Pradeep Singh
Balasubramanian Raman

# The Geometry of Intelligence: Foundations of Transformer Networks in Deep Learning

Springer

# Studies in Big Data

Volume 175

**Series Editor**

Janusz Kacprzyk, Polish Academy of Sciences, Warsaw, Poland

The series "Studies in Big Data" (SBD) publishes new developments and advances in the various areas of Big Data- quickly and with a high quality. The intent is to cover the theory, research, development, and applications of Big Data, as embedded in the fields of engineering, computer science, physics, economics and life sciences. The books of the series refer to the analysis and understanding of large, complex, and/or distributed data sets generated from recent digital sources coming from sensors or other physical instruments as well as simulations, crowd sourcing, social networks or other internet transactions, such as emails or video click streams and other. The series contains monographs, lecture notes and edited volumes in Big Data spanning the areas of computational intelligence including neural networks, evolutionary computation, soft computing, fuzzy systems, as well as artificial intelligence, data mining, modern statistics and Operations research, as well as self-organizing systems. Of particular value to both the contributors and the readership are the short publication timeframe and the world-wide distribution, which enable both wide and rapid dissemination of research output.

The books of this series are reviewed in a single blind peer review process.

Indexed by SCOPUS, EI Compendex, SCIMAGO and zbMATH.

All books published in the series are submitted for consideration in Web of Science.

Pradeep Singh · Balasubramanian Raman

# The Geometry
# of Intelligence: Foundations
# of Transformer Networks
# in Deep Learning

Pradeep Singh
Department of Computer Science
and Engineering
IIT Roorkee
Roorkee, India

Balasubramanian Raman
Department of Computer Science
and Engineering
IIT Roorkee
Roorkee, India

If disposing of this product, please recycle the paper.

*To the Boundless Frontiers of Thought and Technology*

# Preface

The advent of Transformer networks has not merely advanced the field of artificial intelligence; it has redefined the landscape entirely. As the current state-of-the-art across various domains—including natural language processing, computer vision, time series forecasting, and signal analysis—Transformers have demonstrated an unparalleled ability to model complex patterns, understand intricate relationships, and deliver breakthrough performance. The self-attention mechanism at the heart of these models allows them to capture dependencies within data in ways that traditional architectures could never achieve, making Transformers the backbone of modern AI research and applications.

Despite the widespread adoption and success of Transformers, much of the literature remains focused on their practical implementation, often overlooking the deep mathematical structures that enable their effectiveness. This gap presents a significant opportunity for researchers and practitioners alike. Understanding the mathematical foundations of Transformers is essential for those who seek to push the boundaries of what these models can achieve. A solid mathematical understanding equips us to innovate, optimize, and potentially discover the next generation of models that will build on the success of Transformers.

*The Geometry of Intelligence: Foundations of Transformer Networks in Deep Learning* is crafted for those who wish to explore the profound theoretical underpinnings of Transformer networks. This book is intended for researchers, academics, and advanced practitioners who aspire to grasp the elegant mathematical principles that make Transformers work. By focusing exclusively on the theoretical aspects, we aim to provide readers with a deep and thorough understanding of the geometry, symmetry, and intelligence encoded within these models, without the distractions of implementation details, for which a plethora of resources already exist.

The structure of this book reflects our commitment to a comprehensive and rigorous exploration of the mathematics of Transformers. We begin with essential mathematical preliminaries before delving into detailed, domain-specific explorations of how Transformers operate in various abstract spaces. Each chapter is designed to present the mathematical formulations, theoretical insights, and analyses that reveal the true power and potential of Transformer models. The pursuit of a

mathematical understanding of Transformers is more than an intellectual endeavor; it is a journey toward the future of artificial intelligence. By uncovering the principles that drive their success, we open the door to new possibilities and innovations that can extend the capabilities of these models even further. We hope this book serves as a valuable resource on your path to mastering the intricate and fascinating world of Transformers.

Roorkee, India                                                                  Pradeep Singh
August 2024                                              Balasubramanian Raman

# Contents

# About the Authors

**Dr. Pradeep Singh** earned his Ph.D. and Master's degrees from the Indian Institute of Technology (IIT) Delhi, specializing in Dynamical Systems, and a Bachelor's degree in Data Science from IIT Madras. Currently, he is a Post-doctoral Researcher and Principal Investigator at the Machine Intelligence Lab within the Department of Computer Science and Engineering at IIT Roorkee, where he is actively engaged in research at the intersection of Geometric Deep Learning, Neuro-symbolic AI, and Dynamical Systems. His research is supported by the National Post Doctoral Fellowship (N-PDF) from the Science and Engineering Research Board (SERB), Department of Science and Technology. Dr. Singh has earned multiple accolades, including All India Rank 1 in IIT GATE 2020, IIT JAM 2015, and CSIR NET 2019, along with prestigious fellowships such as the National Board for Higher Mathematics (NBHM) Masters, Doctoral, and Post Doctoral Fellowships from the Department of Atomic Energy, India. In 2019, he was one of only two researchers nationwide to be awarded the esteemed Shyama Prasad Mukherjee (SPM) Doctoral Fellowship in Mathematics by the Council of Scientific and Industrial Research, India.

**Dr. Balasubramanian Raman** (Senior Member, IEEE) is Professor (HAG) and Head of the Department of Computer Science & Engineering at IIT Roorkee, where he also serves as iHUB Divyasampark Chair Professor and holds a joint appointment with the Mehta Family School of Data Science & AI. He earned his Ph.D. from IIT Madras (2001). Dr Raman's research—spanning machine learning, computer vision, image/video processing, and pattern recognition—has yielded more than 250 peer-reviewed publications and a Google-Scholar h-index of 50. He has held post-doctoral and visiting appointments at Rutgers University, the University of Missouri-Columbia, Osaka Metropolitan University, Curtin University, and the University of Cyberjaya. He has supervised 32 Ph.D. scholars to completion—with another 17 in progress—and regularly directs large, multi-institutional research programmes supported by leading national agencies and industry partners. Dr Raman is an inventor of an Indian patent for a real-time fog-removal imaging system. He has co-authored a book—Deep Learning Through the Prism of Tensors (Springer Nature Singapore, 2024). His honours include the DST BOYSCAST Fellowship, two IIT Roorkee Outstanding Teacher Awards, the Ramkumar Prize for Outstanding Teaching and Research, and the ICPC Coach Award, the last recognising his teams' top-50 finishes at the ACM ICPC World Finals.

# Abbreviations

| | |
|---|---|
| AI | Artificial Intelligence |
| ANN | Artificial Neural Network |
| BERT | Bidirectional Encoder Representations from Transformers |
| BLEU | Bilingual Evaluation Understudy |
| CBOW | Continuous Bag of Words |
| CNN | Convolutional Neural Network |
| ELMo | Embeddings from Language Models |
| FFT | Fast Fourier Transform |
| GAN | Generative Adversarial Network |
| GPT | Generative Pre-trained Transformer |
| GPU | Graphics Processing Unit |
| LSTM | Long Short-Term Memory |
| MFCC | Mel-Frequency Cepstral Coefficients |
| ML | Machine Learning |
| NLP | Natural Language Processing |
| PCA | Principal Component Analysis |
| RNN | Recurrent Neural Network |
| RoBERTa | Robustly optimized BERT approach |
| STFT | Short-Time Fourier Transform |
| SVD | Singular Value Decomposition |
| SWIN | Shifted Window Transformer |
| TF-IDF | Term Frequency-Inverse Document Frequency |
| TPU | Tensor Processing Unit |
| VAE | Variational Autoencoder |
| ViT | Vision Transformer |
| WER | Word Error Rate |
| XLNet | Generalized autoregressive pretraining for language understanding |

# Chapter 1
# Foundations of Representation Theory in Transformers

## 1.1 Introduction

The study of transformers, particularly in the context of natural language processing and machine learning, has revolutionized the way we understand and process data. The central theme of this revolution is the concept of representation: how data, whether it be words in a sentence, pixels in an image, or nodes in a graph, is transformed into a mathematical structure that a machine can manipulate. Understanding this process from a mathematical perspective requires us to delve into the theory of vector spaces and linear algebra, which form the backbone of representation theory.

In transformer models, data is represented as vectors in high-dimensional spaces. These vectors capture not only the intrinsic properties of the data but also the relationships and interactions between different data points. The operations performed on these vectors, such as attention mechanisms and linear transformations, rely heavily on the principles of vector spaces. Therefore, a deep understanding of vector spaces, subspaces, and bases is essential for comprehending how transformers encode and manipulate information. Moreover, transformers leverage the concept of symmetry—a principle deeply rooted in group theory and representation theory. Symmetry allows us to understand how certain transformations, such as rotations or translations, affect the data representations. By exploring these symmetries within the framework of vector spaces, we can gain insights into the invariances and equivariances that make transformer models so powerful.

This chapter serves as a foundation for the mathematical framework that underpins transformers. We begin by introducing the fundamental concepts of vector spaces and linear algebra. These concepts will not only provide the necessary tools to analyze and understand transformer models but will also reveal the deep connections between geometry, symmetry, and intelligence. As we progress, we will see how the abstract mathematical notions introduced here manifest in the practical operations of transformers, setting the stage for more advanced topics in later chapters. Through mathematical exploration, we aim to build an intuition for how transformers operate,

grounded in the precise language of vector spaces and their transformations. This will enable us to appreciate the elegance and power of transformers from a purely mathematical standpoint, where the focus is not on implementation or code, but on the underlying mathematical structures that drive these models.

### 1.1.1   Vector Spaces

A vector space (also called a linear space) is a fundamental concept in mathematics that encapsulates the notions of addition and scalar multiplication. Formally, a vector space $V$ over a field $\mathbb{F}$ (such as $\mathbb{R}$ or $\mathbb{C}$) is a set equipped with two operations: vector addition and scalar multiplication. These operations satisfy the following axioms:

1. Closure under addition: For all $\mathbf{u}, \mathbf{v} \in V$, the sum $\mathbf{u} + \mathbf{v} \in V$.
2. Associativity of addition: For all $\mathbf{u}, \mathbf{v}, \mathbf{w} \in V$, $(\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$.
3. Existence of additive identity: There exists an element $\mathbf{0} \in V$ such that $\mathbf{v} + \mathbf{0} = \mathbf{v}$ for all $\mathbf{v} \in V$.
4. Existence of additive inverses: For each $\mathbf{v} \in V$, there exists an element $-\mathbf{v} \in V$ such that $\mathbf{v} + (-\mathbf{v}) = \mathbf{0}$.
5. Commutativity of addition: For all $\mathbf{u}, \mathbf{v} \in V$, $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$.
6. Closure under scalar multiplication: For all $\alpha \in \mathbb{F}$ and $\mathbf{v} \in V$, $\alpha \mathbf{v} \in V$.
7. Distributivity of scalar multiplication with respect to vector addition: For all $\alpha \in \mathbb{F}$ and $\mathbf{u}, \mathbf{v} \in V$, $\alpha(\mathbf{u} + \mathbf{v}) = \alpha \mathbf{u} + \alpha \mathbf{v}$.
8. Distributivity of scalar multiplication with respect to scalar addition: For all $\alpha, \beta \in \mathbb{F}$ and $\mathbf{v} \in V$, $(\alpha + \beta)\mathbf{v} = \alpha \mathbf{v} + \beta \mathbf{v}$.
9. Associativity of scalar multiplication: For all $\alpha, \beta \in \mathbb{F}$ and $\mathbf{v} \in V$, $\alpha(\beta \mathbf{v}) = (\alpha \beta)\mathbf{v}$.
10. Existence of multiplicative identity: For every $\mathbf{v} \in V$, $1 \cdot \mathbf{v} = \mathbf{v}$, where 1 is the multiplicative identity in $\mathbb{F}$.

These axioms ensure that vector spaces generalize the concept of Euclidean spaces to potentially infinite dimensions and to fields other than $\mathbb{R}$. For instance, the space $\mathbb{R}^n$ of all $n$-tuples of real numbers forms a vector space over $\mathbb{R}$, where vector addition and scalar multiplication are defined component-wise. Another important example is the space of all continuous functions from $\mathbb{R}$ to $\mathbb{R}$, denoted $C(\mathbb{R}, \mathbb{R})$, which is also a vector space under pointwise addition and scalar multiplication.

