitation in some settings where access to high-performance computing resources may be limited or costly. Additionally, scaling AI models to handle large datasets or real-time prediction applications may also pose challenges in terms of computational efficiency and scalability.

Addressing these limitations is crucial for the responsible and effective use of AI in biodegradation prediction and ensuring that the models are reliable, transparent, and compliant with relevant regulations and ethical standards.

### **10.6.2 FUTURE DIRECTIONS OF AI IN BIODEGRADATION PREDICTION**

The field of AI in biodegradation prediction is constantly evolving, and there are several potential future directions that can be explored. Some of these directions include:

1. **Improved Data Collection and Integration:** Enhancing data collection methods for biodegradation prediction, including more comprehensive and diverse datasets, can improve the accuracy and robustness of AI models. This could involve leveraging advanced data acquisition techniques, such as high-throughput screening, omics technologies (genomics, proteomics, metabolomics), and sensor-based monitoring, to generate more data on biodegradation rates, environmental conditions, and chemical properties. Integrating data from different sources and domains, such as environmental, biological, and chemical data, can also provide a more holistic understanding of biodegradation processes and improve model performance.
2. **Advancements in Machine Learning Algorithms:** There is ongoing research in developing advanced machine learning algorithms specifically tailored for biodegradation prediction. This includes developing novel algorithms, such as deep learning architectures, recurrent neural networks, and graph-based models, that can capture complex relationships and patterns in biodegradation data. Additionally, integrating domain-specific knowledge, such as enzymatic mechanisms, metabolic pathways, and microbial interactions, into machine learning algorithms can enhance their predictive capabilities for biodegradation prediction.
3. **Model Interpretability and Explainability:** Improving the interpretability and explainability of AI models for biodegradation prediction is an important future direction. This can involve developing methods to interpret and explain the predictions of complex AI models, such as deep learning models, to gain insights into the decision-making process and increase trust and acceptance by stakeholders. Techniques such as explainable AI (XAI) and model-agnostic explanations can be explored to provide interpretable and transparent results for biodegradation prediction models.
4. **Integration of Multi-Omics Data:** This, which includes genomics, proteomics, metabolomics, and other high-dimensional data, can provide valuable information about the functional and metabolic activities of microorganisms involved in biodegradation processes. Integrating multi-omics data with AI models can enable a more comprehensive understanding of biodegradation pathways, metabolic interactions, and microbial communities, leading to more accurate and robust predictions.
5. **Model Transferability and Scalability:** Developing AI models that are transferable and scalable across different environmental conditions,



chemical compounds, and biodegradation pathways is an important future direction. This could involve developing transfer learning techniques, where models trained on one biodegradation scenario can be fine-tuned for another scenario with limited data, or developing ensemble models that combine multiple models to improve prediction accuracy and robustness. Additionally, optimizing AI models for computational efficiency and scalability can enable real-time or large-scale applications of biodegradation prediction in practical settings.

6. **Integration with Experimental Validation:** Experimental validation of biodegradation prediction models is crucial for their reliability and accuracy. Future directions may involve integrating AI models with experimental validation methods, such as lab-scale biodegradation experiments, microbial culturing, and molecular biology techniques, to validate and refine the predictions. This can help bridge the gap between computational predictions and real-world biodegradation processes, improving the practical applicability of AI models.
7. **Ethical and Societal Considerations:** Ethical and societal considerations will continue to be important in the future development and use of AI in biodegradation prediction. Ensuring fairness, transparency, and accountability in AI models, addressing issues such as bias and discrimination, and incorporating ethical guidelines and regulations into the development and deployment of AI models for biodegradation prediction will be crucial in shaping the future direction of this field.
8. **Real-Time Monitoring and Prediction:** Developing real-time monitoring and prediction systems using AI can enable timely detection and prediction of biodegradation processes in the environment, facilitating proactive environmental management and pollution mitigation strategies.
9. **Industry Applications:** The application of AI in biodegradation prediction can have significant implications for various industries, such as pharmaceuticals, chemicals, agriculture, and waste management. Future research can focus on developing industry-specific applications of AI for biodegradation prediction, tailored to the unique needs and challenges of different sectors.
10. **Academia, Industry, and Regulatory Agencies Collaboration:** Collaboration between academia, industry, and regulatory bodies can speed up the research and implementation of artificial intelligence in biodegradation prediction. Future research can focus on developing

such cooperation to ensure that AI models are produced and implemented in accordance with regulatory standards, industrial needs, and environmentally sustainable practices.

There are several potential future directions for AI in biodegradation prediction. One direction is the incorporation of multi-modal data, such as chemical, environmental, and biological data, to improve the accuracy and reliability of biodegradation prediction models. Another avenue of investigation is the use of ensemble approaches, which entail mixing numerous AI models to improve prediction accuracy.

Integration of AI with other emerging technologies such as the Internet of Things (IoT) and big data analytics is another potential future avenue. IoT sensors, for example, might be used to collect real-time data on environmental factors affecting biodegradation rates. This data could be analyzed using big data analytics techniques to identify patterns and relationships that can inform biodegradation prediction models.

Overall, there are many exciting possibilities for the future of AI in biodegradation prediction. By incorporating new data sources, exploring new techniques, and integrating with other technologies, AI has the potential to dramatically increase our capacity to forecast biodegradation rates and make intelligent environmental management decisions.

## **10.7 CONCLUSION**

AI potentially increase our capacity to anticipate biodegradation rates. To construct accurate and trustworthy biodegradation prediction models, several AI techniques such as machine learning, deep learning, and data mining can be applied. However, there are also challenges and limitations to using AI in biodegradation prediction, such as data availability, model robustness, and ethical considerations. Performance evaluation is an important aspect of developing and deploying AI-based biodegradation prediction models. It entails evaluating the accuracy and dependability of the model's predictions using unseen data. Interpreting AI models for biodegradation prediction can be challenging due to issues such as model explain ability and transparency. There are several potential future directions for AI in biodegradation prediction, such as incorporating multi-modal data, exploring ensemble methods, and integrating AI with other emerging technologies like IoT and big data analytics.

Biodegradation prediction using AI involves the use of machine learning and data-driven approaches to predict the biodegradation potential of chemicals, pollutants, or other substances in the environment. AI models can provide accurate and fast predictions of biodegradation rates, pathways, and microbial communities, which can have significant applications in environmental risk assessment, pollution mitigation, and sustainable waste management.

AI model performance must be evaluated, and commonly used metrics such as accuracy, precision, recall, and F1-score can aid in this process. Challenges in interpreting AI models for biodegradation prediction include model explain ability and transparency, as complex models may lack interpretability, making it difficult to understand the decision-making process. Data availability, model robustness, and ethical considerations are among the limitations of AI in biodegradation prediction that must be addressed for the reliable and responsible application of AI in this sector.

Potential future directions of AI in biodegradation prediction include improved data collection and integration, advancements in machine learning algorithms, model interpretability and explain ability, integration of multi-omics data, model transferability and scalability, integration with experimental validation, and addressing ethical and societal considerations.

In nutshell, artificial intelligence has the potential to transform biodegradation prediction by increasing accuracy, efficiency, and sustainability in environmental management practices. However, careful consideration of model limitations, interpretability, and ethical implications is crucial for responsible and reliable use of AI in this field.

## **KEYWORDS**

- **artificial intelligence**
- **artificial neural networks**
- **biological oxygen demand**
- **convolutional neural networks**
- **deep learning**
- **gated recurrent unit**
- **gradient boosting machines**
- **internet of things**

## REFERENCES

1. Zheng, Y. K., Qiao, X. G., Miao, C. P., Liu, K., Chen, Y. W., Xu, L. H., & Zhao, L. X. (2016). Diversity, distribution and biotechnological potential of endophytic fungi. *Annals of Microbiology*, 66, 529–542. <https://doi.org/10.1007/s13213-015-1153-7>.
2. Lan, T. D., Huong, D. T. T., & Trang, C. T. T. (2012). Assessment of natural resources use for sustainable development–DPSIR framework for case studies in Hai Phong and Nha Trang, Vietnam. In *Proceedings of the 12th International Coral Reef Symposium* (Vol. 7, July).
3. Patowary, R., Patowary, K., Kalita, M. C., Deka, S., Borah, J. M., Joshi, S. J., Zhang, M., Peng, W., Sharma, G., Rinklebe, J., & Sarma, H. (2022). Biodegradation of hazardous naphthalene and cleaner production of rhamnolipids–Green approaches of pollution mitigation. *Environmental Research*, 209, 112875. <https://doi.org/10.1016/j.envres.2022.112875>.
4. Bátori, V., Åkesson, D., Zamani, A., Taherzadeh, M. J., & Sárvári Horváth, I. (2018). Anaerobic degradation of bioplastics: A review. *Waste Management*, 80, 511–523. <https://doi.org/10.1016/j.wasman.2018.09.040>.
5. Dung, D. T., Cuong, N. H., Duy, D. K., & Huong, D. T. (2021). Orientation to sustainable and climate adapted agriculture with advanced technology in Industry 4.0 in Vietnam. *VNU Journal of Science: Economics and Business*, 37(1). <https://doi.org/10.25073/2588-1108/vnueab.4431>.
6. Singh, A. K., Bilal, M., Iqbal, H. M. N., & Raj, A. (2021). Trends in predictive biodegradation for sustainable mitigation of environmental pollutants: Recent progress and future outlook. *Science of The Total Environment*, 770, 144561. <https://doi.org/10.1016/J.SCITOTENV.2020.144561>.
7. Zumstein, M. T., Narayan, R., Kohler, H. P. E., McNeill, K., & Sander, M. (2019). Dos and do nots when assessing the biodegradation of plastics. *Environmental Science & Technology*, 53(17), 9967–9969. <https://doi.org/10.1021/ACS.EST.9B04513>.
8. Madsen, E. L. (1991). Determining in situ biodegradation. *Environmental Science & Technology*, 25(10), 1662–1673. <https://doi.org/10.1021/ES00022A001>.
9. Taniguchi, I., Yoshida, S., Hiraga, K., Miyamoto, K., Kimura, Y., & Oda, K. (2019). Biodegradation of PET: Current status and application aspects. *ACS Catalysis*, 9(5), 4089–4105. <https://doi.org/10.1021/ACSCATAL.8B05171>.
10. Nielsen, P. H., Bjerg, P. L., Nielsen, P., Smith, P., & Christensen, T. H. (1996). In situ and laboratory determined first-order degradation rate constants of specific organic compounds in an aerobic aquifer. *Environmental Science & Technology*, 30(1), 31–37. <https://doi.org/10.1021/ES940722O>.
11. Jin, H., & Ma, Q. (2021). Impacts of permafrost degradation on carbon stocks and emissions under a warming climate: A review. *Atmosphere*, 12(11), 1425. <https://doi.org/10.3390/ATMOS12111425>.
12. Khan, M. (2013). Biodegradable waste to biogas: Renewable energy option for the Kingdom of Saudi Arabia. *International Journal of Innovation, Management and Technology*, 4(1), 101–113. Retrieved from: researchgate.net (accessed on 25 July 2024).
13. Butnariu, M., & Bonciu, E. (2022). Biodegradable waste: Renewable energy source. In *Environmental Biotechnology* (pp. 159–199). <https://doi.org/10.1201/9781003277279-7/biodegradable-waste-renewable-energy-source-monica-butnariu-elena-bonciu>.

14. Butnariu, M. (2022). Biodegradable waste: Renewable energy source. Retrieved from: Taylor & Francis Online
15. Katinas, V., Marčiukaitis, M., Perednis, E., & Dzenajavičienė, E. F. (2019). Analysis of biodegradable waste use for energy generation in Lithuania. *Renewable and Sustainable Energy Reviews*, 101, 559–567. <https://doi.org/10.1016/J.RSER.2018.11.022>.
16. Sintim, H. Y., Bandopadhyay, S., English, M. E., Bary, A. I., DeBruyn, J. M., Schaeffer, S. M., Miles, C. A., Reganold, J. P., & Flury, M. (2019). Impacts of biodegradable plastic mulches on soil health. *Agriculture, Ecosystems & Environment*, 273, 36–49. <https://doi.org/10.1016/J.AGEE.2018.12.002>.
17. Wackett, L. P. (1999). Predicting biodegradation. *Environmental Microbiology*, 1(2), 119–124. <https://doi.org/10.1046/j.1462-2920.1999.00029.x>.
18. Finley, S. D., Broadbelt, L. J., & Hatzimanikatis, V. (2009). Computational framework for predictive biodegradation. *Biotechnology and Bioengineering*, 104(6), 1086–1097. <https://doi.org/10.1002/BIT.22489>.
19. Joutey, N. T., Bahafid, W., Sayel, H., & El Ghachtouli, N. (2013). Biodegradation: Involved Microorganisms and Genetically Engineered Microorganisms. In *Biodegradation*; Chamy, R., Rosenkranz, F., (eds.), *IntechOpen: Rijeka*, <https://doi.org/10.5772/56194>.
20. Loonen, H., Lindgren, F., Hansen, B., Karcher, W., Niemelä, J., Hiromatsu, K., Takatsuki, M., Peijnenburg, W., Rorije, E., & Struijs, J. (1999). Prediction of biodegradability from chemical structure: Modeling of ready biodegradation test data. *Environmental Toxicology and Chemistry*, 18(8), 1763–1768. <https://doi.org/10.1002/ETC.5620180822>.
21. Mohapatra, B. R., Dinardo, O., Gould, W. D., & Koren, D. W. (2011). 6.54 – Molecular Aspects of Microbial Dissimilatory Reduction of Radionuclides: A Review A2 – Moo-Young, Murray BT – *Comprehensive Biotechnology (Second Edition)*, 709–718. <https://doi.org/https://doi.org/10.1016/B978-0-08-088504-9.00402-5>.
22. Beulke, S., Dubus, I. G., Brown, C. D., & Gottesbüren, B. (2000). Simulation of pesticide persistence in the field on the basis of laboratory data—A review. *Journal of Environmental Quality*, 29(5), 1371–1379. <https://doi.org/10.2134/JEQ2000.00472425002900050001X>.
23. Dd, S. S., Blockeel, H., Kompare, B., Kramer, S., Pfahringer, B., & Laer, V. (2021). *Experiments in Predicting Biodegradability*. Springer.
24. Lv, S., Li, Y., Zhao, S., & Shao, Z. (2024). Biodegradation of Typical Plastics: From Microbial Diversity to Metabolic Mechanisms. *International Journal of Molecular Sciences*. Multidisciplinary Digital Publishing Institute (MDPI) January 1, <https://doi.org/10.3390/ijms25010593>.
25. Kamali, M., Appels, L., Yu, X., Aminabhavi, T. M., & Dewil, R. (2021). Artificial intelligence as a sustainable tool in wastewater treatment using membrane bioreactors. *Elsevier*, 417, 128070.
26. Blockeel, H., Džeroski, S., Kompare, B., Kramer, S., & Van Laer, B. P. W. (2004). *Experiments in Predicting Biodegradability*. Taylor & Francis. <https://doi.org/10.1080/08839510490279131>.
27. He, L., Bai, L., Dionysiou, D. D., Wei, Z., Spinney, R., Chu, C., Lin, Z., & Xiao, R. (2021). Applications of computational chemistry, artificial intelligence, and machine learning in aquatic chemistry research. *Elsevier*, 426, 121810.
28. Boutra, B., Sebti, A., & Trari, M. (2022). Response surface methodology and artificial neural network for optimization and modeling the photodegradation of organic pollutants in water. *International Journal of Environmental Science and Technology*, 19(11), 11263–11278. <https://doi.org/10.1007/S13762-021-03875-1>.

29. Boutra, B., & Sebti, A., (2022). Response surface methodology and artificial neural network for optimization and modeling the photodegradation of organic pollutants in water. *Science, MTIJ of E., Volume 19, Issue 2, Springer.*
30. Yang, M., Chen, L., Wang, J., Msigwa, G., Osman, A. I., Fawzy, S., Rooney, D. W., & Yap, P. S. (2023). Circular economy strategies for combating climate change and other environmental issues. *Environmental Chemistry Letters*, 21(1), 55–80. <https://doi.org/10.1007/S10311-022-01499-6>.
31. Nishant, R., Kennedy, M., & Corbett, J. (2020). Artificial intelligence for sustainability: Challenges, opportunities, and a research agenda. *Elsevier*, 53, 102104.
32. Moshood, T. D., Nawanir, G., & Mahmud, F. (2022). Sustainability of biodegradable plastics: A review on social, economic, and environmental factors. *Critical Reviews in Biotechnology*, 42(6), 892–912. <https://doi.org/10.1080/07388551.2021.1973954>.
33. Kompare, B. (1998). Estimating environmental pollution by xenobiotic chemicals using QSAR (QSBR) models based on artificial intelligence. In *Water Science and Technology* (Vol. 37). [https://doi.org/10.1016/S0273-1223\(98\)00231-5](https://doi.org/10.1016/S0273-1223(98)00231-5).
34. Raymond, J. W., Rogers, T. N., Shonnard, D. R., & Kline, A. A. (2001). A review of structure-based biodegradation estimation methods. *Elsevier*, 84, 189–215.
35. Klopman, G., Wang, S., & Balthasar, D. M. (1992). Estimation of aqueous solubility of organic molecules by the group contribution approach. Application to the study of biodegradation. *Journal of Chemical Information and Computer Sciences*, 32(5), 474–482. <https://doi.org/10.1021/CI00009A013>.
36. Kishore, S. C., Perumal, S., Atchudan, R., Alagan, M., Sundramoorthy, A. K., & Lee, Y. R. (2022). A critical review on artificial intelligence for fuel cell diagnosis. *Catalysts*, 12(7), 743. <https://doi.org/10.3390/catal12070743>.
37. Ren, L., Cui, J., & Sun, Y. (2017). Multi-bearing remaining useful life collaborative prediction: A deep learning approach. *Journal of Computational Systems Engineering. Elsevier*, 43, 248–256.
38. Kamali, M., Appels, L., Yu, X., Aminabhavi, T. M., & Dewil, R. (2021). Artificial intelligence as a sustainable tool in wastewater treatment using membrane bioreactors. *Elsevier*, 417, 128070.
39. Baker, J. R., Gamberger, D., Mihelcic, J., & Sabljic, A. (2004). Evaluation of artificial intelligence-based models for chemical biodegradability prediction. *MDPI*, 9(12), 989–1004.
40. Jaworska, J. S., Boethling, R. S., & Howard, P. H. (2003). Recent developments in broadly applicable structure-biodegradability relationships. *Environmental Toxicology and Chemistry*, 22(8), 1710–1723. <https://doi.org/10.1897/01-302>.
41. Howard, P. H., Stiteler, W. M., Meylan, W. M., Hueber, A. E., Beauman, J. A., Larosche, M. E., & Boethling, R. S. (1992). Predictive model for aerobic biodegradability developed from a file of evaluated biodegradation data. *Environmental Toxicology and Chemistry*, 11(5), 593–603. <https://doi.org/10.1002/ETC.5620110502>.
42. Pavan, M., & Worth, A. P. (2008). Review of estimation models for biodegradation. *QSAR & Combinatorial Science*, 27(1), 32–40. <https://doi.org/10.1002/QSAR.200710117>.
43. Nash, W., Drummond, T., & Degradation, N. B. M. (2018). A review of deep learning in the study of materials degradation. *Nature*.
44. Roberts, H., Cowls, J., Morley, J., Taddeo, M., Wang, V., & Floridi, L. (2021). The Chinese approach to artificial intelligence: An analysis of policy, ethics, and regulation. *Philosophical Studies Series*, 144, 47–79. [https://doi.org/10.1007/978-3-030-81907-1\\_5](https://doi.org/10.1007/978-3-030-81907-1_5).

45. Bell, J. (2022). What is machine learning? In *Machine Learning and the City: Applications in Architecture and Urban Design* (pp. 209–216). [https://doi.org/10.1007/978-3-319-18305-3\\_1](https://doi.org/10.1007/978-3-319-18305-3_1).
46. Alpaydin, E. (2021). *Machine Learning*. MIT Press.
47. What is machine learning? *Wiley Online Library*. J. B. M. L. & C. A., 2022.
48. Na Chu, Yong Jiang, Qinjun Liang, Panpan Liu, Donglin Wang, Xueming Chen, Daping Li, Peng Liang, Raymond Jianxiong Zeng, & Yifeng Zhang. (2023). Electricity-Driven Microbial Metabolism of Carbon and Nitrogen: A Waste-to-Resource Solution. *Environmental Science & Technology*, 57(11), 4379–4395. <https://doi.org/10.1021/acs.est.2c07588>.
49. Mahesh, B. (2019). Machine learning algorithms—A review. *ResearchGate*, 9(1). <https://doi.org/10.21275/ART20203995>.
50. Ahmad, Z., Zhong, H., Mosavi, A., Sadiq, M., Saleem, H., Khalid, A., Mahmood, S., & Nabipour, N. (2020). Machine learning modeling of aerobic biodegradation for azo dyes and hexavalent chromium. *Mathematics*, 8(6), 913. <https://doi.org/10.3390/MATH8060913>.
51. Schulze-Kremer, S., & King, R. D. (1992). IPSA—Inductive protein structure analysis. *Protein Engineering, Design and Selection*, 5(5), 377–390. <https://doi.org/10.1093/PROTEIN/5.5.377>.
52. Saiakhov, R. D., Stefan, L. R., & Klopman, G. (2000). Multiple computer-automated structure evaluation model of the plasma protein binding affinity of diverse drugs. *Perspectives in Drug Discovery and Design*, 19(1), 133–155. <https://doi.org/10.1023/A:1008723723679>.
53. Wellawatte, G. P., Gandhi, H. A., Seshadri, A., & White, A. D. (2023). A perspective on explanations of molecular prediction models. *Journal of Chemical Theory and Computation*, 19(8), 2149–2160. <https://doi.org/10.1021/ACS.JCTC.2C01235>.
54. Boethling, R. S., Gregg, B., Frederick, R., Gabel, N. W., Campbell, S. E., & Sabljic, A. (1989). Expert systems survey on biodegradation of xenobiotic chemicals. *Ecotoxicology and Environmental Safety*, 18(3), 252–267. [https://doi.org/10.1016/0147-6513\(89\)90019-5](https://doi.org/10.1016/0147-6513(89)90019-5).
55. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press. <https://doi.org/10.4258/hir.2016.22.4.351>.
56. Alzubaidi, L., Zhang, J., Humaidi, A. J., Al-Dujaili, A., Duan, Y., Al-Shamma, O., Santamaría, J., Fadhel, M. A., Al-Amidie, M., & Farhan, L. (2021). Review of Deep Learning: Concepts, CNN Architectures, Challenges, Applications, Future Directions. *J Big Data*, 8(1), 53. <https://doi.org/10.1186/s40537-021-00444-8>.
57. Hao, X., Zhang, G., & Ma, S. (2016). Deep learning. *International Journal of Semantic Computing*, 10(3), 417–439. <https://doi.org/10.1142/S1793351X16500045>.
58. LeCun, Y., Bengio, Y., & Nature, G. H. (2015). Deep learning. *Nature*.
59. Clustering, D. L. H. (2020). Deep learning. *ICPME*, <https://www.icpme.us>.
60. Kamilaris, A. (2018). Deep learning in agriculture: A survey. *Elsevier*.
61. Bengio, Y., Goodfellow, I., & Courville, A. (2017). *Deep Learning*. MIT Press.
62. Deng, L., & Processing, D. Y. F. and trends® in signal. (2014). Deep learning: Methods and applications. *NowPublishers*. <https://doi.org/10.1561/20000000039>.
63. Vaz, J. M., & Balaji, S. (2021). Convolutional neural networks (CNNs): Concepts and applications in pharmacogenomics. *Molecular Diversity*, 25(3), 1569–1584. <https://doi.org/10.1007/S11030-021-10225-3>.

64. Silva, R. (2020). A novel approach to condition monitoring of the cutting process using recurrent neural networks. *Sensors*, 20(16), 4493. <https://doi.org/10.3390/s20164493>.
65. Hand, D. (2001). Data mining. In *Encyclopedia of Environmetrics*. <https://doi.org/10.1002/9780470057339.VAD002>.
66. Cabena, P., Hadjinian, P., Stadler, R., Verhees, J., & Zanasi, A. (1998). *Discovering Data Mining: From Concept to Implementation*. Prentice Hall.
67. Hand, D. J. (2007). Principles of data mining. *Drug Safety*, 30(7), 621–622. <https://doi.org/10.2165/00002018-200730070-00010>.
68. Chen, M., Han, J., & Yu, P. S. (1996). Data mining: An overview from a database perspective. *IEEE Transactions on Knowledge and Data Engineering*. Retrieved from: *IEEE Xplore*.
69. Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. (2005). *Practical Machine Learning Tools and Techniques*. Morgan Kaufmann.
70. Polishchuk, P. (2017). Interpretation of quantitative structure-activity relationship models: Past, present, and future. *Journal of Chemical Information and Modeling*, 57(11), 2618–2639. <https://doi.org/10.1021/ACS.JCIM.7B00274>.
71. Muhammad, U., & Uzairu, A. (2018). Review on: Quantitative structure activity relationship (QSAR) modeling. *International Journal of Advances in Applied Sciences*, 4(5), 2488–9849.
72. Dudek, A., & Arodz, T. (2006). Computational methods in developing quantitative structure-activity relationships (QSAR): A review. *Journal of Computational Chemistry & High-Performance Computing*. Retrieved from *Ingenta Connect*
73. Bouarab-Chibane, L., Forquet, V., Lantéri, P., Clément, Y., Léonard-Akkari, L., Oulahal, N., Degraeve, P., & Bordes, C. (2019). Antibacterial properties of polyphenols: Characterization and QSAR (quantitative structure-activity relationship) models. *Frontiers in Microbiology*, 10, 829. <https://doi.org/10.3389/fmicb.2019.00829>.
74. Esposito, E. X., Hopfinger, A. J., & Madura, J. D. (2004). Methods for applying the quantitative structure-activity relationship paradigm. In *Methods in Molecular Biology* (Vol. 275, pp. 131–214). <https://doi.org/10.1385/1-59259-802-1:131>.
75. Mansouri, K., Ringsted, T., Ballabio, D., Todeschini, R., & Consonni, V. (2013). Quantitative structure-activity relationship models for ready biodegradability of chemicals. *Journal of Chemical Information and Modeling*, 53(4), 867–878. <https://doi.org/10.1021/CI4000213>.
76. *Ensemble Methods in Machine Learning*. (2000). Workshop, T. D. C. S. F. I., & MCS. Springer.
77. Kazienko, P., & Lughofer, E. (2013). Hybrid and ensemble methods in machine learning. *Journal of Universal Computer Science*. Special Issue. *academia.edu*.
78. Sagi, O., & Rokach, L. (2018). Ensemble learning: A survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 8(4). <https://doi.org/10.1002/WIDM.1249>.
79. Kunapuli, G. (2023). *Ensemble Methods for Machine Learning*. Springer.
80. Rezai, B., & Allahkarami, E. (2021). Application of neural networks in wastewater degradation process for the prediction of removal efficiency of pollutants. *Elsevier*, 2021, 75–93.
81. Ballabio, D., Biganzoli, F., Todeschini, R., & Consonni, V. (2017). Qualitative consensus of QSAR ready biodegradability predictions. *Toxicology and Environmental Chemistry*, 99(7–8), 1193–1216. <https://doi.org/10.1080/02772248.2016.1260133>.



82. Khan, S. (2018). A review on the application of deep learning in system health management. *Mathematics and Statistics*, 107, 241–265.
83. Navidpour, A. H., Hosseinzadeh, A., Huang, Z., Li, D., & Zhou, J. L. (2022). Application of machine learning algorithms in predicting the photocatalytic degradation of perfluorooctanoic acid. *Catalysis Reviews Science and Engineering*, 66(2), 687–712. <https://doi.org/10.1080/01614940.2022.2082650>.
84. He, L., Bai, L., Dionysiou, D., & Wei, Z. (2021). Applications of computational chemistry, artificial intelligence, and machine learning in aquatic chemistry research. *Computational and Environmental Engineering*, 426, 131810.
85. Padma, K. (2022). Application of artificial intelligence to detect and recover contaminated soil: An overview. In *Advances in Bioremediation and Phytoremediation for Sustainable Soil Management* (pp. 417–427). Springer. [https://doi.org/10.1007/978-3-030-89984-4\\_26](https://doi.org/10.1007/978-3-030-89984-4_26).
86. Bao, Q., Zhang, Z., Luo, H., & Tao, X. (2023). Evaluating and modeling the degradation of PLA/PHB fabrics in marine water. *Polymers*, 15(1). <https://doi.org/10.3390/POLYM15010082>.
87. Padma, K. R., & Don, K. R. (2022). Application of artificial intelligence to detect and recover contaminated soil: An overview. In *Advances in Bioremediation and Phytoremediation for Sustainable Soil Management* (pp. 417–427). Springer. [https://doi.org/10.1007/978-3-030-89984-4\\_26](https://doi.org/10.1007/978-3-030-89984-4_26).
88. Singh, A., Bilal, M., & Iqbal, H. (2021). Trends in predictive biodegradation for sustainable mitigation of environmental pollutants: Recent progress and future outlook. *Advanced Research in Environmental Sciences and Technology*, 770, 144561.
89. Khan, S. (2018). A review on the application of deep learning in system health management. *Mathematics and Statistics*, 107, 241–265.
90. Houssein, E. H., Hosney, M. E., Oliva, D., Ortega-Sánchez, N., Mohamed, W. M., & Hassaballah, M. (2021). Drug design and discovery: Theory, applications, open issues, and challenges. In *Studies in Computational Intelligence* (Vol. 967, pp. 337–358). [https://doi.org/10.1007/978-3-030-70542-8\\_15](https://doi.org/10.1007/978-3-030-70542-8_15).
91. Kar, S., & Leszczynski, J. (2020). Open access in silico tools to predict the ADMET profiling of drug candidates. *Expert Opinion on Drug Discovery*, 15(12), 1473–1487. <https://doi.org/10.1080/17460441.2020.1798926>.
92. Eli-Chukwu, N. C. (2019). Applications of artificial intelligence in agriculture: A review. [pdfs.semanticscholar.org/9\(4\), 4377–4383](https://pdfs.semanticscholar.org/9(4)/4377-4383).
93. Padma, K. R., & Don, K. R. (2022). Application of artificial intelligence to detect and recover contaminated soil: An overview. In *Advances in Bioremediation and Phytoremediation for Sustainable Soil Management* (pp. 417–427). Springer. [https://doi.org/10.1007/978-3-030-89984-4\\_26](https://doi.org/10.1007/978-3-030-89984-4_26).
94. Chivenge, P. P., Murwira, H. K., Giller, K. E., Mapfumo, P., & Six, J. (2007). Long-Term Impact of Reduced Tillage and Residue Management on Soil Carbon Stabilization: Implications for Conservation Agriculture on Contrasting Soils. *Soil Tillage Res*, 94(2), 328–337. <https://doi.org/10.1016/J.STILL.2006.08.006>.
95. Haripriyan, U., Gopinath, K. P., Arun, J., & Gobccvarthanan, M. (2022). Bioremediation of organic pollutants: A mini review on current and critical strategies for wastewater treatment. *Archives of Microbiology*, 204(5). <https://doi.org/10.1007/S00203-022-02907-9>.
96. Zacharof, A. (2004). Stochastic modeling of landfill processes incorporating waste heterogeneity and data uncertainty. *Waste Management & Research*, 24(3), 241–250.

97. Rezai, B. (2021). Application of neural networks in wastewater degradation process for the prediction of removal efficiency of pollutants. *Environmental and Sustainable Waste Management*, 2021, 75–93.
98. Restitution, M. N. P. and E. (2018). Modeling applications in environmental bioremediation studies. In *Environmental Bioremediation Studies* (pp. 143–160). Springer. [https://doi.org/10.1007/978-981-13-1187-1\\_7](https://doi.org/10.1007/978-981-13-1187-1_7).
99. Onesios, K. M., Yu, J. T., & Bouwer, E. J. (2009). Biodegradation and removal of pharmaceuticals and personal care products in treatment systems: A review. *Biodegradation*, 20(4), 441–466. <https://doi.org/10.1007/S10532-008-9237-8>.
100. Malviya, A., & Jaspal, D. (2021). Artificial intelligence as an upcoming technology in wastewater treatment: A comprehensive review. *Environmental Technology Reviews*, 10(1), 177–187. <https://doi.org/10.1080/21622515.2021.1913242>.
101. Nyika, J., & Dinka, M. O. (2022). A Mini-Review on Wastewater Treatment through Bioremediation towards Enhanced Field Applications of the Technology. AIMS Environmental Science. *American Institute of Mathematical Sciences*, pp 403–431. <https://doi.org/10.3934/environsci.2022025>.
102. Jaworska, J., & Bianchini, R. (2003). Recent developments in broadly applicable structure–biodegradability relationships. *Environmental Toxicology and Chemistry*, 22(8), 1710–1723. <https://doi.org/10.1897/01-302>.
103. Backhaus, T., & Faust, M. (2012). Predictive environmental risk assessment of chemical mixtures: A conceptual framework. *Environmental Science & Technology*, 46(5), 2564–2573. <https://doi.org/10.1021/ES2034125>.
104. Tabbussum, R., & Dar, A. Q. (2021). Performance evaluation of artificial intelligence paradigms—Artificial neural networks, fuzzy logic, and adaptive neuro-fuzzy inference system for flood prediction. *Environmental Science and Pollution Research*, 28(20), 25265–25282. <https://doi.org/10.1007/S11356-021-12410-1>.
105. Heys, K. A., Shore, R. F., Pereira, M. G., Jones, K. C., & Martin, F. L. (2016). Risk Assessment of Environmental Mixture Effects. *RSC Adv.*, 6(53), 47844–47857. <https://doi.org/10.1039/C6RA05406D>.
106. Tian, H., Wang, Z., Zhu, T., Yang, C., Shi, Y., & Sun, Y. (2021). Degradation prediction and products of polycyclic aromatic hydrocarbons in soils by highly active bimetal/AC-activated persulfate. *ACS ES&T Engineering*, 1(8), 1183–1192. <https://doi.org/10.1021/ACSESTENG.1C00063>.
107. Biotechnology, Z. U.-M. (2022). Big data and computational advancements for the next generation of microbial biotechnology. *NCBI*. Retrieved from: <https://ncbi.nlm.nih.gov> (accessed on 25 July 2024).
108. Fajobi, M., & Lasode, O. (2022). Effect of biomass co-digestion and application of artificial intelligence in biogas production: A review. *Energy Sources, Part A: Recovery, Utilization, and Environmental Effects*, 44(2), 5314–5339. <https://doi.org/10.1080/15567036.2022.2085823>.
109. Kishore, S. C., Perumal, S., Atchudan, R., Alagan, M., Sundramoorthy, A. K., & Lee, Y. R. (2022). A critical review on artificial intelligence for fuel cell diagnosis. *Catalysts*, 12(7), 743. <https://doi.org/10.3390/catal12070743>.
110. Shahsavar, M. M., Akrami, M., Gheibi, M., Behzadian, K., & Fathollahi-Fard, A. M. (2022). A smart framework for supplying the biogas energy in green buildings using an integration of response surface methodology, artificial intelligence, and Petri net. *Elsevier*.