Vector spaces provide the language for expressing geometric and algebraic properties in a unified way. They serve as the foundation for many areas of mathematics, including the theory of linear transformations, which plays a crucial role in understanding the architecture of transformers in deep learning (see Fig. 1.1). The concepts of symmetry and invariance, which are central to representation theory, are naturally expressed in terms of vector spaces and their linear transformations.

**Fig. 1.1** A simple network of neurons illustrating the flow of data through layers represented as vector spaces, highlighting the role of linear transformations in capturing geometric and algebraic properties essential for deep learning

### Subspaces and Bases

A subspace $W$ of a vector space $V$ is a subset of $V$ that is itself a vector space under the operations of addition and scalar multiplication inherited from $V$. Formally, a subset $W \subseteq V$ is a subspace if it satisfies the following conditions:

1. The zero vector of $V$ is in $W$, i.e., $\mathbf{0} \in W$.
2. $W$ is closed under addition: For all $\mathbf{u}, \mathbf{v} \in W$, $\mathbf{u} + \mathbf{v} \in W$.
3. $W$ is closed under scalar multiplication: For all $\alpha \in \mathbb{F}$ and $\mathbf{v} \in W$, $\alpha \mathbf{v} \in W$.

The concept of a subspace is crucial in many areas of mathematics, particularly in linear algebra and geometry. For example, the set of all solutions to a homogeneous system of linear equations is a subspace of $\mathbb{R}^n$. This subspace, known as the solution space or null space, encapsulates the degrees of freedom within the system.

A basis of a vector space $V$ is a set of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \ldots, \mathbf{v}_n\}$ in $V$ that is linearly independent and spans $V$. A set of vectors is said to be linearly independent if no vector in the set can be expressed as a linear combination of the others. The span of a set of vectors is the set of all possible linear combinations of those vectors. Mathematically, a set $\{\mathbf{v}_1, \mathbf{v}_2, \ldots, \mathbf{v}_n\}$ is a basis for $V$ if:

1. The vectors are linearly independent: $\sum_{i=1}^{n} \alpha_i \mathbf{v}_i = \mathbf{0}$ implies $\alpha_i = 0$ for all $i$.
2. The vectors span $V$: For every $\mathbf{v} \in V$, there exist scalars $\alpha_1, \alpha_2, \ldots, \alpha_n$ such that $\mathbf{v} = \sum_{i=1}^{n} \alpha_i \mathbf{v}_i$.

The number of vectors in a basis is called the dimension of the vector space. In the context of transformers, the notion of a basis can be seen as a way of encoding the essential features of data in a minimal, yet complete, form. The concept of basis extends naturally to function spaces, where Fourier and wavelet bases are particularly relevant in understanding the representations used in machine learning models.

For instance, in $\mathbb{R}^3$, the standard basis is $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$, where $\mathbf{e}_i$ is the vector with a 1 in the $i$th position and 0 elsewhere. Any vector in $\mathbb{R}^3$ can be uniquely expressed as a linear combination of these basis vectors. This idea generalizes to infinite-dimensional spaces, where bases may consist of an infinite number of vectors, such as the set of functions $\{1, x, x^2, x^3, \ldots\}$ in the vector space of polynomials.

Understanding subspaces and bases is fundamental for grasping the structure of vector spaces, which in turn is crucial for analyzing the linear algebraic operations underlying the mechanisms of transformers. These concepts also pave the way for exploring more advanced topics like eigenvectors, eigenspaces, and their role in simplifying linear transformations—a theme that recurs in the study of symmetries and invariances within transformer models.

### 1.1.2   Linear Transformations

Linear transformations, also known as linear maps, are the backbone of linear algebra and play a pivotal role in understanding the geometric and algebraic structure of vector spaces. A linear transformation $T$ from a vector space $V$ over a field $\mathbb{F}$ to another vector space $W$ over the same field is a function $T : V \rightarrow W$ that preserves vector addition and scalar multiplication. Formally, for all $\mathbf{u}, \mathbf{v} \in V$ and $\alpha \in \mathbb{F}$, the map $T$ satisfies the following properties:

$$T(\mathbf{u} + \mathbf{v}) = T(\mathbf{u}) + T(\mathbf{v})$$

$$T(\alpha\mathbf{u}) = \alpha T(\mathbf{u}).$$

These properties ensure that the transformation $T$ respects the linear structure of the vector space $V$, making $T$ an essential tool for analyzing and manipulating vectors in $V$. The concept of linear transformations generalizes the idea of matrix multiplication and encompasses many operations in both pure and applied mathematics, including the transformations used in neural networks and, specifically, in the architecture of transformers.

The geometric interpretation of a linear transformation is that it maps lines to lines and preserves the origin. This geometric perspective is crucial in understanding how linear transformations can be used to represent symmetries and invariances, which are

fundamental in the design of intelligent systems such as transformers. For example, in computer vision, linear transformations can represent rotations, translations, and scalings, which are important for recognizing objects regardless of their orientation or size.

**Kernels and Images**

Given a linear transformation $T : V \rightarrow W$, two subspaces of particular interest are the kernel and the image of $T$. The kernel of $T$, denoted $\ker(T)$, is the set of all vectors in $V$ that are mapped to the zero vector in $W$:

$$\ker(T) = \{\mathbf{v} \in V \mid T(\mathbf{v}) = \mathbf{0}\}.$$

The kernel is a subspace of $V$ because it satisfies the conditions for a subspace: it contains the zero vector, is closed under vector addition, and is closed under scalar multiplication. The dimension of the kernel, known as the nullity of $T$, provides important information about the linear dependence among the vectors in $V$. A transformation with a trivial kernel (i.e., $\ker(T) = \{\mathbf{0}\}$) is injective (one to one), meaning that $T$ preserves the distinctness of vectors, which is crucial in applications where uniqueness is important, such as encoding information.

The image of $T$, denoted $\text{Im}(T)$, is the set of all vectors in $W$ that can be expressed as $T(\mathbf{v})$ for some $\mathbf{v} \in V$:

$$\text{Im}(T) = \{\mathbf{w} \in W \mid \mathbf{w} = T(\mathbf{v}) \text{ for some } \mathbf{v} \in V\}.$$

The image is a subspace of $W$, and its dimension is called the rank of $T$. The rank of a transformation provides a measure of how much of the vector space $W$ is "covered" by the transformation $T$. The Rank-Nullity Theorem, a fundamental result in linear algebra, relates the dimensions of the kernel and image of a linear transformation:

**Theorem 1.1**  (Rank-Nullity Theorem) *Let V and W be vector spaces over a field* $\mathbb{F}$*, and let* $T : V \rightarrow W$ *be a linear transformation. Then the dimension of the vector space V is equal to the sum of the rank of T (the dimension of the image of T) and the nullity of T (the dimension of the kernel of T). Formally,*

$$dim(V) = rank(T) + nullity(T),$$

*where* $rank(T) = \dim(Im(T))$ *is the dimension of the image of T, and* $nullity(T) = \dim(\ker(T))$ *is the dimension of the kernel of T.*

This theorem encapsulates the idea that the dimension of the domain $V$ of a linear transformation is partitioned into two parts: the dimension of the image

(which corresponds to the effective output of the transformation) and the dimension of the kernel (which corresponds to the loss of information or degeneracy in the transformation) [4, 29].

The concepts of kernel and image are intimately related to symmetry and geometry. For example, if a linear transformation represents a symmetry operation, such as a rotation or reflection, the kernel may represent directions that remain invariant under the transformation (e.g., the axis of rotation), while the image represents the space that is affected by the transformation. These ideas are directly applicable to the design of transformer models, where understanding the flow of information and the preservation of structure under transformations is key to building intelligent systems.

**Matrix Representations**

Every linear transformation $T : V \to W$ can be represented by a matrix once bases for $V$ and $W$ are chosen. If $\{\mathbf{v}_1, \mathbf{v}_2, \ldots, \mathbf{v}_n\}$ is a basis for $V$ and $\{\mathbf{w}_1, \mathbf{w}_2, \ldots, \mathbf{w}_m\}$ is a basis for $W$, then the matrix representation of $T$, denoted $[T]$, is the $m \times n$ matrix whose $j$th column is the vector of coordinates of $T(\mathbf{v}_j)$ with respect to the basis $\{\mathbf{w}_i\}$:

$$[T] = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix},$$

where $T(\mathbf{v}_j) = \sum_{i=1}^{m} a_{ij}\mathbf{w}_i$. The matrix $[T]$ encapsulates all the information about the linear transformation $T$ and allows us to perform computations with $T$ in a straightforward manner using matrix arithmetic.

The matrix representation of a linear transformation provides a bridge between abstract linear algebra and concrete numerical methods. For instance, the composition of two linear transformations corresponds to the multiplication of their respective matrices. This correspondence is crucial in many applications, including those in machine learning, where complex operations are often broken down into sequences of linear transformations, each represented by a matrix.

Matrices also have a natural geometric interpretation. For example, in $\mathbb{R}^2$, a rotation matrix

$$R(\theta) = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}$$

represents a rotation by an angle $\theta$ about the origin. Similarly, a reflection matrix

$$M = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

represents a reflection across the x-axis. These geometric transformations are critical in understanding symmetries in physical systems and in the design of artificial intelligence models that must be invariant to such transformations.

In the context of transformers, matrix representations are used extensively to describe the operations within the network. For example, the weights of the layers in a transformer can be viewed as matrices that transform input vectors into output vectors. Understanding the properties of these matrices, such as their eigenvalues and eigenvectors, can provide deep insights into the behavior of the model, including its ability to capture symmetries and invariances in the data.

Moreover, the study of matrices and their transformations is closely related to the study of eigenvalues and eigenvectors, which describe the directions in which a transformation acts as a simple scaling operation. These concepts are fundamental in many areas of applied mathematics, including quantum mechanics, where the eigenvectors of an operator represent the possible states of a system, and in machine learning, where they can be used to identify the most important features or components of the data. Understanding linear transformations, their kernels, images, and matrix representations are essential for analyzing and designing transformer architectures. These concepts provide the mathematical framework for exploring how information is processed and transformed within the model, how symmetries are preserved or broken, and how the geometry of the input space is manipulated to produce intelligent behavior.

### *1.1.3 Eigenvalues and Eigenvectors*

Eigenvalues and eigenvectors are fundamental concepts in linear algebra that provide deep insights into the structure of linear transformations [54]. Given a linear transformation $T : V \rightarrow V$ on a vector space $V$, an eigenvector $\mathbf{v} \in V$ is a non-zero vector that is scaled by $T$ by a scalar factor known as the eigenvalue. Formally, $\mathbf{v}$ is an eigenvector of $T$ corresponding to the eigenvalue $\lambda \in \mathbb{F}$ if

$$T(\mathbf{v}) = \lambda \mathbf{v}.$$

This equation indicates that the action of $T$ on $\mathbf{v}$ does not change its direction but only its magnitude. The scalar $\lambda$ captures the amount of scaling, and the vector $\mathbf{v}$ reveals the direction along which this scaling occurs. To find the eigenvalues of a linear transformation, one must solve the characteristic equation, which is derived from the determinant of the matrix representation of $T$ minus $\lambda$ times the identity matrix:

$$\det([T] - \lambda I) = 0.$$

The solutions to this polynomial equation in $\lambda$ give the eigenvalues, and for each eigenvalue, the corresponding eigenvectors are found by solving the system of linear equations:

$$([T] - \lambda I)\mathbf{v} = \mathbf{0}.$$

Eigenvalues and eigenvectors are pivotal in many areas of mathematics, including the study of symmetries, stability analysis, and the decomposition of linear transformations into simpler, more interpretable components. In the context of transformers, eigenvalues and eigenvectors help us understand how information is propagated and transformed through the layers, providing insights into the model's ability to capture complex patterns and symmetries in data.

The geometric interpretation of eigenvectors is closely tied to the concepts of symmetry and invariance. For instance, in quantum mechanics, the eigenvectors of an operator represent the possible states of a system, each associated with a specific measurement outcome (the eigenvalue). Similarly, in the analysis of mechanical systems, eigenvectors describe the principal directions of vibration, while eigenvalues indicate the frequencies. These interpretations are directly applicable to machine learning models, where understanding the directions in which data is most naturally transformed can lead to better feature extraction and dimensionality reduction.