111. Elsayed, M., Abomohra, A., Ai, P., & D. W. B. (2018). Of rice straw by sequential fermentation and anaerobic digestion for bioethanol and/or biomethane production: Comparison of structural properties and energy output. *Elsevier*.
112. Adikaram, K., Hussein, M., & M. E. T. S. W. (2014). Outlier detection method in linear regression based on sum of arithmetic progression. *Hindawi*. Retrieved from: <https://hindawi.com> (accessed on 25 July 2024).
113. Magesh, S., Niveditha, V. R., Rajkumar, P. S., & Radha Rammohan, S. (2020). Pervasive computing in the context of COVID-19 prediction with AI-based algorithms. *International Journal of Pervasive Computing and Communications*, 16(5), 477–487. <https://doi.org/10.1108/IJPCC-07-2020-0082>.
114. Sun, L., Shang, Z., Xia, Y., Bhowmick, S., & Nagarajaiah, S. (2020). Review of bridge structural health monitoring aided by big data and artificial intelligence: From condition assessment to damage detection. *Journal of Structural Engineering*, 146(5). [https://doi.org/10.1061/\(ASCE\)ST.1943-541X.0002535](https://doi.org/10.1061/(ASCE)ST.1943-541X.0002535).
115. Vision, M. N. J. of M. I. and. (2004). A variational approach to remove outliers and impulse noise. *Journal of Mathematical Imaging and Vision*, 20, 99–120. <https://doi.org/10.1023/B:JMIV.0000011920.58935.9c>.
116. Kaaya, I., Lindig, S., Weiss, K. A., Virtuani, A., Sidrach de Cardona Ortin, M., & Moser, D. (2020). Photovoltaic lifetime forecast model based on degradation patterns. *Progress in Photovoltaics: Research and Applications*, 28(10), 979–992. <https://doi.org/10.1002/PIP.3280>.
117. Sambasivan, N., Kapania, S., Highfill, H., Akrong, D., Paritosh, P., & Aroyo, L. (2021). “Everyone wants to do the model work, not the data work”: Data cascades in high-stakes AI. *ACM Digital Library*. <https://doi.org/10.1145/3411764.3445518>.
118. Mansouri, K., Ringsted, T., Ballabio, D., Todeschini, R., & Consonni, V. (2013). Quantitative structure-activity relationship models for ready biodegradability of chemicals. *Journal of Chemical Information and Modeling*, 53(4), 867–878. <https://doi.org/10.1021/CI4000213>.
119. Baker, J. R., Gamberger, D., Mihelcic, J. R., & Sabljic, A. (2004). Evaluation of artificial intelligence-based models for chemical biodegradability prediction. *Molecules: A Journal of Synthetic Chemistry and Natural Product Chemistry*, 9(12), 989. <https://doi.org/10.3390/91200989>.
120. Sauban, S., & G. A. (2020). A comprehensive review of database resources in chemistry. *Eclética Química*, 45(3), 57–68. <https://doi.org/10.26850/1678-4618EQJ.V45.3.2020.P57-68>.
121. Anser, M. K., Usman, M., Sharif, M., Bashir, S., Malik, M., Shabbir, S., Ghulam, M., Khan, Y., Lydia, M., & Lopez, B. (2022). The dynamic impact of renewable energy sources on environmental economic growth: Evidence from selected Asian economies. *Environmental Science and Pollution Research*, 1(3), 3. <https://doi.org/10.1007/s11356-021-17136-8>.
122. Adebayo, T. S. (2022). Environmental consequences of fossil fuel in Spain amidst renewable energy consumption: A new insight from the wavelet-based Granger causality approach. *Journal of Sustainable Development*, 29(7), 579–592. <https://doi.org/10.1080/13504509.2022.2054877>.
123. Liu, X., Lu, D., Zhang, A., Liu, Q., & Jiang, G. (2022). Data-driven machine learning in environmental pollution: Gains and problems. *Environmental Science & Technology*, 56(4), 2124–2133. <https://doi.org/10.1021/ACS.EST.1C06157>.

124. Thanh Noi, P., & Kappas, M. (2018). Comparison of random forest, k-nearest neighbor, and support vector machine classifiers for land cover classification using Sentinel-2 imagery. *Sensors (Basel)*, 18(1). <https://doi.org/10.3390/S18010018>.
125. Lee, M., & Min, K. (2022). A comparative study of the performance for predicting biodegradability classification: The quantitative structure–Activity relationship model vs the graph convolutional network. *ACS Omega*, 7(4), 3649–3655. <https://doi.org/10.1021/ACSOMEGA.1C06274>.
126. Hernández, M. A., & Stolfo, S. J. (1998). Real-world data is dirty: Data cleansing and the merge/purge problem. *Data Mining and Knowledge Discovery*, 2(1), 9–37. <https://doi.org/10.1023/A:1009761603038>.
127. García, S., Ramírez-Gallego, S., Luengo, J., Benítez, J. M., & Herrera, F. (2016). Big data preprocessing: Methods and prospects. *Big Data Analytics*, 1(1). <https://doi.org/10.1186/S41044-016-0014-0>.
128. May, C. R., Cummings, A., Girling, M., Bracher, M., Mair, F. S., May, C. M., Murray, E., Myall, M., Rapley, T., & Finch, T. (2018). Using normalization process theory in feasibility studies and process evaluations of complex healthcare interventions: A systematic review. *Implementation Science*, 13(1). <https://doi.org/10.1186/S13012-018-0758-1>.
129. Dalkin, S. M., Hardwick, R. J. L., Highton, C. A., & Finch, T. L. (2021). Combining realist approaches and normalization process theory to understand implementation: A systematic review. *Implementation Science Communications*, 2(1). <https://doi.org/10.1186/S43058-021-00172-3>.
130. McEvoy, R., Ballini, L., Maltoni, S., O'Donnell, C. A., Mair, F. S., & MacFarlane, A. (2014). A qualitative systematic review of studies using the normalization process theory to research implementation processes. *Implementation Science*, 9(1). <https://doi.org/10.1186/1748-5908-9-2>.
131. May, C. R., Mair, F., Finch, T., MacFarlane, A., Dowrick, C., Treweek, S., Rapley, T., Ballini, L., Ong, B. N., Rogers, A., Murray, E., Elwyn, G., Légaré, F., Gunn, J., & Montori, V. M. (2009). Development of a theory of implementation and integration: Normalization process theory. *Implementation Science*, 4(1). <https://doi.org/10.1186/1748-5908-4-29>.
132. Huang, K., & Zhang, H. (2022). Classification and regression machine learning models for predicting aerobic ready and inherent biodegradation of organic chemicals in water. *Environmental Science & Technology*, 56(17), 12755–12764. <https://doi.org/10.1021/ACS.EST.2C01764>.
133. Naser, M. Z., & Alavi, A. H. (2021). Error metrics and performance fitness indicators for artificial intelligence and machine learning in engineering and sciences. *Architecture, Structures and Construction*. <https://doi.org/10.1007/S44150-021-00015-8>.
134. Xu Wang, Liang He, Lulu Xu, Zhongshou Liu, Yao Xiong, Weiqi Zhou, Hang Yao, Yangping Wen, Xiang Geng, & Ruimei Wu. (2023). Intelligent analysis of carbendazim in agricultural products based on a ZSHPC/MWCNT/SPE portable nanosensor combined with machine learning methods. *Analytical Methods*, 15(5), 562–571. <https://doi.org/10.1039/D2AY01779B>.
135. Jiang, S., Liang, Y., Shi, S., Wu, C., & Shi, Z. (2023). Improving Predictions and Understanding of Primary and Ultimate Biodegradation Rates with Machine Learning Models. *Science of The Total Environment*, 904, 166623. <https://doi.org/10.1016/J.SCITOTENV.2023.166623>.

136. Yin, H., Lin, C., Tian, Y., & Yan, A. (2022). Prediction and structure-activity relationship analysis on ready biodegradability of chemical using machine learning method. *Chemical Research in Toxicology*, 36(4). <https://doi.org/10.1021/ACS.CHEMRESTOX.2C00330>.
137. Zhao, Y., Zhang, C., Zhang, Y., Wang, Z., & Li, J. (2020). A review of data mining technologies in building energy systems: Load prediction, pattern identification, fault detection, and diagnosis. *Energy and Built Environment*, 1(2), 149–164. <https://doi.org/10.1016/J.ENBENV.2019.11.003>.
138. Bunke, H., & Kandel, A. (2002). *Hybrid Methods in Pattern Recognition*.
139. Caldas, R. D., Rodrigues, A., Gil, E. B., Rodrigues, G. N., Vogel, T., & Pelliccione, P. (2020). A hybrid approach combining control theory and AI for engineering self-adaptive systems. *ACM Digital Library*, 9–19. <https://doi.org/10.1145/3387939.3391595>.
140. Talreja, N., Chauhan, D., & Ashfaq, M. (2023). Photo-Antibacterial Activity of Two-Dimensional (2D)-Based Hybrid Materials: Effective Treatment Strategy for Controlling Bacterial Infection. *Antibiotics (Basel)*, 12(2), 398. doi: 10.3390/antibiotics12020398. PMID: 36830308; PMCID: PMC9952232.
141. Huang, K., & Zhang, H. (2022). Classification and regression machine learning models for predicting aerobic ready and inherent biodegradation of organic chemicals in water. *Environmental Science & Technology*, 56(17), 12755–12764. <https://doi.org/10.1021/ACS.EST.2C01764>.
142. Yushu Cheng, Kai Zhang, Kuan Huang, & Huichun Zhang (2024). Meta-Analysis and Machine Learning Models for Anaerobic Biodegradation Rates of Organic Contaminants in Sediments and Sludge. *Environmental Science & Technology*, 58(29), 12976–12988. <https://doi.org/10.1021/acs.est.4c01033>.
143. Ecosystem, A. T. O. (2016). Adoption of machine learning techniques in ecology and earth science. *One Ecosystem*. <https://doi.org/10.3897/oneeco.1.e8621>.
144. Thomas, R., & Uminsky, D. (2020). The problem with metrics is a fundamental problem for AI. *Unpublished Manuscript*.
145. Chen Zhao, Xinyue Xu, Hongmei Chen, Fengwen Wang, Penghui Li, Chen He, Quan Shi, Yuanbi Yi, Xiaomeng Li, Siliang Li, & Ding He. (2023). Exploring the Complexities of Dissolved Organic Matter Photochemistry from the Molecular Level by Using Machine Learning Approaches. *Environmental Science & Technology*, 57(46), 17889–17899. <https://doi.org/10.1021/acs.est.3c00199>.
146. Haobo Wang, Wenjia Liu, Jingwen Chen, & Zhongyu Wang. (2023). Applicability Domains Based on Molecular Graph Contrastive Learning Enable Graph Attention Network Models to Accurately Predict 15 Environmental End Points. *Environmental Science & Technology*, 57(44), 16906–16917. <https://doi.org/10.1021/acs.est.3c03860>.
147. Sommer, S., & K., S. (2007). Three data mining techniques to improve lazy structure–activity relationships for noncongeneric compounds. *Journal of Chemical Information and Modeling*, 47(6), 2035–2043.
148. Lee, J., Im, J., Kim, U., & Löffler, F. E. (2016). A data mining approach to predict in situ detoxification potential of chlorinated ethenes. *Environmental Science & Technology*, 50(10), 5181–5188. <https://doi.org/10.1021/ACS.EST.5B05090>.

# Computer-Based Technologies for Prediction of Biodegradation

KUMARI NEHA,<sup>1</sup> KALICHARAN SHARMA,<sup>2</sup> and SHARAD WAKODE<sup>1</sup>

<sup>1</sup>*Department of Pharmaceutical Chemistry, Delhi Institute of Pharmaceutical Sciences and Research, DPSR University, New Delhi, India*

<sup>2</sup>*Department of Pharmaceutical Chemistry, Delhi Pharmaceutical Sciences and Research University, New Delhi, India*

---

### ABSTRACT

Artificial Intelligence (AI) techniques have been increasingly utilized in progressive years for the prediction of biodegradation due to their capability to process enormous datasets and retrieve pertinent data. The prediction of biodegradation is vital for the development of ecologically friendly products and the planning of waste management strategies. AI-based models have shown great promise in forecasting biodegradability and evaluating the potential toxicants related to the flushing of substances into the atmosphere. The expenditure of AI methods in predicting biodegradation involves the creation of models that can examine a variety of environmental parameters and forecast the pace and degree of biodegradation under different circumstances. Machine learning (ML) and artificial neural networks (ANNs) are two of the most often utilized AI techniques in biodegradation prediction. ML is a subdivision of AI that provides processors the skill to acquire from documents/facts and develop over time. ML algorithms can be trained on massive datasets of biodegradation information to identify patterns and links between environmental conditions and biodegradation outcomes based on their chemical structure and other pertinent properties. ANNs are a form of machine learning algorithm that imitates the composition and operation of

the human brain. ANNs are particularly useful in forecasting biodegradation because they can manage huge, complicated datasets and recognize non-linear correlations between inputs and outcomes. ANNs have been used to predict the biodegradability of many different molecules, such as insecticides, medications, and industrial chemicals.

In addition to ML and ANN, Other AI tools, including fuzzy logic, genetic algorithms, and expert systems, have also been used to forecast biodegradation. Compared to conventional methods, using AI techniques to forecast biodegradation has several benefits. Large datasets can be analysed by AI-based models to extract pertinent data, enabling more precise predictions. The complicated and non-linear interactions between environmental variables and biodegradation outcomes can also be handled by AI approaches. Additionally, AI-based models are particularly helpful for tracking environmental conditions since they can be trained to spot patterns and make predictions in real time.

## **11.1 INTRODUCTION**

In current years, the use of artificial intelligence (AI) in the area of biodegradation prediction has gained significant attention due to its potential to accelerate the discovery of new, sustainable materials [1]. Biodegradation is the breaking of molecules of organic matter by biological microbes into less toxic substances that can be taken up by the ecosystem without damage. This process is an essential part of the natural cycle of carbon, nitrogen, and other elements in the environment. The process of biodegradation occurs when microbes such as algae, bacteria, and fungi ingest carbon-based compounds as a spring of energy and nutrients [2]. During this process, the organic molecules are broken down into uncomplicated compounds such as water, carbon dioxide, and minerals. Biodegradation engages in a perilous part in preserving the wellbeing of ecosystems by removing harmful pollutants and recycling nutrients. It is also an important process for the degradation of carbon-based matter in discarded water management and composting [3]. Biodegradation is critical in determining the biological and ecological risk of organic compounds created by the modern chemical and pharmaceutical industries that are unable to keep up by the rate at which these molecules are manufactured, creating attentions about their final destiny if discharged into the environment [4]. However, predicting the biodegradability of a given molecule is a complex task that requires considering a multitude of factors such as chemical structure, functional groups, and environmental conditions.

Historically, experimental methods have been used to study biodegradation, but they have limitations, such as requiring large amounts of time and resources to perform, as well as difficulties in obtaining accurate and reliable data. The latest developments in machine learning, in specifically deep learning, have unbolted up new potentials for biodegradation prediction. With the advent of machine learning methods and the availability of large datasets/records, AI has appeared as an encouraging tool for predicting the biodegradability of molecules with high accuracy [5]. This has imperative inferences for an eclectic range of uses, from drug discovery to the design of ecologically friendly resources. In this context, the use of AI for predicting biodegradation represents an influential tool for advancing the progress of sustainable technologies and reducing the environmental impact of chemical compounds.

In addition, *in silico* methods such as quantitative structure-activity relationship (Q)SARs have been proposed as a means of supplementing and enhancing experimental data [6]. These models use molecular descriptors and other chemical properties to predict biodegradability, providing a complementary approach to experimental methods. The ability to predict biodegradation accurately is crucial for environmental risk assessment and management, as it allows for the identification of hazardous chemicals and the development of effective waste management strategies.

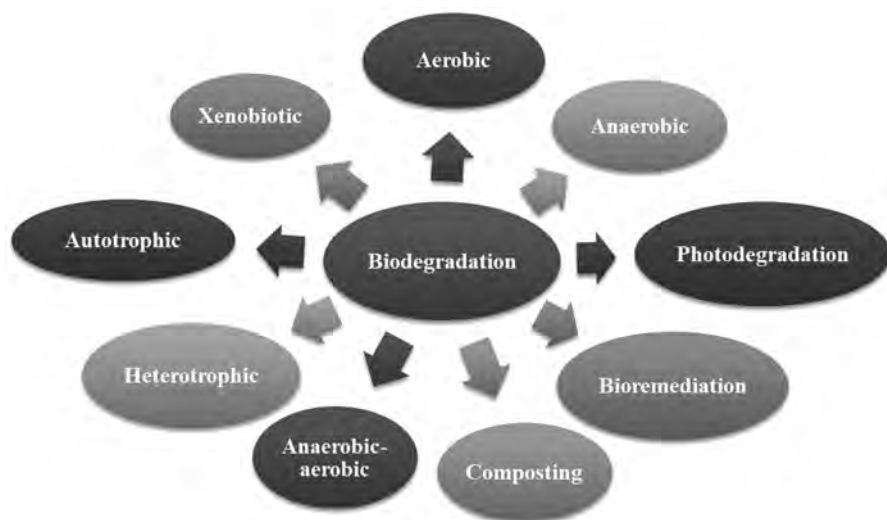
Knowledge on a substance's potential for breakdown is necessary for an effective risk evaluation because biological degradation is the predominant mechanism for the transformation of the environment of most compounds. However, there are few ways to forecast a compound's ability to degrade because this is contingent upon both the chemical's design and the environment to which it is exposed to [7]. Both the biological and ecological aspects of the breakdown mechanism based on discussion of quantitative structure-biodegradability relationship models (QSBRs) are highlighted in the current chapter's [8]. Investigations on the microbiological features of biodegradation and the techniques for determining biodegradability are provided. The latest developments in biodegradation modeling, such as contributions to computerized biodegradability predicting systems, are examined. The process of validating several recently created models for evaluating risk and ecological effect in marine and terrestrial systems is addressed [9]. Microbiology, sciences of the environment, biological technology, and bioremediation processes are all active fields of study. The report will be crucial for policymakers as they evaluate the present form of acquaintance on the breakdown of chemicals.



### 11.1.1 TYPES OF BIODEGRADATIONS

Several types of biodegradations (Figure 11.1) can occur depending on the behavior of the organic compound and the ecological situations. Some of the common types of biodegradations are:

1. **Aerobic Biodegradation:** This type of biotransformation happens in the existence of oxygen. Aerobic microorganisms utilize oxygen to digest organic molecules, generating carbon dioxide and water as by-products [10].
2. **Anaerobic Biodegradation:** This occurs in the absence of oxygen. In this process, bacteria consume alternative electron acceptors like sulfate, nitrate, or carbon dioxide to degrade organic compounds [11].
3. **Photo-Degradation:** It involves the degradation of organic compounds through exposure to sunlight or other forms of radiation. The energy from the radiation can break down the chemical bonds of the carbon-based molecules, leading to their degradation [12].
4. **Bioremediation:** It is a type of biodegradation that comprises the usage of microbes to clear-out contaminated environments. Microorganisms are introduced to contaminated soil or water, where they metabolize the organic compounds and reduce their concentrations to safe levels [13].
5. **Composting:** It is a form of biodegradation that involves the controlled degradation of organic matter such as food waste or plant material. Microorganisms in the composting pile break down the organic matter, producing a nutrient-rich soil amendment [14].
6. **Anaerobic-Aerobic:** This respiration is a type of biodegradation process that occurs in environments where oxygen availability is fluctuating or limited. In this process, microorganisms can switch between anaerobic and aerobic respiration depending on the availability of oxygen [15].
7. **Heterotrophic:** This type of biodegradation comprises the use of organic substances as a source of energy and carbon by heterotrophic microorganisms.
8. **Autotrophic:** This type of biodegradation comprises the use of organic substances as a source of energy and carbon by autotrophic microorganisms [16].
9. **Xenobiotic Biodegradation:** This type of biodegradation involves the breakdown of synthetic or man-made organic compounds, such as pesticides or plastics, by microorganisms.



**FIGURE 11.1** Various types of biodegradations.

### **11.1.2 MECHANISM OF BIODEGRADATION**

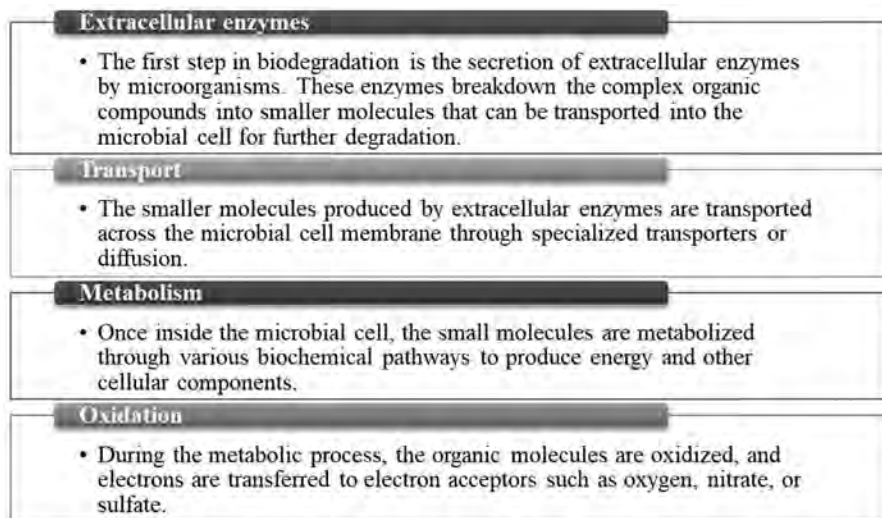
The mechanism of biodegradation involves several steps that are carried out by microorganisms in the environment. These steps can vary reliant on the category of compound being degraded, the ecological conditions, and the microbial community involved [17]. However, some general processes that occur during biodegradation shown in Figure 11.2.

### **11.1.3 FACTORS AFFECTING THE RATE OF BIODEGRADATION**

The biotransformation of organic compounds is exaggerated by a variety of elements that can influence the rate and extent of degradation. These factors include:

1. **Chemical Structure:** The chemical design of the organic compound is one of the most significant features affecting its biodegradability. Compounds with complex structures, such as aromatic compounds and halogenated compounds, are generally more resistant to biodegradation than simpler compounds.
2. **Molecular Weight:** This of the compound can also affect its biodegradability. Compounds with high molecular weights are often

more difficult for microorganisms to metabolize, and therefore, they degrade more slowly.



**FIGURE 11.2** Stages of biodegradations.

- 3. Solubility:** This of the compound in water and other solvents can also impact its biodegradability. Compounds that are more soluble in water are typically more easily degraded, as they can be more readily transported into microbial cells.
- 4. Environmental Conditions:** Ecological surroundings such as temperature, oxygen availability, pH, and nutrient availability can have a significant impact on biodegradation. For example, optimal temperature and pH levels are required for the growth and actions of microbes involved in breakdown of molecules [18].
- 5. Microbial Community:** The occurrence of microbial community in the habitat can also distress the rate and extent of biotransformation. Different microorganisms have varying metabolic capabilities, and the occurrence of explicit microbial species can promote or inhibit the biodegradation of certain compounds.

Overall, the biodegradability of a compound is influenced by a complex interplay of various factors, and a better understanding of these factors can help in designing added sustainable and ecologically friendly substances.

## 11.2 BACKGROUND

Predicting biodegradation accurately is critical for analyzing the environmental factors associated with chemicals, calculating toxicity, and establishing appropriate waste management techniques. Recent breakthroughs in machine learning, particularly deep learning, have unlocked new-fangled avenues for predicting biodegradation. These methods use massive datasets of chemical structures and biodegradation outcomes to discover patterns and make predictions. Deep learning algorithms are capable of learning complicated correlations between chemical characteristics and biodegradability, outperforming earlier prediction methods [19].

The absence of first-class experimental facts is a major difficulty in biodegradation prediction. *In-silico* methodologies, namely QSAR (quantitative structure-activity relationship) models, have been presented as a way to augment and improve experimental data. These models predict biodegradability using molecular descriptors and other chemical parameters, giving a supplementary approach to experimental approaches [20]. To summarize, biodegradation prediction is an important area of research having implications for environmental risk assessment, toxicology, and waste management. Recent breakthroughs in deep learning and *In-silico* technologies offer interesting options for enhancing our understanding of biodegradability and its environmental consequences, as well as the creation of effective remediation measures.

## 11.3 MATERIAL AND METHODS

A regular internet browser could be utilized for accessing both websites indicated in the protocols. Java applets are used for executing some of its functionalities. You might require altering the settings in the web browser or downloading additional applications to execute these applets. The *In-silico* methodology is commonly utilized in the design of biotransformation and bioremediation experiments. In addition, artificial neural networks (ANN), machine learning (ML), and genetic algorithms (GA)-based systems predict decomposition, toxic interactions, and ecosystem fate. Moreover, current improvements in QSAR modeling, algorithms, and specific biodegradation forecast systems with distinguishing features have been discussed. The protocols that follow describe how to use one technologies of specific type, namely EAWAG-PPS (formerly UM-PPS) and SVM Biodegradability predictor.

Both use chemical formulas as initial response, but they provide slightly unlike and harmonizing evidence on their utmost expected ecological result. Other biodegradation prediction software is Biowin, Catabol-301C model, and META-PC is also utilized. In addition to experimental testing and computer modeling, databases such as the University of Minnesota biocatalysis/biodegradation (UM-BBD), Plastics Microbial Biodegradation Database, AromaDeg, ONDB (Organonitrogen Degradation Database), MetaCyc, OxDBase (biodegradative oxygenases database), PBT Profiler, and Bionemo (biodegradation network-molecular biology database) are accessible [21].

#### **11.4 ROLE OF AI ALGORITHMS AND TECHNIQUES IN BIODEGRADATION**

AI tools have revealed pronounced ability in the prediction of biodegradation. These tools include various fuzzy logic, artificial neural networks (ANNs), genetic algorithms, machine learning (ML) algorithms, and expert systems. Here are some examples of how these tools are being used in the field of biodegradation prediction:

1. **Machine Learning (ML) algorithms:** These algorithms can be trained on large datasets of biodegradation information to recognize patterns and relationships between environmental factors and biodegradation outcomes. Based on the molecular arrangement and other physical parameters, ML algorithms were employed to foresee the biodegradability of organic molecules. The results showed that the ML algorithms outperformed traditional methods in predicting biodegradation [22].
2. **Artificial Neural Networks (ANNs):** These are particularly useful in predicting biodegradation because they can handle large and complex datasets and identify non-linear relationships between factors and outcomes. For example, in a study conducted on ANNs, which were used to predict the biodegradability of five common pesticides. The fallouts showed that the ANN model was capable to accurately foresee the biodegradability of the pesticides based on their physicochemical properties [23, 24].
3. **Fuzzy Logic:** It is a mathematical technique that can handle uncertain or vague information and is particularly useful in dealing with complex environmental factors. Fuzzy logic was used to expect the

biodegradability of organic composites in the presence of multiple pollutants. The outcomes showed that the fuzzy logic model was capable to accurately forecast the biodegradability of the compounds under different environmental conditions [25].

4. **Genetic Algorithms:** These are optimization techniques that can search for the best combination of environmental factors to maximize biodegradation. These algorithms were utilized to optimize the biotransformation of polycyclic aromatic hydrocarbons (PAHs) in polluted loam. The results showed that the genetic algorithm was able to identify the optimal combination of environmental factors to maximize PAH biodegradation [26].
5. **Expert Systems:** These are computer software package that can replicate the decision-making method of human experts and can be used to predict biodegradation grounded on a number of directions. An expert system was developed to envisage the biodegradability of carbon-based compounds based on their molecular arrangement and other physicochemical properties [27].

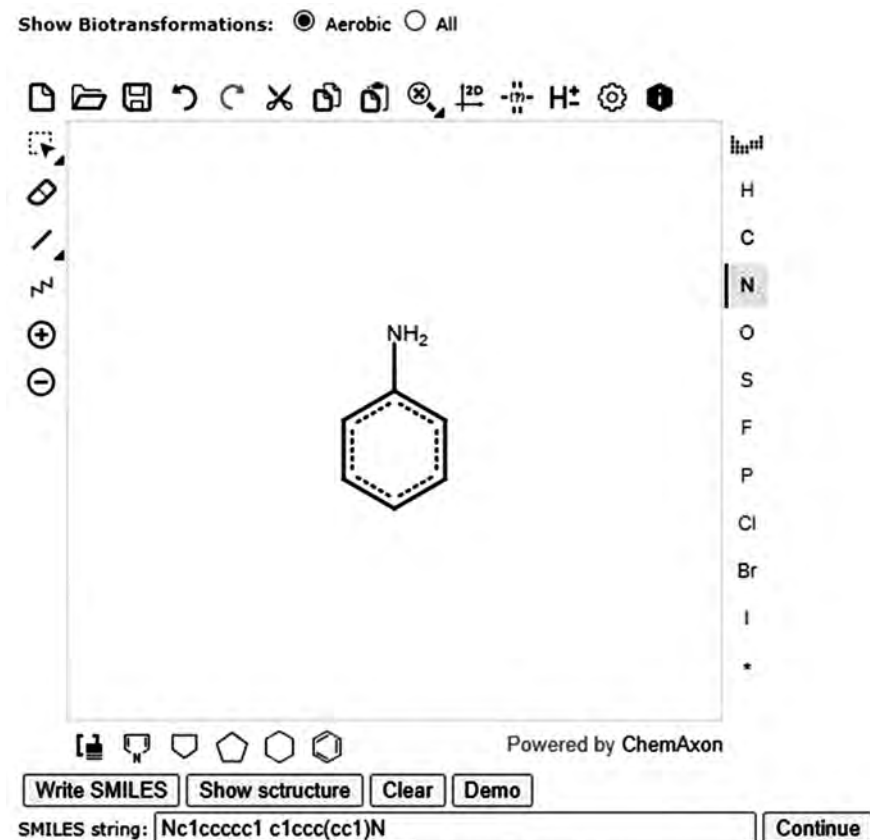
## 11.5 PREDICTION OF BIODEGRADATION USING AI TOOLS

Biodegradation prediction tools are important for estimating the potential conservational impact of chemicals and compounds. These tools use various methods, including machine learning algorithms and biodegradation pathways analysis, to predict the biodegradability of a substance. Biodegradability prediction tools can provide valuable information on the potential for a chemical to break down in the environment and the likelihood of it persisting and gathering in the environment [28, 29]. This information is vital for measuring the environmental hazard of a substance and for guiding decisions related to its use and regulation. With the growing concern over environmental pollution, the development of accurate and reliable biodegradation prediction tools is becoming increasingly important to protect the health and well-being of both human and ecological communities. Some of the tools are discussed in subsections.

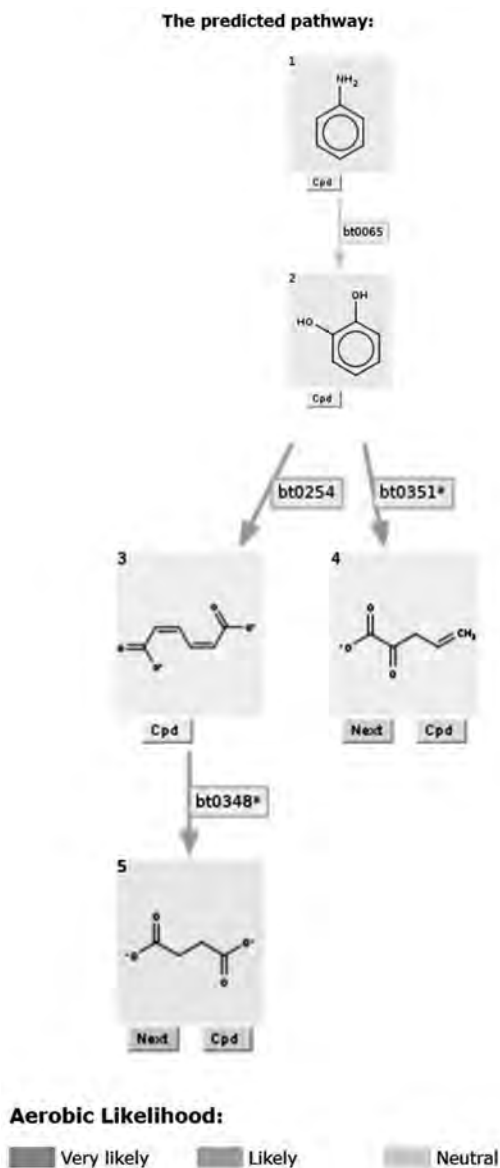
### 11.5.1 EAWAG-PPS

EAWAG-PPS is a project focused on studying the presence and behavior of pharmaceuticals and their byproducts in the environment. It is carried

out in conjunction with healthcare firms, treatment plants, and government organizations by the Swiss Federal Institute of Aquatic Science and Technology. This program can be opened at the following address: <http://eawag-bbd.ethz.ch/predict/>. It has open access web page, which uses chemical structure or Simplified Molecular Input Line Entry System (SMILES) format as input file to run the process. Insert the SMILES string encoding of your molecule in the appropriate checkbox. If otherwise, simply “draw” the molecular framework in the given molecular editor and then select “Write SMILES” to automatically produce the SMILES and click “continue” (Figure 11.3(a)). After a while, a representation of a portion of your compound’s expected biodegradative pathway appears (Figure 11.3(b)). A color scale indicates the aerobic likelihood of the



(a)



(b)

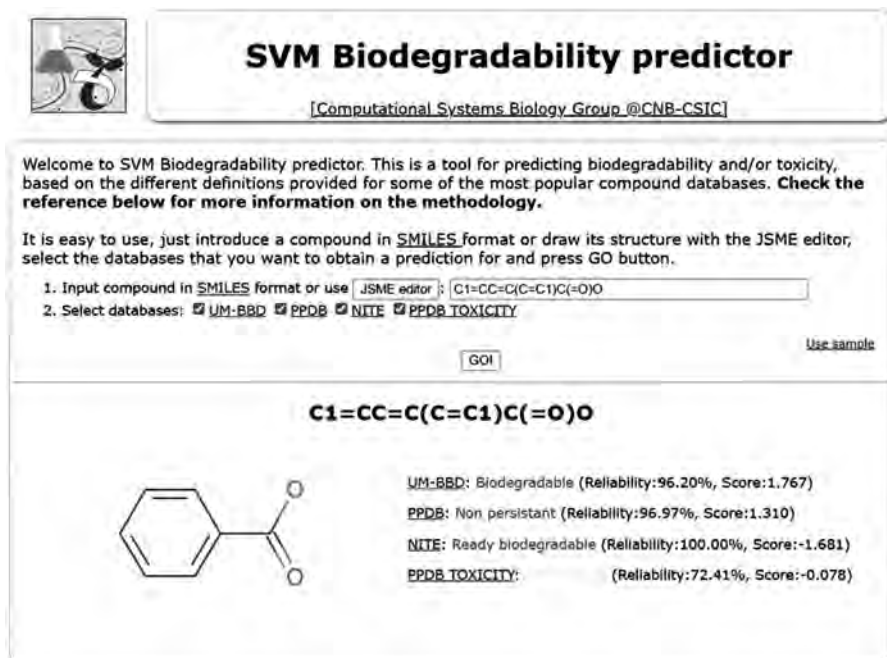
**FIGURE 11.3** Screenshots of the EAWAG-PPS system when forecasting the biodegradation paths for aniline. (a) The input form comprises a molecular editor to “draw” the structure of the input molecule. And (b) in the forecasted biodegradative routes, molecules existing in the EAWAG-PPS have a “Cpd” key to go to the corresponding pages of compound data and routes, while non-end compounds have a next key to recover the single downstream biodegradative paths starting with initial one.



various transformation processes in this illustration. The codes (such as bt0065) provide immediate links to UM-BBD pages that contain full material about the reaction's guidelines. The main objective of the EAWAG-PPS program is to create and test technologies that can remove pharmaceuticals and their byproducts from wastewater, and assess the potential hazards connected with their presence in the ecosystem. Additionally, the project aims to increase public awareness about the issue of pharmaceuticals in the environment and promote sustainable practices within the pharmaceutical industry [30]. The EAWAG-PPS project has had a noteworthy impact on our consideration of the environmental impact of pharmaceuticals. The project has headed to the expansion of new technologies and methodologies for removing pharmaceuticals from wastewater and has helped inform policy decisions regarding the regulation of pharmaceuticals in the environment.