### Spectral Theorem

The spectral theorem is a cornerstone result in linear algebra, particularly in the study of normal operators on finite-dimensional vector spaces [25, 38]. It provides a powerful tool for analyzing linear transformations by decomposing them into simpler components. For a linear operator $T : V \rightarrow V$ on a finite-dimensional vector space $V$ over $\mathbb{C}$, the spectral theorem states that $T$ is diagonalizable if and only if it is normal, meaning $TT^* = T^*T$, where $T^*$ is the conjugate transpose of $T$. The spectral theorem asserts that a normal operator $T$ can be represented as

$$T = U \Lambda U^*,$$

where $U$ is a unitary matrix (satisfying $U^*U = UU^* = I$), and $\Lambda$ is a diagonal matrix whose entries are the eigenvalues of $T$. This result generalizes the diagonalization of symmetric matrices to a broader class of matrices and operators, allowing for a more profound understanding of their action.

The spectral theorem has profound implications in various domains, including quantum mechanics, where it underpins the spectral decomposition of observables, and in functional analysis, where it plays a critical role in the study of compact operators on Hilbert spaces. In the context of transformers, the spectral theorem can be used to analyze the stability and behavior of the layers by examining their spectral properties.

The geometric interpretation of the spectral theorem is that the transformation $T$ can be understood as a rotation (given by $U$) followed by a scaling along orthogonal

directions (given by $\Lambda$). This decomposition reveals the intrinsic symmetries of the transformation, which are crucial for understanding how transformers maintain or break symmetry during the processing of data. In machine learning, these symmetries often correspond to invariances in the data, such as rotational or translational invariance, which the model must learn to recognize and exploit.

## Diagonalization

Diagonalization is the process of finding a basis for a vector space $V$ in which the matrix representation of a linear transformation $T$ is diagonal. Formally, a matrix $A$ is said to be diagonalizable if there exists an invertible matrix $P$ such that

$$P^{-1}AP = D,$$

where $D$ is a diagonal matrix. The columns of $P$ are the eigenvectors of $A$, and the diagonal entries of $D$ are the corresponding eigenvalues. Diagonalization simplifies the analysis of linear transformations by reducing them to scaling operations along the eigenvectors, which are the principal directions of the transformation.

The process of diagonalization reveals the underlying structure of the transformation and allows for the decomposition of complex transformations into simpler, more interpretable components. This decomposition is particularly useful in the study of dynamical systems, where the long-term behavior of the system can often be understood by examining the eigenvalues of the transformation matrix. If the eigenvalues are all distinct and non-zero, the transformation is not only diagonalizable but also invertible, meaning that the system's behavior can be fully described by its eigenvalues and eigenvectors.

In the context of transformers, diagonalization can provide insights into the behavior of layers and their ability to propagate information. For instance, in the self-attention mechanism, understanding the eigenvalues and eigenvectors of the attention matrix can reveal how information is weighted and combined across different parts of the input sequence. This understanding is crucial for designing models that effectively capture the hierarchical and symmetric structures in data.

The process of diagonalization is closely related to the concepts of symmetry and invariance. A diagonal matrix is invariant under the change of basis represented by its eigenvectors, meaning that the transformation acts independently along each of these directions. This invariance is a key feature in many intelligent systems, where the ability to recognize and exploit symmetry can lead to more efficient and robust models. In particular, transformers, which are designed to process data with complex, hierarchical structures, can benefit greatly from the insights provided by diagonalization and spectral analysis.

## 1.2   Group Theory and Symmetries

### 1.2.1   Basic Concepts of Group Theory

Group theory is a branch of abstract algebra that studies algebraic structures known as groups. Groups are fundamental in understanding the concept of symmetry in mathematics, physics, and various other disciplines, including the study of neural networks and transformers. A group $G$ is a set equipped with a binary operation $\cdot$ (often called multiplication) that combines any two elements $a$ and $b$ in $G$ to form another element $a \cdot b$ in $G$. The group operation must satisfy the following axioms:

1. Closure: For all $a, b \in G$, the product $a \cdot b \in G$.
2. Associativity: For all $a, b, c \in G$, $(a \cdot b) \cdot c = a \cdot (b \cdot c)$.
3. Identity Element: There exists an element $e \in G$ such that for all $a \in G$, $e \cdot a = a \cdot e = a$. This element $e$ is called the identity element of the group.
4. Inverse Element: For each $a \in G$, there exists an element $a^{-1} \in G$ such that $a \cdot a^{-1} = a^{-1} \cdot a = e$, where $e$ is the identity element.

These axioms define the abstract structure of a group and provide a framework for studying the symmetry and invariance properties of various mathematical objects. Groups can be finite or infinite, and they can arise in many different contexts, from the symmetries of geometric shapes to the permutations of a set, and even to the automorphisms of algebraic structures.

   The study of groups is motivated by the need to understand symmetry in a structured way. For example, the set of all rotations of a regular polygon forms a group under the operation of composition, reflecting the symmetrical properties of the shape. This group encapsulates the geometric symmetries of the polygon and allows for a deep understanding of how these symmetries interact.

**Groups, Subgroups, and Cosets**

A subgroup $H$ of a group $G$ is a subset of $G$ that is itself a group under the operation of $G$. Formally, $H \subseteq G$ is a subgroup if it satisfies the following conditions:

1. The identity element of $G$ is in $H$.
2. $H$ is closed under the group operation: For all $h_1, h_2 \in H$, $h_1 \cdot h_2 \in H$.
3. $H$ is closed under taking inverses: For all $h \in H$, $h^{-1} \in H$.

   Subgroups are important because they inherit the algebraic structure of the larger group and often reveal the internal symmetries of the group. For example, the set of all rotations of a cube that leave a given face fixed forms a subgroup of the full rotation group of the cube. This subgroup reflects the symmetry of the cube with respect to that particular face.

   Given a subgroup $H$ of a group $G$, one can form cosets of $H$ in $G$. A left coset of $H$ in $G$ is a set of the form $gH = \{g \cdot h \mid h \in H\}$ for some $g \in G$. Similarly, a right

coset is a set of the form $Hg = \{h \cdot g \mid h \in H\}$. Cosets partition the group $G$ into disjoint subsets, each of which is a translation of the subgroup $H$ by some element of $G$. The number of distinct cosets of $H$ in $G$ is called the index of $H$ in $G$ and is denoted by $[G : H]$.

The concept of cosets is crucial in understanding the structure of groups and their subgroups. For instance, Lagrange's theorem, a fundamental result in group theory, states that the order of any finite subgroup $H$ of a finite group $G$ divides the order of $G$. This theorem can be understood through the lens of cosets, as it implies that the number of elements in $G$ is equal to the number of elements in $H$ multiplied by the number of cosets of $H$. Cosets also play a significant role in the construction of quotient groups, which are groups formed by "collapsing" a normal subgroup $N$ of $G$ to the identity element. The quotient group $G/N$ consists of the cosets of $N$ in $G$, and it inherits the group structure from $G$. Quotient groups are essential in understanding how larger groups can be decomposed into simpler components, a theme that is central to the study of symmetry and intelligence in mathematical systems.

In the context of machine learning, groups, subgroups, and cosets can be used to model the symmetries present in data. For example, in computer vision, the symmetries of an object under rotation and translation can be represented by a group, and the subgroup structure can reveal invariant features of the object. Transformers, with their ability to capture hierarchical structures, can be seen as exploiting these symmetries to efficiently process and transform data.

**Group Homomorphisms**

A group homomorphism is a function between two groups that preserves the group structure. Formally, if $G$ and $H$ are groups, a homomorphism from $G$ to $H$ is a function $\phi : G \to H$ such that for all $a, b \in G$,

$$\phi(a \cdot b) = \phi(a) \cdot \phi(b),$$

where $\cdot$ denotes the group operation in $G$ and in $H$. The property of preserving the group operation means that the image of a product under the homomorphism is the product of the images. This preservation of structure makes homomorphisms a central concept in the study of groups, as they allow for the comparison and classification of groups based on their structural similarities.

A homomorphism that is both injective (one to one) and surjective (onto) is called an isomorphism. If there exists an isomorphism between two groups $G$ and $H$, then $G$ and $H$ are said to be isomorphic, denoted $G \cong H$. Isomorphic groups are structurally identical, meaning that they have the same group-theoretic properties, though their elements may be different. This concept is fundamental in understanding how different mathematical objects can exhibit the same symmetries.

The kernel of a group homomorphism $\phi : G \to H$ is the set of elements in $G$ that map to the identity element in $H$:

$$\ker(\phi) = \{g \in G \mid \phi(g) = e_H\},$$

where $e_H$ is the identity element in $H$. The kernel is a normal subgroup of $G$, and the First Isomorphism Theorem states that the image of $\phi$ is isomorphic to the quotient group $G/\ker(\phi)$:

$$G/\ker(\phi) \cong \mathrm{Im}(\phi).$$

This theorem provides a deep connection between homomorphisms, kernels, and quotient groups, and it plays a crucial role in the classification and analysis of groups.

Group homomorphisms are not just abstract constructs; they have concrete applications in many areas of mathematics and science. For example, in geometry, the symmetry group of a shape can be related to the symmetry group of a different shape through a homomorphism, revealing how one set of symmetries can be mapped onto another. In physics, homomorphisms between symmetry groups can describe how different physical systems are related, such as how the symmetries of a molecule are related to the symmetries of its constituent atoms.

In the context of transformers and machine learning, group homomorphisms play a crucial role in modeling how symmetries in the input data are preserved or transformed by the network. To understand this, we can consider a scenario where the input data exhibits certain symmetries that can be described by a group $G$. The elements of $G$ represent symmetry operations, such as rotations, translations, or permutations, that can be applied to the data.

Let $X$ denote the input space, which could be a space of images, sequences, or any other structured data. The group $G$ acts on $X$ via a group action $\alpha : G \times X \to X$, where for each $g \in G$ and $x \in X$, the action $\alpha(g, x)$ describes the transformed version of $x$ under the symmetry operation $g$. In mathematical terms, this means that for each fixed $g \in G$, the map $\alpha_g : X \to X$ defined by $\alpha_g(x) = \alpha(g, x)$ is a transformation of $X$.

In machine learning, the data $X$ is often embedded into a higher-dimensional feature space $V$ via a representation map $\phi : X \to V$. This embedding can be seen as a preprocessing step where the raw data $x \in X$ is transformed into a feature vector $\phi(x) \in V$, which is then fed into the neural network. If the data $X$ has an inherent symmetry described by the group $G$, it is desirable for the embedding $\phi$ to respect this symmetry. Specifically, if $\alpha_g(x)$ is the transformed data under the symmetry $g$, we would like the embedding of the transformed data $\phi(\alpha_g(x))$ to be related to the embedding of the original data $\phi(x)$ in a way that reflects the action of $g$.

This leads to the concept of a group homomorphism associated with the neural network layers. Let $\rho : G \to \mathrm{GL}(V)$ be a representation of the group $G$ on the vector space $V$, where $\mathrm{GL}(V)$ denotes the group of invertible linear transformations on $V$.

The map $\rho(g)$ represents the action of the group element $g$ on the feature space $V$. Ideally, we want the embedding $\phi$ to satisfy the following equivariance condition:

$$\phi(\alpha_g(x)) = \rho(g)\phi(x) \quad \text{for all } g \in G \text{ and } x \in X.$$

This equation states that applying a group transformation $g$ to the input data $x$ and then embedding the result into the feature space $V$ should be equivalent to first embedding $x$ into $V$ and then applying the linear transformation $\rho(g)$ to the resulting feature vector. The map $\rho$ is a group homomorphism because it preserves the group structure:

$$\rho(g_1 g_2) = \rho(g_1)\rho(g_2) \quad \text{for all } g_1, g_2 \in G.$$

The layers of a transformer can then be viewed as a sequence of homomorphisms $T_1, T_2, \ldots, T_n$, each mapping the feature space $V_i$ of the $i$th layer to the feature space $V_{i+1}$ of the next layer. Each transformation $T_i : V_i \to V_{i+1}$ is expected to preserve the group action, meaning that for each layer $i$ there exists a homomorphism $\rho_i : G \to \text{GL}(V_i)$ such that

$$T_i(\rho_i(g)\mathbf{v}) = \rho_{i+1}(g)T_i(\mathbf{v}) \quad \text{for all } \mathbf{v} \in V_i \text{ and } g \in G.$$

This condition ensures that the symmetry described by the group $G$ is preserved throughout the layers of the network. In other words, the action of the group on the data at any layer is consistent with the action of the group on the data at previous and subsequent layers.

To capture this symmetry in practice, one often designs the architecture of the transformer so that each layer $T_i$ naturally respects the group action. For example, in the case of rotational symmetries, the feature spaces $V_i$ can be chosen to be spaces of spherical harmonics, which are representations of the rotation group SO(3). The layers $T_i$ are then constructed to act on these spaces in a way that is consistent with the rotational symmetry.

Understanding these transformations through the lens of group theory leads to more robust and interpretable models. The equivariance condition ensures that the network processes data in a way that respects the symmetries present in the input, leading to better generalization to unseen data that exhibits the same symmetries. For instance, if the input data consists of images, and the group $G$ represents the group of rotations, then a network that respects the rotational symmetry will perform consistently regardless of the orientation of the input images. Moreover, this approach allows for a more efficient representation of the data, as the model can leverage the symmetry to reduce the complexity of the learning task. Instead of learning separate representations for each possible transformation of the data, the model can learn a single representation that is equivariant under the group action. This not only reduces the number of parameters needed but also leads to a more interpretable model, as the learned features align with the underlying symmetries of the data.

## *1.2.2   Representation Theory of Finite Groups*

Representation theory is a powerful tool that allows us to study abstract algebraic structures, such as groups, by representing their elements as matrices and their operations as matrix multiplication. This approach provides a concrete way to analyze the symmetries of a system, which is particularly relevant in areas such as physics, chemistry, and computer science. In the context of transformers, understanding group representations is crucial for analyzing how symmetries in data are captured and processed by the model.

**Group Representations**

A representation of a group $G$ on a vector space $V$ over a field $\mathbb{F}$ is a homomorphism $\rho : G \to \mathrm{GL}(V)$, where $\mathrm{GL}(V)$ is the group of all invertible linear transformations of $V$. In other words, a representation is a way of associating each element $g \in G$ with an invertible matrix $\rho(g)$ in such a way that the group operation is preserved:

$$\rho(g_1 g_2) = \rho(g_1)\rho(g_2)$$

for all $g_1, g_2 \in G$. The vector space $V$ is called the representation space of $\rho$, and the dimension of $V$ is called the degree of the representation. If $V$ is finite-dimensional, the representation $\rho$ can be described by a set of $n \times n$ matrices, where $n$ is the dimension of $V$.

The study of group representations allows us to understand the structure of a group by examining how it acts on vector spaces. For example, consider the cyclic group $C_n = \langle g \rangle$ of order $n$. A representation of $C_n$ is given by associating the generator $g$ with a matrix $\rho(g)$ such that $\rho(g^n) = I$, where $I$ is the identity matrix. One simple representation of $C_n$ is the one-dimensional representation where $\rho(g) = e^{2\pi i/n}$, which corresponds to a rotation by $2\pi/n$ in the complex plane.

Group representations are particularly important in understanding symmetries in geometry and physics. For instance, the rotation group SO(3) has representations that describe how objects in three-dimensional space can be rotated. These representations are used extensively in quantum mechanics, where the symmetries of a system are represented by the group of rotations, and the states of the system correspond to vectors in a representation space.

In the context of machine learning and transformers, group representations can be used to model how data is transformed as it passes through the layers of the network. For example, the self-attention mechanism in transformers can be interpreted as a representation of a permutation group, where the elements of the group correspond to different ways of permuting the input sequence. By understanding the representation theory of the group, we can gain insights into how the model captures the symmetries and invariances in the data.

**Character Theory**

Character theory is a branch of representation theory that focuses on the trace of the matrices associated with group elements. The character of a representation $\rho : G \to$ $GL(V)$ is a function $\chi : G \to \mathbb{F}$ defined by

$$\chi(g) = \text{Tr}(\rho(g)),$$

where $\text{Tr}(\rho(g))$ denotes the trace of the matrix $\rho(g)$. The character $\chi$ encodes important information about the representation, such as its degree (given by $\chi(e)$, where $e$ is the identity element of $G$) and how the representation decomposes into irreducible components.

One of the key results in character theory is that characters are class functions, meaning they are constant on conjugacy classes of the group. In other words, if $g_1$ and $g_2$ are conjugate in $G$ (i.e., there exists $h \in G$ such that $g_2 = hg_1h^{-1}$), then

$$\chi(g_1) = \chi(g_2).$$

This property greatly simplifies the study of representations, as it reduces the problem to understanding the characters on a finite number of conjugacy classes.

Characters play a crucial role in the classification of representations. The orthogonality relations for characters provide a powerful tool for determining whether two representations are equivalent and for decomposing a given representation into irreducible components. For a finite group $G$ and two irreducible characters $\chi$ and $\psi$, the orthogonality relation states that

$$\frac{1}{|G|} \sum_{g \in G} \chi(g)\overline{\psi(g)} = \delta_{\chi\psi},$$

where $|G|$ is the order of the group, and $\delta_{\chi\psi}$ is the Kronecker delta, which is 1 if $\chi = \psi$ and 0 otherwise. This relation implies that the characters of different irreducible representations are orthogonal, and it provides a method for finding the multiplicities of irreducible components in a given representation.

Character theory also has deep connections to the geometry of the underlying space on which the group acts. For example, the characters of the rotation group $SO(3)$ are related to spherical harmonics, which describe the symmetries of functions on the sphere. These symmetries are exploited in many areas of physics and mathematics, including the analysis of atomic orbitals and the study of vibrations in mechanical systems.

In the context of transformers, character theory can be used to analyze how the model captures the symmetries in the data. For example, if the data exhibits a certain group symmetry, the character of the representation associated with the model can reveal how this symmetry is preserved or broken as the data is processed. This analysis can lead to a deeper understanding of the model's behavior and its ability to generalize to new data.

## 1.2.3 Applications to Transformers

The application of group theory and representation theory to transformer architectures provides a profound understanding of how these models capture and utilize symmetries in data. Transformers, by their design, exploit various symmetries to achieve invariance and equivariance in tasks such as natural language processing, image recognition, and other domains where hierarchical and structured data are prevalent. By exploring the symmetries inherent in these architectures, we can develop deeper insights into how transformers generalize across different tasks and how their performance can be enhanced through mathematical principles.

### Symmetries in Transformer Architectures

Symmetry plays a central role in the design and functioning of transformer architectures. A key symmetry in transformers is the permutation symmetry of the input sequences. In a standard transformer model, the input data, such as a sequence of words in a sentence, is processed without any inherent order bias, thanks to the self-attention mechanism. The self-attention mechanism treats all positions in the input sequence symmetrically, allowing the model to focus on different parts of the sequence based on the learned attention weights.

Mathematically, this permutation symmetry can be understood through the lens of group theory. Let $S_n$ denote the symmetric group on $n$ elements, which represents all possible permutations of an $n$-element sequence. The transformer model is invariant under the action of this group, meaning that if $\sigma \in S_n$ is a permutation of the input sequence indices, the output of the transformer remains consistent with this permutation, modulo the learned attention weights:

$$\text{Transformer}(\sigma \cdot \mathbf{x}) = \sigma \cdot \text{Transformer}(\mathbf{x}),$$

where $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ is the input sequence, and $\sigma \cdot \mathbf{x} = (x_{\sigma(1)}, x_{\sigma(2)}, \ldots, x_{\sigma(n)})$ represents the permuted sequence. This symmetry ensures that the model does not privilege any specific order of the input sequence, making it particularly effective in tasks where the order of elements can vary or where relationships between elements are more important than their positions.

Another important symmetry in transformer architectures is related to the attention mechanism itself. The self-attention mechanism computes a weighted sum of the input features, where the weights are determined by the learned attention scores. These attention scores can be viewed as elements of a matrix, and the transformation induced by the attention mechanism can be seen as a linear operator acting on the input features. The symmetry here lies in the fact that the attention mechanism is equivariant to the input features: it applies the same operation to all elements of the sequence, preserving the overall structure of the data.

In more advanced transformer architectures, such as those used in image processing or graph-based models, additional symmetries may be present. For example, in vision transformers (ViTs), rotational and translational symmetries of the input images can be exploited to improve model performance. By incorporating these symmetries into the architecture, the model can become more robust to variations in the input data, leading to better generalization across different tasks and datasets.

### Invariant Subspaces in Attention Mechanisms

In the context of attention mechanisms, invariant subspaces play a crucial role in understanding how information is processed and retained throughout the layers of a transformer model. An invariant subspace under a linear operator $T : V \rightarrow V$ is a subspace $W \subseteq V$ such that for every vector $\mathbf{w} \in W$, $T(\mathbf{w}) \in W$. In the context of transformers, the linear operator can represent the action of the attention mechanism on the feature space, and the invariant subspaces correspond to directions in the feature space that are preserved under the attention transformation.

To formalize this, let $A$ represent the attention matrix derived from the self-attention mechanism, and let $\mathbf{x}$ be the input vector representing the features at a particular position in the sequence. The attention mechanism can be viewed as applying a linear transformation $T_A$ to $\mathbf{x}$, where

$$T_A(\mathbf{x}) = A\mathbf{x}.$$

An invariant subspace $W$ under $T_A$ satisfies

$$T_A(\mathbf{w}) = A\mathbf{w} \in W \quad \text{for all } \mathbf{w} \in W.$$

Invariant subspaces are significant because they represent directions in the feature space that remain unchanged by the attention mechanism. These directions can correspond to important features or patterns in the data that the model has learned to recognize and preserve throughout the layers. For example, in natural language processing, certain syntactic or semantic relationships between words may correspond to invariant subspaces under the attention transformation, allowing the model to maintain these relationships as the information propagates through the network.

The study of invariant subspaces is closely related to the eigenvalues and eigenvectors of the attention matrix $A$. The eigenvectors corresponding to eigenvalues with absolute value 1 define invariant subspaces that are neither amplified nor diminished by the attention mechanism. These eigenvectors can be interpreted as the principal directions along which the attention mechanism operates, preserving certain features while allowing others to be modified.

This concept can be further extended to analyze the stability and convergence properties of the transformer model. If the attention mechanism consistently preserves certain invariant subspaces, the model is more likely to retain important information throughout its layers, leading to better performance and generalization. On the other

hand, if the attention mechanism disrupts these invariant subspaces, the model may struggle to maintain coherence in its output, leading to poorer performance.

## 1.3  Metric Spaces and Topological Preliminaries

### 1.3.1  Definition of Metric Spaces

A metric space is a set equipped with a function that defines a notion of distance between any two elements. Formally, a metric space is a pair $(X, d)$, where $X$ is a set, and $d : X \times X \to \mathbb{R}$ is a function (called a metric or distance function) satisfying the following properties for all $x, y, z \in X$:

1. Non-negativity: $d(x, y) \geq 0$ (the distance between any two points is non-negative).
2. Identity of indiscernibles: $d(x, y) = 0$ if and only if $x = y$ (the distance between two distinct points is positive).
3. Symmetry: $d(x, y) = d(y, x)$ (the distance is symmetric with respect to its arguments).
4. Triangle inequality: $d(x, z) \leq d(x, y) + d(y, z)$ (the direct distance between two points is always less than or equal to the sum of the distances through a third point).

These axioms encapsulate the essential properties of distance in a geometric space and provide the foundation for many topological concepts. The metric $d$ induces a topology on $X$, which allows us to define concepts such as continuity, convergence, and compactness. The importance of metric spaces in the study of attention mechanisms lies in the ability to measure distances between representations (such as word embeddings or image features) and to understand how these distances influence the behavior of the model.

In the context of machine learning and transformers, the metric space $(X, d)$ can represent the space of input data (such as the space of word embeddings in natural language processing or the space of image features in computer vision). The metric $d$ measures how similar or different two data points are, which is critical for tasks such as clustering, classification, and information retrieval. For example, in a language model, the distance between two word embeddings might reflect the semantic similarity between the corresponding words, influencing the attention mechanism's ability to focus on relevant parts of the input sequence.