### **11.5.2 SVM BIODEGRADABILITY PREDICTOR**

The SVM (Support Vector Machine) Biodegradability Predictor is a computational tool that uses machine-learning algorithms to foretell the biodegradability of chemical composites. This software can be visited at the given address: <http://csbg.cnb.csic.es/BiodegPred>. The system requires SMILES string compound information as an input. To execute the chosen forecasters on the input framework, click "Go." The outcome page includes a depiction of the organic arrangement restored from the input SMILES (to ensure its accuracy) as well as the results of the predictors used (Figure 11.4). The character of the predictions is shown using a color code (green as biodegradable/nontoxic; red as recalcitrant/toxic). The score of the predictor and its related dependability are also shown. The reliability standards assigned to respectively score were derived from a test set of substances with recognized fates, and they reflect the fraction of molecules in the test set with a given score or greater that were properly prophesied. This has helped researchers and industries in designing more sustainable products, falling the effect of chemicals on the habitat, and complying with environmental regulations. The SVM Biodegradability Predictor is based on a large database of biochemical structures and their corresponding biodegradability data. The tool uses this database to train a machine learning model that can predict the biodegradability of new chemical compounds based on their structural features [31].



**SVM Biodegradability predictor**  
[Computational Systems Biology Group @CNB-CSIC]

Welcome to SVM Biodegradability predictor. This is a tool for predicting biodegradability and/or toxicity, based on the different definitions provided for some of the most popular compound databases. **Check the reference below for more information on the methodology.**

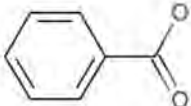
It is easy to use, just introduce a compound in SMILES format or draw its structure with the JSME editor, select the databases that you want to obtain a prediction for and press GO button.

1. Input compound in SMILES format or use JSME editor :

2. Select databases: ☒ UM-BBD ☒ PPDB ☒ NITE ☒ PPDB TOXICITY

[Use sample](#)

**C1=CC=C(C=C1)C(=O)O**



UM-BBD: Biodegradable (Reliability:96.20%, Score:1.767)

PPDB: Non persistent (Reliability:96.97%, Score:1.310)

NITE: Ready biodegradable (Reliability:100.00%, Score:-1.681)

PPDB TOXICITY: (Reliability:72.41%, Score:-0.078)

**FIGURE 11.4** Screenshot of the SVM biodegradability predictor being used to forecast benzoic acid's destiny in the environment.

### 11.5.3 CATALOGIC

CATALOGIC is a software suite of OASIS (<http://oasis.lmc.org/products/software/catalogic.aspx>) designed for the evaluation of ecological fate and eco-toxicity endpoints. In these models, the delineation of the biological degradation pathways for an input molecule is associated on a set of catabolic modifications that have been “weighted” with information from experiments on biodegradative destiny gathered from records and with other factors like the “biological oxygen demand.”

A database pertaining to the environmental destiny like abiotic and biotic degradation, bioaccumulation, and acute marine noxiousness of chemicals, the system has features for managing endpoint and metabolic data, running models in batches and individually, predicting the acute aquatic toxicity of particular metabolites, and providing interactive help. It is a tool used in ecological harm estimation to predict the behavior and potential impacts of

chemicals in the environment. It employs various models and algorithms to estimate important parameters such as biodegradability, bioaccumulation, and toxicity of chemicals. These predictions aid in evaluating the potential environmental fate of a substance such as its persistence, mobility, and potential to bioaccumulate in organisms [32, 33].

By integrating information on the physicochemical properties of a chemical with data on environmental conditions, CATALOGIC enables users to assess the environmental risk associated with its use. This software suite can be valuable in decision-making processes related to chemical management, regulatory compliance, and the development of safer and more sustainable products [34]. Overall, CATALOGIC serves as a useful tool for environmental risk assessors, researchers, and regulators to evaluate the environmental fate and ecotoxicity of chemicals, ultimately contributing to informed decision-making and environmental protection.

#### **11.5.4 BIOWIN**

Biowin is a computer-based model used for envisaging the biodegradability of carbon-based chemicals in the environment. It is a fragment of the US EPA's (Environmental Protection Agency) EPI Suite software, which is widely used in environmental risk assessment and chemical management. It is paid software having annual, perpetual and academic license. This program can be obtained from EnviroSim Associates Ltd., Canada, website (<https://envirosim.com/>) [35]. The combined triggered sludge/anaerobic digestion (AS/AD) model, is used in this program. In addition, BioWin provides three models to assess settling and methods of separation, namely a point separation model, an ideal separation model, and a flux based model. The Biowin model is based on a set of mathematical equations that contains numerous factors namely molecular structure of the organic, its solubility, and its partition coefficient. The model uses these parameters to predict the potential for biodegradation of the chemical in different environmental conditions [36].

The Biowin has been extensively validated against experimental data and has been shown to accurately predict the biodegradability of substrates/elements under diverse ecological surroundings. The use of the Biowin model has helped to identify chemicals that might pose a hazard to the environment and has contributed to the development of sustainable chemical management.

### **11.5.5 PBT PROFILER**

The PBT Profiler model is a web-based tool (part of the US Environmental Protection Agency (EPI) Suite software) used to estimate the persistence, bioaccumulation, and toxicity (PBT) of organic chemicals in the environment. It is a thorough and current database with details on a variety of organic substances and their potential ecological effects. The process considers investigational data from databases in the QSAR Toolbox and output from (Q)SAR models (<https://repository.qsartoolbox.org/Tools/List/Profilers>). In accordance with Annex XIII of REACH Regulation (EC) No 1907/2006 and ECHA Guidance on Data Needs and Chemical Safety evaluation Chapter R.11: PBT/vPvB evaluation, the fallouts of this examination is founded on solitary verge scores for P/vP, B/vB, and T (ENV).

The PBT Profiler is centered on a set of mathematical equations that take into account various parameters such as the chemical's molecular structure, its environmental fate, and its toxicological properties. The model uses these parameters to forecast the potential diligence, bioaccumulation, and harmfulness of the organic chemicals under diverse ecological circumstances. It is an important tool for measuring the potential environmental impact of organic chemicals. It has been extensively validated against experimental data and has been shown to accurately predict the PBT characteristics of a wide range of organic chemicals under different environmental conditions [37]. The use of the PBT Profiler has helped to identify organic chemicals that could pose a danger to the ecosystem and has contributed to the development of sustainable chemical management practices.

### **11.5.6 META-PC**

Meta-PC is a computer-based tool used for predicting the tenacity and long-term fortune of compounds in the ecosystem. It is a slice of the US EPA's (Environmental Protection Agency) EPI Suite software, which is widely utilized in environmental risk valuation and chemical treatment. It is a product of MultiCASE Inc. (<https://www.multicase.com/>). Data about 1,467 Mammalian Metabolism, 505 Aerobic Microbial Biotransformation, 344 Anaerobic Microbial Biodegradation, and 1,193 responses to photo degradation can be found in Meta-PC. The Meta-PC model is also subjected on a set of mathematical equations that take into account various parameters such as the chemical's molecular structure, its environmental fate, and its

toxicological properties. The model uses these parameters to predict the potential for the persistence of the chemical in different environmental conditions [38].

The Meta-PC model is an significant utensil for evaluating the potential ecological impact of substances. It has been extensively validated against experimental data and has been shown to accurately predict the persistence of a varied range of chemicals under diverse environmental conditions. The use of the Meta-PC model has helped to identify chemicals that might pose harm to the habitat and has contributed to the development of sustainable chemical management practices.

## **11.6 AI-BASED BIODEGRADATION DATABASES**

In addition to experimental testing and computer modeling, biodegradation database is a repository of knowledge regarding the capacity of different compounds to be broken down by biological processes. It often contains information on the rates, processes, and byproducts of various substances when they are subjected to biological agents like bacteria or enzymes. The databases are covered in subsections.

### **11.6.1 UNIVERSITY OF MINNESOTA BIOCATALYSIS/BIODEGRADATION (UMBB)**

The UMBB databank is an online resource that delivers information on bacterial biodegradation pathways for environmental contaminants. It is a comprehensive and up-to-date databank that comprises evidence on the biodegradation of a wide range of organic chemicals. It contains details on more than 1,300 chemicals, 900 enzymes, around 1,500 reactions, and about 543 entries for microorganisms. The UM-PPS now has 249 biotransformation principles that were resulting from reactions discovered in the UM-BBD and academic research. Data from the UM-BBD are becoming more widely available. The public chemical databases PubChem and ChemSpider now receive compound data from the UM-BBD. At ETH Zürich, a novel mirror website of the UM-BBD, UM-PPS and UM-BPT is being created to increase the speed and dependability of online access from any location. The information is organized in a user-friendly format that allows researchers and practitioners to quickly access relevant data [39]. It has been extensively

used in the sphere of ecological science to assess the potential environmental impact of chemicals and to identify microorganisms that can be used in bioremediation processes.

The databank is regularly modernized with new information and is widely regarded as a reliable and authoritative source of information on biodegradation pathways. Its use has contributed to the development of sustainable chemical management practices and has helped to decrease the ecological influence of organic chemicals. Overall, the UMBBD is an important resource for environmental scientists and practitioners and vital part in advancing the field of environmental science.

### **11.6.2 PLASTICS MICROBIAL BIODEGRADATION DATABASE**

The Plastics Microbial Biodegradation Database is an internet-based tool which gives details on microbiological plastic decomposition. The database includes information on the microbes and enzymes responsible for biotransformation of plastics, as well as the metabolic intermediates that are produced during the degradation process. In this repository, 79 genes and 949 interactions between microbes and plastics were carefully compiled and validated through literature searches. Additionally, through the TrEMBL component of the UniProt database, above 8,000 automatically annotated enzyme successions were taken as may be incorporated in the biodegradation of plastics. It is a comprehensive and up-to-date databank that holds information on the biodegradation of a varied range of plastics under different ecological situations. It also offers data on the physical and chemical possessions of different types of plastics that influence their biodegradability [40]. The Plastics Microbial Biodegradation Database is a valuable tool for environmental risk assessment and plastic waste management. It has been extensively used to validate the potential environmental impact of plastic waste and to identify microorganisms that can be used in bioremediation processes.

The databank is regularly re-equipped with innovative information and is widely regarded as a reliable and authoritative source of information on the biodegradation of plastics. Its use has contributed to the expansion of sustainable plastic leftover management practices and has helped to diminish the plastic waste. Overall, the Plastics Microbial Biodegradation Database is an important resource for environmental scientists and practitioners for plastic waste management.

### **11.6.3 ONDB (ORGANONITROGEN DEGRADATION DATABASE)**

The Organonitrogen Degradation Database (ONDB) is a web-based service that contains knowledge regarding the microbe-mediated breakdown of organic nitrogen-containing molecules in the environment. It is a comprehensive and up-to-date database that contains information on the biodegradation of an extensive range of organonitrogen molecules. The ONDB was created to compile data on the price, chemistry, and biodegradability of frequently used organonitrogen compounds. It offers information about the enzymatic processes, microbial processes, and routes involved in the breakdown of organonitrogen compounds such as urea, methylenediurea. The database includes information on the bacteria and enzymes allied in the biotransformation of organonitrogen compounds, as well as the metabolic intermediates that are produced during the degradation process. It also offers physical and chemical data of different types of organonitrogen compounds that influence their biodegradability [41]. The ONDB is a helpful tool for ecological toxicant evaluation and biochemical management. It has been extensively used in the sphere of ecological science to measure the latent ecological impression of organonitrogen compounds and to identify microorganisms that can be employed in bioremediation techniques.

The databank is regularly upgraded with innovative information and is widely regarded as a reliable and authoritative source of information on the biodegradation of organonitrogen compounds. Its use has contributed to the development of sustainable chemical management practices and has helped to lessen the conservatoire impact of organic nitrogen-containing compounds [42].

### **11.6.4 BIONEMO (BIODEGRADATION NETWORK-MOLECULAR BIOLOGY DATABASE)**

Bionemo is a comprehensive and up-to-date online resource that offers evidence on the microbial biodegradation of organic molecules. Bionemo store manually organized data on genes and proteins essential in the breakdown. The database contains details on the sequence, domains, and structures of proteins as well as the sequence, regulatory fundamentals, and transcription components of genes, whenever this is possible. Consuming data from available studies, Bionemo was created by manually linking sequence databank entries to biodegradation reactions. For biodegradation genes, data on transcription elements and their control were also taken from

the pre-evaluated works and connected to the basic biochemical network. Bionemo now has 324 reactions' sequencing data as well as data on above 100 promoters and 100 transcription factors' transcription regulation. The details in the Bionemo database is reachable via a website hosted on a server, and a PostgreSQL dumps of every record is also obtainable for downloading. It is a prevailing tool that combines molecular biology and bioinformatics to predict the biotransformation pathways of carbon-based compounds in the ecosystem. The database contains information on the enzymes, genes, and biochemical pathways connected in the biodegradation of a comprehensive variety of carbon-based compounds, covering contaminants along with natural substances [43]. It has been extensively used in the sphere of ecological science to evaluate the possible ecological impact of organic molecules and to identify microorganisms that can be employed in bioremediation techniques.

The Bionemo database can be accessed by researchers, environmental scientists, and practitioners who need information on the biodegradation pathways of organic compounds. Its use has contributed to the development of sustainable chemical management practices and has helped to lessen the environmental impact of carbon-based compounds. Overall, Bionemo is an important resource for the field of habitat science and plays a crucial role in advancing our acquaintance of the microbial biodegradation of carbon-based compounds in the atmosphere [44].

### **11.6.5 METACYC**

MetaCyc is a comprehensive and up-to-date databank of biochemical pathways metabolic from all spheres of life and enzymes found in living organisms. It contains information on over 3,105 metabolic pathways, including pathways involved in energy generation, biosynthesis, degradation, and detoxification. Primary and secondary metabolic trails, in addition to related metabolites, processes, enzymes, and genes, are all included in the MetaCyc database. MetaCyc aims to compile a comprehensive catalogue of all known metabolic processes. It is regularly updated with new information and is widely regarded as a reliable and authoritative source of data on biochemical pathways and enzymes. It is used by investigators in a massive fields, including biochemistry, genetics, and biotechnology. One of the key features of MetaCyc is its ability to utilize an enzyme database to support metabolic engineering and metabolomics research is aided by a metabolite database. This information can be used to identify potential drug targets,



develop new biocatalysts, and understand the metabolic processes involved in diseases [45].

MetaCyc is an imperative utensil for bioinformatics and computational biology. It is used to model and simulate metabolic pathways, predicts the effects of genetic mutations on metabolism, and design new metabolic engineering strategies. Overall, MetaCyc is an essential resource for researchers and practitioners in the fields of biochemistry, genetics, and biotechnology. Its use has contributed to the development of new drugs, biocatalysts, and metabolic engineering strategies, and has advanced our understanding of biochemical pathways and bacteria in living organisms [46].

#### **11.6.6 OXDBASE (BIODEGRADATIVE OXYGENASES DATABASE)**

OxDBase is a comprehensive and up-to-date database that provides information on biodegradative oxygenases. These are enzymes that play a crucial part in the dilapidation of an eclectic variety of organic compounds in the atmosphere. The database contains information on the structure, function, and substrate specificity of biodegradative oxygenases. It covers an extensive choice of oxygenases, counting monooxygenases, dioxygenases, and peroxidases [47]. There are two distinct search engines, one for the monooxygenases and the other for the dioxygenases directory. A single entry for a specific enzyme includes information about the enzyme's general name, reaction, family, subfamily, structure, gene link, and literature reference. The records also include links to numerous additional databases, namely KEGG, ENZYME, BRENDA, and UM-BBD, which provide extensive additional information. Currently, the repository has information pertaining to more than 235 oxygenases.

One of the key features of OxDBase is its ability to provide information on the genes that encode biodegradative oxygenases. This information can be used to identify microorganisms that are skilled of breaking specific organic compounds and to design new bioremediation strategies. OxDBase is an important tool for ecological danger estimation and compound management. It is used to measure the possible ecological impression of organic molecules and to identify microorganisms that can be used in bioremediation processes. Overall, OxDBase is an essential resource for researchers and practitioners in the fields of environmental science, biotechnology, and pharmacology. Its use has contributed to the expansion of new bioremediation approaches, the identification of potential drug targets, and our understanding of the role of biodegradative oxygenases in the habitat.

### **11.6.7 AROMADEG**

The Aromatic Hydrocarbon Degradation Database (AHDD) is a novel, complete online resource that delivers data on the degradation of aromatic hydrocarbons by aerobic microorganisms. Aromatic hydrocarbons are a class of carbon-based molecules that contain one or more aromatic rings. The AHDD database contains information on the genes, enzymes, and metabolic paths involved in the degradation of an enormous variety of aromatic hydrocarbons, including pollutants and natural compounds [48]. AromaDeg enables inquiry and facts extraction of new metagenomic or metatranscriptomic, genomic groups of data and is founded on phylogenetic studies of the protein sequences of important catabolic protein groups and of enzymes with known functions. The database is regularly restructured with novel information and is widely regarded as a reliable and authoritative source of material on the microbial dilapidation of aromatic hydrocarbons. It is a valuable tool for environmental hazard valuation and compound management, which has been extensively used in the sphere of ecological science to assess the potential environmental impact of aromatic hydrocarbons and to identify microorganisms that can be used in bioremediation processes. The AHDD database can be accessed by researchers, environmental scientists, and practitioners who need information on the degradation pathways of aromatic hydrocarbons. Its use has contributed to the development of sustainable chemical management practices and has helped to reduce the environmental impact of aromatic hydrocarbons [49]. AromaDeg uses phylogenetic tree building and the amalgamation of experimental findings to produce greater accuracy interpretations of novel scientific information pertaining to aerobic aromatic breakdown routes, addressing the shortcomings of homology-based protein activity forecasting. Overall, the AHDD database is an vital resource for the arena of environmental discipline and plays a crucial role in advancing our knowledge of the contagious degradation of aromatic hydrocarbons in the environment.

## **11.7 APPLICATION OF AI IN CHEMICAL SCIENCES**

Artificial intelligence (AI) has evolved into a necessary instrument in the discipline of chemical sciences, providing extraordinary forecasts for analysis and revolution. Artificial intelligence (AI) has become a necessary instrument in the field of chemical sciences, providing extraordinary prospects for study and development. AI has made tremendous progress in drug development,

synthesis of chemicals, design of materials, and process management through the integration of machine learning, large data analysis, and predictive modeling.

AI algorithms in drug discovery can evaluate massive volumes of chemical and biological data to identify possible drug targets and forecast the efficacy of novel drug candidates. This method allows researchers to speed up the drug development process while lowering the costs involved with clinical trials. By predicting reaction outcomes and adjusting reaction circumstances, AI may also optimize chemical reactions. This can result in considerable gains in process efficiency, as well as a reduction in waste and energy usage.

## **11.8 LIMITATIONS**

Using AI technologies to predict biodegradation has its limits. The complexity of environmental factors that can affect biodegradability, the lack of measurement standardization, the variability of microbial populations, and the constrained application of AI models for predicting biodegradability are a few of the main limitations [50]. For AI-based models to predict biodegradation with greater precision and applicability, it is critical to address these limitations.

## **11.9 CONCLUSIONS**

In conclusion, the prediction of biodegradation is a complex task that requires the analysis of various environmental factors and microbial interactions. Predicting biodegradability can aid in the identification of dangerous substances and assist the creation of safer alternatives that are less persistent and toxic. AI techniques such as ML, ANN, fuzzy logic, genetic algorithms, and expert systems have shown great promise in predicting biodegradability and assessing the potential threats allied with the discharge of constituents into the habitat. The use of AI-based models for predicting biodegradation has several advantages over traditional methods, including the ability to analyze large datasets, handle complicated interactions, and make real-time forecasts. As AI tools/methods continue to develop, they are probable to show a progressively vital part in predicting biodegradation and ensuring the sustainability of our environment.

## KEYWORDS

- **artificial intelligence**
- **artificial neural network**
- **biodegradation**
- **bioremediation chemical structure**
- **deep learning**
- **machine learning**
- **quantitative structure-activity relationship**

## REFERENCES

1. Singh, A. K., Bilal, M., Iqbal, H. M., & Raj, A. (2021). Trends in predictive biodegradation for sustainable mitigation of environmental pollutants: Recent progress and future outlook. *Science of The Total Environment*, 770, 144561. <https://doi.org/10.1016/j.scitotenv.2020.144561>.
2. Eskander, S., & Saleh, H. E. D. (2017). Biodegradation: Process mechanism. *Environmental Science & Engineering*, 8(8), 1–31.
3. Yan, Y., & Wang, X. (2022). Integrated energy view of wastewater treatment: A potential of electrochemical biodegradation. *Frontiers of Environmental Science & Engineering*, 16, 52. <https://doi.org/10.1007/s11783-021-1486-3>.
4. Sharma, S., Saxena, S., Mudgil, B., & Vats, S. (2022). Advances in biodegradation and bioremediation of environmental pesticide contamination. In *Biological Approaches to Controlling Pollutants* (pp. 79–106). Woodhead Publishing. <https://doi.org/10.1016/B978-0-12-824316-9.00009-4>.
5. Moussa, Z., Darwish, D. B., Aldahe, S. S., & Saber, W. I. (2021). Innovative artificial-intelligence-based approach for the biodegradation of feather keratin by *Bacillus paramycoides*, and cytotoxicity of the resulting amino acids. *Frontiers in Microbiology*, 12, 731262. <https://doi.org/10.3389/fmicb.2021.731262>.
6. Dimitrov, S., Pavlov, T., Nedelcheva, D., Reuschenbach, P., Silvani, M., Bias, R., Comber, M., Low, L., Lee, C., Parkerton, T., & Mekenyan, O. (2007). A kinetic model for predicting biodegradation. *SAR and QSAR in Environmental Research*, 18(5–6), 443–457. <https://doi.org/10.1080/10629360701429027>.
7. Singh, A., & Ward, O. P. (Eds.). (2004). *Biodegradation and Bioremediation* (Vol. 2). Springer Science & Business Media.
8. Acharya, K., Werner, D., Dolfing, J., Barycki, M., Meynet, P., Mrozik, W., Komolafe, O., Puzyn, T., & Davenport, R. J. (2019). A quantitative structure-biodegradation relationship (QSBR) approach to predict biodegradation rates of aromatic chemicals. *Water Research*, 157, 181–190. <https://doi.org/10.1016/j.watres.2019.03.086>.

9. Wackett, L. P. (2004). Prediction of microbial biodegradation: An annotated selection of world wide web sites relevant to the topics in environmental microbiology. *Environmental Microbiology*, 6(3), 313–313. <https://doi.org/10.1111/j.1462-2920.2004.00606.x>.
10. Zambrano, M. C., Pawlak, J. J., Daystar, J., Ankeny, M., Goller, C. C., & Venditti, R. A. (2020). Aerobic biodegradation in freshwater and marine environments of textile microfibers generated in clothes laundering: Effects of cellulose and polyester-based microfibers on the microbiome. *Marine Pollution Bulletin*, 151, 110826. <https://doi.org/10.1016/j.marpolbul.2019.110826>.
11. Tomei, M. C., Mosca Angelucci, D., Clagnan, E., & Brusetti, L. (2021). Anaerobic biodegradation of phenol in wastewater treatment: Achievements and limits. *Applied Microbiology and Biotechnology*, 105, 2195–2224. <https://doi.org/10.1007/s00253-021-11182-5>.
12. Jaiswal, K. K., Dutta, S., Banerjee, I., Pohrmen, C. B., Singh, R. K., Das, H. T., Dubey, S., & Kumar, V. (2022). Impact of aquatic microplastics and nanoplastics pollution on ecological systems and sustainable remediation strategies of biodegradation and photodegradation. *Science of The Total Environment*, 806, 151358. <https://doi.org/10.1016/j.scitotenv.2021.151358>.
13. Sharma, I. (2020). Bioremediation techniques for polluted environment: Concept, advantages, limitations, and prospects. In *Trace Metals in the Environment—New Approaches and Recent Advances*. IntechOpen.
14. Meena, A. L., Karwal, M., Dutta, D., & Mishra, R. P. (2021). Composting: Phases and factors responsible for efficient and improved composting. *Agriculture and Food: e-Newsletter*, 1, 85–90.
15. Cai, Q., Shi, C., Yuan, S., & Tong, M. (2023). Integrated anaerobic–aerobic biodegradation of mixed chlorinated solvents by electrolysis coupled with groundwater circulation in a simulated aquifer. *Environmental Science and Pollution Research*, 30, 31188–31201. <https://doi.org/10.1007/s11356-022-24377-8>.
16. Yan, R., Wang, Y., Li, J., Wang, X., & Wang, Y. (2022). Determination of the lower limits of antibiotic biodegradation and the fate of antibiotic resistant genes in activated sludge: Both nitrifying bacteria and heterotrophic bacteria matter. *Journal of Hazardous Materials*, 425, 127764. <https://doi.org/10.1016/j.jhazmat.2021.127764>.
17. Torgbo, S., & Sukyai, P. (2020). Biodegradation and thermal stability of bacterial cellulose as biomaterial: The relevance in biomedical applications. *Polymer Degradation and Stability*, 179, 109232. <https://doi.org/10.1016/j.polymdegradstab.2020.109232>.
18. Pischedda, A., Tosin, M., & Degli-Innocenti, F. (2019). Biodegradation of plastics in soil: The effect of temperature. *Polymer Degradation and Stability*, 170, 109017. <https://doi.org/10.1016/j.polymdegradstab.2019.109017>.
19. Brintha, V. P., Rekha, R., Nandhini, J., Sreekaarthick, N., Ishwaryaa, B., & Rahul, R. (2020). Automatic classification of solid waste using deep learning. In L. Kumar, L. Jayashree, & R. Manimegalai (Eds.), *Proceedings of International Conference on Artificial Intelligence, Smart Grid and Smart City Applications* (pp. 83). Springer. [https://doi.org/10.1007/978-3-030-24051-6\\_83](https://doi.org/10.1007/978-3-030-24051-6_83).
20. Astuto, M. C., Di Nicola, M. R., Tarazona, J. V., Rortais, A., Devos, Y., Liem, A. D., Kass, G. E., Bastaki, M., Schoonjans, R., Maggiore, A., & Charles, S. (2022). In silico methods for environmental risk assessment: Principles, tiered approaches, applications, and future perspectives. In E. Benfenati (Ed.), *In Silico Methods for Predicting Drug Toxicity* (Vol. 2425, pp. 23). Humana. [https://doi.org/10.1007/978-1-0716-1960-5\\_23](https://doi.org/10.1007/978-1-0716-1960-5_23).

21. Pazos, F., & de Lorenzo, V. (2015). Biodegradation prediction tools. In T. McGenity, K. Timmis, & B. Nogales Fernández (Eds.), *Hydrocarbon and Lipid Microbiology Protocols* (pp. 651–665). Springer. [https://doi.org/10.1007/8623\\_2015\\_87](https://doi.org/10.1007/8623_2015_87).
22. Ahmad, Z., Zhong, H., Mosavi, A., Sadiq, M., Saleem, H., Khalid, A., Mahmood, S., & Nabipour, N. (2020). Machine learning modeling of aerobic biodegradation for azo dyes and hexavalent chromium. *Mathematics*, 8(6), 913. <https://doi.org/10.3390/math8060913>.
23. Rashtbari, S., & Dehghan, G. (2021). Biodegradation of malachite green by a novel laccase-mimicking multicopper BSA-Cu complex: Performance optimization, intermediates identification and artificial neural network modeling. *Journal of Hazardous Materials*, 406, 124340. <https://doi.org/10.1016/j.jhazmat.2020.124340>.
24. Priya, P. P., Jenit, A., & Sharma, N. K. (2023). Predictive biodegradation of multiple toxic pollutants in bioreactors treating real wastewater using ANN and GP. In *IOP Conference Series: Earth and Environmental Science* (Vol. 1130, No. 1, p. 012040). IOP Publishing. <https://doi.org/10.1088/1755-1315/1130/1/012040>.
25. Ruan, J., Chen, X., Huang, M., & Zhang, T. (2017). Application of fuzzy neural networks for modeling of biodegradation and biogas production in a full-scale internal circulation anaerobic reactor. *Journal of Environmental Science and Health, Part A*, 52(1), 7–14. <https://doi.org/10.1080/10934529.2016.1221216>.
26. Huang, M., Zhang, T., Ruan, J., & Chen, X. (2017). A new efficient hybrid intelligent model for biodegradation process of DMP with fuzzy wavelet neural networks. *Scientific Reports*, 7(1), 41239. <https://doi.org/10.1038/srep41239>.
27. Schomburg, I., Jeske, L., Ulbrich, M., Placzek, S., Chang, A., & Schomburg, D. (2017). The BRENDA enzyme information system—From a database to an expert system. *Journal of Biotechnology*, 261, 194–206. <https://doi.org/10.1016/j.jbiotec.2017.04.020>.
28. Arora, P. K., & Shi, W. (2010). Tools of bioinformatics in biodegradation. *Reviews in Environmental Science and Biotechnology*, 9, 211–213. <https://doi.org/10.1007/s11157-010-9211-x>.
29. Zhan, Z., Li, L., Tian, S., Zhen, X., & Li, Y. (2017). Prediction of chemical biodegradability using computational methods. *Molecular Simulation*, 43(13–16), 1277–290. <https://doi.org/10.1080/08927022.2017.1328556>.
30. Chibwe, L., Titaley, I. A., Hoh, E., & Simonich, S. L. M. (2017). Integrated framework for identifying toxic transformation products in complex environmental mixtures. *Environmental Science & Technology Letters*, 4(2), 32–43. <https://doi.org/10.1021/acs.estlett.6b00455>.
31. Garcia-Martin, J. A., Chavarría, M., de Lorenzo, V., & Pazos, F. (2020). Concomitant prediction of environmental fate and toxicity of chemical compounds. *Biology Methods and Protocols*, 5(1), bpaa025. <https://doi.org/10.1093/biomethods/bpaa025>.
32. Dimitrov, S., Kamenska, V., Walker, J. D., Windle, W., Purdy, R., Lewis, M., & Mekenyan, O. (2004). Predicting the biodegradation products of perfluorinated chemicals using CATABOL. *SAR and QSAR in Environmental Research*, 15(1), 69–82. <https://doi.org/10.1080/1062936032000169688>.
33. Dimitrov, S., Nedelcheva, D., Dimitrova, N., & Mekenyan, O. (2010). Development of a biodegradation model for the prediction of metabolites in soil. *Science of the Total Environment*, 408(18), 3811–3816. <https://doi.org/10.1016/j.scitotenv.2010.02.008>.
34. Sakuratani, Y., Yamada, J., Kasai, K., Noguchi, Y., & Nishihara, T. (2005). External validation of the biodegradability prediction model CATABOL using data sets of existing

- and new chemicals under the Japanese Chemical Substances Control Law. *SAR and QSAR in Environmental Research*, 16(5), 403–431. <https://doi.org/10.1080/10659360500320289>.
35. Tang, W., Li, Y., Yu, Y., Wang, Z., Xu, T., Chen, J., Lin, J., & Li, X. (2020). Development of models predicting biodegradation rate rating with multiple linear regression and support vector machine algorithms. *Chemosphere*, 253, 126666. <https://doi.org/10.1016/j.chemosphere.2020.126666>.
  36. Balakrishnan, A., Kanchinadham, S. B. K., & Kalyanaraman, C. (2020). Assessment on biodegradability prediction of tannery wastewater using EPI Suite BIOWIN model. *Environmental Monitoring and Assessment*, 192, 732. <https://doi.org/10.1007/s10661-020-08661-z>.
  37. Cassani, S., & Gramatica, P. (2015). Identification of potential PBT behavior of personal care products by structural approaches. *Sustainable Chemistry and Pharmacy*, 1, 19–27. <https://doi.org/10.1016/j.scp.2015.10.002>.
  38. Sedykh, A., Saiakhov, R., & Klopman, G. (2001). META V: A model of photodegradation for the prediction of photoproducts of chemicals under natural-like conditions. *Chemosphere*, 45(6–7), 971–981. [https://doi.org/10.1016/S0045-6535\(01\)00007-8](https://doi.org/10.1016/S0045-6535(01)00007-8).
  39. Seema, K., Kakoli, D., & Veena, G. (2014). Computational analysis of biodegradation pathways for chlorpyrifos using EAWAG-Biocatalysis/Biodegradation Database Pathway Prediction System. *World Journal of Pharmaceutical Research*, 3(9), 1071–1082.
  40. Gan, Z., & Zhang, H. (2019). PMBD: A comprehensive plastics microbial biodegradation database. *Database*. <https://doi.org/10.1093/database/baz119>.
  41. Robinson, S. L., Biernath, T., Rosenthal, C., Young, D., Wackett, L. P., & Martinez-Vaz, B. M. (2021). Development of the organonitrogen biodegradation database: Teaching bioinformatics and collaborative skills to undergraduates during a pandemic. *Journal of Microbiology & Biology Education*, 22(1), ev22i1–2351. <https://doi.org/10.1128/jmbe.v22i1.2351>.
  42. Arora, P. K., Kumar, A., Srivastava, A., Garg, S. K., & Singh, V. P. (2022). Current bioinformatics tools for biodegradation of xenobiotic compounds. *Frontiers in Environmental Science*, 1499. <https://doi.org/10.3389/fenvs.2022.980284>.
  43. Arora, P. K., & Bae, H. (2014). Integration of bioinformatics to biodegradation. *Biological Procedure Online*, 16, 8. <https://doi.org/10.1186/1480-9222-16-8>.
  44. Nigam, S., & Sinha, S. (2023). Bioinformatics and its contribution to bioremediation and genomics: Recent trends and advancement. In *Genomics Approach to Bioremediation: Principles, Tools, and Emerging Technologies* (pp. 455–466). <https://doi.org/10.1002/9781119852131.ch24>.
  45. Caspi, R., Billington, R., Keseler, I. M., Kothari, A., Krummenacker, M., Midford, P. E., Ong, W. K., Paley, S., Subhraveti, P., & Karp, P. D. (2020). The MetaCyc database of metabolic pathways and enzymes—A 2019 update. *Nucleic Acids Research*, 48(D1), D445–D453. <https://doi.org/10.1093/nar/gkz862>.
  46. Karp, P. D., Riley, M., Paley, S. M., & Pellegrini-Toole, A. (2002). The MetaCyc database. *Nucleic Acids Research*, 30(1), 59–61. <https://doi.org/10.1093/nar/30.1.59>.
  47. Arora, P. K., Kumar, M., Chauhan, A., Raghava, G. P. S., & Jain, R. K. (2009). OxDBase: A database of oxygenase's involved in biodegradation. *BMC Research Notes*, 2, 67. <https://doi.org/10.1186/1756-0500-2-67>.
  48. Duarte, M., Jauregui, R., Vilchez-Vargas, R., Junca, H., & Pieper, D. H. (2014). AromaDeg, a novel database for phylogenomic of aerobic bacterial degradation of aromatics. *Database*. <https://doi.org/10.1093/database/bau118>.