### Examples of Metric Spaces

The concept of a metric space is quite general and can be instantiated in various ways, depending on the nature of the set $X$ and the metric $d$. Some common examples of metric spaces include

1. Euclidean Space: The most familiar example of a metric space is the $n$-dimensional Euclidean space $\mathbb{R}^n$ equipped with the Euclidean metric. For any two points $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ and $\mathbf{y} = (y_1, y_2, \ldots, y_n)$ in $\mathbb{R}^n$, the Euclidean distance is given by

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}.$$

This distance corresponds to the straight-line distance between the points $\mathbf{x}$ and $\mathbf{y}$ in $\mathbb{R}^n$, and it satisfies all the properties of a metric. The Euclidean space is fundamental in geometry and is the setting for many problems in machine learning, including clustering and classification.

2. Discrete Metric Space: For any set $X$, the discrete metric is defined by

$$d(x, y) = \begin{cases} 0 & \text{if } x = y, \\ 1 & \text{if } x \neq y \end{cases}.$$

This metric treats all distinct points as being at the same distance from each other, regardless of their nature. The discrete metric space is useful in theoretical considerations where we want to analyze properties that hold in the most extreme case of separation between points.

3. Manhattan Distance (Taxicab Metric): In $\mathbb{R}^n$, another important metric is the Manhattan distance, defined by

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{n} |x_i - y_i|.$$

This metric reflects the distance traveled along grid-like paths, such as the streets in a city laid out in a grid pattern. The Manhattan distance is widely used in machine learning, particularly in problems where the grid structure is inherent, such as image processing with pixel grids.

4. Cosine Similarity as a Metric: In the context of high-dimensional vector spaces, such as those used in word embeddings, the cosine similarity is often used as a measure of similarity between vectors. While cosine similarity itself is not a metric (because it does not satisfy the triangle inequality), the related cosine distance defined by

$$d(\mathbf{x}, \mathbf{y}) = 1 - \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

can be used to construct a valid metric, under certain conditions. This metric is particularly useful in natural language processing, where the direction of vectors (rather than their magnitude) is often more informative.

5. Hamming Distance: For a set of strings of equal length over a fixed alphabet, the Hamming distance between two strings is the number of positions at which the corresponding symbols differ. Mathematically, for two strings $s_1$ and $s_2$ of length $n$,

$$d(s_1, s_2) = \sum_{i=1}^{n} \delta(s_1[i], s_2[i]),$$

where $\delta(a, b) = 1$ if $a \neq b$, and 0 otherwise. The Hamming distance is widely used in coding theory and information retrieval, where it measures the difference between binary strings or other discrete sequences.

These examples illustrate the versatility of metric spaces in capturing different notions of distance and similarity across various domains. In machine learning, choosing the appropriate metric for the problem at hand is crucial for the success of the model, as the metric determines how the model perceives the relationships between data points and how it learns to focus on relevant features.

### Convergence and Completeness

Convergence and completeness are fundamental concepts in the study of metric spaces, providing a framework for understanding the behavior of sequences and the structure of the space itself. A sequence $\{x_n\}$ in a metric space $(X, d)$ is said to converge to a point $x \in X$ if, for every $\epsilon > 0$, there exists an integer $N$ such that for all $n \geq N$, $d(x_n, x) < \epsilon$. In other words, the elements of the sequence eventually get arbitrarily close to $x$, and $x$ is called the limit of the sequence, denoted $\lim_{n \to \infty} x_n = x$.

Convergence in metric spaces generalizes the notion of convergence in Euclidean spaces to more abstract settings. For example, in the space of continuous functions, convergence can describe the behavior of function sequences as they approach a limiting function. This concept is crucial in machine learning, particularly in the analysis of algorithms that iteratively update models to minimize a loss function. Understanding the conditions under which these sequences converge ensures that the algorithm behaves predictably and leads to an optimal solution.

A metric space $(X, d)$ is said to be complete if every Cauchy sequence in $X$ converges to a point in $X$. A sequence $\{x_n\}$ is a Cauchy sequence if, for every $\epsilon > 0$, there exists an integer $N$ such that for all $m, n \geq N$, $d(x_m, x_n) < \epsilon$. Completeness ensures that the space has no "holes" or "gaps," meaning that the space is sufficiently "large" to contain the limits of all converging sequences. The Completeness Theorem can be stated as follows:

**Theorem 1.2** (Completeness of Metric Spaces) *Let $(X, d)$ be a metric space. Then $X$ is complete if and only if every Cauchy sequence in $X$ has a limit in $X$.*

This theorem is fundamental in functional analysis, where the completeness of function spaces under various norms determines whether certain methods, such as

Fourier series expansions or integral transforms, can be applied. In the context of transformers and neural networks, completeness plays a role in ensuring that optimization algorithms, such as gradient descent, behave well, particularly when working in high-dimensional spaces where the geometry of the space can become complex.

For instance, in training deep neural networks, the parameter space is often modeled as a high-dimensional metric space, and the completeness of this space ensures that sequences of iteratively updated parameters converge to a well-defined limit, representing the optimal set of parameters that minimizes the loss function. Without completeness, the optimization process might diverge or oscillate without reaching a stable solution.

## 1.3.2  Topology of Metric Spaces

Topology provides a framework for understanding the structure and properties of metric spaces beyond mere distances. In this section, we delve into the fundamental topological concepts within metric spaces that are crucial for studying the geometry, symmetry, and continuity properties that are often leveraged in the design of intelligent systems like transformers.

### Open and Closed Sets

In the context of a metric space $(X, d)$, the notions of open and closed sets are fundamental in defining the topology induced by the metric. An open set in $X$ is a subset $U \subseteq X$ such that for every point $x \in U$, there exists a radius $\epsilon > 0$ for which the open ball centered at $x$ with radius $\epsilon$ is entirely contained within $U$. Formally, this means

$$\text{For each } x \in U, \text{ there exists } \epsilon > 0 \text{ such that } B(x, \epsilon) = \{y \in X \mid d(x, y) < \epsilon\} \subseteq U.$$

Open sets can be thought of as generalizations of the concept of intervals in real analysis. In Euclidean spaces, for example, the set of all points within a certain distance from a given point forms an open ball, which is a prototypical example of an open set. Open sets are the building blocks of the topology on $X$, as they satisfy the following properties:

1. The union of any collection of open sets is open.
2. The intersection of a finite number of open sets is open.
3. The empty set $\emptyset$ and the entire space $X$ are open.

These properties allow us to define a topology on $X$ as the collection of all open sets. This topology governs the notions of convergence, continuity, and compactness within the space.

Conversely, a closed set in $X$ is defined as a subset $C \subseteq X$ whose complement $X \setminus C$ is open. Equivalently, $C$ can be characterized by the property that it contains all its limit points, i.e., if a sequence $\{x_n\}$ in $C$ converges to a point $x \in X$, then $x$ must also belong to $C$. Mathematically, $C$ is closed if:

For every convergent sequence $\{x_n\}$ in $C$ with $\lim_{n \to \infty} x_n = x$, we have $x \in C$.

Closed sets in Euclidean space include familiar examples such as closed intervals $[a, b]$, finite sets, and the whole space $\mathbb{R}^n$. Closed sets are also important in defining the concept of continuity and in understanding the structure of metric spaces, as they possess properties analogous to those of open sets:

1. The intersection of any collection of closed sets is closed.
2. The union of a finite number of closed sets is closed.
3. The empty set $\emptyset$ and the entire space $X$ are closed.

In many applications, including those in machine learning and neural networks, the distinction between open and closed sets plays a critical role in defining neighborhoods, boundaries, and the behavior of functions within the space.

### Continuity and Compactness

Continuity is a key concept that relates the structure of metric spaces to the behavior of functions defined on them. A function $f : (X, d_X) \to (Y, d_Y)$ between two metric spaces is said to be continuous if it preserves the notion of "closeness" between points. Formally, $f$ is continuous at a point $x \in X$ if, for every $\epsilon > 0$, there exists a $\delta > 0$ such that

$$\text{If } d_X(x, x') < \delta, \text{ then } d_Y(f(x), f(x')) < \epsilon.$$

This definition, known as the $\epsilon$-$\delta$ definition of continuity, ensures that small changes in the input $x$ result in small changes in the output $f(x)$. Intuitively, a continuous function does not exhibit any "jumps" or "breaks," making it essential for analyzing the smoothness and stability of transformations within machine learning models.

A function $f : X \to Y$ is continuous on the entire space $X$ if it is continuous at every point in $X$. In the language of topology, $f$ is continuous if the preimage of every open set in $Y$ is an open set in $X$. That is,

$$f^{-1}(V) \text{ is open in } X \text{ whenever } V \text{ is open in } Y.$$

This topological definition of continuity connects directly with the earlier definition in metric spaces and generalizes to more abstract settings.

Compactness is another fundamental topological property that has deep implications in analysis and geometry. A subset $K \subseteq X$ of a metric space $(X, d)$ is said to be compact if every open cover of $K$ has a finite subcover. An open cover of $K$ is a collection of open sets $\{U_\alpha\}_{\alpha \in A}$ such that $K \subseteq \bigcup_{\alpha \in A} U_\alpha$, and a finite subcover is a finite subcollection $\{U_{\alpha_1}, U_{\alpha_2}, \ldots, U_{\alpha_n}\}$ that still covers $K$.

Compactness can be understood as a generalization of the concept of finiteness. In Euclidean spaces, the Heine–Borel theorem provides a useful characterization: a subset of $\mathbb{R}^n$ is compact if and only if it is closed and bounded. This characterization is particularly important in analysis, where compact sets exhibit properties such as:

1. Every sequence in a compact set has a convergent subsequence (sequential compactness).

2. Continuous functions on compact sets attain their maximum and minimum values (extreme value theorem).

In the context of machine learning and attention mechanisms, compactness is important because it ensures the existence of solutions and the stability of algorithms. For example, when optimizing a loss function over a compact set of parameters, one can guarantee that a minimum exists and that the optimization process will converge to a well-defined solution. Moreover, in the analysis of neural networks, compactness can be used to study the behavior of functions representing the network layers, ensuring that the transformation of data through the network preserves essential properties.

Continuity and compactness also intersect in the study of the stability and generalization of machine learning models. For instance, the Arzelà–Ascoli theorem provides conditions under which a family of continuous functions is relatively compact, meaning that any sequence within the family has a subsequence that converges uniformly to a continuous function. This result has applications in understanding how neural networks approximate functions and how well they generalize to unseen data.

### 1.3.3   Mappings Between Metric Spaces

Mappings between metric spaces play a crucial role in understanding how structures and properties are preserved or transformed. These mappings can reveal deep insights into the geometry and symmetry of spaces, and are foundational to many concepts in analysis, geometry, and machine learning. In the context of transformers and attention mechanisms, understanding these mappings helps us explore how information is processed, how distances between data points are preserved or altered, and how convergence to fixed points is achieved.

**Isometries and Contractions**

An isometry is a mapping between metric spaces that preserves distances. Formally, if $(X, d_X)$ and $(Y, d_Y)$ are two metric spaces, a function $f : X \to Y$ is called an isometry if for all $x_1, x_2 \in X$,

$$d_Y(f(x_1), f(x_2)) = d_X(x_1, x_2).$$

This property means that the structure of $X$ is preserved under $f$; the distances between points in $X$ remain unchanged when mapped to $Y$. Isometries are important in geometry because they reflect symmetries of the space. For instance, in Euclidean space $\mathbb{R}^n$, rotations, translations, and reflections are examples of isometries, as they preserve the Euclidean distance between points.

Isometries are not only significant in pure mathematics but also in the design of machine learning models, particularly in ensuring that certain geometric properties of the data are preserved as it is transformed through different layers of a neural network. For example, in applications like image processing, it is often desirable that the network's transformations preserve the spatial relationships between pixels, which is a type of isometry.

A contraction mapping is a stronger notion where the distances between points are not only preserved but actually reduced. A function $f : X \to X$ on a metric space $(X, d)$ is called a contraction if there exists a constant $0 \leq c < 1$ such that for all $x_1, x_2 \in X$,

$$d(f(x_1), f(x_2)) \leq c \cdot d(x_1, x_2).$$

The constant $c$ is called the Lipschitz constant of the contraction. The significance of contraction mappings lies in their tendency to bring points closer together, which has profound implications for the existence and uniqueness of fixed points—points $x^* \in X$ such that $f(x^*) = x^*$.