49. Bargiela, R., Gertler, C., Magagnini, M., Mapelli, F., Chen, J., Daffonchio, D., Golyshin, P. N., & Ferrer, M. (2015). Degradation network reconstruction in uric acid and ammonium amendments in oil-degrading marine microcosms guided by metagenomic data. *Frontiers in Microbiology*, 6, 1270. <https://doi.org/10.3389/fmicb.2015.01270>.
50. Dick, C., Rey, S., Boschung, A., Miffon, F., & Seyfried, M. (2016). Current limitations of biodegradation screening tests and prediction of biodegradability: A focus on fragrance substances. *Environmental Technology & Innovation*, 5, 208–224. <https://doi.org/10.1016/j.eti.2016.03.002>.





# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## **PART III**

### **Application of Expert Systems and AI in Fault Diagnosis and Structure Representation**



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## CHAPTER 12

---

# Exploring the Range of Knowledge-Based Prediction Applications in Chemistry

ROHINI N. SHELKE,<sup>1</sup> LAXMI G. KATHAWATE,<sup>1</sup>  
DATTATRAYA N. PANSARE,<sup>2</sup> ANIKET P. SARKATE,<sup>3</sup> AJIT K. DHAS,<sup>2</sup>  
PRAVIN N. CHAVAN,<sup>4</sup> SHAILEE V. TIWARI,<sup>5</sup> DEEPAK K. LOKWANI,<sup>6</sup>  
and SHIVRAJ N. MAWALE<sup>7</sup>

<sup>1</sup>*Department of Chemistry, Radhabai Kale Mahila Mahavidyalaya, Maharashtra, India*

<sup>2</sup>*Department of Chemistry, Deogiri College, Aurangabad, Maharashtra, India*

<sup>3</sup>*Department of Chemical Technology, Dr. Babasaheb Ambedkar Marathwada University, Aurangabad, Maharashtra, India*

<sup>4</sup>*Department of Chemistry, Doshi Vakil Arts College and G. C. U. B. Science & Commerce College, Goregaon, Raigad, Maharashtra, India*

<sup>5</sup>*Department of Pharmaceutical Chemistry, Shri Ramkrishna Paramhans College of Pharmacy, Hasnapur, Parbhani, Maharashtra, India*

<sup>6</sup>*Rajarashi Shahu College of Pharmacy, Buldana, Maharashtra, India*

<sup>7</sup>*Department of Pharmaceutical Chemistry, AISSMS College of Pharmacy, Pune, Maharashtra, India*

---

## ABSTRACT

Predicting the outcomes of organic transformations is a vital and difficult undertaking in the field of molecular synthesis. The fusion of machine learning and chemical expertise presents a distinctive and potent approach for generating predictions in synthesis. This comprehensive analysis delves into

---

Artificial Intelligence for Chemical Sciences: Concepts, Models, and Applications. Shrikaant Kulkarni, Shashikant Bhandari, Dushyant Varshney, & P. William (Eds.)

© 2025 Apple Academic Press, Inc. Co-published with CRC Press (Taylor & Francis)

the most recent embedding techniques and model designs that facilitate the development of machine learning models capable of reliably predicting yield and selectivity in molecular synthesis. Recent advancements in the utilization of machine learning (ML) techniques in chemistry have showcased the potential for data-driven forecasting of synthesis efficiency. The integration of digitization and ML modeling plays a pivotal role in fully harnessing experimental data's potential and accurately predicting performance and selectivity. Multiple studies have emphasized the importance of integrating chemical knowledge into ML models, enhancing their capacity to make predictions that surpass human capabilities. This succinct analysis provides an overview of state-of-the-art techniques and model designs in forecasting synthetic presentations, with a focus on the effective integration of chemical knowledge into machine learning as of June 2022. By incorporating strategies from organic synthesis and chemical information, our objective is to furnish chemists with a roadmap and inspiration for digitizing and automating organic chemistry principles. Chemists rely on their domain expertise to predict reaction efficiencies, considering factors such as reactant properties, molecular-level reaction mechanisms, optimal steps for rates and selectivity, and the quantum chemical basis of desired performance. This knowledge significantly enhances prediction accuracy, but remains a daunting task even for seasoned experts in synthesis and catalysis. Computational chemistry, which encompasses chemistry-based software, has yielded various applications, including chemical design, automated reaction synthesis, analysis of spectral data, and molecular docking.

## **12.1 INTRODUCTIONS**

This subject delves into the range of knowledge-based prediction applications in the field of chemistry, with a specific focus on the utilization of artificial intelligence (AI) models. AI has gained recognition as a field within computer science that aims to develop software capable of intelligent computations akin to those performed by the human brain. It encompasses a range of methods, tools, and systems designed to simulate human knowledge acquisition, logical reasoning, and problem-solving abilities. The development of AI can be broadly classified into two main types. The first category involves methods and systems that leverage human expertise and decision-making rules, such as expert systems. The second category encompasses approaches that seek to replicate the functioning of the human brain, such as artificial neural networks (ANNs). Expert systems operate under the premise

that logical reasoning can be organized by gathering logical propositions and executing logical transformations on them. These systems offer distinct advantages in medical analysis and diagnostic problem-solving [1, 2], as they provide guidance for prediction and decision-making under various environmental conditions. Machine learning algorithms are frequently utilized in computer-aided drug discovery [3–5]. The considerable attention in recent times due to its capacity to automatically extract features from input data and capture intricate input-output relationships [6, 7]. In the field of drug discovery, deep learning has experienced a renewed interest, leading to the emergence of innovative modeling methodologies and applications [8–12].

Over half a century ago, the QSAR/QSPR modeling field first emerged [13]. Its significant impact on drug discovery is evident from successful predictions of biological activity and considerations such as ADMET [14–17]. This discourse encompasses various subjects, including quantitative structure-activity and property relationships, structure-based modeling, also biochemical production calculation. Currently, there is an increasing focus on deep learning applications and the future potential of AI in drug discovery. AI has greatly propelled computer-assisted medicine detection. AI is actively utilized in fully automated synthesis planning within the field of chemistry, employing tools like LHASA for retrosynthetic planning and leveraging reaction templates and sub-molecular patterns to represent changes in atomic connectivity [18]. Additionally, the use of AI technology to enhance toxicity prediction models is an emerging concept, showing promise in achieving scientific consensus and facilitating practical applications.

## **12.2 OVERVIEWS OF POTENTIAL APPLICATIONS**

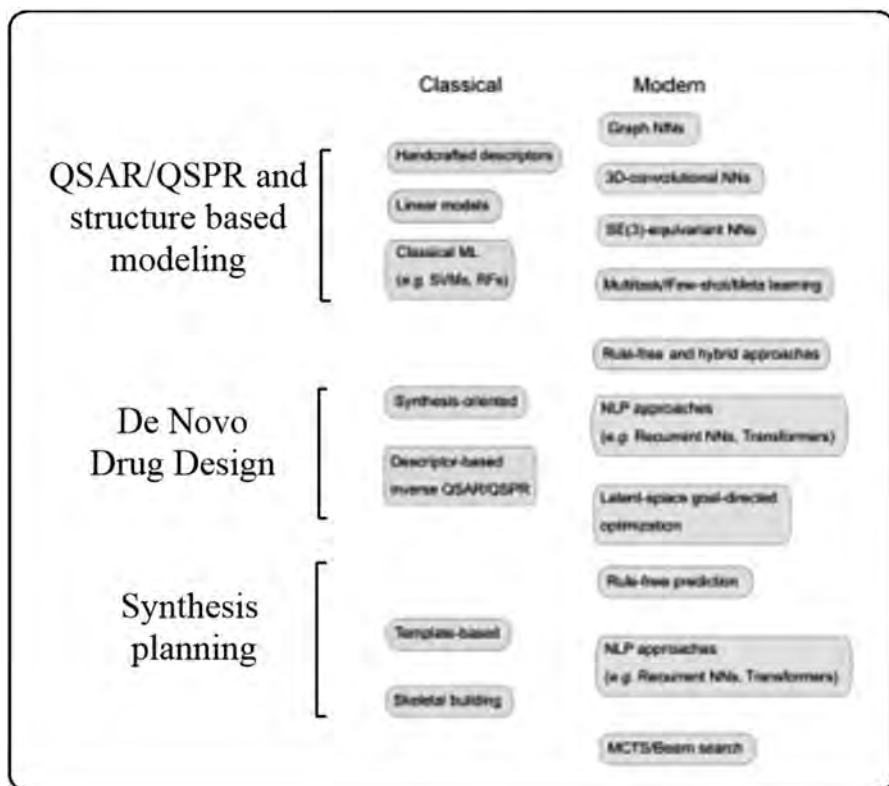
The ANN approach is stimulated by the operative of the hominoid mind in processing data, making it distinct from traditional statistical methods. It offers promising modeling techniques, particularly for datasets with non-linear relationships commonly encountered in pharmaceutical advancements. Additionally, they can utilize multiple training algorithms. Artificial neural networks do not demand knowledge of the data source for model design, but they do require extensive training sets due to their approximate weights. Moreover, ANNs can integrate experimental and literature-based information to solve problems. The field of pharmaceutical research is experiencing growth in the utilization of ANNs [19–23]. In the medical field, monitoring collective networks through traditional reaction surface methodology (RSM) can be applied. Unconfirmed feature-removing

linkages provide an alternative to principal component analysis (PCA) by mapping input datasets onto a two-dimensional space. Non-adaptive web record datasets can restructure their patterns in the presence of noisy designs, enabling their application in pattern recognition tasks. The potential applications of ANN processes in pharmacological science are extensive, ranging from interpretation data collection (quality control in medicinal research), biological target identification (molecular targets), and dosage form optimization in developed processes, to *in vivo* correlation through biotechnology and clinical pharmacy. ANN, or Artificial Neural Network, is utilized for pattern recognition and analysis of analytical data, optimization of pharmaceutical production, quantitative structure-property relationship (QSPR) and pharmacokinetics, protein structure and function prediction, and molecular modeling. Artificial intelligence is also employed in molecular model-based and ligand-based virtual screening, physicochemical and ADMET property estimation, as well as drug replication. It is used in automated synthesis planning, retrosynthesis, and molecular simulations. The impact of minor particles on the action of therapeutic proteins, the identification of radiographic images related to diseases, although the success rates may vary [24]. Neural networks in machine learning provide an optimal approach for biomedical and biological research due to their ability to process large datasets. In biomedical applications, artificial intelligence and computer-based modeling are essential for drawing conclusions and gaining insights beyond human capabilities.

Currently, the rapid generation of biomedical datasets using high-speed data throughput technologies has opened up opportunities to revolutionize biotechnology and pharmacy through the use of device education. Drug toxicity refers to harmful effects on living organisms, such as cells, caused by the actions or metabolic processes of certain compounds [25].

### **12.3 THE REVIEW EXAMINES THE USE OF MOLECULAR EMBEDDING TECHNIQUES AND THEIR APPLICATIONS IN PREDICTING REACTION PERFORMANCE**

It explores the advancements made by deep learning methods in cheminformatics, surpassing traditional approaches in various aspects. Figure 12.1 visually represents these concepts. Furthermore, the chapter emphasizes the existing limits of (AI) in one-to-one of these fields then speculates on its potential impact on the future of computer-aided medicinal innovation [26, 27].



**FIGURE 12.1** Illustrates the shift from conventional to modern techniques.

## 12.4 THE APPLICATION OF ARTIFICIAL INTELLIGENCE (AI) HAS GREATLY ENHANCED STRUCTURE-BASED VIRTUAL SCREENING (SBVS)

Leading to significant improvements in optimizing the hit detection workflow. By leveraging the powerful learning and generalization capabilities of AI methods, the effectiveness and probability of high-throughput screening (HTS) have been increased, while also reducing costs [28, 29]. Within the virtual screening pipeline, AI has played a crucial role in various aspects. Structure-based virtual screening (SBVS) encompasses the utilization of structural data derived from the binding between ligands and targets to effectively identify and rank potential ligands from the vast pool of available chemical compounds. The application of molecular docking methodologies in SBVS has been extensively employed since the 1980s [30].

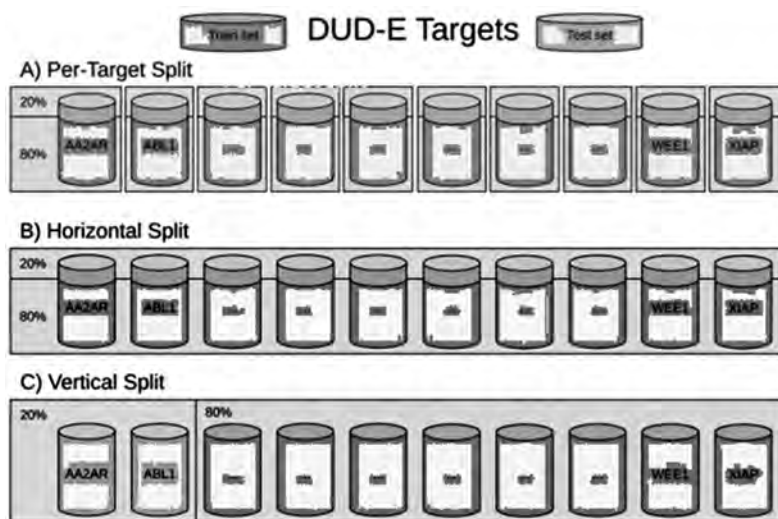


While initial studies faced computational limitations with fully flexible ligands and a restricted set of targets, advancements in hardware have made this methodology more accessible and widely applicable [31]. The initial stage of docking involves positioning molecules from a library onto the target's binding site, considering their steric and physiochemical properties. Subsequently, a scoring function is utilized to evaluate the generated binding poses and predict the energetic interaction between the target and the ligand. The selection of an appropriate scoring method has been a subject of extensive discussion. Traditional approaches have relied on experimental methods utilizing drive models or practical characteristics for pose calculation. However, recent advancements have focused on knowledge-based data mining [32]. After the counting step, a post-processing stage is conducted to rank candidate results. Several docking software options now provide vigorous and precise sample procedures for pose grouping, such as matching algorithms [33] that map ligands to the target's active site based on molecular size, as well as the Incremental Construction Algorithm [34].

The active site of the ligand was accurately positioned using a systematic approach that incorporated multiple techniques. Random conformation changes, facilitated by Monte Carlo (MC) algorithms [35], were employed to generate various poses, while the molecular dynamics (MD) algorithm [36] facilitated the movement of individual atoms within the goal in the presence of surrounding atoms. In the field of structure bases virtual broadcast aided by AI. A notable application involved the development of exact goal counting model known as SVM-SP [37]. This model utilized mechanical ligand docking [38] to create protein-ligand complexes, and SVMGen [39], the scoring model derived from numerical pairwise possibilities of these docked pairs, formed an overall scoring function. Furthermore, a successful SBVS approach called MIEC-SVM [40], which effectively differentiates amongst vigorous and sluggish particles, have demonstrated promising performance in various reflective computer-generated educations [41, 42]. The education conducted by Sun et al. [43] focused on optimizing the construction of the MIEC-SVM model.

Extensive research has been conducted on the practicality of Support Vector Machine (SVM) methods in post-docking analysis. Leong and colleagues [44] proposed a collaborative docking scheme that utilizes a combinatorial approach to choice required positions plus determine the required empathy. Similarly, Yan et al. [45] presented PLEIC-SVM, a post-docking classification method that employs thumbprint to represent observed interactions in each complex. Furthermore, Rodri Gez-Perez et al. [46] developed SVM-based workflows that integrate the powers and limits of both constructions based practical broadcast techniques. The scientist Xie

and coworkers [47] joint SVM to classify potential inhibitors of c-Met tyrosine kinase from vast compound library. Meslamani et al. [48] successfully applied RF-Score296 in various structure-based virtual screening (SBVS) applications. They evaluated the performance of RF-Score against different targets and created RF-Score-VS, a practical scoring function trained on a comprehensive collection of enhanced decoy data [49]. To assess the quality of the model, a stratified 5-fold cross-validation approach was employed (Figure 12.2). The results indicated that the top 1% of compounds ranked by RF-Score-VS achieved a virtual hit rate of 56%, while the best traditional scoring function yielded only 16% hits. These findings illustrate the significant improvement in virtual screening performance achievable through the utilization of machine learning techniques.

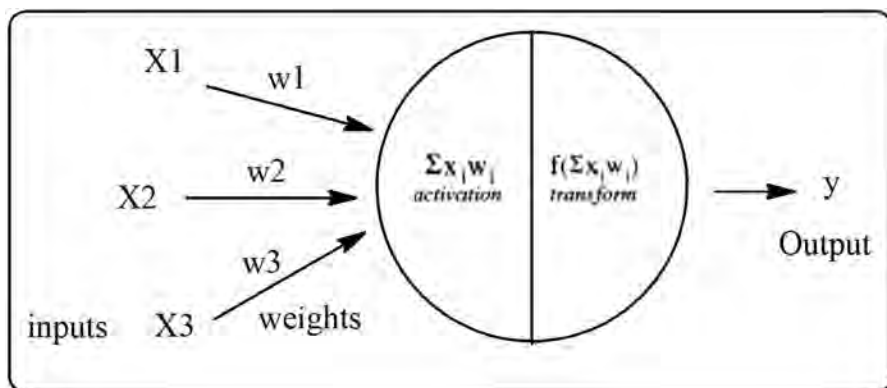


**FIGURE 12.2** A summary of the quality valuation approach [49].

Source: Reprinted from Ref.[49]. <https://creativecommons.org/licenses/by/4.0/>

An artificial neural network (ANN), also known as a connectionist model, is a computational framework comprising synthetic number of known as Processing Elements (PE) [50, 51]. These PEs process information by adjusting weighted inputs, employing transmission functions, and generating outputs. ANNs simulate the information processing abilities of the human brain, acquiring knowledge by identifying patterns and relationships in data rather than relying on explicit programming. They learn and enhance their performance through experience and training. While ANNs

are capable of handling vast amounts of data and occasionally provide highly accurate predictions, they do not exhibit human-like intelligence. Therefore, it may be more appropriate to describe them as computer intellect rather than intelligence in the traditional sense (Figure 12.3).



**FIGURE 12.3** Model of an artificial neuron.

#### 12.4.1 MACHINE KNOWLEDGE

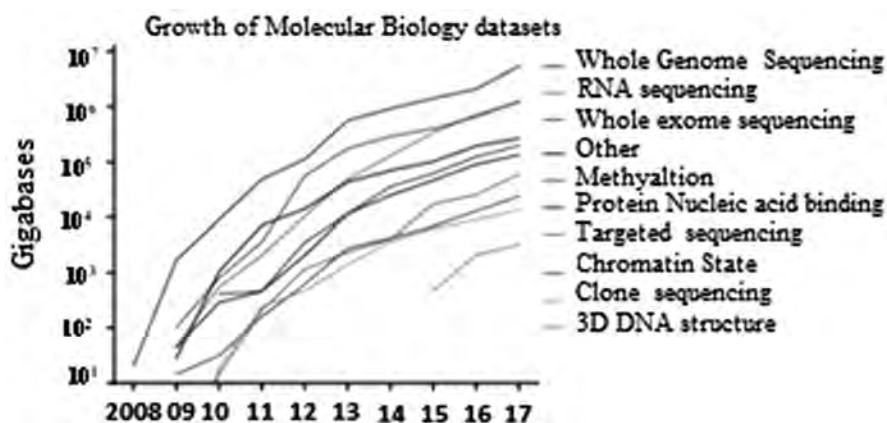
Machine Knowledge involves a broad range of techniques used to extract meaningful insights from data without explicit programming. To illustrate, in computer vision, an image can be used as input, and the desired outcome is incline of identified matters. The model's performance is evaluated using a damage purpose, which measures the deviation from the desired outcome. A model with improperly set parameters will produce numerous errors and exhibit high training error. However, an effective training process will optimize these parameters to minimize errors and align predictions with the training labels. Once the training phase is finished, the model cab be deployed to handle novel contribution scenarios. Learning methods are characterized by their ability to learn mathematical.

#### 12.4.2 APPLICATIONS OF QSAR MODELING IN MEDICINE DISCOVERY

The process of drug detection comprises two primary elements: goal documentation and the identification of mixtures that exhibit desired on the goal effects while minimizing off-goal things. While deep neural networks

(DNNs) analyzing molecular effects can aid in target identification, our focus lies on the prediction of the biological activity of potential drugs towards a specific target. However, QSAR modeling encounters a significant challenge due to the complex structural nature of proteins, which exceeds that of small molecules. The estimate of a proteins tertiary construction solely based on it is polypeptide order remainders an unexplained problematic. Current deep learning efforts in protein structure calculation primarily concentrate on the prediction of secondary structure using persistent neural networks [52, 53], as predicting tertiary structure solely from sequence is currently not feasible. To address this complexity, one approach is to utilize ligand-based models that solely consider small molecules and do not involve protein modeling. These descriptors contain information about the quantity and arrangement of atom and various function group in the molecule. Alternatively, molecular descriptors can be derived mechanically from biochemical structures consuming models such as autoencoders [54] or dynamically learned using specific neural network architectures [55, 56].

A modified method utilizes a deep learning framework that incorporates multi-task learning to train a unified model for diverse proteins. The team achieved success by employing a combination of single-task and multi-task deep neural networks (Figure 12.4) [57, 58]. Building upon this multi-task strategy, the authors further expanded it to encompass more than 200 targets, demonstrating continuous progress without reaching saturation, even as additional targets and tasks were included. However, together one and



**FIGURE 12.4** The emergence and impact, highlights, remarkable surge in bio-medicinal data through the example of genomic sequence data (measured in giga bases) across different assays (it is important to note the utilization of a logarithmic scale, lastly, subsection).

multiple task ligand bases model suffer from a significant limitation—they are unable to accurately predict proteins that were not included in the training set, as these models do not explicitly represent proteins [59]. Consequently, this creates a disadvantage for proteins that require precise predictions since they have limited available data. To overcome this challenge, a promising approach has emerged that incorporates three-dimensional convolutional networks to directly model protein structure, alongside small molecule structure. This empowers deep neural networks to make generalizable predictions for entirely new proteins, without relying solely on experimental biological activity data [60].

## **12.5 RESULT AND DISCUSSION**

With substantial investments of time and capital, coupled with technological advancements, it is expected that the aforementioned obstacles can be overcome in the near future. Consequently, various factors that have hindered drug discovery endeavors for the past five decades are likely to undergo further advancements soon. The prospects for utilizing AI in repurposing drugs show significant promise. Given the exorbitant costs and high failure rates associated by unique, chemically inventive drug research and improvement, exploring new applications for current medicines can recover the danger outline and provide opportunities for abandoned healing categories. Historically, medical reproducibility has relied on practical approaches grounded in clinical observations. This abundance offers an opportunity to adopt a groundbreaking approach that integrates these diverse data sources into a more cohesive and nonstop stream of income and pharmacological visions. Existing computational high-throughput methods face similar challenges to their experimental counterparts, including high rates of false positives and an imbalance between helpful and bad information. In the coming years, it is expected that virtual screening technologies utilizing deep learning will emerge, either as replacements or complements to traditional screening methods, thereby improving the efficiency and success rates of the screening process. Currently, lead optimization represents one of the most complex and multifaceted stages in drug development, traditionally relying on the expertise and ruling of therapeutic chemist. Key challenges revolve around defining the characteristics of an effective drug, promoting favorable ADMET possessions, and instantaneously optimizing these possessions alongside desirable on-target activity. It is crucial to acknowledge the interdependence and equal importance of these factors. By

harnessing the power of AI, we can strive to optimize all these parameters concurrently, refining our QSAR model and facilitating faster identification of safe compounds for synthesis. Ongoing research actively supports the development of automated synthesis methods.

## 12.6 CONCLUSION

AI has showcased its potential in various domains of drug discovery. Although it is not a panacea, it should be embraced as a supportive tool for experts in their respective parts and specialisms during the medicine growth. The implementation of AI in exact sectors within the business is still in its early stages. We shouldn't expect overnight revolutionary transformations due to the cautious and gradual nature of the drug discovery process, which prioritizes risk management and accountable development of innovative scientific approaches for the benefit of patients and shareholders. Nevertheless, by integrating these mindsets, there is a significant opportunity to enhance efficiency in certain aspects of drug discovery. This integration enables researchers to allocate their time and attention to different challenges by assigning routine tasks to a combination of AI and automation. Additionally, AI can offer valuable visions to experienced researchers based on its extensive "experiential recall," providing a fundamentally new approach. However, the introduction of these innovations will inevitably face obstacles and duplicated efforts. Despite these challenges, there is no doubt that AI will drive changes in some drug innovation processes, promoting the exploration and advancement of new drugs.

## KEYWORDS

- **adversarial autoencoder**
- **artificial intelligence**
- **artificial neural network**
- **machine learning**
- **molecular embedding**
- **natural language processing**
- **neural network**
- **reaction efficiency prediction**

## REFERENCES

1. Heckerman, D. E., & Shortliffe, E. H. (1992). Artificial Intelligence in medicine. *Artificial Intelligence in Medicine*, 4(1), 35–52.
2. Jimison, H. B., Fagan, L. M., Shachter, R. D., & Shortliffe, E. H. (1992). Artificial Intelligence in medicine. *Artificial Intelligence in Medicine*, 4(2), 191–205.
3. Vamathevan, J., Clark, D., Czodrowski, P., et al. (2019). Applications of machine learning in drug discovery and development. *Nature Reviews Drug Discovery*, 18(6), 463–477.
4. Lavecchia, A. (2015). Machine-learning approaches in drug discovery: Methods and applications. *Drug Discovery Today*, 20(3), 318–331.
5. Lo, Y. C., Rensi, S. E., Tornø, W., et al. (2018). Machine learning in cheminformatics and drug discovery. *Drug Discovery Today*, 23(8), 1538–1546.
6. Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural Networks*, 61, 85–117.
7. Yang, X., Wang, Y., Byrne, R., et al. (2019). Concepts of artificial intelligence for computer-assisted drug discovery. *Chemical Reviews*, 119(18), 10520–10594.
8. Schneider, G. (2019). Mind and machine in drug design. *Nature Machine Intelligence*, 1(3), 128–130.
9. Wu, Z., Ramsundar, B., Feinberg, E., et al. (2018). MoleculeNet: A benchmark for molecular machine learning. *Chemical Science*, 9(2), 513–530.
10. Feinberg, E. N., Sur, D., Wu, Z., et al. (2018). PotentialNet for molecular property prediction. *ACS Central Science*, 4(11), 1520–1530.
11. Kearnes, S., McCloskey, K., Berndl, M., et al. (2016). Molecular graph convolutions: Moving beyond fingerprints. *Journal of Computer-Aided Molecular Design*, 30(8), 595–608.
12. Gilmer, J., Schoenholz, S. S., Riley, P. F., et al. (2017). Neural message passing for quantum chemistry. *arXiv preprint [cs.LG]*.
13. Hansch, C., Maloney, P. P., Fujita, T., et al. (1962). Correlation of biological activity of phenoxyacetic acids with Hammett substituent constants and partition coefficients. *Nature*, 194(4824), 178–180.
14. Göller, A., Kuhnke, L., Montanari, F., et al. (2020). Bayer's in silico ADMET platform: A journey of machine learning over the past two decades. *Drug Discovery Today*, 25(9), 1702–1709.
15. Winiwarter, S., Ahlberg, E., Watson, E., et al. (2018). In silico ADME in drug design—Enhancing the impact. *ADMET & DMPK*, 6(1), 15–33.
16. Beresford, A. P., Segall, M., & Tarbit, M. H. (2004). In silico prediction of ADME properties: Are we making progress? *Current Opinion in Drug Discovery & Development*, 7(1), 36–42.
17. Norinder, U., & Bergström, C. A. S. (2006). Prediction of ADMET properties. *Chem MedChem*, 1(9), 920–937.
18. Pensak, D. A., & Corey, E. J. (1977). Computer-Assisted Organic Synthesis. In *ACS Symposium Series* (Vol. 61, pp. 1–32). <https://doi.org/10.1021/bk1977-0061.ch001>.
19. Hussain, A. S., Xuanqiang, Y., & Johnson, R. D. (1991). *Pharm. Research*, 8, 1248–1252.
20. Murtoniemi, E., Merkkä, P., Kinnunen, P., Leiviska, K., & Yliruusi, J. (1994). *International Journal of Pharmaceutics*, 110(2), 101–108.
21. Gasperlin, M., Tusar, L., Tusar, M., Kristl, J., & Smid-Korbar, J. (1998). *International Journal of Pharmaceutics*, 168(2), 243–254.

22. Takayama, K., Fujikawa, M., & Nagai, T. (1999). *Pharmaceutical Research*, 16(1), 1–6.
23. Achanta, A. S., Kowalski, J. G., & Rhodes, C. T. (1995). *Drug Development and Industrial Pharmacy*, 21(2), 119–155.
24. Agatonovic-Kustrin, S., & Beresford, R. (2000). Basic concepts of artificial neural network (ANN) modeling and its application in pharmaceutical research. *Journal of Pharmaceutical and Biomedical Analysis*, 22(5), 717–727.
25. Rumelhart, D. E., & McClelland, J. L. (1986). *Parallel Distributed Processing: Explorations In the Microstructure of Cognition* (Vol. 1: Foundations). Cambridge, MA: MIT Press.
26. Gawehn, E., Hiss, J. A., & Schneider, G. (2016). Deep learning in drug discovery. *Molecular Informatics*, 35(1), 3–14.
27. Zhang, L., Tan, J., Han, D., et al. (2017). From machine learning to deep learning: Progress in machine intelligence for rational drug discovery. *Drug Discovery Today*, 22(11), 1680–1685.
28. Schneider, G. (2010). Virtual screening: An endless staircase? *Nature Reviews Drug Discovery*, 9(4), 273–276.
29. Scior, T., Bender, A., Tresadern, G., Medina-Franco, J. L., Martínez-Mayorga, K., Langer, T., Cuanalo-Contreras, K., & Agrafiotis, D. K. (2012). Recognizing pitfalls in virtual screening: A critical review. *Journal of Chemical Information and Modeling*, 52(5), 867–881.
30. Kuntz, I. D., Blaney, J. M., Oatley, S. J., Langridge, R., & Ferrin, T. E. (1982). A geometric approach to macromolecule-ligand interactions. *Journal of Molecular Biology*, 161(2), 269–288.
31. Lyu, J., Wang, S., Balius, T. E., Singh, I., Levit, A., Moroz, Y. S., O’Meara, M. J., Che, T., Algaa, E., Tolmachova, K., et al. (2019). Ultra-large library docking for discovering new chemotypes. *Nature*, 566(7743), 224–229.
32. Sottriffer, C. A., Gohlke, H., & Klebe, G. (2002). Docking into knowledge-based potential fields: A comparative evaluation of DrugScore. *Journal of Medicinal Chemistry*, 45(9), 1967–1970.
33. Norel, R., Fischer, D., Wolfson, H. J., & Nussinov, R. (1994). Molecular surface recognition by a computer vision-based technique. *Protein Engineering, Design & Selection*, 7(1), 39–46.
34. Rarey, M., Kramer, B., Lengauer, T., & Klebe, G. (1996). A fast flexible docking method using an incremental construction algorithm. *Journal of Molecular Biology*, 261(3), 470–489.
35. Liu, M., & Wang, S. (1999). MCDock: A Monte Carlo simulation approach to the molecular docking problem. *Journal of Computer-Aided Molecular Design*, 13(6), 435–451.
36. Brooks, B. R., Brucoleri, R. E., Olafson, B. D., States, D. J., Swaminathan, S., & Karplus, M. (1983). CHARMM: A program for macromolecular energy, minimization, and dynamics calculations. *Journal of Computational Chemistry*, 4(2), 187–217.
37. Li, L., Khanna, M., Jo, I., Wang, F., Ashpole, N. M., Hudmon, A., & Meroueh, S. O. (2011). Target-specific support vector machine scoring in structure-based virtual screening: Computational validation, in vitro testing in kinases, and effects on lung cancer cell proliferation. *Journal of Chemical Information and Modeling*, 51(4), 755–759.
38. Xu, D., Li, L., Zhou, D., Liu, D., Hudmon, A., & Meroueh, S. O. (2017). Structure-based target-specific screening leads to small-molecule CaMKII inhibitors. *ChemMedChem*, 12(5), 660–677.



39. Xu, D., & Meroueh, S. O. (2016). Effect of binding pose and modeled structures on SVMGen and GlideScore enrichment of chemical libraries. *Journal of Chemical Information and Modeling*, 56(6), 1139–1151.
40. Hou, T., Zhang, W., Case, D. A., & Wang, W. (2008). Characterization of domain–peptide interaction interface: A case study on the Amphiphysin-1 SH3 domain. *Journal of Molecular Biology*, 376(5), 1201–1214.
41. Ding, B., Li, N., & Wang, W. (2013). Characterizing binding of small molecules. II. Evaluating the potency of small molecules to combat resistance based on docking structures. *Journal of Chemical Information and Modeling*, 53(5), 1213–1222.
42. Li, N., Ainsworth, R. I., Wu, M., Ding, B., & Wang, W. (2016). MIECSVM: Automated pipeline for protein peptide/ligand interaction prediction. *Bioinformatics*, 32(7), 940–942.
43. Sun, H., Pan, P., Tian, S., Xu, L., Kong, X., Li, Y., Li, D., & Hou, T. (2016). Constructing and validating high-performance MIEC-SVM models in virtual screening for kinases: A better way for actives discovery. *Scientific Reports*, 6, 24817.
44. Leong, M. K., Syu, R. G., Ding, Y. L., & Weng, C. F. (2017). Prediction of N-methyl-D-aspartate receptor GluN1-ligand binding affinity by a novel SVM-pose/SVM-score combinatorial ensemble docking scheme. *Scientific Reports*, 7, 40053.
45. Yan, Y., Wang, W., Sun, Z., Zhang, J. Z., & Ji, C. (2017). Protein–ligand empirical interaction components for virtual screening. *Journal of Chemical Information and Modeling*, 57(8), 1793–1806.
46. Rodríguez-Perez, R., Vogt, M., & Bajorath, J. (2017). Influence of varying training set composition and size on support vector machine-based prediction of active compounds. *Journal of Chemical Information and Modeling*, 57(4), 710–716.
47. Xie, Q. Q., Zhong, L., Pan, Y. L., Wang, X. Y., Zhou, J. P., Diwu, L., Huang, Q., Wang, Y. L., Yang, L. L., Xie, H. Z., et al. (2011). Combined SVM-based and docking-based virtual screening for retrieving novel inhibitors of c-Met. *European Journal of Medicinal Chemistry*, 46(8), 3675–3680.
48. Meslamani, J., Bhajun, R., Martz, F., & Rognan, D. (2013). Computational profiling of bioactive compounds using a target-dependent composite workflow. *Journal of Chemical Information and Modeling*, 53(9), 2322–2333.
49. Wojcikowski, M., Ballester, P. J., & Siedlecki, P. (2017). Performance of machine-learning scoring functions in structure-based virtual screening. *Scientific Reports*, 7, 46710.
50. Zupan, J., & Gasteiger, J. (1992). *Analytica Chimica Acta*, 248(1), 1–30.
51. Zurada, J. M. (1992). *Introduction to Artificial Neural Systems*. PWS Publishing Company.
52. Baldi, P., Brunak, S., Frasconi, P., Soda, G., & Pollastri, G. (1999). Exploiting the past and the future in protein secondary structure prediction. *Bioinformatics*, 15(11), 937–946.
53. Pollastri, G., Przybylski, D., Rost, B., & Baldi, P. (2002). Improving the prediction of protein secondary structure in three and eight classes using recurrent neural networks and profiles. *Proteins*, 47(2), 228–235.
54. Tan, J., Ung, M., Cheng, C., & Greene, C. S. (2015). Unsupervised feature construction and knowledge extraction from genome-wide assays of breast cancer with denoising autoencoders. *Pacific Symposium on Biocomputing*, 20, 132–143.
55. Duvenaud, D., et al. (2015). Convolutional networks on graphs for learning molecular fingerprints. *Preprint*. <https://arxiv.org/abs/1509.09292>.
56. Kearnes, S., McCloskey, K., Berndl, M., Pande, V., & Riley, P. (2016). Molecular graph convolutions: Moving beyond fingerprints. *Journal of Computer-Aided Molecular Design*, 30(8), 595–608.

57. Dahl, G. E., Jaitly, N., & Salakhutdinov, R. (2014). Multi-task neural networks for QSAR predictions. *Preprint*. <https://arxiv.org/abs/1406.1231>.
58. Ma, J., Sheridan, R. P., Liaw, A., Dahl, G. E., & Svetnik, V. (2015). Deep neural nets as a method for quantitative structure-activity relationships. *Journal of Chemical Information and Modeling*, 55(12), 263–274.
59. Ramsundar, B., et al. (2015). Massively multitask networks for drug discovery. *Preprint*. <https://arxiv.org/abs/1502.02072>.
60. Wallach, I., Dzamba, M., & Heifets, A. (2015). AtomNet: A deep convolutional neural network for bioactivity prediction in structure-based drug discovery. *Preprint*. <https://arxiv.org/abs/1510.02855>.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## CHAPTER 13

---

# Fault Diagnosis of Chemical Process Plant Using Artificial Intelligence

TEJAS TEKAWADE, R. B. DHUMALE, and P. B. MANE

*AISSMS Institute of Information Technology, Pune, Maharashtra, India*

---

### ABSTRACT

Fault diagnosis is an essential aspect of process plant management, and it is critical to detect and diagnose faults accurately and promptly to prevent downtime and reduce the risk of accidents. Artificial Intelligence (AI) techniques, such as machine learning and deep learning, have shown great potential in fault diagnosis for chemical process plants. This chapter presents an overview of the use of AI in fault diagnosis for chemical process plants. The first part of the chapter discusses the importance of fault diagnosis in chemical process plants, the challenges associated with traditional fault diagnosis methods and role of AI in improving fault diagnosis. The second part of the chapter provides an overview of AI techniques, including machine learning and deep learning, comparison of AI techniques with traditional fault diagnosis methods and their Advantages and limitations in fault diagnosis. The third part of the chapter describes a case study of fault diagnosis using AI for a chemical process plant. The study demonstrates how a combination of machine learning algorithms, such as random forest and support vector machine, and deep learning techniques, such as convolutional neural networks and long short-term memory, can be used to diagnose faults accurately and promptly. The results show that the AI-based approach achieved a high accuracy rate of fault diagnosis and reduced the downtime of the process plant. The final part of the chapter discusses the potential benefits and limitations of using AI for fault diagnosis in chemical process plants.

The chapter concludes that AI-based fault diagnosis has great potential in improving the efficiency and reliability of chemical process plants, but its success depends on the quality and quantity of data collected and the accuracy of the AI model. Continuous monitoring and updating of the AI model is necessary to ensure its effectiveness.

### **13.1 INTRODUCTION**

Condition-based maintenance (CBM) is a maintenance strategy that focuses on monitoring the actual condition of equipment or assets to determine when maintenance activities should be performed [1]. Rather than relying on fixed schedules or time-based maintenance, CBM uses real-time data and various monitoring techniques to identify the optimal time for maintenance. The benefits of condition-based maintenance include increased equipment uptime, reduced maintenance costs, improved safety, and extended equipment life. By detecting issues in advance and intervening only when necessary, CBM allows organizations to allocate maintenance resources more efficiently and minimize disruptions to operations [2].

Fault detection, fault diagnosis, and fault prognosis are key components of condition-based maintenance (CBM) that aid in identifying and addressing potential issues with equipment or assets. Fault detection involves the process of identifying deviations or anomalies from normal operating conditions. Once a fault is detected, fault diagnosis is performed to determine the underlying cause of the deviation. Prognosis is the process of estimating the remaining useful life (RUL) or the time to failure of the equipment. The combination of fault detection, diagnosis, and prognosis enables CBM systems to provide timely and accurate information about the health and condition of equipment [3]. This allows maintenance teams to take proactive measures, such as scheduling maintenance actions in advance, ordering necessary spare parts, and allocating resources efficiently.

Advanced technologies, such as Machine Learning, Artificial Intelligence (AI), and Data Analytics, play a crucial role in fault detection, diagnosis, and prognosis in CBM. These technologies enable the automated analysis of large volumes of data, the identification of complex patterns, and the development of predictive models. As a result, CBM systems can continuously learn and improve their fault detection, diagnosis, and prognosis capabilities over time, leading to more accurate and reliable maintenance decisions.

### **13.1.1 IMPORTANCE OF FAULT DIAGNOSIS IN CHEMICAL PROCESS PLANTS**

Fault diagnosis shows an important role in chemical process plants as it helps identify and resolve problems that can effect the safety, efficiency, and productivity of the plant. Faults in chemical process plants can lead to hazardous situations, such as leaks, equipment failures, or process deviations. Fault diagnosis helps in identifying potential safety hazards and allows for timely interventions to mitigate risks. Fault diagnosis helps identify deviations from normal operating conditions and allows for prompt corrective actions. Early fault detection allows plant operators to take the necessary precautions to safeguard workers, prevent accidents, and guarantee that safety requirements are being followed. By addressing faults, operators can optimize process parameters, improve efficiency, and minimize waste or resource consumption, leading to cost savings and improved overall plant performance [4]. Maintenance teams are able to minimize the impact on production continuity and prevent expensive downtime by identifying faults and foreseeing failures. This allows them to arrange repairs or replacements during planned shutdowns or maintenance windows. Fault diagnosis relies on data analysis techniques, including statistical analysis, pattern recognition, and machine learning algorithms. These techniques enable plant operators to make informed decisions based on real-time and historical data. By leveraging data-driven insights, operators can identify recurring patterns, trends, or correlations that contribute to faults and use this knowledge to optimize operations, improve maintenance strategies, and enhance overall plant performance [5].

Fault diagnosis in chemical process plants is essential for maintaining safe and efficient operations. It allows for proactive maintenance, reduces downtime, optimizes resource utilization, and supports data-driven decision making. By investing in robust fault diagnosis systems and practices, chemical process plants can enhance safety, productivity, and profitability while minimizing risks and operational disruptions [6].

### **13.1.2 CHALLENGES ASSOCIATED WITH TRADITIONAL FAULT DIAGNOSIS METHODS**

Traditional fault diagnosis methods in chemical process plants often face several challenges that can hinder their effectiveness. These methods rely on a limited number of sensors to collect data from the process equipment.

This limited sensor coverage may not capture all the relevant parameters or provide a comprehensive understanding of the equipment's behavior. As a result, certain faults or anomalies may go undetected, leading to incomplete or inaccurate fault diagnosis [7].

These methods often involve manual analysis of collected data by human experts. This manual process is time-consuming, subjective, and prone to human error. It requires extensive domain knowledge and expertise to identify patterns, correlations, and fault signatures from the data. Manual analysis becomes increasingly challenging with the growing complexity and volume of data generated by modern process plants. Some traditional fault diagnosis methods rely on periodic or offline data collection and analysis. This delayed approach to fault diagnosis means that faults or deviations may not be detected in real-time. By the time a fault is diagnosed, it may have already caused significant damage or disruption to the process. Real-time monitoring is crucial for proactive fault detection and timely intervention [8].

Chemical process plants generate large volumes of multivariate data from various sensors and equipment. Traditional fault diagnosis methods may struggle to effectively handle and analyze such complex and high-dimensional data. The interdependencies and interactions among different variables can be challenging to capture and interpret using conventional techniques, limiting the accuracy and reliability of fault diagnosis [9].

These methods often focus solely on data analysis without considering the broader contextual information surrounding the process plant. Factors such as operating conditions, environmental conditions, and historical data may provide valuable insights for fault diagnosis [10]. Ignoring this contextual information can result in false positives or false negatives during fault diagnosis.

Traditional fault diagnosis methods are often designed for specific equipment or process configurations. Adapting these methods to new equipment or process variations can be time-consuming and resource-intensive. Additionally, traditional methods may struggle to scale up to handle large-scale process plants with multiple interconnected systems and complex process dynamics [11].

There is a growing trend towards leveraging advanced technologies such as machine learning, data analytics, and real-time monitoring systems for fault diagnosis in chemical process plants. These technologies can address the limitations of traditional methods by enabling automated data analysis, handling complex multivariate data, providing real-time insights, and incorporating contextual information for more accurate and efficient fault diagnosis.

### **13.1.3 ROLE OF ARTIFICIAL INTELLIGENCE IN IMPROVING FAULT DIAGNOSIS**

AI has emerged as a powerful tool for improving fault diagnosis in various industries, including chemical process plants. AI techniques, such as machine learning, data analytics, and expert systems, can significantly enhance fault diagnosis capabilities. AI algorithms can analyze large volumes of data generated by sensors and process equipment in real-time. Machine learning techniques, such as anomaly detection and pattern recognition, can automatically identify abnormal behavior or fault signatures without relying on predefined rules or thresholds. This automated data analysis enables early fault detection, reducing the time and effort required for manual analysis [12].

AI techniques excel in handling complex and high-dimensional multivariate data. They can capture complex relationships and interdependencies among different variables in the process plant, facilitating more accurate fault diagnosis. AI models can detect subtle correlations and interactions that traditional methods may overlook, improving the understanding of fault propagation and root causes [13].

AI-based fault diagnosis can incorporate predictive maintenance capabilities. By analyzing historical data, sensor trends, and equipment degradation patterns, AI models can predict the remaining useful life (RUL) of critical components or equipment. This enables proactive maintenance planning, scheduling repairs or replacements before failures occur, minimizing downtime, and optimizing maintenance resources [14].

AI algorithms can learn from historical data and identify fault patterns that are indicative of specific issues. These patterns can be complex and nonlinear, making them challenging for traditional methods to detect. AI-based fault diagnosis can recognize recurring patterns and associate them with known fault types, aiding in accurate diagnosis and reducing false positives or false negatives [15].

AI models can integrate contextual information, such as operating conditions, environmental factors, and maintenance history, into the fault diagnosis process. This contextual understanding helps in generating more accurate diagnoses by considering the broader context of the process plant. For example, AI models can identify how variations in operating conditions impact fault behavior or assess the impact of external factors on equipment performance.

AI-based fault diagnosis systems can continuously learn and improve over time. They can adapt to changing process conditions, equipment variations,



and evolving fault patterns. By incorporating feedback from maintenance actions and outcomes, AI models can refine their fault diagnosis capabilities, enhancing accuracy and reliability [16].

AI plays a pivotal role in improving fault diagnosis by enabling automated data analysis, handling complex multivariate data, predicting failures, recognizing fault patterns, considering contextual information, and continuously learning from feedback. By leveraging AI techniques, chemical process plants can enhance their fault diagnosis capabilities, reduce downtime, improve safety, and optimize maintenance strategies.

## **13.2 ARTIFICIAL INTELLIGENCE TECHNIQUES FOR FAULT DIAGNOSIS IN CHEMICAL PROCESS PLANTS**

AI techniques offer a range of powerful tools for fault diagnosis in chemical process plants. AI techniques can be complemented with data preprocessing, feature selection, and optimization techniques to further enhance fault diagnosis performance. The selection of the appropriate AI technique depends on the specific requirements of the chemical process plant, the available data, and the nature of the faults to be diagnosed. It is essential to have a well-curated dataset and appropriate training and validation procedures to develop accurate and reliable fault diagnosis models (FDMs).

### **13.2.1 OVERVIEW OF AI TECHNIQUES**

AI techniques can be applied to various stages of fault diagnosis, including fault detection, fault classification, and prognosis. They require historical data, sensor measurements, and equipment information to train and develop accurate models. Data preprocessing, feature selection, and model optimization techniques are often used to improve the performance of AI-based fault diagnosis systems.

Machine Learning algorithms learn patterns from historical data to make predictions or identify anomalies. Supervised learning algorithms can classify faults based on labeled data, while unsupervised learning algorithms can detect anomalies and deviations from normal behavior. Machine Learning techniques such as decision trees, support vector machines, and random forests are commonly applied for fault diagnosis in chemical process plants.

Artificial Neural Networks are inspired by the structure and functioning of the human brain. They consist of interconnected nodes or “neurons” that

process and propagate information. ANNs can learn complex patterns in data and are effective for fault diagnosis tasks. They can detect anomalies, classify faults, and predict equipment failures by learning from historical data. Expert systems combine human expertise with rule-based reasoning to diagnose faults. They consist of a knowledge base that stores rules and an inference engine that applies these rules to analyze data and make diagnoses. Expert systems can incorporate domain knowledge, process rules, and fault signatures to provide accurate diagnoses and recommend appropriate actions. Deep learning is a subset of machine learning that utilizes deep neural networks with multiple layers to learn intricate patterns and representations. Techniques such as Convolutional Neural Networks and Recurrent Neural Networks can be applied for fault diagnosis in chemical process plants. Deep learning models excel in processing large volumes of data and can automatically extract relevant features, aiding in fault detection, classification, and prognosis.

Ensemble methods combine multiple AI models to improve fault diagnosis accuracy. By aggregating the predictions of individual models, ensemble methods can make collective decisions and reduce the risk of overfitting. Techniques such as random forests, gradient boosting, and stacking can be employed as ensemble methods for fault diagnosis, incorporating diverse perspectives and improving overall performance.

### ***13.2.2 COMPARISON OF TRADITIONAL AND ARTIFICIAL INTELLIGENCE-BASED FAULT DIAGNOSIS METHODS***

AI-based fault diagnosis methods offer superior accuracy, real-time monitoring, predictive maintenance, adaptability, and continuous learning capabilities compared to traditional methods. Comparison of AI-based fault diagnosis methods and traditional methods in the chemical industry presented in Table 13.1.

## **13.3 CASE STUDY OF AI-BASED FAULT DIAGNOSIS METHODS IN CHEMICAL PROCESS PLANT**

Examples of how AI-based fault diagnostic techniques were applied in a chemical processing facility to enhance problem detection and maintenance procedures are discussed in this section. The case study demonstrates the benefits and effectiveness of AI techniques in a real-world industrial setting.

**TABLE 13.1** Comparison of AI-based Fault Diagnosis Methods and Traditional Methods in the Chemical Industry

Criteria	AI-based Fault Diagnosis Methods	Traditional Methods
Accuracy and reliability	Provide accurate and reliable fault detection and diagnosis based on data analysis.	Reliance on manual analysis or predefined rules, may be less accurate and reliable.
Real-time monitoring and early detection	Enable real-time monitoring and early detection of faults.	Delayed fault detection due to manual inspections or periodic monitoring.
Handling complex data	Effective at analyzing complex and high-dimensional data, capturing intricate relationships and interactions.	May struggle to handle complex data sets and capture underlying relationships.
Predictive maintenance and prognostics	Integrate predictive maintenance capabilities, predict remaining useful life (RUL) of equipment.	Focus on reactive or preventive maintenance without considering specific equipment conditions.
Adaptability and scalability	Adaptable and scalable to different equipment configurations and process variations.	May require extensive customization and struggle to scale up or adapt to new systems.
Continuous learning and improvement	Can continuously learn and update models based on new data, allowing for ongoing optimization.	Reliance on static rules or heuristics without easy updates or refinements.
Integration of contextual information	Ability to incorporate contextual information, such as operating conditions and maintenance history, into fault diagnosis.	May overlook or underutilize contextual information in the diagnosis process.

### 13.3.1 MACHINE LEARNING-BASED FAULT DIAGNOSIS

FDM based on an optimized long short-term memory network (LSTMN) is suggested in Ref. [17]. The link to shaping the ideal number of hidden layer nodes using a repetitive technique based on the LSTMN is improved because the number of hidden layer nodes in the LSTMN greatly influences the diagnosis outcome. To increase the accuracy of diagnosing chemical method faults, the LSTMN is then optimized. The findings of the simulation testing of the Tennessee Eastman (TE) chemical process confirm that the optimized LSTM network outperforms the BP neural network, the multi-layer perceptron approach, as well as the unique LSTMN in terms of performance in identifying chemical process faults. From the recurrent neural network (RNN), the LSTM neural network was created. A neural network called the RNN is utilized to process the sequence data. When RNN is learning, the concept of time is incorporated, unlike the regular ANN. The hidden layer and the output layer of conventional neural networks, such back propagation neural networks, are fully connected, but the hidden layer's neurons are not. The neurons in the hidden layer of the RNN have a feedback mechanism, creating a closed-loop structure in the layer. In order to realize the transmission of information before and after, it can be stretched, forming a time series. The RNN's structural layout is shown in Figure 13.1.

In the single FDMs, the testing data's normal and fault data are distributed at random. The testing data is categorized and diagnosed after the fault diagnostic model has been trained, and the diagnosis information is then generated.

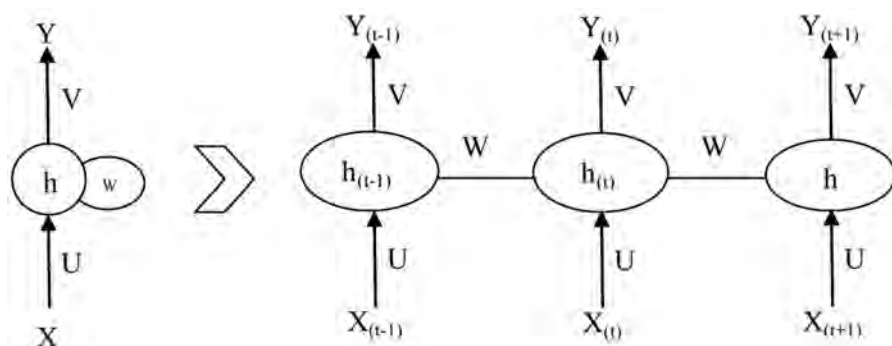


FIGURE 13.1 RNN's structural layout.

The starting settings for the experiment's model's parameters are input size 41, output size 1, batch size 6, and time step 1. The model is trained by iterating until the least diagnostic mistake is found, and then the testing data is forecasted.

The viability of the model for multi-fault diagnosis of the chemical process is examined based on the aforementioned single fault diagnosis experiment. Setting Fault 6 and Fault 8 as inputs to the TE process simulation produced training data and test data. Similar to the single-fault experimental technique, the testing data is set as the model's input after training to produce a fault diagnosis result.

The fault diagnostic model shows outstanding performance when compared to the original LSTMN, the BPNN approach, and the multi-layer perceptron method in the fault diagnosis of chemical processes. Meanwhile, the outcome of the TE chemical process' simulation experiment confirms that the optimized LSTMN has a beneficial impact on both single-fault and multiple fault detection.

### **13.3.2 DEEP LEARNING-BASED FAULT DIAGNOSIS**

Due to its efficiency in processing the highly nonlinear and strongly correlated industrial process data, deep learning networks have lately been used for FDMs. With layer-wise feature compression in conventional deep networks, the valuable information in the raw input can be filtered. This won't help with the later fault classification fine-tuning phase. The extended deep belief network (EDBN), which combines the raw data with hidden features as inputs to each extended restricted Boltzmann machine (ERBM) during the pre-training phase, is offered as a solution to this issue. The dynamic properties of process data are then taken into account while building a dynamic EDBN-based fault classifier. Finally, the Tennessee Eastman (TE) methodology for fault classification is used to evaluate the performance of the suggested method [18].

According to the information bottleneck theory, there would be less and less information that is significant between the extracted deep features and the raw data as the number of neural network layers rises. Therefore, numerous valuable information in the raw data may typically be lost in high layers during the layer-wise compression process used by the majority of existing deep networks. An extended deep belief network (EDBN) is suggested as a solution to this issue in order to adequately capture the important information in the raw data, which is stacked by numerous ERBMs. The raw input data

can take part in the entire compression process by being used as extra inputs to the visible layer of each ERBM for pre-training. As a result, the recovered deep features have a close relationship to the original data, where any potentially relevant information has been completely withheld. When compared to current techniques, EDBN can repeatedly extract useful information from raw data and can offer deep compressed representations that are highly connected with the original data. Additionally, for classification tasks, EDBN can obtain improved accuracy and a reduced false positive rate. Figure 13.2 illustrates the EDBN structure. The average rate of fault diagnosis with the EDBN model is 94.31%, up 0.42% from the DBN. Therefore, compared to the original DBN, EDBN can excerpt extra useful characteristics from raw data for subsequent fault classification performance, showing tremendous promise for defect diagnosis in chemical processes.

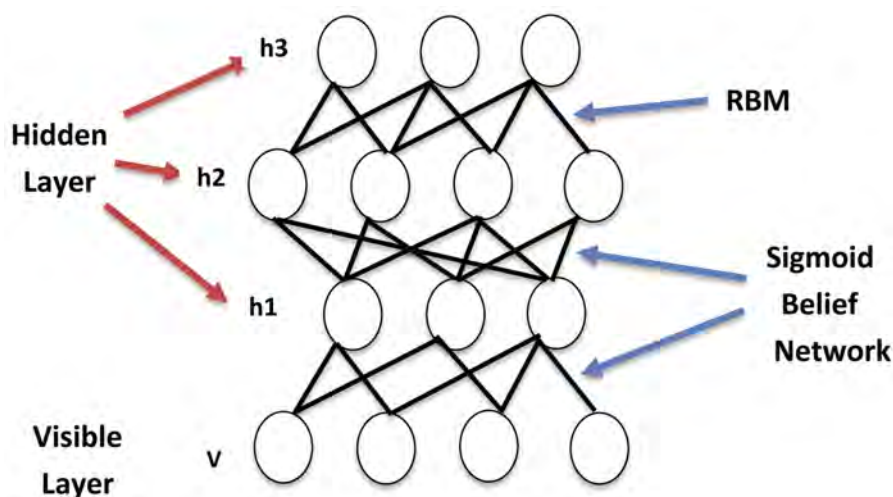


FIGURE 13.2 Extended deep belief network structure.

### 13.3.3 INTEGRATION OF MACHINE LEARNING AND DEEP LEARNING TECHNIQUES FOR FAULT DIAGNOSIS

Integrating machine learning and deep learning techniques offers a data-driven approach to fault diagnosis, enabling accurate and automated identification of faults. The integration of machine learning and deep learning techniques allows for a hybrid approach to fault diagnosis in the chemical process industry.

This integration combines the strengths of both techniques, leveraging the interpretability of machine learning algorithms and the ability of deep learning models to handle complex patterns and high-dimensional data.

Deep expert systems and neural networks, two fundamentally dissimilar AI approaches, are combined in the operator advisory system in Ref. [19]. As a first level filter, a diagnostic method based on the hierarchical application of neural networks is utilized to identify defects frequently found in chemical processing plants. The deep knowledge expert system examines the data and validates the diagnosis or suggests alternate remedies after the neural networks have localized the flaws inside the process. The object-oriented knowledge base of the model-based expert system includes information about the structure and operation of the plant. The diagnostic approach may deal with many defects, noisy process sensor readings, and unique or previously unrecognized faults. An example of the operator advisory system is presented utilizing a multi-column distillation facility.

The designed diagnostic system demonstrated effective diagnostic performance in a range of situations, including the presence of unique problems and sensor noise.

## **13.4 ADVANTAGES AND LIMITATIONS OF AI-BASED FAULT DIAGNOSIS IN THE CHEMICAL INDUSTRY**

### **13.4.1 ADVANTAGES**

AI-based fault diagnosis methods, the chemical industry can benefit from increased operational efficiency, reduced downtime, optimized maintenance strategies, improved safety, and enhanced overall plant performance. These advantages contribute to cost savings, improved product quality, and increased competitiveness in the market.

- 1. Improved Accuracy:** AI techniques, such as machine learning and deep learning, can analyze large volumes of data and identify complex patterns. This results in more accurate fault detection and diagnosis compared to traditional methods, which may rely on manual analysis or predefined rules.
- 2. Real-Time Monitoring:** AI-based fault diagnosis methods enable real-time monitoring of process data. They can continuously analyze data and detect anomalies or deviations from normal behavior promptly. This allows for immediate action and reduces the risk of prolonged faults or operational disruptions.

3. **Early Fault Detection:** AI techniques can detect faults at an early stage, even before they escalate into major issues. By identifying subtle anomalies or deviations in process data, AI-based methods can provide early warning signals, allowing operators to take preventive measures and avoid equipment failures or costly downtime.
4. **Predictive Maintenance:** AI-based fault diagnosis methods incorporate predictive maintenance capabilities. By analyzing historical data and equipment degradation patterns, these methods can predict the remaining useful life (RUL) of critical components or equipment. This enables proactive maintenance planning, reduces unplanned downtime, and optimizes maintenance resources.
5. **Handling Complex Data:** The chemical industry generates vast amounts of complex and high-dimensional data from sensors, instruments, and control systems. AI techniques, such as neural networks and deep learning, excel at handling such data and can capture intricate relationships and dependencies among variables. This enhances the accuracy of fault diagnosis in complex chemical processes.
6. **Adaptability and Scalability:** AI-based fault diagnosis methods are adaptable and scalable to different equipment configurations, process variations, and operating conditions. They can be trained on specific datasets and then applied to similar systems, making them applicable to a wide range of chemical process plants. This flexibility allows for broader utilization and transferability of FDMs.
7. **Continuous Learning and Improvement:** AI-based methods have the ability to continuously learn and improve over time. They can update their models based on new data and feedback, allowing for ongoing optimization and enhanced fault diagnosis performance. This adaptability ensures that the fault diagnosis methods remain effective as process conditions change or new fault patterns emerge.

### 13.4.2 LIMITATIONS AND CHALLENGES

While AI-based fault diagnosis methods offer numerous advantages in the chemical industry, it is important to consider their limitations. Some of the key limitations and challenges include:

1. **Data Availability and Quality:** AI models heavily rely on high-quality and relevant data for training and inference. However, obtaining comprehensive and reliable data in the chemical industry



can be challenging. Limited or incomplete data, sensor failures, and data inconsistencies can impact the performance of AI models and hinder accurate fault diagnosis.

2. **Need for Labeled Training Data:** Supervised machine learning techniques require labeled training data, where faults are identified and labeled. Collecting a diverse and representative dataset with labeled fault instances can be time-consuming and costly. The availability of labeled data for rare or complex faults may be particularly limited, leading to challenges in training accurate models.
3. **Model Interpretability:** Deep learning models, such as neural networks, are often considered black boxes, making it difficult to interpret the reasoning behind their predictions. In the chemical industry, where explainability is crucial for decision-making and troubleshooting, the lack of interpretability in AI models can pose challenges.
4. **Scalability and Generalization:** AI models trained on specific datasets and conditions may struggle to generalize to new environments or handle unseen fault scenarios. Adapting AI models to different chemical processes or scaling them for large-scale industrial applications may require additional efforts in retraining or fine-tuning the models.
5. **Computational Resources and Complexity:** AI-based fault diagnosis methods, especially deep learning models, can be computationally demanding and require substantial computing resources. Training complex models with large datasets may necessitate high-performance computing infrastructure, which can be costly for some organizations. Additionally, the complexity of AI models may require skilled personnel and specialized expertise for their development and maintenance.
6. **Continuous Learning and Adaptability:** Fault conditions in chemical processes can evolve over time due to changes in operating conditions, equipment degradation, or process modifications. AI models need to adapt and continuously learn from new data to stay effective. Implementing mechanisms for online learning and real-time model updates can be challenging and require careful monitoring and validation.
7. **Integration with Existing Systems:** Integrating AI-based fault diagnosis methods into existing plant infrastructure and systems can be complex. Ensuring compatibility, data integration, and real-time

communication between the AI models and control systems may require additional investments and coordination.

Despite these limitations, ongoing research and advancements in AI technologies, along with efforts to address these challenges, continue to improve the effectiveness and applicability of AI-based fault diagnosis in the chemical industry. It is essential to carefully consider these limitations and assess the feasibility and suitability of AI-based solutions for specific applications and operational contexts.

### 13.5 AI-BASED FAULT DIAGNOSIS SOFTWARE USED IN THE CHEMICAL PROCESS INDUSTRY

There are several AI-based fault diagnosis software solutions that are commonly used in the chemical industry. Some notable examples include:

1. **Aspen Mtell:** This, developed by AspenTech, is a predictive maintenance software that utilizes AI and machine learning algorithms to identify and diagnose equipment failures in real-time. It provides early warning notifications, root cause analysis, and recommended actions to prevent unplanned downtime and optimize maintenance strategies.
2. **Seeq:** It is a process analytics software that combines AI and machine learning techniques to analyze time-series data from various sources, including sensors, historians, and process databases. It helps identify patterns, anomalies, and potential faults in the data, enabling proactive fault detection and diagnosis.
3. **Cognite Data Fusion:** It is an industrial data operations and contextualization platform that incorporates AI and machine learning capabilities. It collects, integrates, and analyzes data from diverse sources to enable advanced fault diagnosis, anomaly detection, and predictive maintenance in the chemical industry.
4. **GE Digital APM:** GE Digital's asset performance management (APM) software utilizes AI and machine learning algorithms to optimize asset performance and predict equipment failures. It enables real-time monitoring, fault diagnosis, and proactive maintenance strategies for chemical process plants.
5. **ABB Ability™ Asset Health for Process Industries:** ABB's Asset Health for Process Industries is an AI-powered software solution that

combines advanced analytics, fault detection, and diagnosis capabilities. It provides real-time insights into equipment health, predicts failures, and recommends appropriate maintenance actions for the chemical industry.

6. **Siemens Mindsphere:** It is an open IoT operating system that incorporates AI and analytics tools. It enables real-time monitoring, fault detection, and diagnosis of assets and processes in the chemical industry, supporting predictive maintenance and optimization of operations.

These software solutions leverage AI and machine learning techniques to analyze large volumes of data, detect anomalies, and provide actionable insights for fault diagnosis in the chemical industry. It is important to note that the selection of a specific software solution depends on the specific requirements, infrastructure, and objectives of each chemical process plant.

## 13.6 FUTURE DIRECTIONS AND CONCLUSION

Future directions and the potential for AI-based fault diagnosis in chemical process plants are promising. Here are some areas of development and potential advancements:

1. **Advanced Data Analytics:** AI-based fault diagnosis can benefit from advancements in data analytics techniques. Incorporating advanced analytics methods, such as anomaly detection, outlier analysis, and pattern recognition, can enhance the accuracy and effectiveness of fault diagnosis systems.
2. **Deep Reinforcement Learning:** The integration of deep reinforcement learning techniques with fault diagnosis can enable the development of intelligent systems that learn from continuous interactions with the environment. These systems can optimize fault detection and diagnosis processes by dynamically adapting their strategies based on feedback and rewards.
3. **Edge Computing and IoT Integration:** The deployment of AI-based fault diagnosis systems at the edge, closer to the data source, can enable real-time monitoring and decision-making. Integration with IoT devices and sensors can provide a wealth of real-time data, enhancing the accuracy and timeliness of fault diagnosis.
4. **Explainable AI:** Addressing the interpretability challenge is crucial for gaining trust and acceptance of AI-based fault diagnosis systems

in the chemical industry. Developing techniques and approaches that provide explanations for the model's decisions can enhance transparency and facilitate human understanding and validation.

5. **Hybrid Models:** Integrating different AI techniques, such as combining machine learning with physics-based models or expert systems, can leverage the strengths of each approach. Hybrid models can provide more accurate and robust fault diagnosis by combining data-driven learning with domain knowledge and rules.
6. **Proactive Maintenance and Predictive Analytics:** AI-based fault diagnosis can be integrated with predictive maintenance strategies to enable proactive maintenance actions. By analyzing historical data and identifying early warning signs, AI models can predict impending faults and recommend preventive measures, reducing downtime and optimizing maintenance schedules.

These future directions and potential advancements demonstrate the wide-ranging possibilities for AI-based fault diagnosis in chemical process plants. Continued research, collaboration between academia and industry, and advancements in AI technologies will contribute to the development of more robust, efficient, and intelligent fault diagnosis systems in the chemical industry.

In conclusion, AI-based fault diagnosis has emerged as a promising approach in the chemical industry, offering numerous advantages in terms of accuracy, efficiency, and safety. By leveraging advanced machine learning and deep learning techniques, these systems can effectively analyze complex process data, detect faults, and classify them with high precision. However, it is important to recognize that implementing AI-based fault diagnosis in the chemical industry comes with certain challenges and limitations, such as data availability, interpretability, scalability, and integration with existing systems. Despite these challenges, ongoing research and advancements in AI technologies continue to address these limitations and push the boundaries of fault diagnosis in the chemical industry. Future directions include advanced data analytics, integration of edge computing and IoT, explainable AI, hybrid models, proactive maintenance, collaborative fault diagnosis, and continuous learning and adaptation. Overall, AI-based fault diagnosis has the potential to significantly enhance operational efficiency, reduce downtime, and improve safety in chemical process plants. As the field continues to evolve, it is crucial to strike a balance between the benefits and limitations of AI-based approaches, ensuring the development of robust and reliable fault diagnosis systems that can effectively support decision-making and maintenance actions in the chemical industry.