Contractions are particularly important in the analysis of iterative algorithms, where a sequence of approximations $\{x_n\}$ is generated by repeatedly applying a contraction mapping. The property that contractions reduce distances ensures that the sequence converges to a unique fixed point, as guaranteed by the Banach fixed-point theorem.

**The Banach Fixed-Point Theorem**

The Banach fixed-point theorem, also known as the contraction mapping theorem, is a fundamental result in metric space theory that provides conditions under which a contraction mapping on a complete metric space has a unique fixed point. The theorem is stated as follows:

**Theorem 1.3** (Banach Fixed-Point Theorem) *Let $(X, d)$ be a complete metric space, and let $f : X \to X$ be a contraction mapping with Lipschitz constant $c$ (i.e., $d(f(x_1), f(x_2)) \leq c \cdot d(x_1, x_2)$ for all $x_1, x_2 \in X$ with $0 \leq c < 1$). Then*

1. *$f$ has a unique fixed point $x^* \in X$ such that $f(x^*) = x^*$.*
2. *For any initial point $x_0 \in X$, the sequence $\{x_n\}$ defined by $x_{n+1} = f(x_n)$ converges to $x^*$ as n tends to infinity.*

The Banach fixed-point theorem is not only a cornerstone of analysis but also has far-reaching applications in various fields, including numerical analysis, differential equations, and machine learning. Its utility lies in the fact that it guarantees convergence to a solution under mild conditions, making it an invaluable tool for proving the existence and uniqueness of solutions in iterative processes.

In the context of machine learning, the Banach fixed-point theorem can be applied to analyze the convergence properties of training algorithms, particularly in the optimization of neural networks. For example, when training a model, one often seeks to minimize a loss function by iteratively updating the model's parameters. If the update rule can be modeled as a contraction mapping, the Banach fixed-point theorem guarantees that the sequence of parameter updates will converge to a unique set of optimal parameters, provided the metric space of parameters is complete.

Furthermore, in the design of attention mechanisms within transformers, the concept of fixed points and contractions can be used to ensure that the iterative process of refining attention weights leads to stable and consistent results. For instance, in some advanced attention mechanisms, iterative refinement of weights may be employed to achieve better focus on relevant parts of the input sequence. By ensuring that these refinements act as contractions, one can guarantee convergence to a stable set of attention weights, thereby improving the robustness and performance of the model.

The Banach fixed-point theorem also plays a role in understanding the stability of recurrent neural networks (RNNs), where the network's output is fed back as input in a recursive manner. Ensuring that the transformation applied at each step is a contraction helps in preventing the explosion or vanishing of gradients, thus maintaining the stability of the learning process.

## 1.4 Mathematical Foundations of Attention

### 1.4.1 Attention as a Mapping

Attention mechanisms are at the heart of modern transformer architectures, providing a powerful method for dynamically focusing on relevant parts of the input data. Mathematically, attention can be understood as a mapping that transforms input features into output features by assigning different weights to different parts of the input. This mapping is central to the model's ability to capture dependencies and relationships across different elements in the input sequence.

**Formulation of Attention as a Function**

Let $X = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n\}$ be a set of input vectors, where each $\mathbf{x}_i \in \mathbb{R}^d$ represents a feature vector in a $d$-dimensional space. The goal of the attention mechanism is to produce a weighted combination of these input vectors, emphasizing those that are most relevant to a particular context or query. This process can be mathematically formulated as a function $A : \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d}$, where the first argument $Q$ represents the query matrix, composed of $n$ query vectors; the second argument $K$ represents the key matrix, composed of $n$ key vectors; and the third argument $V$ represents the value matrix, composed of $n$ value vectors.

The attention function $A(Q, K, V)$ is typically defined by first computing the attention scores as a similarity measure between the query vectors and the key vectors. A common choice for this similarity measure is the scaled dot product, given by

$$S_{ij} = \frac{\langle \mathbf{q}_i, \mathbf{k}_j \rangle}{\sqrt{d_k}},$$

where $\mathbf{q}_i$ and $\mathbf{k}_j$ are the $i$th query vector and $j$th key vector, respectively, $d_k$ is the dimensionality of the key vectors, and $\langle \mathbf{q}_i, \mathbf{k}_j \rangle$ denotes the dot product. The scaling factor $\frac{1}{\sqrt{d_k}}$ is introduced to mitigate the effect of increasing dimensionality, ensuring that the magnitude of the dot products remains stable as $d_k$ grows.

The attention scores $S_{ij}$ are then normalized using the softmax function to produce the attention weights:

$$\alpha_{ij} = \frac{\exp(S_{ij})}{\sum_{k=1}^{n} \exp(S_{ik})}.$$

These attention weights $\alpha_{ij}$ indicate the relative importance of the $j$th input vector to the $i$th query vector. The final output of the attention mechanism is computed as a weighted sum of the value vectors:

$$\mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{v}_j,$$

where $\mathbf{v}_j$ is the $j$th value vector, and $\mathbf{z}_i$ is the output vector corresponding to the $i$th query vector. The complete attention mapping can be expressed as

$$A(Q, K, V) = \mathrm{Softmax} \left( \frac{Q K^{\top}}{\sqrt{d_k}} \right) V.$$

This formulation highlights the role of attention as a mapping that transforms the input sequence $X$ into a new sequence $Z = \{\mathbf{z}_1, \mathbf{z}_2, \ldots, \mathbf{z}_n\}$ by focusing on the most relevant components of the input based on the learned similarity structure.

**Properties of Attention Mappings**

The attention mapping $A : \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d}$ possesses several key properties that are essential for understanding its behavior and effectiveness in capturing relationships within the data. These properties can be analyzed in terms of geometry, symmetry, and stability.

1. Linearity in Value Vectors: The attention mechanism is linear with respect to the value vectors $V$. Given a linear combination of value matrices $V_1$ and $V_2$ with corresponding scalars $\alpha$ and $\beta$, the attention mapping satisfies

$$A(Q, K, \alpha V_1 + \beta V_2) = \alpha A(Q, K, V_1) + \beta A(Q, K, V_2).$$

This linearity property allows the attention mechanism to combine different sets of value vectors in a controlled manner, making it suitable for tasks where different sources of information need to be aggregated.

2. Invariance Under Permutations: The attention mechanism is invariant under simultaneous permutations of the input sequence. Let $\sigma$ be a permutation of the indices $\{1, 2, \ldots, n\}$. Then, for any permutation $\sigma$, we have

$$A(Q^\sigma, K^\sigma, V^\sigma) = A(Q, K, V)^\sigma,$$

where $Q^\sigma$, $K^\sigma$, and $V^\sigma$ denote the permuted query, key, and value matrices, respectively. This invariance reflects the fact that the attention mechanism does not impose any fixed order on the input sequence, allowing it to flexibly capture dependencies across different parts of the sequence regardless of their positions.

3. Commutativity of Attention Heads: In multi-head self-attention, the order in which the attention heads are applied does not affect the final output. This commutativity property can be expressed as

$$\bigoplus (Z_{\sigma(1)}, Z_{\sigma(2)}, \ldots, Z_{\sigma(h)}) W_O = \bigoplus (Z_1, Z_2, \ldots, Z_h) W_O$$

for any permutation $\sigma$ of the indices $\{1, 2, \ldots, h\}$. This property allows the model to process different attention heads in parallel, contributing to the efficiency and scalability of transformers.

4. Boundedness and Stability: The attention mapping is bounded and stable under small perturbations of the input vectors. Let $\delta Q$, $\delta K$, and $\delta V$ be small perturbations of the query, key, and value matrices. Then the change in the output of the attention mapping is bounded by

$$\|A(Q + \delta Q, K + \delta K, V + \delta V) - A(Q, K, V)\| \le C(\|\delta Q\| + \|\delta K\| + \|\delta V\|)$$

for some constant $C$. This property ensures that the attention mechanism is robust to small changes in the input, making it resilient to noise and variations in the data.

5. Non-negativity and Probability Interpretation: The attention weights $\alpha_{ij}$ are non-negative and sum to 1 for each query vector, i.e.:

$$\alpha_{ij} \geq 0 \quad \text{and} \quad \sum_{j=1}^{n} \alpha_{ij} = 1.$$

This property gives the attention mechanism a probabilistic interpretation, where the weights $\alpha_{ij}$ can be viewed as probabilities that determine the importance of each value vector in constructing the output vector. This probabilistic nature is crucial for tasks that require a soft selection of relevant information from the input, such as language modeling and translation.

## 1.4.2  The Geometry of High-Dimensional Spaces

As attention mechanisms and transformers often operate in high-dimensional spaces, understanding the geometric properties of these spaces is crucial for analyzing how these models function and why they are effective. The study of high-dimensional spaces introduces unique challenges and opportunities that are not present in lower dimensions, particularly concerning the phenomena known as the curse of dimensionality and the concentration of measure.

### Curse of Dimensionality

The curse of dimensionality refers to the various phenomena that arise when working in high-dimensional spaces that make intuitive concepts from low-dimensional spaces fail to generalize. These effects become particularly pronounced as the dimensionality $d$ of the space increases, leading to challenges in computation, data analysis, and learning.

One of the most significant manifestations of the curse of dimensionality is the rapid increase in the volume of the space as the dimension increases. Consider the unit hypercube in $\mathbb{R}^d$, defined as the set

$$C = \{\mathbf{x} \in \mathbb{R}^d \mid 0 \leq x_i \leq 1, \, i = 1, 2, \ldots, d\}.$$

The volume $V_d(C)$ of this hypercube is simply $1^d = 1$, regardless of the dimension $d$. However, consider a hypersphere inscribed within this hypercube with radius $r = \frac{1}{2}$. The volume of this hypersphere, given by $V_d(S)$, is

$$V_d(S) = \frac{\pi^{d/2}}{\Gamma\left(\frac{d}{2} + 1\right)} r^d = \frac{\pi^{d/2}}{\Gamma\left(\frac{d}{2} + 1\right)} \left(\frac{1}{2}\right)^d,$$

where $\Gamma$ is the Gamma function, which generalizes the factorial function. As $d$ increases, the volume of the hypersphere decreases exponentially compared to the volume of the hypercube. This result implies that in high-dimensional spaces, most of the volume of the hypercube is concentrated near its corners, rather than being uniformly distributed throughout the interior.

This phenomenon has profound implications for machine learning and data analysis:

1. Sparse Data: In high-dimensional spaces, data points tend to be sparsely distributed, making it difficult to find close neighbors or meaningful clusters. As a result, algorithms that rely on proximity or density, such as k-nearest neighbors or clustering algorithms, may perform poorly without appropriate dimensionality reduction techniques.

2. Distance Metrics: The relative distances between points in high-dimensional spaces become less informative. For instance, in high dimensions, the difference between the maximum and minimum pairwise distances between random points becomes negligible, leading to the "distance concentration" effect. Mathematically, if $\mathbf{x}_1$, $\mathbf{x}_2$ are random points in $\mathbb{R}^d$ with a fixed distribution, the ratio of the distance between these points to the mean distance approaches 1 as $d$ increases:

$$\lim_{d \to \infty} \frac{\|\mathbf{x}_1 - \mathbf{x}_2\|}{\mathbb{E}[\|\mathbf{x}_1 - \mathbf{x}_2\|]} = 1.$$

This effect makes it challenging to distinguish between points based on distance alone, requiring the use of alternative measures or embeddings to capture meaningful relationships.

3. Dimensionality Reduction: To combat the curse of dimensionality, various dimensionality reduction techniques, such as Principal Component Analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE), are often employed. These methods seek to project high-dimensional data onto a lower-dimensional subspace where the data's essential structure and relationships are preserved.

In the context of transformers and attention mechanisms, the curse of dimensionality underscores the importance of carefully designing the model's architecture to manage the complexity that arises from high-dimensional input spaces. For instance, the self-attention mechanism inherently mitigates some of these challenges by focusing on specific subsets of the input data, thereby reducing the effective dimensionality that the model needs to process.