## KEYWORDS

- artificial intelligence
- chemical process plant
- decision-making
- deep learning
- fault diagnosis
- hybrid models
- machine learning

## REFERENCES

1. Ali, A., & Abdelhadi, A. (2022). Condition-based monitoring and maintenance: State of the art review. *Applied Sciences*, 12(2), 688. <https://doi.org/10.3390/app12020688>.
2. Rojek, I., Jasiulewicz-Kaczmarek, M., Piechowski, M., & Mikołajewski, D. (2023). An artificial intelligence approach for improving maintenance to supervise machine failures and support their repair. *Applied Sciences*, 13(8). <https://doi.org/10.3390/app13084971>.
3. Yodo, N., Afrin, T., Yadav, O. P., Wu, D., & Huang, Y. (2023). Condition-based monitoring as a robust strategy towards sustainable and resilient multi-energy infrastructure systems. *Sustainable and Resilient Infrastructure*, 8(sup1), 170–189. <https://doi.org/10.1080/23789689.2022.2134648>.
4. Tripathi, V., Chattopadhyaya, S., Mukhopadhyay, A. K., Saraswat, S., Sharma, S., Li, C., Rajkumar, S., & Georgise, F. B. (2022). A novel smart production management system for the enhancement of industrial sustainability in Industry 4.0. *Mathematical Problems in Engineering*, 2022. <https://doi.org/10.1155/2022/6424869>.
5. Khalid, S., Song, J., Raouf, I., & Kim, H. S. (2023). Advances in fault detection and diagnosis for thermal power plants: A review of intelligent techniques. *Mathematics*, 11(8). <https://doi.org/10.3390/math11081767>.
6. Liao, M., Lan, K., & Yao, Y. (2022). Sustainability implications of artificial intelligence in the chemical industry: A conceptual framework. *Journal of Industrial Ecology*, 26(1), 164–182. <https://doi.org/10.1111/jiec.13214>.
7. Torabi, N., Gunay, H. B., O'Brien, W., & Moromisato, R. (2022). A holistic sequential fault detection and diagnostics framework for multiple zone variable air volume air handling unit systems. *Building Services Engineering Research and Technology*, 43(5), 605–625. <https://doi.org/10.1177/01436244221097827>.
8. Ren, J., Xu, C., Wang, J., Zhang, J., Mao, X., & Shen, W. (2023). An edge-fog-cloud computing-based digital twin model for prognostics health management of process manufacturing systems. *CMES – Computer Modeling in Engineering and Sciences*, 135(1), 599–618. <https://doi.org/10.32604/cmes.2022.022415>.

9. Patil, P. (n.d.). A comparative study of different time series forecasting methods for predicting traffic flow and congestion levels in urban networks. *Unpublished manuscript*, 1–20.
10. Rafati, A., Shaker, H. R., & Ghahghahzadeh, S. (2022). Fault detection and efficiency assessment for HVAC systems using non-intrusive load monitoring: A review. *Energies*, 15(1). <https://doi.org/10.3390/en15010341>.
11. Nguyen, T. N., & Nielsen, P. (2023). The dynamics of information system development in developing countries: From mutual exclusion to hybrid vigor. *Electronic Journal of Information Systems in Developing Countries*, 2022, 1–21. <https://doi.org/10.1002/isd2.12266>.
12. Tkach, V., Kudin, A., KEBande, V. R., Baranovskyi, O., & Kudin, I. (2023). Non-pattern-based anomaly detection in time-series. *Unpublished Manuscript*, 1–25.
13. Arinez, J. F., Chang, Q., Gao, R. X., Xu, C., & Zhang, J. (2020). Artificial intelligence in advanced manufacturing: Current status and future outlook. *Journal of Manufacturing Science and Engineering, Transactions of the ASME*, 142(11), 1–16. <https://doi.org/10.1115/1.4047855>.
14. Qin, Y., Cai, N., Gao, C., Zhang, Y., Cheng, Y., & Chen, X. (2022). Remaining useful life prediction using temporal deep degradation network for complex machinery with attention-based feature extraction. *arXiv preprint*. <http://arxiv.org/abs/2202.10916>.
15. Abubakar, A., Almeida, C. F. M., & Gemignani, M. (2021). Review of artificial intelligence-based failure detection and diagnosis methods for solar photovoltaic systems. *Machines*, 9(12), 328. <https://doi.org/10.3390/machines9120328>.
16. Bhuiyan, M. R., & Uddin, J. (2023). Deep transfer learning models for industrial fault diagnosis using vibration and acoustic sensors data: A review. *Vibration*, 6(1), 218–238. <https://doi.org/10.3390/vibration6010014>.
17. Han, Y., Ding, N., Geng, Z., Wang, Z., & Chu, C. (2020). An optimized long short-term memory network-based fault diagnosis model for chemical processes. *Journal of Process Control*, 92, 161–168. <https://doi.org/10.1016/j.jprocont.2020.06.005>.
18. Wang, Y., Pan, Z., Yuan, X., Yang, C., & Gui, W. (2020). A novel deep learning-based fault diagnosis approach for chemical process with extended deep belief network. *ISA Transactions*, 96, 457–467. <https://doi.org/10.1016/j.isatra.2019.07.001>.
19. Becraft, W., Lee, P., & Newell, R. (1991). Integration of neural networks and expert systems for process fault diagnosis. *Proceedings of the 12th International Conference on Artificial Intelligence*, 832–837. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.77.3731&rep=rep1&type=pdf> (accessed on 25 July 2024).



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

## CHAPTER 14

---

# Structure Representation Techniques and Applications in Cheminformatics

DEEP V. SHAH<sup>1</sup> and PRASHANT S. KHARKAR<sup>2</sup>

<sup>1</sup>*Indian Institute of Science Education and Research, Pune, Maharashtra, India*

<sup>2</sup>*Department of Pharmaceutical Sciences and Technology, Institute of Chemical Technology, Matunga, Mumbai, Maharashtra, India*

---

### ABSTRACT

Chemical information management and its effective use depends on many factors including the information storage and retrieval systems, tools and techniques for encoding and decoding generated/stored data, computational power, property predictions models and many others. The age-old practices in chemical information management have metamorphosed into the modern-day AI/ML dominant strategies. Structure representation, either in 1D, 2D, or higher dimensions, have significantly affected many disciplines including life and material sciences. Starting from simple linear notations, the structure representation now encompasses complex molecular graph theory and graphical representations, along with molecular fingerprints and others. These complicated forms of structure representation are more machine-friendly and at times, poorly understood by humans. The present chapter summarizes the progress made in structure representations with reference to chemical information management, with particular emphasis on AI/ML-driven strategies, which have revolutionized the way chemical information is stored, accessed and retrieved. The future perspectives in various application areas such as reaction representation, and computer-assisted structure elucidation are briefly discussed.

---

Artificial Intelligence for Chemical Sciences: Concepts, Models, and Applications. Shrikaant Kulkarni, Shashikant Bhandari, Dushyant Varshney, & P. William (Eds.)

© 2025 Apple Academic Press, Inc. Co-published with CRC Press (Taylor & Francis)



## 14.1 INTRODUCTION

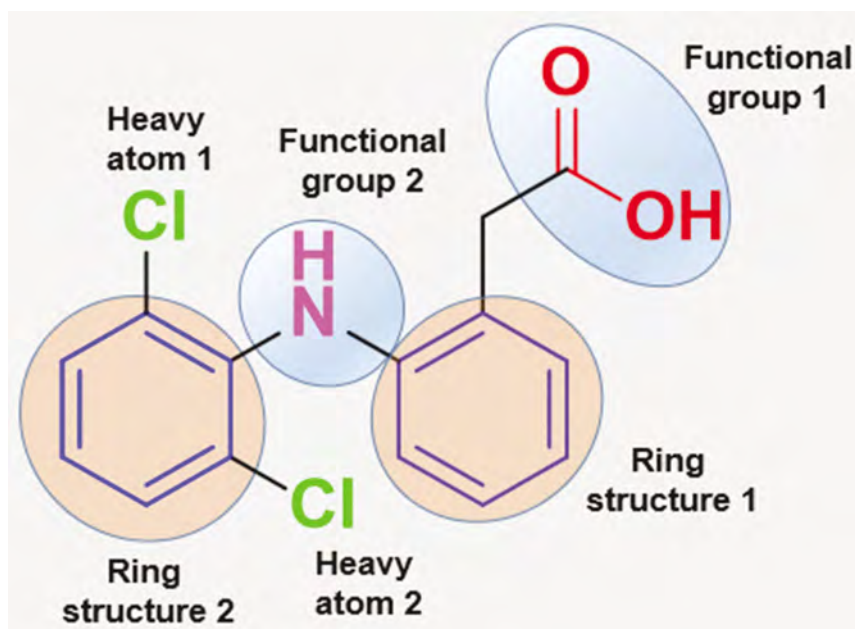
In the present era, managing chemical information is crucial in various fields of Life, Pharmaceutical, Chemical and Material Sciences, to name a few. The domain experts, i.e., Chemical Scientists, play crucial part in searching, organizing, researching and disseminating chemical information to various stakeholders including specialists or experts as well as general public. The pace with which chemical information is generated is humongous and overwhelming. Thousands of research articles, reports, conference proceedings are produced every week and are accessed by researchers, students, industry professionals, regulators and many others across the globe. Storing, retrieving and accessing the historical chemical information is central to chemical research. Of course, these tasks are critically dependent on the softwares, tools, search engines, property calculators and what not. Large chemical databases such as Chemical Abstracts Service (CAS) Registry [1], PubChem [2], hold very large amount of chemical information. For example, as on date, PubChem, world's largest collection of freely accessible chemical information, contains 116 million compounds along with 36 million literature records and 42 million patents [3].

The major tasks in chemical information management include:

1. Access to chemical information (search, read, understand and interpret) by technically competent audience.
2. Chemical database search including chemical structure and patent information retrieval.
3. Present chemical information in readable form or abstract summary.
4. Develop tools and softwares for chemical information storing, retrieving, accessing and processing.
5. Devise domain-specific (e.g., chemical, pharmaceutical, material, etc.,) strategies and processes for chemical information management.

Majority of these tasks are highly dependent on the conversion of chemical information in machine-readable form for easy storage and on-the-go access. Chemical structures, especially of organic compounds, which form the basis of our very existence, are a major part of the primary chemical information that we understand, from which secondary or derived chemical information such as molecular properties, can be generated. Each chemical can be represented in many ways such as molecular formula, empirical formula, if any, molecular structure, isomeric structures, stereoisomers, if applicable, and higher-dimensional formats. Historically, the molecular structure representation has undergone metamorphosis since early 19<sup>th</sup> century [4, 5]. If a

chemist is told about a chemical, say, diclofenac, the most logical thing that comes to her mind is a molecular structure (Figure 14.1), a two-dimensional representation wherein the chemical, in this case diclofenac, is comprised of a collection of atoms, bonds, rings, functional groups, and peculiar way in which they are connected in a unique way; some atoms are implicit, such as nonpolar hydrogens. The chemical information, represented by a molecular structure, can be further used for many different purposes. These include designing chemical synthetic routes, predicting molecular, electronic and steric properties, degradation pathways, studying interaction of the ligand (represented by the structure) with its therapeutic target, metabolic pathways, environmental persistence and many others.



**FIGURE 14.1** Molecular structural representation of diclofenac—a drug.

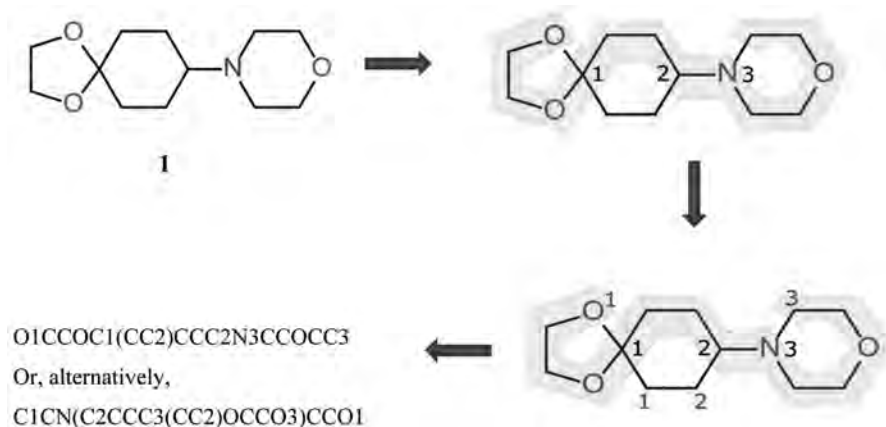
For successful applications of the chemical structure, its representation in human-readable form (Figure 14.1) is quite appropriate. For the chemical information management tasks outlined above, the computer must be able to convert the structural representation from human-to-machine-readable form and back. For advanced applications such as property calculations, the machine-readable form will be further subjected to mathematical manipulations, wherein the chemical will be represented in a complicated way, which

may not truly be understood by the humans, and it is perfectly fine. The task is executed by the computer and interpretable output is generated at the end. The exchange of information between the human input and the machine output is at the core of the chemical information management. Few excellent reviews highlight the progress made over several decades or even centuries, leading to our in-depth understanding of the structural representation [6–8]. The present chapter discusses the state-of-the-art and future perspectives on structural presentation in the era dominated by artificial intelligence (AI), machine- and deep-learning (ML/DL).

## 14.2 METHODS OF STRUCTURE (OR MOLECULE) REPRESENTATION

### 14.2.1 LINEAR NOTATIONS

Numerous representations of molecules in the pre-machine readable formats were introduced in mid-20<sup>th</sup> century [9]. One of the most popular even today, representations include SMILES (Simplified Molecular Input Line Entry System), which came into existence in mid-1960s to early 1970s [10]. It comprised of a one-dimensional (1D), textual representation of a structure as a string of letters and numbers. The grammar of SMILES notations is quite complicated. The interested readers can refer to earlier texts on the subject. It is beyond the scope of this chapter and thus, not discussed further. However, Figure 14.2 demonstrates the SMILES notations for a simple organic molecule, **1**.



**FIGURE 14.2** Illustration of derivation of SMILES notation for a simple organic molecule.

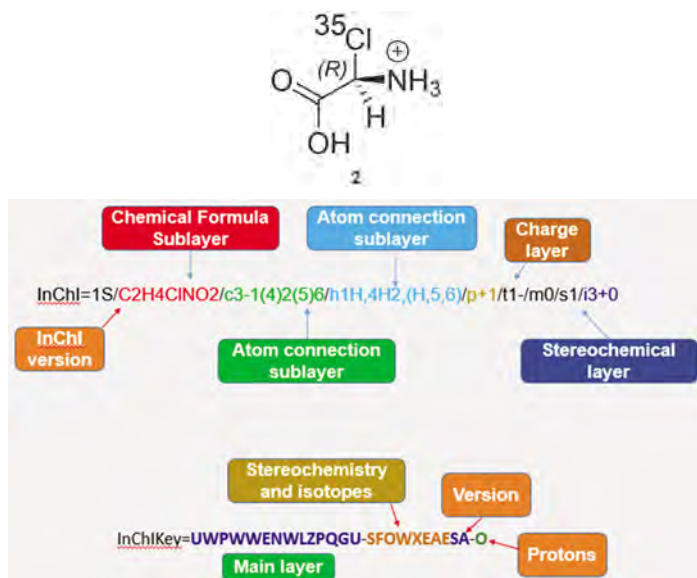
In the first step, the longest chain of connected atoms is located (highlighted in green), followed by numbering the atoms representing the cyclic systems 1–3. Next, the remaining atoms are marked as per the atom they are connected to in the main chain to yield SMILES, moving from left to right. Due to its compact, 1D form, compared to original 2D form of molecular structure, SMILES represent a machine-friendly, yet human-readable format, with a possibility and convenience of incorporating isotopes and stereochemistry, providing a complete and accurate depiction of molecules. It automatically saves disk space while saving very large amount of structural information, especially chemical databases. This obviously relieves the strain on retrieval systems, e.g., during (sub)structure and patent database searches. The multiple representations of same molecule (Figure 14.2) can be problematic when searching for unique structures or molecules. Nonetheless, SMILES is still a favored method of structure representation for chemical information management.

A second form of linear notations include InChI, i.e., International Chemical Identifier, a standardized and unique textual representation of chemical structures [11]. It was developed by the International Union of Pure and Applied Chemistry (IUPAC) and aimed at providing a universal and machine-readable format for conveying molecular information. InChI is designed to address some of the limitations of SMILES and other chemical notations, such as the lack of uniqueness, difficulties in handling stereochemistry, and challenges with interconvertibility. Unlike SMILES, InChI is a linear string of characters that fully represents the molecular structure, including stereochemistry and isotopes, without the need for 2D drawings. It is not only human-readable but also easily interpretable by computers, making it suitable for large chemical databases and automated processing. On the other hand, InChI Key, or InChIKey, is a hashed version of the full InChI representation. It serves as a fixed-size and unique identifier for a specific chemical structure. The InChI Key is designed for quick and efficient searching and indexing in chemical databases, as well as for facilitating the retrieval of specific compounds without the need to store the entire InChI string.

The InChI Key is typically 27 character-long and provides a compact alphabetical code for the InChI. It is generated by applying a cryptographic hash function to the full InChI string, ensuring that the same structure always produces the same InChI Key. As a result, chemical databases can index compounds based on their InChI Keys, enabling fast and reliable searches.

The differences between InChI and InChI Key include (i) InChI is a full and complete representation of the chemical structure, while the InChI Key is a shortened and hashed version of the InChI, designed solely for

identification and fast searching purposes; (ii) InChI can be quite lengthy, especially for complex molecules, as it represents the entire structure in a linear format. In contrast, the InChI Key is fixed at 27 characters, providing a compact and consistent identifier for each chemical compound; (iii) InChI is human-readable to some extent, as it uses standard chemical symbols and conventions. However, its length and complexity may make it less intuitive for direct interpretation by humans. In contrast, the InChI Key is not designed for human interpretation and serves as a machine-friendly identifier; and (iv) InChI serves as a universal and unique identifier for chemical structures, facilitating data exchange and storage in chemical databases, while the InChI Key is specifically designed for efficient searching and indexing of chemical compounds in databases. Figure 14.3 depicts InChI and InChIKey representations for a chiral molecule **2**.



**FIGURE 14.3** Components of InChI and InChIKey for a charged amino acid.

SMARTS (SMILES Arbitrary Target Specification) is an extension of the SMILES notation, designed to enable pattern matching and substructure searching in chemical databases and computational chemistry applications [12]. Developed by Jonathan W. Frey and Christopher J. O. Wyeth in the early 1990s, SMARTS provides a powerful and flexible language for defining complex chemical patterns and querying chemical structures for specific features.

### Basics of SMARTS:

1. **Atoms and Bonds:** In SMARTS, atoms are represented by the same symbols used in SMILES (e.g., C for carbon, O for oxygen, etc.). Additionally, one can use atom class notation (e.g., [C] for any carbon atom) to define generic atom types. Bonds are represented similarly to SMILES (e.g., '-' for single bonds, '=' for double bonds, '#' for triple bonds).
2. **Wildcards:** SMARTS allows the use of wildcards to represent any atom or any bond. The wildcard for any atom is '\*', and for any bond is '~.' This is useful when the specific atom or bond type is not crucial to the search pattern.
3. **Ring Constraints:** SMARTS enables the definition of ring constraints to identify specific ring sizes or types within the target molecule.
4. **Logical Operators:** SMARTS supports logical operators such as AND ('&'), OR ('|'), and NOT ('!'). These operators are used to combine multiple atom or bond patterns to create more complex queries.
5. **Quantifiers:** SMARTS allows the use of quantifiers to specify the number of times a pattern should be repeated. For example, 'n' is used to represent any positive integer, '\*' for zero or more occurrences, '+' for one or more occurrences, and '?' for zero or one occurrence.

SMARTS are a good way to describe the type of molecule rather than a specific molecule. SMARTS take in all the notation of SMILES and add onto them using the aforementioned logical operations, substructure querying and environmental descriptions. For example, to describe the notion of a heteroatom containing six-membered rings, one will have to specify every single possible case, which is not practically possible every time. However, in SMARTS notation, a simple: '[!C&!c&\$( \*1~\*~\*~\*~\*~\*~1)]'.

Some of the useful SMARTS notations to query molecules can be found below:

1. C, N, O, etc., can be used, but [#<Atomic Number>] is also a supported way to specify an atom.
2. Aromatic atoms are shown with a lowercase symbol, like c, n, o.
3. Bonds are specified using the same notation as that of SMILES, with the addition of '~,' which represents any bond and '@,' which represents any ring bond.

4. The conditions can be specified using '!' (not), '||' (or), '&' (and) or ';' (and-low precedence in evaluation).
5. Atoms can be further specified using their chirality, number of bonds, implicit and explicit hydrogens and so on.
6. Environments can be mentioned using the '\$' operator. For example, \$( \*1 ~ \* ~ \* ~ \* ~ \* ~ \*1 ) represents an environment of 6 membered rings with any type of atoms and bonds in between them.

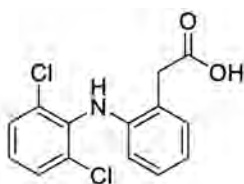
#### 14.2.1.1 APPLICATIONS OF SMARTS

1. One of the primary applications of SMARTS is substructure searching. Researchers can use SMARTS patterns to define specific chemical motifs they are interested in and then search databases or collections of chemical structures to identify molecules containing those substructures. This is invaluable in drug discovery, where scientists may look for molecules with specific functional groups or pharmacophores.
2. SMARTS is useful for filtering and selecting molecules based on specific criteria. For example, cheminformatics methods may use SMARTS patterns to exclude undesirable compounds or select molecules with certain properties for further analysis.

SYBYL Line Notation (SLN) is yet another linear notation to represent chemical structures, molecular fragments, reactions, formulations, molecular queries, reaction queries as well as virtual of *in silico* libraries [13]. Just like other linear notations, SLN was too inspired by SMILES but eventually differed substantially from its predecessor. Interestingly, SLN does not make any assumption pertaining to atomic valence and leaves the choice for representing a particular cohort of bonds, e.g., aromatic bonds, with the user. Aromaticity is treated as a property to atoms, as per SMILES conventions, while, SLN treats it as a property of bonds. This enables SLN to take the best of SMILES and SMARTS. The SLN syntax appeared in the literature in late 1990s. Figure 14.4 illustrates the SLN of few representative molecules along with the SMILES to give the user a basic understanding of the difference between the two. More details on the structural atom attributes, specification of stereochemistry, structural bond attributes, connection table attributes can be found elsewhere [13].

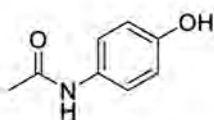
SLN is also favorably used for specifying reactions, reaction queries using atom and bond attributes. SLN is proposed as a preferred method for chemical structure encoding over computer networks. Several software

applications use SLN, e.g., SYBYL. SLN's capability for storing chemical structures in relational databases is particularly appreciated. SLN also uses a concept of macro-atoms, which represent atomic groups such as amino acids. For example, {Ala: NHC[s]l]H(CH<sub>3</sub>)C=O<v=1,9>} and {Gly:NHCH<sub>2</sub>C=O<v=1,6>} for alanine and glycine peptide residues. Overall, SLN represents a unique notation which is a hybrid of SMILES and SMARTS. The intriguing feature of SLN is the representation of a combinatorial SLN, representing multiple structures in one SLN. An illustrative example is shown in Figure 14.4. An interesting application of SLN were reported lately for chemical information management [14].



SLN: OC(CC[4]=C(NC[12])(=C(Cl)C=CC=C@12Cl)C=CC=C@5)=O

SMILES: OC(CC1=C(NC2=C(Cl)C=CC=C2Cl)C=CC=C1)=O



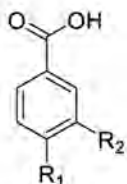
SLN: O=C(C)NC[3]=CC=C(O)C=C@4

SMILES: O=C(C)NC1=CC=C(O)C=C1



SLN: O[0]CC[2](CNC@3)C@1

SMILES: O1CC2(CNC2)C1



R<sub>1</sub> = Me, Et, MeOCH<sub>2</sub>

R<sub>2</sub> = Cl, OH, COOMe, CF<sub>3</sub>

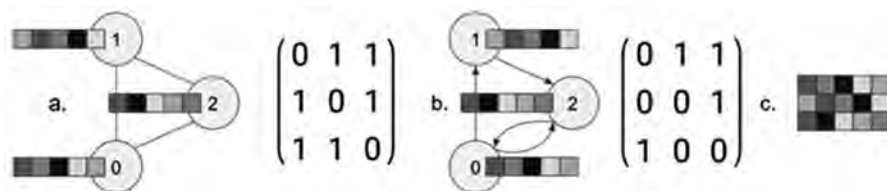
cSLN: C[1](C(C=O)OH):CH:C(G1):C(G2):CH:CH:@1\{G1:CH3<v=1>|CH2CH3<v=1>|CH2OCH3<v=1>}\{G2:Cl,OH,CH3OC(=O)<v=6>,FC(F)F<v=2>}

**FIGURE 14.4** Comparison of SLN and SMILES for select few molecules and the representative example of a combinatorial SLN.



### 14.2.2 GRAPH REPRESENTATIONS

Apart from linear strings (SMILES) and notations (SMARTS), molecules can also be represented as a graph. This way of representation borrows concepts and formalism from graph theory which describes a graph mathematically as  $G = (V, E)$ , where  $V$  is an ordered collection of all the nodes (in this case, set of all atoms present in a molecule), and  $E$  describes the connectivity between those nodes (in this case, set of all bonds connecting the atoms represented by  $V$ ) in the graph [15]. These connections may be weighted, i.e., have a particular number associated with them, or just represent a connection. There are two types of graphs, directed and undirected (Figure 14.5). The fundamental difference between them is that a connection between two nodes is mutual in an undirected graph. From a cheminformatics perspective, a molecule can be described as an undirected graph, with weighted connections, where the weights describe the bond order of the edge. Each node also has a node feature associated to it, which generally is a bit string of atomic descriptors like one hot encoded vectors representing type of atom (C, O, N, etc.), formal charge, number of hydrogens, hybridization of the atom, number of lone pairs, number of radical electrons and so on. The edges, too, can each have a feature matrix associated with them though, it is less commonly used. The most common deep learning (DL) method used for the purposes of cheminformatics is a Graph Convolution Network (GCN) [16]. The purpose of a GCN is to conduct message passing on each layer and scaling using weights which are learned. The applications for these are wide-ranging but outside the scope of this chapter.



**FIGURE 14.5** (a) Directed and (b) undirected graphs along with their adjacency matrices; and their (c) node feature.

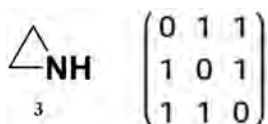
1. **Atomic Features:** The atomic features are of the utmost importance from an AI/ML perspective. They are the input for the algorithms and what will be predicted if that is the aim of the algorithm. These feature vectors are then concatenated to one another in a set predefined order

(the order of atomic index), and that is what will form the Node feature matrix, an  $n \times F$  matrix where  $n$  is the number of atoms and  $F$  is the length of the feature vector per atom. Keeping in mind the same order of atoms, a connectivity index is also generated. This connectivity matrix is an  $n \times n$  matrix which is not weighted. This connectivity matrix is also called the adjacency matrix [17], which has use cases even outside of AI/ML [18]. The adjacency matrix is what dictates the message passing of node feature vectors and their aggregation. PyTorch, a popular ML library, does not support an adjacency matrix as an input to the layer but rather takes in a tensor of size  $2e \times 2$  where  $e$  is the total number of connections in an undirected graph, and each entry in the tensor shows the starting index of the connection and receiving index (here an undirected graph is thought of as a directed graph). The atomic features can be atom describing or discriminating. For example, if the prediction to be made is at the entire graph level, having features that describe the atom are more useful, thus including things like electronegativity and hybridization of the atom helps, however if you are running a semi-supervised algorithm and the goal is to predict the type of node, just a one hot encoded vector describing the type of atom should be enough, as the node level predictions can directly be processed through a softmax function, which is an activation function that scales numbers or logits into probabilities.

- 2. Adjacency Matrices and Their Construction:** An adjacency matrix is a mathematical representation of a graph where the graph's nodes (vertices) are represented as rows and columns, and the presence or absence of edges (connections) between the nodes is denoted by entries in the matrix. In cheminformatics, adjacency matrices are used to represent molecular structures as graphs, with atoms as nodes and chemical bonds as edges. These matrices have found applications in cataloguing chemical reactions [19]. In cheminformatics, molecules can be represented as graphs, where atoms are nodes, and chemical bonds are edges connecting these nodes. An adjacency matrix provides an efficient way to represent these molecular graphs in a matrix format, capturing the connectivity information between atoms. An interesting fact about adjacency matrices of a molecule (actually, of all undirected graphs) is that the matrix is always symmetric.

To construct the adjacency matrix for a molecule, you create a square matrix with dimensions equal to the number of atoms in the molecule. Each row and column of the matrix correspond to an atom in the molecule. The

entry at row  $i$  and column  $j$  ( $A_{ij}$ ) is 1 if there is a chemical bond between atoms  $i$  and  $j$ , and 0 if there is no bond. You can also use  $A_{ij}$  equal to the bond type (2 for a double bond, 3 for a triple bond). Note: the ordering of atoms is crucial and should be consistent throughout. For example, consider a simple cyclic structure, aziridine (3), with three atoms ( $C_1CN_1$ ). The adjacency matrix of this would be as in Figure 14.6.



**FIGURE 14.6** Adjacency matrix for aziridine (3).

*Note:* Where the order of the atoms is [ $'C'$ ,  $'C'$ ,  $'N'$ ] (it does not matter here, but it does in general).

2. **Feature Vectors:** In ML, a feature vector is a numerical representation of data instance used as input for training models for making predictions [16]. It is a fundamental concept in various ML tasks, including classification, regression, clustering, and more. A feature vector combines different features, each representing a specific characteristic or property of the data, into a single, fixed-length vector. Few interesting characteristics of Feature Vectors include:
  - i. **Numerical Representation:** Feature vectors are composed of numerical values, which can be real numbers, integers, or binary values. This numeric representation allows the data to be processed and used in machine learning algorithms.
  - ii. **Fixed Length:** Feature vectors have a fixed length, meaning each data instance is represented by the same number and types of features. This uniformity is essential for ensuring compatibility and consistency when training machine learning models so that there is meaningful data on which learning can take place.
  - iii. **Feature Selection and Engineering:** The process of creating feature vectors involves feature selection and engineering. Feature selection involves choosing relevant and informative features that best describe the data, while feature engineering may involve transforming, scaling, or combining existing features to improve model performance.

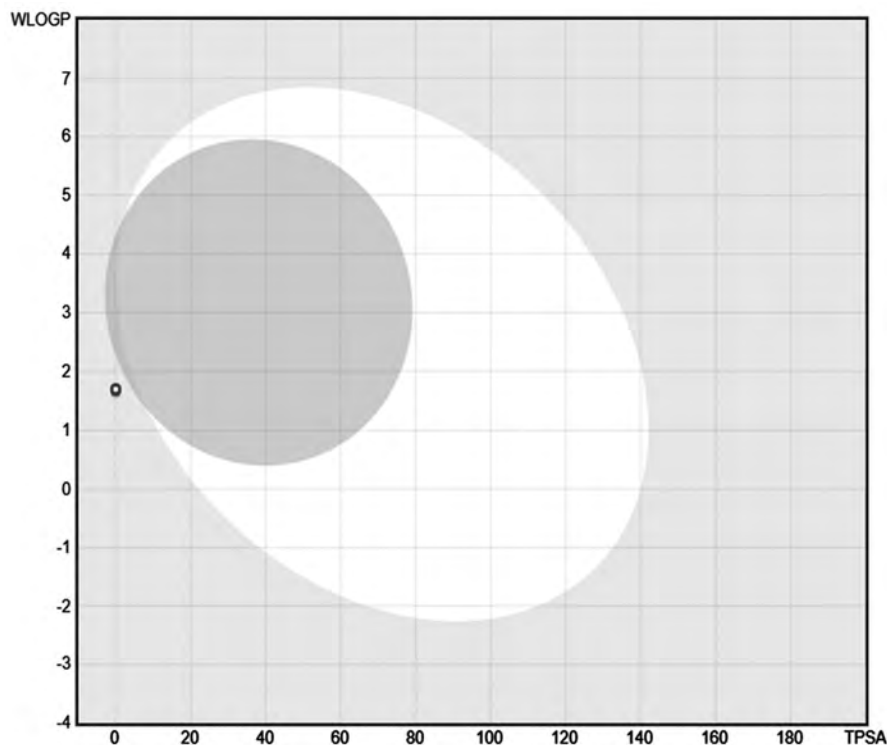
Feature vectors are, in a way, embeddings that represent molecules which can be used to train a Dense Neural Network [20], Recurrent Neural Networks [21] and much more. Feature

vectors can be of multiple types, they can either be learned embeddings through Variational Autoencoders (VAE) [22] that have been trained on small organic molecules, or they can be a chain of mathematical and/or physicochemical molecular descriptors. The feature vectors can also be fingerprints which are vectors that describe the presence or absence of some particular substructure. Each will be discussed in detail in the coming subsections. However, one should make note that the components of a feature vector are not just restricted to molecules but also to a mixture of compounds. For example, a feature vector generated in the way mentioned before for both solute and solvent can be concatenated together and fed into a Deep Neural Network (DNN) algorithm so as to generate results about the solute that are solvent specific.

- 3. Learned Embeddings:** A VAE, by Design, work on the principle that there exists a hidden distribution of data that the model tries to learn. It is composed of two parts, an encoder and a decoder. The encoder tries to encode the molecule into an  $n$ -dimensional Gaussian function centered around  $\mu$  with a standard deviation of  $\sigma$ , both of which are  $n$ -dimensional vectors describing the Gaussian curve in  $n$ -D. The decoder in turn, uses these generated embeddings to recreate the molecules. The algorithm is trained to generate the outputs of the decoder as closely as possible to the inputs of the encoder, processed in a sequential manner. Unlike a Generative Adversarial Network (GAN), here we can directly correlate the molecules to their embeddings and thus, have some meaningful representation of the molecule.

Feature vectors can also be in the form of a vector with descriptors of the molecule used to give 2D, 3D, and/or charge based information about the molecule. These descriptors can be either mathematically calculated like topological descriptors (Balaban Index [23], Weiner Index [24]), geometric descriptors (surface area, volume), electro-topological descriptors (dipole moment), etc., or they can be experimentally determined values like Melting point, Log P (partition coefficient), LogD (distribution coefficient), and ALogS (aqueous solubility). Popular cheminformatics libraries like RDKit [25] or software's such as Dragon [26], Molecular Operating Environment (MOE) [27] have functionalities that help calculate mathematical descriptors, and Python APIs like Leruli [28] help provide queries to a database of experimental information. A lot of web-scrappers can also be employed to help in data-mining, however, their use is debated. Services like swissADME

[29] also provide the user with a lot of useful information when queried with molecules. One of the useful things they generate is a boiled egg diagram. It is a plot between the total polar surface area (TPSA) on the x-axis and WLogP values on the y-axis (Figure 14.7). The graph has two regions of interest, namely the white region, which shows all the possible combinations of values for which the molecule is absorbed by our gastrointestinal tract, while the inner yellow portion (yolk) shows the values which are blood-brain-barrier permeable. This is a good way to visualize a molecule given an embedding. One can train an algorithm to arrive at these two values using the embedding or perhaps use previously documented algorithms that predict these so as to visualize these molecules in the 2D plane. One can also use principle component analysis to reduce dimensions and visualize the molecules. It can provide help in debugging cases and also help understand the algorithm better.



**FIGURE 14.7** Boiled egg diagram for benzene.

*Source:* Generated using SwissADME.

### 14.2.3 FINGERPRINTS

Molecular fingerprints are essential tools used in cheminformatics and computational chemistry to encode the structural information of molecules into a format that can be efficiently compared and analyzed. These fingerprints are numerical representations of molecular structures, enabling researchers to perform similarity searches, virtual screening, and machine learning algorithms for drug discovery, clustering, and other cheminformatics applications [30]. Molecular fingerprints encode molecular information based on the presence or absence of specific substructural features or chemical fragments within a molecule. Different types of fingerprints employ distinct algorithms and methodologies to generate these representations. Here are some commonly used molecular fingerprints:

1. **Morgan (Circular) Fingerprints:** This, also known as circular fingerprints or ECFP (Extended-Connectivity Fingerprints) [31], are one of the most widely used fingerprinting methods. They are based on the concept of topological substructures. The algorithm works by iteratively generating circular atom neighborhoods around each atom in the molecule up to a specified radius. The presence or absence of these substructures within the molecule is recorded in the fingerprint.
2. **MACCS (Molecular ACCess System) Keys:** These keys are a fixed-length fingerprinting system developed by Molecular Design Limited (MDL). It uses a predefined set of 166 structural keys (0 or 1) to represent various substructural patterns. MACCS keys are particularly useful for similarity searching and have been employed in many virtual screening studies [32].
3. **Daylight Fingerprints [33]:** This, also known as substructure keys, are based on encoding substructures and patterns in a molecule. It uses a linear notation for the representation of these patterns. Each substructure in the molecule is hashed to generate a unique identifier, which is then used to create the fingerprint.
4. **Topological Torsion Fingerprints [34]:** These capture the topological features of molecules by encoding patterns of bonded atom triplets. The presence or absence of these torsion patterns is used to generate the fingerprint.
5. **Pharmacophore Fingerprints [35]:** These are based on the concept of pharmacophore modeling. These fingerprints encode the essential features responsible for a molecule's biological activity, such as hydrogen bond donors, acceptors, aromatic rings, etc.

Each type of fingerprint has its strengths and weaknesses, depending on the specific use case and the desired level of molecular information captured. Researchers choose the most appropriate fingerprint type based on their specific research objectives and the characteristics of the molecular data being analyzed. In practice, multiple fingerprint types are often used in combination to leverage their complementary information for improved performance in cheminformatics applications. These fingerprints employ various techniques to generate the vector. However, many of them can be employed and customized by using RDKit and SMARTS substructure search methods to generate a vector, which can be used as a molecular descriptor for learning, filtering, sorting and even for storage. One of the advantages of using a custom fingerprint is that it cannot just be a mixture of above-mentioned standardized fingerprints but can also encode information that one feels could be crucial for accurate learning, such as functional groups, essential drug-molecule substructures like morpholine or pyridine.

#### 14.2.4 MISCELLANEOUS

1. **Unconventional Descriptors:** When training a ML algorithm, often providing all the information at hand is a wise choice. Though this can lead to underfitting sometimes, there are other methods to address that issue. No single representation is complete, and thus many a time, people prefer to use multiple representations of molecules in conjunction. For example, while MACCS fingerprints give a good idea of the molecule and its substructures, using them together with molecular descriptors has often proven to give more accurate models.
2. **GeoGNN:** A general trend has been noticed in the representations discussed so far is that they generally aim to encode the spatial arrangements of the atoms or substructures. However, they lose out on the information of bond angles. According to the graph representation of the molecules we have seen so far, it is not possible to encode bond angles, as edge features include bond length or bond order, while the node features include atomic information. Thus, a new graph has to be made, where the nodes represent pair of atoms (a bond), and the edges represent the presence of one common atom. Let us call this new graph a 'Meta' graph. Now we can encode bond angle, angle strain and many more properties about the angle between any three atoms. This way, we can either run a parallel network that processes those angles and gives us an output or uses meta-message

passing, where information from another graph is incorporated into the molecular graph during the message passing and aggregation step, like what has been described by Fang et al. [36].

3. **Persistent Homologies:** Imagine one has a set of data points scattered around in space, and she wants to understand the shapes and patterns that exist in this data. One way to do this is by using a mathematical tool called ‘persistent homology’ [37]. It is a technique from a branch of mathematics called ‘topology,’ which is all about understanding the properties of shapes and spaces. The field of Topology is a vast ocean, so diving into it is beyond the scope of this chapter. To apply persistent homology, one first creates a family of shapes (generally a Gaussian kernel or a spherical kernel) around the data points. These shapes are like little bubbles that grow and shrink, covering the data points. As one changes the size of the bubbles, some features, like holes (like a doughnut), will appear, change shape and size, and eventually disappear. The key idea of persistent homology is to track these features as we change the size of the bubbles. One keeps track of when features appear and how long they stick around at different scales. Persistent homology helps to identify the most persistent features, which are likely to be more meaningful and representative of the underlying shape in the data.

### 14.3 SUMMARY

The present discussion provided a comprehensive overview of the molecular representation techniques in cheminformatics and drug discovery. It highlighted the significance of molecular representation in various applications, including AI/ML, data storage and drug development. Next, the generic methods like molecular formulas and their drawbacks, were explored, leading to the introduction of SMILES as a more informative and concise representation. The advantages and disadvantages of SMILES were outlined, followed by an explanation of InChI and InChI Key as standardized textual representations with unique identifiers for chemical structures. The chapter delved into the detailed methodology of determining SMILES and InChI representations from molecular structures. The discussion extended to SMARTS, an extension of SMILES used for substructure searching and pattern matching in chemical databases. It also covered graph representation of molecules through node feature matrices, adjacency matrices and their significance in cheminformatics. The use and generation of feature vectors



through the use of descriptors, fingerprints and a mixture of both was mentioned. The concepts of VAE used in learning embeddings that could perhaps give a more useful representation computationally. Additionally, it introduced various mathematical descriptors, boiled egg diagrams, and the utilization of fingerprints for molecular encoding. The chapter concluded by explaining the concept of persistent homologies and its application in understanding the patterns and shapes present in molecular data. Overall, the chapter provided a thorough exploration of the diverse methods for representing molecules in cheminformatics.

#### **14.4 CONCLUSIONS AND FUTURE PERSPECTIVES**

The molecular/structural representation has evolved over fairly large period to suit the need of present-day technologies such as AI, ML and DL. Simple linear notations provided the simplicity and convenience in terms of memory requirements and retrieval ease. These were superseded by complex, at times, proprietary representations to overcome some of the limitations of their predecessors, which made them more robust and machine-readable. Increased computational power catalyzed their evolution even further. The graph theory formalism increased the capabilities of the already used structure representation methods, delving into core areas of AI/ML, further helping the expansion of their applicability domain. Various graph representation methods have been incorporated in modern-day DNN tools. The chemist's e-toolkit now contains various routines based on newer structure representations, few of which are specific for the intended purpose, e.g., property predictions. These tools overcome the inherent limitations of various chemometric methods of old times.

The use of many structure representations in reaction representation has helped the scientific community in general and chemists in particular, for disseminating the enormous chemical information which is generated in real-time on massive scale. The resulting searches of the open- and patent literature is way much faster than one could imagine. Thanks to these AI/ML compatible structure representation modalities! Another obvious application area of structure representation is computer-assisted structure elucidation, which has made substantial progress in last two decades. The representation of stereochemistry in complex structures such as natural products, can now be dealt with confidence with the help of modern structural representations. The Markush Structure enumeration is one sought-after application area in chemical information management. The use of AI/ML coupled with

newer AI/ML-compatible structure representations have made the Markush enumeration easier and accurate. The obvious explosion in the number of known chemicals and the utility of the machine-friendly structure representations have contributed in managing chemical information effectively. Future innovations in this area will further fuel the growth of chemical information management discipline.

## KEYWORDS

- **artificial intelligence**
- **deep learning**
- **deep neural network**
- **graph convolution network**
- **InChI key**
- **machine learning**
- **molecular structure**
- **simplified molecular input line entry system**
- **SMILES arbitrary target specification**

## REFERENCES

1. Chemical Abstracts Service. (2007). *National Historic Chemical Landmark*. Retrieved from: <https://www.acs.org/education/whatischemistry/landmarks/cas.html> (accessed on 25 July 2024).
2. Kim, S., Chen, J., Cheng, T., Gindulyte, A., He, J., He, S., Li, Q., Shoemaker, B. A., Thiessen, P. A., Yu, B., Zaslavsky, L., Zhang, J., & Bolton, E. E. (2023). PubChem 2023 update. *Nucleic Acids Research*, 51(D1), D1373–D1380. <https://doi.org/10.1093/nar/gkac956>.
3. PubChem. (2023). Retrieved from: <https://pubchem.ncbi.nlm.nih.gov/> (accessed on 25 July 2024)
4. Lawlor, B. (2016). The chemical structure association trust. *Chemistry International*, 38(2), 12–15. <https://doi.org/10.1515/ci-2016-0206>.
5. Wiswesser, W. J. (1968). 107 years of line-formula notations (1861–968). *Journal of Chemical Documentation*, 8(3), 146–150.
6. David, L., Thakkar, A., Mercado, R., & Engkvist, O. (2020). Molecular representations in AI-driven drug discovery: A review and practical guide. *Journal of Cheminformatics*, 12(1), 56. <https://doi.org/10.1186/s13321-020-00460-5>.

7. (a) An, X., Chen, X., Yi, D., Li, H., & Guan, Y. (2022). Representation of molecules for drug response prediction. *Briefings in Bioinformatics*, 23(1), bbab393. <https://doi.org/10.1093/bib/bbab393>. (b) Chuang, K. V., Gunsalus, L. M., & Keiser, M. J. (2020). Learning molecular representations for medicinal chemistry. *Journal of Medicinal Chemistry*, 63(16), 8705–8722. <https://doi.org/10.1021/acs.jmedchem.0c00385>.
8. Santa Maria, J. P., Jr, Wang, Y., & Camargo, L. M. (2023). Perspective on the challenges and opportunities of accelerating drug discovery with artificial intelligence. *Frontiers in Bioinformatics*, 3, 1121591. <https://doi.org/10.3389/fbinf.2023.1121591>.
9. Gelberg, A. (1970). Chemical notations. In A. Kent & H. Lancour (Eds.), *Encyclopedia of Library and Information Science* (Vol. 4, pp. 510–528). Marcel Dekker.
10. Weininger, D. (1988). SMILES, a chemical language and information system: 1. Introduction to methodology and encoding rules. *Journal of Chemical Information and Computer Science*, 28(1), 31–36.
11. Heller, S. R., McNaught, A., Pletnev, I., Stein, S., & Tchekhovskoi, D. (2015). InChI, the IUPAC International Chemical Identifier. *Journal of Cheminformatics*, 7(1), 23. <https://doi.org/10.1186/s13321-015-0095-0>.
12. (2023). Daylight Theory: SMARTS—A language for describing molecular patterns. Retrieved from: <https://www.daylight.com/dayhtml/doc/theory/theory.smarts.html> (accessed on 25 July 2024).
13. Homer, R. W., Swanson, J., Jilek, R. J., Hurst, T., & Clark, R. D. (2008). SYBYL line notation (SLN): A single notation to represent chemical structures, queries, reactions, and virtual libraries. *Journal of Chemical Information and Modeling*, 48(12), 2294–2307. <https://doi.org/10.1021/ci7004687>.
14. Kochev, N., Jeliakova, N., & Tancheva, G. (2021). Ambit-SLN: An open-source software library for processing of chemical objects via SLN linear notation. *Molecular Informatics*, 40(11), e2100027. <https://doi.org/10.1002/minf.202100027>.
15. Kay, E., Bondy, J. A., & Murty, U. S. R. (1977). *Graph Theory with Applications*. Operational Research Quarterly (1970–1977), 28, 237.
16. Zhang, S., Tong, H., Xu, J., et al. (2019). Graph convolutional networks: A comprehensive review. *Computational Social Networks*, 6, 11. <https://doi.org/10.1186/s40649-019-0069-y>.
17. Berinde, Z., & Mădălina, B. (2004). On a matrix representation of molecular structures. *Carpathian Journal of Mathematics*, 20(2), 205–209. <http://www.jstor.org/stable/43996749>.
18. Kim, Y., Jeong, Y., Kim, J., Lee, E. K., Kim, W. J., & Choi, I. S. (2022). MolNet: A chemically intuitive graph neural network for prediction of molecular properties. *Chemistry—An Asian Journal*, 17(16), e202200269. <https://doi.org/10.1002/asia.202200269>.
19. Ismail, I., Chantreau Majerus, R., & Habershon, S. (2022). Graph-driven reaction discovery: Progress, challenges, and future opportunities. *The Journal of Physical Chemistry A*, 126(40), 7051–7069. <https://doi.org/10.1021/acs.jpca.2c06408>.
20. Krasteva, V., Christov, I., Naydenov, S., Stoyanov, T., & Jekova, I. (2021). Application of dense neural networks for detection of atrial fibrillation and ranking of augmented ECG feature set. *Sensors (Basel, Switzerland)*, 21(20), 6848. <https://doi.org/10.3390/s21206848>.
21. Otović, E., Njirjak, M., Kalafatovic, D., & Mauša, G. (2022). Sequential properties representation scheme for recurrent neural network-based prediction of therapeutic peptides. *Journal of Chemical Information and Modeling*, 62(12), 2961–2972. <https://doi.org/10.1021/acs.jcim.2c00526>.

22. Vogt, M. (2022). Using deep neural networks to explore chemical space. *Expert Opinion on Drug Discovery*, 17(3), 297–304. <https://doi.org/10.1080/17460441.2022.2019704>.
23. Thakur, A., Thakur, M., Khadikar, P. V., Supuran, C. T., & Sudele, P. (2004). QSAR study on benzenesulphonamide carbonic anhydrase inhibitors: Topological approach using Balaban index. *Bioorganic & Medicinal Chemistry*, 12(4), 789–793. <https://doi.org/10.1016/j.bmc.2003.10.058>.
24. Mandloi, M., Sikarwar, A., Sapre, N. S., Karmarkar, S., & Khadikar, P. V. (2000). A comparative QSAR study using Wiener, Szeged, and molecular connectivity indices. *Journal of Chemical Information and Computer Sciences*, 40(1), 57–62. <https://doi.org/10.1021/ci980139h>.
25. Lovrić, M., Molero, J. M., & Kern, R. (2019). PySpark and RDKit: Moving towards big data in cheminformatics. *Molecular Informatics*, 38(6), e1800082. <https://doi.org/10.1002/minf.201800082>.
26. DRAGON 7.0. (2024). Retrieved from: <https://chm.kode-solutions.net/pf/dragon-7-0/> (accessed on 25 July 2024).
27. Molecular Operating Environment (MOE). (2024). Retrieved from: <https://www.chemcomp.com> (accessed on 25 July 2024).
28. Leruli. (2024). Retrieved from: <https://leruli.com/> (accessed on 25 July 2024).
29. Daina, A., Michielin, O., & Zoete, V. (2017). SwissADME: A free web tool to evaluate pharmacokinetics, drug-likeness and medicinal chemistry friendliness of small molecules. *Scientific Reports*, 7, 42717. <https://doi.org/10.1038/srep42717>.
30. Xue, L., Godden, J. W., & Bajorath, J. (2003). Mini-fingerprints for virtual screening: Design principles and generation of novel prototypes based on information theory. *SAR and QSAR in Environmental Research*, 14(1), 27–40. <https://doi.org/10.1080/1062936021000058764>.
31. Rogers, D., & Hahn, M. (2010). Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5), 742–754. <https://doi.org/10.1021/ci100050t>.
32. Carracedo-Reboredo, P., Liñares-Blanco, J., Rodríguez-Fernández, N., Cedrón, F., Novoa, F. J., Carballal, A., Maojo, V., Pazos, A., & Fernandez-Lozano, C. (2021). A review on machine learning approaches and trends in drug discovery. *Computational and Structural Biotechnology Journal*, 19, 4538–4558. <https://doi.org/10.1016/j.csbj.2021.08.011>.
33. Daylight Fingerprints. (2024). Retrieved from: <https://www.daylight.com/meetings/summerschool01/course/basics/fp.html>.
34. Koda, P., & Hoksza, D. (2015). Exploration of topological torsion fingerprints. In *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)* (pp. 822–828). Washington, DC, USA. <https://doi.org/10.1109/BIBM.2015.7359792>.
35. Yang, J., Cai, Y., Zhao, K., Xie, H., & Chen, X. (2022). Concepts and applications of chemical fingerprint for hit and lead screening. *Drug Discovery Today*, 27(11), 103356. <https://doi.org/10.1016/j.drudis.2022.103356>.
36. Fang, X., Liu, L., Lei, J., et al. (2022). Geometry-enhanced molecular representation learning for property prediction. *Nature Machine Intelligence*, 4, 127–134. <https://doi.org/10.1038/s42256-021-00438-4>.
37. Townsend, J., Micucci, C. P., Hymel, J. H., et al. (2020). Representation of molecular structures with persistent homology for machine learning applications in chemistry. *Nature Communications*, 11, 3230. <https://doi.org/10.1038/s41467-020-17035-5>.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

# Index

---

## A

- Abiotic
  - degradation, 303
  - pesticide degradation processes, 250
- Academic organizations, 77
- Accurate metabolic processes, 170
- Acetylation, 158, 159
- Acetylcholinesterase inhibitor, 26
- Acute
  - oral toxicity, 179, 182, 202
  - toxicity, 128, 130, 200, 201, 242
- Adaptive neuro-fuzzy inference system (ANFIS), 5, 8, 9, 12
- Adenosine triphosphate, 49
- ADMET predictor, 155, 179, 192, 200–202
- AdmetSAR*, 128, 130, 198
- Adversarial autoencoder, 332
- Aerobic
  - biodegradation, 294
  - microbial biotransformation, 305
  - microorganisms, 294
- Agitation velocity, 10
- Agranulocytosis, 188
- Agrochemicals, 58, 256
- Algorithm, 5, 6, 13, 22, 24, 25, 48, 59–65, 68, 69, 76, 77, 86–92, 104, 105, 117–122, 127, 132–134, 136, 137, 143, 148, 149, 151, 152, 180, 181, 183, 191, 194, 197–201, 209, 210, 216, 225, 238, 248, 250, 251, 253, 254, 257, 258, 262, 264–268, 271, 277, 280, 291, 292, 297–299, 302, 304, 312, 323, 326, 337, 339, 341, 342, 348, 351, 366, 368, 369, 370, 371
  - development, 19
- Ambient variables, 251
- Amino acids, 83, 364
- Amlodipine, 186
- Ampicillin, 48, 49
- Anaerobic
  - biodegradation, 294
  - microbial biodegradation, 305
- Analytical chemical science, 155
- Androgen receptor, 128
- Angiotensin II receptor antagonism, 27
- Animal experimentation, 102, 136
- Anomaly detection, 273, 341, 351, 352
- Anthelmintic activity, 27
- Antibacterial properties, 47
- Antifungal
  - activity, 48
  - movement, 47
  - properties, 27
- Antimalarial activity, 27
- Antimicrobial
  - activity, 27, 46
  - drugs, 27
- Antioxidant response element, 128
- Antitumor effects, 27
- Antiviral properties, 27
- Applicability domain (AD), 194, 209, 238, 239, 304, 374
- Application toolbox, 195
- Area under curve (AUC), 134, 135, 137, 234, 235, 237
- AromaDeg*, 298, 311
- Aromatic hydrocarbon degradation database (AHDD), 299, 311
- Aromaticity, 364
- Artificial intelligence (AI), 1, 3–13, 62, 63, 73–80, 82–94, 99, 104, 109, 141–143, 146–152, 154–156, 168, 169, 209–214, 225, 226, 238–241, 247–261, 263–266, 270–280, 291–293, 298, 299, 306, 311, 312, 319, 322–326, 330, 331, 337, 338, 341–344, 348–353, 357, 360, 366, 367, 373, 374
  - algorithms, 90, 258, 298
  - application, 73, 77, 80, 83, 84, 311
  - assisted inventions, 75, 78
  - biodegradation
    - databases, 306
    - prediction models, 253, 255, 264, 270, 271, 279

driven drug development, 87  
 enabled tools, 13  
 fault diagnosis  
   software, 351  
   systems, 341, 342, 352  
 medication metabolism, 148  
 patient filing, 89  
 powered  
   chatbots, 92  
   technologies, 92  
 systems, 76, 86, 89–91, 93, 94, 148, 150  
 techniques, 10, 13, 142, 146, 148, 152,  
   169, 252, 258, 265, 266, 279, 291, 292,  
   312, 337, 341–343, 348, 349, 353  
 technology, 13, 75, 87, 92, 93, 147, 150, 323  
 tools, 141, 170  
 Artificial neural network (ANN), 3–5, 7,  
   104, 105, 111, 120–122, 143, 169, 170,  
   253, 265, 280, 291, 292, 297, 298, 312,  
   313, 322–324, 327, 332, 342, 343, 345  
 Aryl hydrocarbon receptor (AhR), 128,  
   134–136  
 Aspen Mtell, 351  
*Aspergillus flavus* (AF), 47–49, 104  
 Asset performance management (APM), 351  
 Association rule mining, 265, 269, 270  
 Astemizole, 181  
 AstraZeneca, 88  
 Atomic  
   features, 366  
   level models, 76  
   number, 222  
   partial charges, 225  
 Attention convolutional neural networks  
   (ACNN), 120  
 Augmented  
   diagnostic tool, 12  
   reality (AR), 128, 134–136  
 AutoDock, 24  
 Autoencoders, 329  
 Automated  
   feature learning, 269  
   reaction synthesis, 322  
   retrosynthetic programs, 22  
   synthesis methods, 331  
 Automatic accurate computerized prediction  
   system, 136  
 Autotrophic microorganisms, 294

## B

Bayesian  
   classifiers, 104, 143  
   principle, 125  
 Benchmarks, 12, 144, 213  
 Benigni-Bossa rule base, 196  
 Benzylisoquinoline alkaloids, 153  
 Bile acid conversions, 160  
 Binary  
   array, 220  
   classification, 105, 106, 214, 226, 236, 267  
   values, 368  
 Binding affinity, 25, 27, 28, 129, 253  
 BindingDB, 153  
 Bioaccumulation, 303–305  
 Bioactivation, 186  
 Bioanalyte signal amplification, 80  
 Biochemical  
   characteristics, 144  
   management, 308  
   production calculation, 323  
 Bioconcentration factor, 194  
 Biodegradability, 249, 257, 265, 266, 293, 304  
   predictor, 302, 303  
 Biodegradation, 99, 194, 247–271, 273–280,  
   291–299, 301, 304–309, 312, 313  
   data, 263, 264, 275  
   databases, 306  
   datasets, 264  
   degree, 262, 291  
   forecast systems, 297  
   mechanism, 295  
   modeling, 293  
   network-molecular biology database, 298,  
     308  
   outcomes, 249, 269, 270, 291, 292, 297, 298  
   pathways, 275, 277, 278, 299, 306, 307, 309  
   prediction, 247–255, 257–271, 273–280,  
     291–293, 297–299, 312  
   AI techniques, 265  
   data mining techniques, 269  
   machine learning algorithms, 266, 268  
   models, 251, 253, 254, 259, 262–264,  
     270, 271, 277–279  
   processes, 250–254, 260, 262, 268, 277,  
     278  
   rates, 249–251, 253–263, 269, 274, 275,  
     277, 279, 280

- reactions, 308
- techniques, 298
- types, 294
- Biodegradative oxygenases, 298, 310
- Biodiesel
  - production, 12
  - reactor, 12
  - reactors, 8
- Biofuels, 249, 257
- Biogas, 249
- Bioinformatics, 83, 143, 169, 226, 309, 310
  - Molecular Design Research Center (BMDRC), 197
  - technologies, 169
  - tools, 168
- Biological
  - activity, 91, 104, 107, 185, 197, 254, 323, 329, 330, 371
  - chemical property predictions, 169
  - degradation prediction models, 265
  - evaluation
    - antibacterial activity, 47
    - antifungal activity, 47
  - half-life, 141
  - microbes, 292
  - nervous systems, 121
  - oxygen demand (BOD), 262, 280, 303
  - technology, 293
- Biomarkers, 89, 151
- Bionemo, 298, 308, 309
- Bioremediation, 294
  - chemical structure, 313
  - processes, 293, 307, 310, 311
  - strategies, 251, 310
  - techniques, 308, 309
- Biotechnology, 79, 92, 93, 190, 250, 257, 270, 271, 309, 310, 324
- Biotic degradation, 303
- Biotransformation, 155, 157, 167, 188, 294–297, 306–309
- Biotransformer, 146, 154–158, 160, 167, 169
  - program, 157
- Biowin, 298, 304
- Black boxes, 93, 273, 350
- Blood
  - brain-barrier permeable, 370
  - pressure, 151
- Bogus pharmaceuticals, 91

- Boltzmann distribution, 225
- Bond type, 222, 363, 367
- Bonded atom triplets, 371
- Bone marrow toxicity, 188
- Bootstrap samples, 130
- Bootstrapping method, 132
- Broad spectrum, 13, 199

## C

- Cancer expert system, 195
- Candida albicans (CA), 47, 48
- Carbon, 58, 194, 248, 292, 294, 299, 304, 309, 311, 362
  - compounds, 292, 299, 309
  - dioxide, 248, 292, 294
  - natural cycle, 292
- Carbutamide amine group, 189
- Carcinogenesis, 242
- Carcinogenicity, 128, 130, 180, 194–196, 199–201, 210, 213, 242
- Cardiotoxicity, 179, 181, 202
- Cardiovascular illness, 151, 152
- Catabolic modifications, 303
- Catalysis, 76, 322
- Cell signaling, 160
- ChemAxon, 191
- ChemBioFinder, 153
- Chemical
  - abstracts service (CAS), 358
  - category, read across, and trend analysis, 185
  - compounds, 76, 117–119, 127–129, 136, 180, 182, 183, 190, 197, 198, 200, 201, 203, 254, 261, 262, 269, 274, 275, 278, 293, 302, 325, 362
  - database, 76, 190, 201, 257, 261, 306, 358, 361, 362, 373
  - search, 358
  - datasets, 214
  - descriptors, 118
  - engineering, 13
    - domain, 3
  - ecotoxicity, 304
  - health
    - care and safety, 137
    - human safety, 136
  - industry, 66, 343, 348–353



- information management, 190, 322,  
357–361, 365, 374
  - discipline, 374
- intoxication, 120
- metabolism, 154
- pharmaceutical science, 73
- process plant, 337–343, 349, 351–354
- properties, 256
- reaction optimization, 88
- reaction, 4–7, 21, 57–61, 63, 64, 68, 69,  
90, 153, 312, 367
  - automation, 21
  - prediction, 90
- reactivity, 20, 63
- risk assessment, 184, 248
- safety
  - assessment, 180
  - evaluation, 197, 200
- sciences, 3–5, 94, 311
  - AI impact, 76
- structure, 63, 87, 110, 153, 182, 190,  
195–201, 214, 250, 251, 253, 255, 256,  
260–262, 264, 268, 295, 297, 361, 362,  
364, 373
  - activity relationships, 198, 199
  - descriptors, 104
- synthesis, 20, 90
- test tube artificial intelligence cleaning  
devices, 80
- toxicity, 117, 119, 127, 134–136, 179,  
202, 203
- toxicokinetic endpoints, 225
- toxicology, 118
- transformations, 212
- ChemIDplus database, 182
- Cheminformatics, 220, 225, 238, 240, 242,  
324, 364, 366, 367, 369–374
  - applications, 225, 371, 372
- Chemistry, 57, 58, 60, 73, 74, 83, 107, 147,  
155, 159, 185, 196, 212, 217, 308, 322,  
323
- Chemoinformatics, 183, 191
  - tools, 191
- Chemotherapy agents, 158
- ChemSpider, 153, 261, 306
- Chromosome segregation, 27
- Chronic illness course, 152
- Cispride, 181
- Clinical
  - data, 86
  - decision-making systems, 89
  - pharmacy, 324
  - records, 88
  - trial, 82, 91–93, 102, 110, 136, 142, 180,  
181, 209, 210, 312
    - optimization, 87
    - volunteer prescreening, 87
- Cluster, 126, 127, 217, 218, 265, 269, 270
  - analysis, 104
- Coagulation factor, 27
- Code snippets, 218, 222, 235
- Cognite data fusion, 351
- Combinatorial library screening, 24
- Commensurate strategies, 9
- Commercial
  - chemicals, 159
  - databases, 212, 213
  - medicinal chemistry databases, 212
- Complex
  - compounds, 21
  - mathematics-based equations, 4
  - molecular graph theory, 357
  - molecules, 59
  - organic
    - substances, 248
    - syntheses, 62
- Compound identification, 91
- Comprehensive risk assessment, 202
- Computation, 106
- Computational
  - approaches, 127, 169, 179, 180, 182, 202,  
203
  - chemistry, 60, 87, 107, 322, 362, 370
  - databases, 262
  - drug innovation, 25
  - efficiency, 68, 225, 276, 278
  - high-throughput methods, 330
  - methods, 58, 61, 66–68, 109, 180, 182,  
209, 265
  - modeling methods, 184
  - models, 69, 117, 128, 197–199, 209, 210,  
225
  - organic synthesis, 67
  - resources, 104, 269, 276, 350
  - simulations, 59, 251
  - software and techniques, 61
  - techniques, 146, 240

- successful applications, 66
- tools, 57–61, 63–71, 102, 180, 182, 183, 200
- toxicology, 202
- Computer
  - aided
    - combination applications, 22
    - design (CAD), 57, 62
    - drug design, 20, 49
    - fusion schemes, 21
    - medicinal innovation, 324
    - reaction prediction, 49
    - synthesis planning (CASP), 21, 22, 49, 61, 62, 71
  - assisted
    - organic synthesis (CAOS), 66
    - retrosynthesis planning, 64, 65
    - structure elucidation, 357, 374
  - controlled robots, 21
  - implemented innovations, 74
  - learning programme, 155
  - program, 65, 154
  - related
    - information, 94
    - invention, 74
  - systems, 13, 73, 94, 142
  - vision, 328
- Computing methods, 168
- Concise representation, 373
- Condition-based maintenance (CBM), 338
- Conformal predictors, 106–108
- Conformational exploration, 24
- Conjugation reactions, 158
- Connectivity matrix, 366
- Conventional tools, 13
- Convolutional
  - layers, 241
  - neural network (CNN), 106, 119, 120, 182, 222, 241, 253, 265, 268, 269, 280, 343, 337
- Cost-effectiveness, 147, 254, 255
- Cramer scheme, 196
- Cross-validation, 258
- Custom fingerprint, 372
- Cypreact, 154, 155
- CyProduct forecasts, 154
- Cytochrome P450 (CYP450), 143, 150, 154, 155, 157, 158, 160
  - mediated metabolism, 155
- Cytotoxicity, 128, 130, 198

**D**

- Daphnia magna*, 197
- Data
  - acquisition, 169, 257, 259, 277
  - aggregation, 19
  - analysis, 75, 92, 148, 191, 201, 242, 250, 312, 339–342
  - analytics, 190, 203, 212, 213, 279, 338, 340, 341, 352, 353
  - anonymization, 264
  - augmentation techniques, 258
  - availability and quality, 275, 349
  - centric metabolic prediction strategy, 142
  - collection, 211, 250, 258–260, 277, 280, 324, 340
  - driven
    - decisions, 76, 217
    - forecasting, 322
    - learning, 353
  - generation, 258, 264
  - integration, 200, 264, 265, 350
  - management, 190
  - mining techniques, 104, 141, 142, 145, 146, 200, 254, 265, 266, 269, 270, 279, 326
  - normalization, 258, 264, 265
  - pre-processing, 132, 190, 211, 214, 254, 259, 263–265, 342
    - challenges, 263
  - privacy, 93, 252, 261, 264, 265, 273, 275
  - processing, 189, 211, 213, 214, 216, 263
  - quality, 88, 248, 257, 271
    - control techniques, 258
  - set, 10, 160, 181
  - standards, 154
  - synthesis, 258
  - transformation, 258
- DataMol, 218, 221
- Dataset, 64, 69, 76, 86, 87, 89, 92, 107, 117, 119, 129, 136, 137, 141, 145, 152–154, 184, 195, 201, 209–213, 217, 221, 226, 231, 241, 250, 251, 253–255, 258, 261, 262, 264–266, 269, 271, 273–277, 291–293, 297, 298, 312, 323, 324, 349, 350
  - collection, 209
  - curation, 212
  - preparation, 212
- Daylight fingerprints, 371

- Decision  
  making, 89, 92, 93, 102, 199, 200, 240, 241, 251, 255, 260, 261, 266, 270, 273, 274, 276, 277, 280, 299, 304, 323, 350, 352–354  
  support systems, 90  
  tree, 105, 123–125, 129–132, 137, 183, 196, 231, 232, 249, 253, 265–268, 271, 342
- Deep  
  learning, 80, 93, 101, 102, 119, 120, 144, 145, 149, 170, 192, 201, 211, 212, 224, 234, 238, 241–243, 253, 254, 265, 266, 268, 269, 273, 276, 277, 279, 280, 293, 297, 313, 323, 324, 329, 330, 337, 343, 346–350, 353, 354, 366, 375  
  fault diagnosis, 346  
  frameworks, 192  
  models, 101, 265, 343, 350  
  techniques, 144, 269, 337, 347, 353  
  neural network (DNN), 144, 182, 201, 328–330, 343, 369, 374, 375  
  reinforcement learning, 352  
  standard learning (DL), 147, 148, 169, 181, 211, 360, 366, 374
- Deepchecks, 153
- DeepChem, 211, 218, 221–223, 234, 235, 243  
  library, 221
- DeepTox, 201
- Degradation, 311  
  pathways, 256, 262, 263, 303, 311, 359
- Dendrogram, 127
- Dense neural network, 368
- Density functional theory, 65
- Deoxyribonucleic acid (DNA)  
  gyrase, 27, 28  
  topology, 27  
  transcription, 27
- Derek nexus, 179, 186, 199, 202
- Design novel molecules, 69
- Detection, 150
- Detoxification, 157, 309
- Developmental toxicity, 197
- Diagnostic problem-solving, 323
- Digital humanities quarterly (DHQ), 26–29, 46, 47  
  heterocyclic system, 27
- Dihydropyridine structure, 186
- Dimensionality reduction, 258, 264, 270
- Dioxygenases, 310  
  directory, 310
- Disease pathways, 86
- Distillation columns, 8
- Docking  
  process, 26, 27  
  software options, 326  
  techniques, 26
- Donepezil, 26
- Dose  
  modifications, 89  
  response and time-response models, 185
- Dragon, 191, 369
- Drug  
  analysis, 93  
  design, 25, 70, 86, 141, 198, 210, 239  
  development, 20, 49, 62, 69, 75, 87, 93, 141, 143, 145, 147, 149, 150, 152, 169, 179, 180, 184, 192, 198, 200, 202, 203, 210, 243, 311, 312, 330, 373  
  discovery, 49, 59, 61, 64, 68–71, 75, 76, 85–87, 93, 102, 118, 141–143, 145, 150, 151, 180, 190, 197, 198, 200, 201, 210, 214, 220, 222, 225, 226, 234, 240, 242, 293, 312, 323, 330, 331, 364, 371, 373  
  implications, 69  
  drug interactions, 149, 170  
  effectiveness, 141  
  excretion, 149, 150, 169  
  failures, 210  
  identification, 91  
  legalization, 102  
  metabolism, 141–143, 145, 147, 149  
  production, 49  
  repurposing, 87, 142  
  retrosynthetic analysis, 49  
  safety, 76, 87, 209  
  target  
    identification, 86  
    interactions, 86  
  toxicity, 104, 109, 209, 222, 239, 240, 324
- Dual-energy X-ray absorptiometry (DXA), 151, 170
- DuPont, 25, 85
- Dynamic nature, 250
- Dynamism, 3, 7, 9, 12