**Concentration of Measure**

The concentration of measure phenomenon is another critical aspect of high-dimensional geometry. It refers to the fact that in high-dimensional spaces, most of the mass of a probability distribution tends to concentrate in a small region near

the mean or median of the distribution. This effect is closely related to the isoperimetric inequalities in mathematics, which describe how the "surface area" of a set relates to its "volume" in high dimensions.

To formalize this concept, consider a high-dimensional sphere $S^{d-1}$ in $\mathbb{R}^d$ with radius $r$. As the dimension $d$ increases, the volume of the sphere's equatorial region (i.e., the region within a small distance $\epsilon$ from the equator) becomes overwhelmingly large compared to the volume near the poles. More generally, for a Lipschitz function $f : S^{d-1} \to \mathbb{R}$, the concentration of measure theorem states that for any $\epsilon > 0$:

$$\mathbb{P}(|f(x) - \mathbb{E}[f]| \geq \epsilon) \leq 2 \exp\left(-\frac{C\epsilon^2}{d}\right),$$

where $C$ is a constant that depends on the Lipschitz constant of $f$. This result implies that as the dimensionality $d$ increases, the probability that a function deviates significantly from its expected value decreases exponentially.

The concentration of measure has profound implications for understanding the behavior of functions in high-dimensional spaces:

1. Robustness of Features: In high-dimensional spaces, most points are close to the "average" behavior of a function, meaning that deviations are rare. This property can be advantageous in machine learning, where models trained on high-dimensional data may exhibit robust performance even under slight perturbations of the input.

2. Random Projections: The concentration of measure phenomenon justifies the use of random projections as a dimensionality reduction technique. The Johnson–Lindenstrauss lemma, a direct consequence of measure concentration, states that a small set of points in a high-dimensional space can be embedded into a lower-dimensional space with almost no distortion of pairwise distances. Mathematically, for any set of $n$ points in $\mathbb{R}^d$ and for any $0 < \epsilon < 1$, there exists a linear map $f : \mathbb{R}^d \to \mathbb{R}^k$ with $k = O(\frac{\log n}{\epsilon^2})$ such that

$$(1 - \epsilon)\|\mathbf{x}_i - \mathbf{x}_j\|^2 \leq \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|^2 \leq (1 + \epsilon)\|\mathbf{x}_i - \mathbf{x}_j\|^2$$

for all points $\mathbf{x}_i, \mathbf{x}_j$ in the set. This result is crucial in reducing the dimensionality of data while preserving its essential geometric properties.

3. Implications for Model Design: In the design of transformers, the concentration of measure suggests that high-dimensional feature spaces can be effectively managed by focusing on the most significant components of the data. Attention mechanisms naturally leverage this principle by assigning higher weights to the most relevant parts of the input, thereby concentrating the "measure" of attention on the critical features.

### *1.4.3  Applications in Transformer Architectures*

The mathematical properties of attention mechanisms, especially when viewed through the lens of geometry, symmetry, and metric spaces, offer profound insights into the expressivity and robustness of transformer models. By understanding attention as a mapping within a high-dimensional space, we can explore its implications for model expressivity and its role as a metric-preserving map, which are crucial for the design and analysis of transformer architectures.

**Implications for Model Expressivity**

Expressivity in the context of transformer architectures refers to the model's ability to capture and represent a wide range of functions or patterns within the input data. The self-attention mechanism, which lies at the core of transformers, significantly enhances the model's expressivity by enabling it to dynamically focus on different parts of the input sequence, thereby capturing complex dependencies and interactions.

Mathematically, the expressivity of a model can be related to its capacity to approximate functions within a certain function space. Let $\mathcal{H}$ denote a Hilbert space of functions defined on the input space $X$, with an associated inner product $\langle \cdot, \cdot \rangle$. The self-attention mechanism can be viewed as an operator $\mathcal{A} : \mathcal{H} \to \mathcal{H}$, where for any function $f \in \mathcal{H}$, the operator $\mathcal{A}$ produces a new function $g = \mathcal{A}(f)$ that represents a weighted combination of $f$ over the input sequence. This operator is defined by

$$g(x_i) = \sum_{j=1}^{n} \alpha_{ij} f(x_j),$$

where $\alpha_{ij}$ are the attention weights, computed as

$$\alpha_{ij} = \frac{\exp(\langle \mathbf{q}_i, \mathbf{k}_j \rangle / \sqrt{d_k})}{\sum_{k=1}^{n} \exp(\langle \mathbf{q}_i, \mathbf{k}_k \rangle / \sqrt{d_k})}.$$

The expressivity of the transformer model is closely related to the ability of $\mathcal{A}$ to approximate any target function $g$ within $\mathcal{H}$. This approximation capability is influenced by several factors:

1. Diversity of Attention Heads: In multi-head attention, multiple attention operators $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_h$ are applied in parallel, each with different parameterizations. The combination of these attention heads increases the model's expressivity by enabling it to capture different aspects of the input data simultaneously. Mathematically, the combined output can be expressed as

$$g(x_i) = \bigoplus (\mathcal{A}_1(f)(x_i), \mathcal{A}_2(f)(x_i), \dots, \mathcal{A}_h(f)(x_i)) W_o,$$

where $W_o$ is a learned weight matrix. The presence of multiple attention heads allows the model to approximate more complex functions by effectively increasing the dimensionality of the function space $\mathcal{H}$ that the model can represent.

2. Non-linearity and Depth: The inclusion of non-linear activation functions, such as ReLU, and the stacking of multiple layers of self-attention further enhances the model's expressivity. Each layer of the transformer applies a sequence of linear and non-linear transformations, which, according to the universal approximation theorem, can approximate any continuous function on a compact domain to arbitrary accuracy, provided the network is sufficiently deep and wide.

3. Attention Weight Diversity: The diversity of attention weights $\alpha_{ij}$ across different queries $x_i$ allows the model to focus on different parts of the input sequence for different contexts. This adaptability is crucial for capturing long-range dependencies and interactions within the input, which are essential for tasks such as translation, summarization, and question answering.

The mathematical implications of these factors suggest that transformers, through their attention mechanisms, possess a high degree of expressivity. This allows them to model complex relationships within data that would be challenging for models lacking such flexible and adaptive attention mechanisms.

## Attention as a Metric-Preserving Map

Another critical aspect of attention mechanisms in transformer architectures is their role as metric-preserving maps. In a metric space $(X, d)$, a map $\phi : X \to Y$ is said to preserve the metric if it approximately maintains the distances between points, i.e., for all $x_1, x_2 \in X$, we have

$$d_Y(\phi(x_1), \phi(x_2)) \approx d_X(x_1, x_2).$$

In the context of transformers, we can analyze the attention mechanism as a mapping from the space of input embeddings $\mathbb{R}^{n \times d}$ to an output space $\mathbb{R}^{n \times d}$. The goal is to understand whether and how this mapping preserves the distances (and hence the relationships) between the input vectors.

Consider the attention mapping $A(Q, K, V)$ discussed earlier. We want to explore whether this mapping preserves the metric structure of the input space. To do this, we can analyze the change in distance between two input vectors $\mathbf{x}_1$ and $\mathbf{x}_2$ under the attention transformation. Let $\mathbf{z}_1$ and $\mathbf{z}_2$ be the corresponding output vectors after applying the attention mechanism:

$$\mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{v}_j \quad \text{for } i = 1, 2.$$

The distance between the output vectors $\mathbf{z}_1$ and $\mathbf{z}_2$ can be analyzed using a suitable norm, such as the Euclidean norm:

$$\|\mathbf{z}_1 - \mathbf{z}_2\|_2 = \left\|\sum_{j=1}^{n}(\alpha_{1j}\mathbf{v}_j - \alpha_{2j}\mathbf{v}_j)\right\|_2.$$

If the attention weights $\alpha_{ij}$ do not vary significantly between $\mathbf{x}_1$ and $\mathbf{x}_2$, then the output distance $\|\mathbf{z}_1 - \mathbf{z}_2\|_2$ will be close to the input distance $\|\mathbf{x}_1 - \mathbf{x}_2\|_2$, implying that the attention mechanism approximately preserves the metric. This property is essential for maintaining the relative relationships between input vectors, which is critical for tasks such as language modeling, where preserving the semantic relationships between words is crucial.

Moreover, when the attention mechanism operates as a metric-preserving map, it ensures that the transformed data retains its structural integrity, enabling the model to effectively leverage the geometric properties of the input space. This preservation of structure is particularly important in high-dimensional spaces, where the relationships between data points can be complex and sensitive to perturbations. In cases where the attention mechanism introduces significant non-linearity, the preservation of the metric may not be exact, but the mechanism can still capture meaningful relationships by focusing on the most relevant dimensions of the input space. This selective focus can be understood as a form of adaptive metric preservation, where the attention mechanism dynamically adjusts the metric to emphasize the most informative features of the data.

## 1.5   Tensor Algebra and Notation

Tensor algebra is a powerful mathematical framework that extends linear algebra to higher dimensions, enabling the manipulation of multi-dimensional arrays called tensors. Tensors generalize vectors and matrices and are central to many areas of mathematics, physics, and machine learning. In the context of transformers, tensors are used extensively to represent and process high-dimensional data, such as word embeddings and attention scores.

### 1.5.1   Introduction to Tensors

A tensor is a multi-dimensional array of numerical values that generalizes the concepts of scalars, vectors, and matrices. Formally, a tensor of order $k$ (also known as a $k$-tensor) over a vector space $V$ is an element of the tensor product of $k$ copies of $V$. In other words, a tensor is a multi-linear map that takes $k$ vectors from $V$ and returns a scalar. Let $V$ be a finite-dimensional vector space over a field $\mathbb{F}$ with a basis $\{\mathbf{e}_1, \mathbf{e}_2, \ldots, \mathbf{e}_n\}$. A tensor of order $k$ on $V$ is an element of the space $V^{\otimes k} = V \otimes V \otimes \cdots \otimes V$ (with $k$ factors). If $T$ is a tensor of order $k$, it can be expressed as

$$T = \sum_{i_1,i_2,\dots,i_k} T_{i_1 i_2 \dots i_k}\, \mathbf{e}_{i_1} \otimes \mathbf{e}_{i_2} \otimes \cdots \otimes \mathbf{e}_{i_k},$$

where $T_{i_1 i_2 \dots i_k}$ are the components of the tensor with respect to the chosen basis, and $\otimes$ denotes the tensor product. The components $T_{i_1 i_2 \dots i_k}$ form a $k$-dimensional array (or hypermatrix), where each index $i_j$ ranges from 1 to $n$. The simplest tensors are

1.  Scalars (order 0 tensors): A scalar is a single numerical value and can be thought of as a tensor with no indices, i.e., a tensor of order 0.
2.  Vectors (order 1 tensors): A vector is a one-dimensional array of numbers, represented as $\mathbf{v} = \sum_i v_i \mathbf{e}_i$, where $v_i$ are the components of the vector.
3.  Matrices (order 2 tensors): A matrix is a two-dimensional array of numbers, represented as $M = \sum_{i,j} M_{ij} \mathbf{e}_i \otimes \mathbf{e}_j$, where $M_{ij}$ are the matrix elements.

Higher-order tensors can be visualized as multi-dimensional arrays, but their manipulation requires careful attention to the indices and the rules of tensor algebra.

### Tensor Products and Contractions

The tensor product is a fundamental operation in tensor algebra that allows the construction of higher-order tensors from lower-order ones. Given two tensors $A \in V^{\otimes k}$ and $B \in W^{\otimes l}$, their tensor product $A \otimes B$ is a tensor of order $k + l$ in the space $V^{\otimes k} \otimes W^{\otimes l}$. If $A$ and $B$ have components $A_{i_1 i_2 \dots i_k}$ and $B_{j_1 j_2 \dots j_l}$, respectively, then the components of the tensor product $A \otimes B$ are given by

$$(A \otimes B)_{i_1 i_2 \dots i_k j_1 j_2 \dots j_l} = A_{i_1 i_2 \dots i_k} B_{j_1 j_2 \dots j_l}.$$

The tensor product operation is bilinear, meaning it satisfies the properties:

$$(A + B) \otimes C = A \otimes C + B \otimes C$$

$$A \otimes (B + C) = A \otimes B + A \otimes C$$

$$\alpha(A \otimes B) = (\alpha A) \otimes B = A \otimes (\alpha B)$$

for any scalar $\alpha \in \mathbb{F}$.