**E**

Ecological  
  impression, 308, 310  
  microbial subsystems, 154  
  toxicant evaluation, 308  
Ecosystems, 153, 249  
Ecotoxicity, 194, 304  
Edge computing, 352  
Electronegativity, 367  
Electronic, 76, 257  
  density, 187  
  health records, 87  
  medical records, 92  
  properties, 64  
Electrostatic methods, 225  
Email spam detection, 226  
Empirical biodegradation testing, 252  
Encoding, 216  
Endocrine disruption, 198  
Energy  
  efficiency, 6  
  minimization techniques, 24  
  production, 257  
  storage, 76  
Ensemble methods, 254  
Entacapone, 187  
Environmental  
  conditions, 296  
  degradation, 169  
  fate, 255  
  hazard valuation, 311  
  management strategy, 255  
  monitoring data, 257  
  parameter, 251, 262, 264–266, 268, 291  
  pollutants, 255  
  pollution, 248, 299  
  protection, 203, 255, 257, 304  
    agency, 304, 305  
  remediation methods, 255  
  risk assessment, 197, 256, 261, 275, 280,  
    293, 297, 304, 307  
  safety evaluations, 184  
  science, 150, 248, 250, 270, 307, 310  
  sustainability, 248, 251, 252  
  toxicity, 200  
  variables, 253, 261, 262, 264, 292  
Enzymatic mechanisms, 277  
Epochs, 241  
Equipment variations, 341

## Error

  correction, 258, 264  
  reduction, 20, 122, 148  
*Escherichia coli* (EC), 47, 48, 154, 305  
Estimation programs interface, 194  
Estrogen receptor (ER), 128, 134–136  
  alpha, 128  
Ethical  
  considerations, 273, 275  
  guidelines, 264, 278  
Euclidean distance, 127  
European Chemicals  
  Agency (ECHA), 191, 305  
  Bureau, 129  
Evaluation metrics, 209, 233, 235, 238  
Evolutionary algorithm (EA), 5, 6  
Excretion, 148  
  prediction, 148, 150  
Exocrine pancreatic insufficiency (EPI),  
  194, 304, 305  
Experimental  
  biological degradation data, 257  
  data, 75, 261, 262  
  values, 231  
  genotoxicity assays, 184  
  research methods, 146  
  validation methods, 278  
Expert systems, 299, 319  
Explainable AI (XAI), 277, 352  
Extended  
  deep belief network (EDBN), 346, 347  
  restricted Boltzmann machine (ERBM),  
    346, 347  
Extreme Gradient Boosting (XGBoost),  
  226, 231, 233

**F**

Factor Xa (fXa), 27  
False  
  negatives (FN), 272  
  positives (FP), 105, 145, 236, 272, 330,  
    340, 341  
Fault  
  detection, 3, 9, 13, 338–343, 346, 348,  
    351, 352  
  diagnosis, 3, 9, 12, 13, 337–343, 346–354  
  capabilities, 341, 342  
  methods, 337, 339, 340, 343, 348–350  
  models (FDMs), 342, 345, 346, 349

Feature  
   extraction, 119, 243, 254, 258  
   ranking, 258  
   selection, 248, 254, 258, 264, 265, 270, 271, 342, 368  
   vectors, 368  
 Featurization method, 218, 224, 225  
 Fertilizers, 84, 249, 256  
 Fingerprint (FPs), 105, 120, 218, 220–222, 224, 225, 239, 242, 263, 368, 370–373  
   models, 235, 238  
 Flecainide, 188  
 Flexibility, 255  
 Fluconazole, 48  
 Food  
   ingredients, 159  
   science, 155  
 Forecast medication interactions, 143  
 Foreseeing failures, 339  
 Forest random (FR), 105  
 Freeware application, 157  
*Fusarium oxysporum* (FO), 47, 48  
 Fuzzy  
   diagnostic  
     system, 10  
     tool, 10  
   logic (FL), 3–5, 7–9, 13, 298  
   model (FM), 5  
   predictive control (FMPC), 7

## G

GastroPlus, 192  
 Gated recurrent unit (GRU), 268, 280  
 Gaussian  
   function distribution, 125  
   process, 117, 135, 136  
   type membership function, 10  
 Generalization, 3, 22, 226, 231, 258, 274, 275, 325  
 Generative adversarial network (GAN), 369  
 Genetic, 88, 152, 309, 310  
   algorithm (GA), 3, 5, 6, 10, 11, 13, 297, 299  
   indicators, 89  
 Genotoxicity, 182  
 Genomes, 86, 87, 144, 145  
 Genomic profiles, 89  
 Genotoxic compounds, 182–184

Genotoxicity, 180, 182–184, 198, 199, 210, 213, 242  
   assessment, 184  
   data, 182, 183  
   prediction, 183, 184  
 Geometric  
   properties, 24  
   theorem, 21  
 Gini coefficient, 130  
 Glide, 24, 27  
   binding energy values, 28  
   docking scores, 27  
 Glucuronidation, 158–160  
 Goal documentation, 328  
 Gradient  
   boosting machines (GBM), 105, 232, 267, 280, 343  
   descent learning, 122  
 Graph  
   attention network, 111  
   convolution network, 375  
   convolutional  
     model, 235  
     network (GCN), 111, 226, 234, 242, 263, 366  
   featurization, 224  
   neural network, 111  
   representations, 365  
   structured data, 234, 242  
   theory formalism, 374  
 GraphConv featurization, 222  
 Graphical  
   representation, 222, 357  
   synthesis, 108  
 Greenhouse gas emissions, 153, 249  
 Grid-based ligand docking, 27  
 Gut  
   metabolism, 157, 160  
   microorganisms, 160

## H

Hammett, 21  
 Hamster carcinogenicity, 194  
 Handling complex data, 349  
 Harmonic mean, 236, 273  
 Hazard expert, 111  
 Hazardous chemicals, 58, 59, 293  
 Health sciences, 79

Healthcare informatics, 82  
Heart shock factor response element (HSE),  
128, 134–136  
Hebbian, 122  
learning, 122  
Hepatotoxic epoxide metabolite, 187  
Hepatotoxicity, 102, 130, 179–181, 202  
Hepatox, 181  
Hepatotoxicity, 128  
Heteroatom, 363  
Heterocycle, 26  
Heterotrophic microorganisms, 294  
Heuristic, 10  
Hierarchical  
application, 348  
clustering, 127  
High  
fidelity reaction, 21  
performance computing infrastructure, 350  
throughput  
applications, 147  
computational toxicity, 104  
prediction, 107  
screening, 128, 129, 191  
Homogeneous, 126  
Human  
ether-a-go-go-related gene (hERG), 181,  
214, 217, 218, 234, 238, 242, 243  
gut  
metabolism, 157  
microbial transformation, 154, 157,  
160, 162, 164, 166–168  
intelligence, 73, 74  
knowledge acquisition, 322  
transcriptome data, 110, 111  
Hurdles, 209  
Hybrid  
approaches, 266  
models, 268, 353, 354  
Hybridization, 222, 366, 367  
state, 222  
Hydrogen, 29, 58, 187, 240, 371  
bond  
acceptors, 29  
donors, 29, 371  
Hydrological systems, 153  
Hyper parameter, 132, 242  
tuning, 233, 271

## I

Ibuprofen, 189  
Image recognition, 226, 268  
Iminosugars, 67  
Imipramine, 187  
Immunotoxicity, 128, 130  
*In silico*, 27, 57, 102, 155, 157, 180, 182,  
184, 203, 293, 364  
ADME prediction, 27, 29, 46, 192, 197  
chemico-biology approach, 48  
genotoxicity prediction, 184  
metabolism prediction, 155  
toxicological assessments, 180  
Inborn metabolic  
abnormalities, 148  
errors, 148  
InChI key, 375  
Incremental construction algorithm, 326  
Inductive  
machine learning, 253  
structure protein analysis (IPSA), 253  
Industrial chemicals, 102, 194, 292  
Industry applications, 278  
Inflammatory gastrointestinal illnesses, 159  
Information and communication technology  
(ICT), 83  
Insecticidal properties, 27  
Insecticides, 120, 292  
*In-silico*, 180, 181  
identification, 202  
Instrumental errors, 9  
Insulin levels, 151  
Intellectual  
landscapes, 111  
property, 73, 75, 93, 94  
Intelligent  
computations, 322  
environmental management decisions, 279  
Interactive web-based environment, 242  
Interconvertibility, 361  
International  
chemical identifier, 361  
greenhouse gas emissions, 153  
Union of Pure and Applied Chemistry  
(IUPAC), 160, 163, 164, 166, 361  
Internet of things (IoT), 279, 280, 352, 353  
integration, 352

Interpretability, 5, 209, 238, 240, 241, 248,  
252, 264, 269, 273, 276, 277, 280, 348,  
350, 352, 353

Intuitivism, 5

*In-vitro*, 105, 180

Isomeric structures, 358

## J

Joint Research Centre (JRC), 196

Jupyter notebook, 209, 214, 216, 217, 242

Juxtaposing, 86

## K

K clusters, 126

Kanamycin (KM), 48

Kernel

models, 194

tricks, 126

Key

area, 85

players, 84, 85

Kinetic models, 144, 146, 168

K-means clustering, 249

K-nearest neighbor (K-NN), 126, 194, 238,  
253, 267

Knowledge

prediction models, 213, 239

techniques, 211

## L

Lab-scale biodegradation, 278

Large-scale process plants, 340

Layer-wise compression process, 346

Lazar, 128, 130, 194, 195

Lazy structure, 128

Leadscope predictive tox suite, 199, 200

Learned embeddings, 369

Life sciences, 94, 107

Ligand

methods, 168

receptor interface dynamisms, 24

Linear

activation function, 122

discriminant analysis, 118, 125, 126, 134,  
135

free energy relationships, 21

regression, 117

Lipid-soluble toxicants, 158

Lipinski's rule, 29

Liver

functions, 180

injury, 128, 180

toxicity, 102, 181, 194, 210, 242

knowledge base, 181

prediction, 210

LiverTox, 181

Local interpretable model-agnostic explana-  
tions (LIME), 274

Log transformation, 264

Logical operators, 363

Logistic regression models, 134, 135, 249,  
265, 267, 268

Long short-term memory (LSTM), 268, 345  
network (LSTMN), 345, 346

Loss function, 243

## M

Machine

knowledge, 21, 104, 328

learning (ML), 47, 48, 59, 61–65, 69, 71,  
76, 78, 86, 92–94, 99, 101, 104, 107,  
109–111, 117–124, 126–130, 132, 133,  
136, 137, 141–154, 157, 160, 168–170,  
179, 181–184, 191, 198–200, 202, 203,  
209–214, 218, 221, 222, 224–226, 231,  
233, 234, 238, 240–243, 248–250,  
253, 254, 257, 262, 265–268, 271, 277,  
279, 280, 291–293, 297–299, 302, 312,  
313, 321, 322, 324, 327, 332, 337–343,  
347, 348, 350–354, 357, 360, 366–368,  
371–375

algorithms, 48, 59, 61, 63–65, 69, 76,  
86, 88, 104, 117–121, 123, 126–128,  
132, 133, 136, 137, 148, 149, 152,  
198, 225, 226, 231, 248, 250, 251,  
253, 254, 257, 262, 265–268, 271,  
277, 280, 298, 299, 323, 337, 339,  
348, 351, 368, 371, 372

approaches, 92, 143, 183

fault diagnosis, 345

models, 65, 143, 145, 149, 179, 202,  
250, 266, 322, 368

techniques, 63, 65, 104, 107, 109,  
143–146, 169, 183, 199, 249, 327,  
341, 342, 350–352

- Maestro builder panel, 27
- Mamdani fuzzy model, 5, 10
- Mammalian metabolism, 305
- Man-made organic compounds, 294
- Marine systems, 293
- Matched molecular pair (MMP), 128, 134–136, 212, 239–241
- Materials science, 57, 58, 60, 70
- Mathematical manipulations, 359
- Maximum recommended daily dose (MRDD), 194
- Mean
  - absolute error (MAE), 235, 238
  - square error (MSE), 10
- Measurement standardization, 312
- Mechanistic methods, 4
- Medical
  - advice, 89
  - histories, 88, 89
  - image analysis, 90
  - technology and pharmaceutical, 79
- Medication
  - compliance, 90
  - development, 142, 143, 149
  - metabolism, 148
  - safety data, 87
  - utilization, 48
- Medicinal poisoning estimate, 104
- Metabolic
  - biotransformations, 169
  - circuits, 145
  - data, 145, 152, 157, 303
  - engineering, 142, 144–146, 168, 170, 309, 310
    - application, 149
    - strategies, 145, 310
  - fate prediction site-at-metabolism, 142, 143, 150
  - indicators, 145
  - interactions, 143, 277
  - pathway, 142–148, 150, 151, 167, 169, 277, 309, 310, 359
    - dynamics, 143, 146
    - kinetics, 148
    - prediction, 144
  - pharmacokinetics analysis, 143
  - phenotypes, 145, 146
  - predictions, 145, 154
    - processes, 145
    - switching, 186, 187
    - syndrome, 145, 146, 151
    - transformation sites, 146
- Metabolism, 29, 141–152, 154, 155–160, 168–170, 186, 188, 198, 310
  - limitations, 156
  - prediction, 142, 146–149, 151, 154, 155, 157, 160, 169
  - prognostication, 149
  - types, 157
- Metabolite, 141, 143–145, 149, 150, 152, 154–157, 159, 170, 185, 250, 303, 309
  - formation, 186, 262
  - structure, 141, 157
    - estimation, 143
- Metabolomic data, 168
- MetaCyc, 298, 309, 310
- Meta-PC, 305, 306
- Methodology, 160, 297, 324, 326, 346, 373
- Methylenediurea, 308
- Metrics, 89, 107, 151, 231, 235, 236, 238, 248, 267, 271–273, 280
- Miconazole, 48
- Microbial
  - activity, 256, 275
  - adaptability, 250
  - biodegradation, 308, 309
  - community, 250, 262, 277, 280, 295, 296
  - DNA gyrase, 27
  - metabolism, 155, 160
  - processes, 308
  - variability, 312
- Microbiological
  - activity, 251
  - plastic decomposition, 307
- Microbiology, 293
- Microbiome composition, 160
- Microchem reactor, 12
- Micrococcus luteus (ML), 47
- Microorganisms, 160, 248, 277, 294–296, 306–311
  - diversity, 250
- Microwave process technology, 12
- Mimics human cognitive processes, 142
- Min-max scaling, 264
- Mitochondrial membrane potential, 106, 128
- Model
  - accuracy, 270
  - biases, 209



- building, 211
  - comparison, 271
  - complexity, 269
  - confidence, 239
  - evaluation, 107, 269
  - expert system, 348
  - generalizability, 271
  - interpretability, 273, 277, 350
  - optimization, 271
  - robustness, 271, 275
  - transferability, 277
  - validation, 271, 274
- Molecular**
- ACCess System (MACCS), 371, 372
  - analysis, 241
  - data, 86, 371, 374
  - descriptors, 105, 128, 182, 183, 191, 198, 199, 218, 238, 240, 263, 293, 297, 329, 368, 372
  - design limited (MDL), 179, 182, 202, 371
  - distribution, 217
  - docking, 25–27, 46, 49, 65, 76, 87, 322, 325
  - dynamics (MD), 25, 147, 326
  - embedding, 332
  - featurization, 209, 211, 212, 218, 224
  - fingerprinting technique, 220
  - fingerprints, 21, 218, 240, 242, 253, 268, 357, 370, 371
  - graph structures, 234
  - modeling, 25, 26, 191, 203, 324
  - modification patterns (MMP's), 241
  - operating environment (MOE), 191, 369
  - queries, 364
  - representation techniques, 373
  - shape analysis, 26
  - similarity maps, 240
  - structural resemblance, 23
  - structure, 21, 24, 59, 64, 76, 91, 107, 154, 169, 218, 224, 225, 242, 250, 265, 304, 305, 358, 359, 361, 367, 370, 373, 375
  - weight, 295
- Molecule**, 20, 22, 24–27, 49, 57–63, 65–70, 86, 87, 118, 127, 129, 130, 142, 145, 146, 149, 153–157, 168, 169, 179–181, 183, 187, 190, 202, 203, 218, 222, 224, 225, 230, 234, 236, 239, 242, 248, 250, 292–294, 296, 298, 301, 302, 308–311, 326, 329, 360–365, 367–374
- designing, 26
  - MoleculeNet, 119, 213
  - MolFeat, 211, 218, 221–223
  - Monetization, 84
  - Monoxygenases, 310
  - Monte Carlo (MC), 326
    - de novo ligand generation, 25
    - replications, 24
    - techniques, 25
    - tree search (MCTS), 23
  - Morgan fingerprint, 105, 218, 220–223, 239
  - Morphological measures, 151
  - MultiCASE, 200, 305
  - Multi-column distillation facility, 348
  - Multidisciplinary
    - area, 120
    - cooperation, 248
  - Multi-layer perceptron (MLP), 5, 121, 122, 345, 346
  - Multi-linear regression, 194
  - Multi-modal data, 279
  - Multi-omics data, 277
  - Multiple interconnected systems, 340
  - Multitask deep neural network (MT-DNN), 182
  - Multivariate systems, 104
  - Mutagenicity, 128, 130, 179, 186, 194, 196, 200–202
- N**
- Naïve Bayes (NB), 118, 122, 134–136, 181, 238
  - National
    - Center for Biotechnology Information, 190
    - Institute of Technology and Evaluation (NITE), 261
  - Natural language processing, 232, 332
  - N-dimensional
    - gaussian function, 369
    - vectors, 369
  - Neural network (NN), 3, 5, 7–9, 12, 23, 93, 105, 106, 111, 117, 119–122, 134, 143, 144, 169, 170, 182, 192, 201, 222, 234, 241–243, 253, 265, 267, 268, 271, 277, 280, 291, 297, 298, 313, 322, 323, 327–330, 332, 337, 343, 345, 346, 348–350, 375
  - layers, 243

- models, 192, 234, 241
- systems (NNSs), 9
- Neuromorphic diagnostic tool, 12
- N-hydroxyaniline, 188, 189
- Nifedipine, 186
- Nitrate, 294
- Nitrogen, 26, 58, 84, 292, 308
- Nitroso derivative, 188
- Node feature matrix, 366
- Non-adaptive web record datasets, 324
- Non-bonded interactions, 28
- Non-genotoxic compounds, 184
- Non-mammals, 157
- Nonlinear
  - classification, 105
  - correlations, 150, 255
  - dynamism, 7
  - problems, 13
- Nonlinearity, 3, 4, 6, 9
- Norfloxacin, 25
- Novel
  - chemical processes, 90
  - compounds, 58, 257
  - medications, 149
  - mirror website, 306
  - neuromorphic diagnostic system, 12
- Nuclear receptor
  - effects (NR), 127, 128
  - estrogen receptor (NR-ER), 128, 137
- Numerical
  - pairwise possibilities, 326
  - representation, 368
- Nutrient
  - availability, 251, 296
  - cycling, 248, 249
  - rich soil amendment, 294
- Nutritional slip-ups, 151
- Nutrition-related data, 150

## O

- Object-oriented knowledge base, 348
- Octanol-water partition coefficient, 29
- Omeprazole, 164
- OncoLogic, 111, 195
- Open-source
  - chemoinformatics toolkit, 191
  - databases, 212
  - libraries models, 212

- OpenTox, 191
- Operational
  - disruptions, 339, 348
  - efficiency, 348, 353
- Optimal
  - hyperparameters, 133
  - hyperplane, 126
- Optimization, 3, 6, 7, 11, 13, 76, 86–88, 92, 132, 150, 201, 212, 232, 251, 252, 271, 299, 324, 330, 342, 349, 352
- Oral drugs, 29
- Orally active medication, 29
- Organic, 20–24, 48, 57–66, 67–71, 79, 153, 194, 195, 248–250, 256, 292, 294, 295, 298, 299, 302, 304–310, 321, 322, 358, 360, 368
  - chemicals, 194, 305–307
  - chemistry, 69, 79, 322
  - compounds, 57, 58, 60, 153, 195, 292, 294, 295, 309, 310, 358
  - materials, 248–250
  - nitrogen, 308
  - production, 21, 48
  - synthesis, 20, 23, 57–71, 322
    - challenges, 60
    - implications, 69
  - transformations, 20, 321
  - waste, 248, 249, 256
- Organization for Economic Co-operation and Development (OECD), 107, 179, 195, 202
- Organonitrogen
  - compounds, 308
  - degradation database (ONDB), 298, 308
- Organ-specific toxicities, 198
- Osiris property explorer, 179, 197, 202
- Osteopenia, 151
- Osteoporosis, 151
- Outlier detection, 239
- OxDBase, 298, 310
- Oxygen, 58, 280, 294, 296, 303, 362
  - availability, 294, 296
- Oxygenases database, 298, 310

## P

- Partial least squares (PLS), 254
- Particle swarm optimization, 6
- Partition coefficient, 29, 194, 304, 369

- PASS (Prediction of Activity Spectra for Substances), 196
- Patent
- application, 74, 75, 78–80, 85–87
  - filing trend, 73, 74, 78, 80, 82, 85
- Patient
- management, 93
  - participation, 92
- Pattern recognition, 104, 265, 324, 339, 341, 352
- Performance evaluation, 269–271, 279
- Peroxidases, 310
- Persistence, bioaccumulation, and toxicity (PBT), 298, 305
- Persistent homology, 372, 373
- Personal health data, 151
- Personalized medicine, 82, 88
- therapy optimization, 88
- Pesticides, 58, 147, 249, 250, 256, 294, 298
- Pharmaceutical, 57, 58, 74, 76, 78, 80, 82, 85, 91, 93, 111, 136, 147, 169, 199, 255, 278, 299, 302
- biotechnology industries, 93
  - businesses, 84
  - chemical sector, 82
  - companies, 84, 93, 169, 202
  - compounds, 91
  - corporations, 87, 88
  - development, 24
  - domains, 75
  - identification, 91
  - industry, 25, 57, 62, 75, 78, 86, 93, 142, 202, 255, 292, 302
  - productivity, 142
  - products, 80, 84
  - science, 73–76, 92–94
    - AI impact, 76
  - sector, 74, 76–79, 84, 86
    - AI application, 77
- Pharmacodynamic, 192
- models, 185
- Pharmacokinetic, 29, 49, 141, 143, 185, 192, 324
- models, 185
  - streams, 142
- Pharmacological intervention, 86
- Pharmacology, 310
- Pharmacophore, 27, 364
- fingerprints, 371
  - modeling, 181
- Phase II transformation, 158, 160, 168
- Phenotypes, 145, 146, 150
- Phosphorus, 58
- Photo-degradation, 294
- Physical-chemical properties, 29
- Physicochemical properties, 107, 249, 256, 262, 298, 299, 304
- Phyto microbiome, 150
- Pipeline, 93, 202, 214, 325
- Plant-wide systems, 12
- Plastic, 84, 256, 294, 307
- microbial biodegradation database, 298, 307
  - waste, 307
- Policy-making, 261
- Pollutants, 255, 256, 265, 266, 280, 292, 299, 311
- Pollution mitigation, 248, 278, 280
- Polycyclic aromatic hydrocarbons (PAHs), 299
- Polymers, 196
- Polynomial function, 226
- Population-level statistics, 88
- Potassium channels, 181
- Potential
- atom-to-atom distances, 23
  - toxic effects, 211
- Practical learning and experimentation, 211
- PreADMET, 197, 198
- Precision, 88, 152, 236, 237, 272, 273
- diagnostics, 88
  - medicine, 152
- Preclinical assessment, 102
- Prioritization, 107, 199–201
- Prioritize compounds, 182, 184, 190, 192
- Priority application data, 77
- Private medical data, 93
- Procainamide, 188
- Processing elements (PE), 327
- Product development, 93, 257
- Production parameters, 91
- Prognosis, 152, 338, 342, 343
- Prognostic toxicology, 108
- Programming languages, 190
- Proprietary representations, 374
- Protein
- structures, 86, 93
  - tertiary construction, 329

*Pseudomonas fluorescens* (PF), 47, 48  
PubChem, 153, 181, 190, 261, 306, 358  
Public health, 101, 119, 255, 257  
Pyridine, 372  
Python, 133, 190, 191, 212–214, 369  
    codes, 209, 211  
PyTorch, 192, 214, 367

## Q

Quality analysis, 91  
Quantifiers, 363  
Quantitative structure  
    activity relationship (QSAR), 25, 26, 104,  
        107, 111, 120, 128, 130, 179, 182, 185,  
        186, 191, 192, 195, 197, 200–202, 240,  
        248, 254, 261–263, 293, 297, 305, 313,  
        323, 329, 331  
        modeling, 328  
        toolbox, 130, 191, 202, 305  
    biodegradability relationship (QSBR),  
        293  
    property relationship (QSPR), 120, 323,  
        324  
Quantum chemistry calculations, 147

## R

Radial basis function (RBF), 226  
Random forest (RF), 105, 117–119, 124,  
    125, 127, 129–135, 137, 143, 144, 183,  
    226, 231, 232, 235, 239, 253, 265, 267,  
    327, 337, 342, 343  
    hyper parameters, 132  
    machine learning  
        algorithm, 119  
        method, 127  
    system architecture, 132  
RdKit, 211, 212, 218, 220–222, 224, 239,  
    240, 243  
    graph featurization, 224  
    similar maps, 240  
Reactant compatibility, 63  
Read across, 111, 184  
Real-time  
    data analysis, 90  
    forecasts, 312  
    monitoring, 278, 340, 343, 348, 351, 352  
    prediction, 168, 276

Recall, 236, 237, 272, 273  
Receiver operating characteristic (ROC),  
    118, 133–135, 137, 234, 235, 237, 273  
    curve, 133, 137  
Receptor-ligand complex, 24  
Recruitment, 91  
Rectified linear unit (ReLU), 243  
Recurrent neural network (RNN), 253, 265,  
    268, 269, 343, 345, 368  
Regulatory  
    agencies, 276  
    compliance, 261, 274, 276, 304  
Reinforcement  
    learning, 226, 352  
        algorithm, 120  
        machine algorithm, 120  
Reliable fault diagnosis systems, 353  
Remaining useful life (RUL), 338, 341, 349  
Remote  
    medical treatments, 89  
    monitoring, 89, 92  
    patient monitoring, 90  
Renewable energy  
    production, 248  
    sources, 249  
Reproductive toxicity, 200  
Research-based techniques, 168  
Resource  
    consumption, 339  
    intensive, 340  
Resting metabolic rate, 146  
Retrosynthesis, 21, 22, 65, 324  
    formalization, 22  
Retrosynthetic  
    analysis, 20–23, 48, 49, 66  
    planning, 22, 323  
Revolutionized, 57, 58, 60, 62–69, 152, 210,  
    357  
Revolutionizing organic synthesis, 61, 70  
    potential, 70  
Ring constraints, 363  
Risk assessment, 88, 200, 241, 266  
Robust  
    patent portfolio, 84  
    predictions, 184, 266, 277  
Robustness, 3, 4, 248, 259, 266, 271, 275,  
    277–280  
Rodent, 194

Root mean squared error (RMSE), 235, 238  
Rule-based molecular property assessments,  
212

## S

Sanofi, 85, 88  
Scalability, 350  
Scatter plots, 198  
Scikit-learn, 105, 106, 191, 226  
Seaborn libraries, 212  
Search engines, 310, 358  
Sedentary compounds, 108  
Self-control, 152  
Self-learning capabilities, 7  
Self-management techniques, 89  
Semi-supervised algorithm, 367  
Sentiment analysis, 226  
Separation techniques, 10  
Sertindole, 181  
Set-point tracking, 12  
Shallow  
  architectures, 105  
  knowledge, 105  
SHapley additive exPlanations (SHAP),  
239, 240, 274  
Siemens mindsphere, 352  
Simple  
  linear notations, 374  
  model of the atmospheric radiative  
  transfer (SMARTS), 362–365, 372, 373  
  applications, 364  
Simplified molecular input line entry  
  system (SMILES), 21, 23, 154, 160, 214,  
  218, 220, 221, 224, 234, 263, 300, 302,  
  360–365, 373, 375  
  arbitrary target specification, 362, 375  
Single-fault experimental technique, 346  
Sitaxentan, 102, 103  
Skin sensitization, 194, 199, 201  
Small drug discovery suite, 27  
SmartCyp, 156  
Softmax, 243, 367  
Software, 21, 24, 29, 61, 66, 84, 88, 94, 142,  
  154–157, 160, 167, 169, 179, 180, 182,  
  186, 189–192, 194–203, 213, 298, 299,  
  302–305, 322, 326, 351, 352, 364, 369  
  tools, 88, 154, 156, 189, 202  
Splitting dataset, 127  
*Staphylococcus aureus* (SA), 47, 48  
Statistical methodologies, 264  
Step-wise feeding, 11  
Stereoisomers, 358  
Stereoselective synthesis, 65  
Steric considerations, 23  
Stratified K-fold splitting, 105  
Streamline production methods, 91  
Stress response (SR), 127, 128, 137  
  heat shock response effect (SR-HSE), 128  
Structural alerts, 185–187, 199, 212, 239, 240  
Structure  
  activity relationships (SAR), 102, 195,  
  200, 305  
  drug  
  design, 24, 25, 88  
  strategy, 24, 25  
  representation, 357  
  virtual screening (SBVS), 325–327  
Suboptimal synthetic routes, 59  
Sulfate, 159, 294  
Sulfation, 158–160  
Sulfur, 58  
Supervised  
  learning, 120, 122, 126, 129, 137, 144, 145  
  algorithm, 120, 121, 342  
  method, 129  
  machine learning algorithms, 119  
Support vector  
  machine (SVM), 105, 111, 117, 119, 126,  
  134–137, 143, 181, 183, 226, 230, 231,  
  235, 238, 253, 254, 265, 267, 268, 271,  
  297, 302, 303, 326, 327, 337, 342  
  biodegradability predictor, 302  
  regression (SVR), 254  
Sustainable  
  agricultural practices, 249, 256  
  agriculture, 248  
  chemical  
  design, 248  
  management practices, 305–309, 311  
  environmental management practices, 252  
  materials, 292  
SYBYL line notation (SLN), 364, 365  
Synthesis planning, 65–67, 323, 324  
Synthetic  
  compounds, 29  
  fertilizers, 249  
  pathways, 49, 59, 65, 90  
Syracuse Research Corporation (SRC), 194

**T**

- Takagi-Sugeno (TS), 5, 8
- Target identification, 93, 324, 329
- Tautomer standardization, 105
- TDCcommons (TDC), 213
- T-distributed stochastic neighbor embedding (T-SNE), 217, 218
- Technological development, 75, 88
- Teleconsultations, 90
- Telemedicine, 89, 90
  - development, 89
- Tennessee Eastman process (TE), 12, 345, 346
- TensorFlow, 192, 211, 214, 222
- Terfenadine, 181
- Tetrahymena pyriformis, 197
- Therapeutic
  - applications, 80, 87
  - candidates, 24, 86
  - effectiveness, 86
  - intervention, 152
  - proteins, 324
- Thermodynamic
  - connections, 27
  - interaction, 27, 28
  - interactions, 27
- Three-dimensional
  - characteristics, 105
  - convolutional networks, 330
- Threshold of toxicological concern (TTC), 196
- Time
  - consuming experimental screenings, 20, 57–60, 63, 117, 148, 210, 249, 255, 340, 350
  - series
    - environmental data, 265
    - measurements, 268
- Tolcapone, 187
- Toolbox, 195, 202
- Topological
  - polar surface area, 29
  - torsion fingerprints, 371
- Total polar surface area (TPSA), 46, 369
- Tox21 dataset, 119, 127, 137
- Toxic
  - interactions, 297
  - substances, 242, 292
- Toxicity, 49, 61, 68, 76, 86, 87, 101, 102, 104, 105, 109–111, 117–120, 127–130, 132–137, 141, 150, 158, 159, 179–182, 184–188, 190–192, 194–203, 209–213, 222, 234, 238–243, 297, 303–305, 323, 324
  - assessment, 185, 192, 196, 201, 241
  - processes, 241
  - dataset, 119
  - endpoints, 191, 198–200, 210–213, 242, 243, 303
  - estimation, 197
    - software tool, 182
  - measurements, 105
  - parameters, 180
  - prediction, 101, 102, 104, 110, 111, 117–119, 127, 132, 133, 136, 137, 180, 184, 189–192, 195, 198–203, 211, 212, 222, 239–241, 323
  - software use, 189
  - system, 117, 119, 127, 136, 137
  - profiles, 192, 201, 209
  - training dataset, 117, 136
  - type, 180
- Toxicological
  - data, 195, 210
  - endpoints, 119, 127, 182, 191, 194, 196, 197, 200–202
  - science, 150
- Toxicology prediction models, 150
- Toxicophore, 183
  - identification, 183
- ToxMic rule, 196
- Toxtree, 179, 186, 196, 202
- Traditional
  - fault diagnosis methods, 337, 340
  - methods, 66, 67, 168, 254, 255, 298, 312, 340, 341, 343, 344, 348
  - organic synthesis methods, 58, 70
  - statistical approaches, 169
  - trial-and-error methods, 60
- Training
  - data coverage, 239
  - dataset, 105, 118–120, 123, 125–128, 130, 132, 133, 136, 194, 234, 241
- Transient receptor potential melastatin 2 (TRPM2), 27
- Transparency, 5, 93, 196, 248, 252, 261, 273, 274, 276, 278–280, 353

## Treatment

- choice, 89
- plans, 92

Tree structure-based algorithm, 123

Trial design, 92

Troubleshooting, 350

## True

- negatives (TN), 236, 272
- positives (TP), 236, 272

**U**

Uncertainty factor models, 185

Unconventional descriptors, 372

UniProt database, 307

United States National Library of Medicine, 190

University of Minnesota biocatalysis/  
biodegradation (UM-BBD), 298, 302,  
306, 307, 310  
databank, 306

Un-supervised learning, 120  
procedure, 126

Utilizing material databases, 76

**V**

Vander-Waals interactions, 23, 28

Variational autoencoders (VAE), 368, 369,  
373

Vector machine, 126

Violation, 29

Viral infections, 67

Virtual screening, 76, 86–88, 324, 325, 327,  
330, 371

Virtualization, 125

Visualization, 137, 180, 190, 198, 216  
tools, 137, 180, 198

**W**

Waste, 10, 59, 62, 65, 76, 155, 248–251,  
256, 261, 270, 271, 278, 280, 291, 293,  
294, 297, 307, 312, 339  
management, 248, 250, 251, 256, 261,  
270, 271, 278, 280, 291, 293, 297, 307  
reduction, 251

Water molecules, 26

molecular docking studies, 26

Website, 153–155, 304, 306, 309

Weight, 8, 234, 303, 327, 366

modification, 122

reduction, 152

Weiner index, 369

Wildcards, 363

Workload, 142

**X**

Xenobiotic, 157, 159, 160

biodegradation, 294

metabolism, 157, 159

XGBoost, 143, 226, 231–233, 235

classifier (XGBClassifier), 233

regressor (XGBRegressor), 233

X-ray crystallography, 26

**Y**

Yard waste, 248

Yield, 6, 11, 12, 48, 59–63, 65, 67, 68, 108,  
144, 146, 149, 168, 322, 361

**Z**

Zolmitriptan, 26

Z-score normalization, 264