Tensor contraction is another essential operation in tensor algebra, analogous to matrix multiplication. Contraction reduces the order of a tensor by summing over one or more pairs of indices. Given a tensor $T \in V^{\otimes k}$ with components $T_{i_1 i_2 \dots i_k}$, a contraction over the indices $i_j$ and $i_l$ (where $j \neq l$) results in a tensor of order $k - 2$ with components:

$$S_{i_1 \dots i_{j-1} i_{j+1} \dots i_{l-1} i_{l+1} \dots i_k} = \sum_{i_j} T_{i_1 i_2 \dots i_j \dots i_l \dots i_k}.$$

Contraction is a linear operation and is particularly useful in reducing the complexity of tensors, transforming them into lower-dimensional objects while preserving specific relationships between the indices. For example, consider two tensors $A \in V \otimes W$ and $B \in W^* \otimes U$, where $W^*$ is the dual space of $W$. The contraction of $A$ and $B$ over the common index space $W$ is a tensor $C \in V \otimes U$ with components:

$$C_{iu} = \sum_j A_{ij} B_{ju}.$$

This operation is analogous to matrix multiplication, where the product of matrices corresponds to the contraction of their associated tensors.

In the context of self-attention mechanisms in transformers, tensor products and contractions are used to compute the attention scores, where the query, key, and value tensors interact through tensor operations. These operations allow the model to dynamically adjust the focus on different parts of the input sequence, capturing complex dependencies and patterns.

## 1.5.2  Algebraic Structures in Transformers

The algebraic structures underlying transformers are crucial for understanding how these models manipulate data and achieve their remarkable performance. Among these structures, matrix multiplication, Kronecker products, and tensor factorization play key roles in the architecture and functioning of transformers, particularly in the computation of attention scores and the efficient representation of high-dimensional data. These algebraic tools allow transformers to capture complex dependencies, exploit symmetries, and manage the computational complexity associated with large-scale models.

### The Role of Matrix Multiplication

Matrix multiplication is a fundamental operation in linear algebra and serves as the backbone of many computations within transformers. In the context of transformers, matrix multiplication is used extensively in the linear transformations applied to input data, in the computation of attention scores, and in the propagation of information across layers.

Consider the basic operation in a transformer layer where an input matrix $X \in \mathbb{R}^{n \times d}$, representing a sequence of $n$ input vectors of dimension $d$, is transformed by a weight matrix $W \in \mathbb{R}^{d \times d'}$ to produce an output matrix $Y \in \mathbb{R}^{n \times d'}$:

$$Y = XW.$$

Here, $Y$ is the result of applying the linear transformation defined by $W$ to each vector in the sequence $X$. This operation is crucial for adjusting the dimensionality of the data, allowing the model to project input vectors into different spaces where specific patterns or features may be more easily captured.

In the self-attention mechanism, matrix multiplication plays a central role in the computation of attention scores. Given query $Q$, key $K$, and value $V$ matrices, the attention output is computed as

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right) V.$$

Here, $QK^\top$ is a matrix multiplication operation that computes the dot-product attention scores between queries and keys. The result is an $n \times n$ matrix where each entry represents the similarity between a query vector and a key vector. This matrix is then normalized using the softmax function, and the result is used to weight the value vectors through another matrix multiplication.

Matrix multiplication in transformers also facilitates the blending of information across different parts of the input sequence. By multiplying matrices, the model can combine contributions from various elements of the input, effectively allowing the model to "attend" to multiple parts of the sequence simultaneously. This ability to aggregate information from different sources is essential for tasks that require understanding long-range dependencies, such as translation or summarization.

Moreover, matrix multiplication in transformers is often used in conjunction with other algebraic operations, such as the addition of bias terms or the application of non-linear activation functions, to create more complex and expressive transformations. These operations are key to the model's capacity to learn and represent intricate patterns in the data.

### Kronecker Products and Factorization

The Kronecker product is a powerful tool in tensor algebra that extends the concept of matrix multiplication to higher-dimensional structures. It is particularly useful in the context of transformers for creating structured and scalable representations of data and for factorizing large matrices into more manageable components.

Given two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times q}$, their Kronecker product $A \otimes B$ is defined as a block matrix of dimensions $mp \times nq$ where each block is a scaled version of $B$ by the corresponding entry in $A$:

$$A \otimes B = \begin{pmatrix} a_{11}B & a_{12}B & \dots & a_{1n}B \\ a_{21}B & a_{22}B & \dots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \dots & a_{mn}B \end{pmatrix}.$$

The Kronecker product inherits properties from matrix multiplication, such as distributivity over matrix addition and associativity with respect to matrix multiplication. However, the Kronecker product also introduces new capabilities, such as the ability to model interactions between different dimensions or to create high-dimensional tensor representations from lower-dimensional ones.

In transformers, Kronecker products can be used to efficiently represent and manipulate large matrices or tensors that arise in the model. For example, consider a situation where the model needs to learn interactions between two sets of features, one represented by matrix $A$ and the other by matrix $B$. Instead of explicitly constructing a large matrix to capture all pairwise interactions, the Kronecker product $A \otimes B$ provides a compact and structured way to encode these interactions.

Factorization techniques based on the Kronecker product are also valuable in reducing the computational complexity of transformer models. Matrix factorization aims to decompose a large matrix $M$ into the product of smaller matrices, thereby reducing the number of parameters and the computational cost of matrix operations. A common factorization approach involves expressing $M$ as

$$M \approx A \otimes B,$$

where $A$ and $B$ are smaller matrices. This decomposition allows the model to approximate large matrix operations using a series of smaller, more efficient operations. In the context of transformers, such factorization can be applied to weight matrices, attention score matrices, or other components of the model to improve scalability without sacrificing performance.

One specific application of Kronecker product-based factorization is in multihead attention, where the attention weights for different heads can be factorized to share information across heads while still allowing for individual specialization. This approach leads to a more efficient representation of the model's parameters and can help prevent overfitting by reducing the model's capacity.

### 1.5.3 Self-attention Mechanisms

The central theme is that of computing a weighted combination of input features, where the weights reflect the relevance of each feature to the others. Mathematically, let $X = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n\}$ be a sequence of input vectors, where each $\mathbf{x}_i \in \mathbb{R}^d$ represents a feature vector in a $d$-dimensional space. The self-attention mechanism transforms this input sequence into an output sequence $Z = \{\mathbf{z}_1, \mathbf{z}_2, \ldots, \mathbf{z}_n\}$, where each output vector $\mathbf{z}_i \in \mathbb{R}^d$ is a weighted sum of the input vectors.

To compute the output vectors $\mathbf{z}_i$, the self-attention mechanism uses three matrices: the query matrix $Q$, the key matrix $K$, and the value matrix $V$, each derived from the input matrix $X$. Specifically, the query, key, and value matrices are defined as

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V,$$

where $W_Q, W_K, W_V \in \mathbb{R}^{d \times d_k}$ are learned weight matrices, and $d_k$ is the dimensionality of the key vectors. The output vector $\mathbf{z}_i$ is then computed as

$$\mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{v}_j,$$

where $\alpha_{ij}$ are the attention weights, defined by

$$\alpha_{ij} = \frac{\exp(S_{ij})}{\sum_{k=1}^{n} \exp(S_{ik})}, \quad S_{ij} = \frac{\langle \mathbf{q}_i, \mathbf{k}_j \rangle}{\sqrt{d_k}}.$$

Here, $\mathbf{q}_i$ and $\mathbf{k}_j$ are the $i$th query vector and $j$th key vector, respectively, and $S_{ij}$ represents the similarity score between $\mathbf{q}_i$ and $\mathbf{k}_j$, scaled by the factor $\frac{1}{\sqrt{d_k}}$ to mitigate the effect of large dot products in high-dimensional spaces. The softmax function ensures that the attention weights $\alpha_{ij}$ sum to 1 for each query vector, providing a probabilistic interpretation of the weights.

The self-attention mechanism can be compactly represented in matrix form as

$$Z = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right) V.$$

This formulation highlights the core operations involved in self-attention: the computation of pairwise similarities between queries and keys (via dot products), the normalization of these similarities to obtain attention weights, and the weighted aggregation of value vectors to produce the final output.

## Multi-head Self-attention

The concept of multi-head self-attention extends the basic self-attention mechanism by allowing the model to focus on different aspects of the input sequence simultaneously. In multi-head self-attention, multiple attention mechanisms (or "heads") operate in parallel, each with its own set of learned weights. This enables the model to capture various types of dependencies and relationships within the input data.

Formally, let $h$ denote the number of attention heads. For each head $i$, the corresponding query, key, and value matrices are computed as

$$Q_i = X W_{Q_i}, \quad K_i = X W_{K_i}, \quad V_i = X W_{V_i},$$

where $W_{Q_i}, W_{K_i}, W_{V_i} \in \mathbb{R}^{d \times d_k}$ are the learned weight matrices for the $i$th head. The output of each head is then computed using the self-attention mechanism:

$$Z_i = \text{Softmax}\left(\frac{Q_i K_i^\top}{\sqrt{d_k}}\right) V_i.$$

The outputs of the $h$ attention heads are concatenated and linearly transformed to produce the final output:

$$Z = \bigoplus (Z_1, Z_2, \ldots, Z_h) W_O,$$

where $W_O \in \mathbb{R}^{hd_k \times d}$ is a learned weight matrix that projects the concatenated output back to the original dimensionality $d$.

The advantage of multi-head self-attention lies in its ability to capture diverse patterns in the data. Each head can attend to different parts of the input sequence or focus on different types of relationships, such as short-range versus long-range dependencies. This parallel processing of multiple attention heads enhances the model's expressivity and allows it to better capture the complexity of the input data.

## 1.6    Matrix Calculus in Self-attention

Matrix calculus is a vital tool in the optimization and analysis of neural networks, including transformer models that rely heavily on self-attention mechanisms. The ability to compute gradients, Jacobians, and Hessians efficiently is crucial for training these models, as these quantities guide the optimization process, ensuring that the model learns to represent the underlying data patterns effectively.

### *1.6.1    Differentiation of Matrix Functions*

Matrix calculus extends the principles of calculus to functions that take matrices as inputs and produce matrices as outputs. In the context of self-attention mechanisms, differentiating matrix functions is essential for backpropagation, the algorithm used to compute gradients and update model parameters during training.

Consider a matrix-valued function $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{p \times q}$, which maps an input matrix $X \in \mathbb{R}^{m \times n}$ to an output matrix $Y = f(X) \in \mathbb{R}^{p \times q}$. The derivative of $f$ with respect to $X$ is a fourth-order tensor that describes how each element of the output matrix $Y$ changes with respect to each element of the input matrix $X$. However, in practice, matrix calculus often focuses on specific cases, such as the gradient, Jacobian, or Hessian, which provide more manageable representations of these derivatives.

**Gradient Computation Techniques**

The gradient of a scalar-valued function $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ with respect to a matrix $X$ is a matrix of partial derivatives, often denoted by $\nabla_X f$ or $\frac{\partial f}{\partial X}$. Each entry of the

gradient matrix $\nabla_X f$ is given by

$$(\nabla_X f)_{ij} = \frac{\partial f}{\partial X_{ij}}.$$

In self-attention mechanisms, gradients are required to update the weight matrices $W_Q$, $W_K$, $W_V$, which are involved in the transformations of the query, key, and value matrices. To compute these gradients, one typically relies on matrix calculus rules that extend scalar calculus to matrix operations.

Consider a common matrix function $f(X) = \text{tr}(AX)$, where $A \in \mathbb{R}^{m \times n}$ is a fixed matrix, and $\text{tr}(\cdot)$ denotes the trace of a matrix. The gradient of $f$ with respect to $X$ is given by

$$\nabla_X f(X) = A^\top.$$

This result follows from the linearity of the trace operator and the fact that the trace of a product of matrices is invariant under cyclic permutations, i.e., $\text{tr}(AB) = \text{tr}(BA)$. For a more complex function such as the Frobenius norm of a matrix, $f(X) = \|X\|_F^2 = \text{tr}(X^\top X)$, the gradient is computed as

$$\nabla_X \|X\|_F^2 = 2X.$$

These gradient computation techniques are used in the optimization of self-attention mechanisms to minimize the loss function during training. For instance, the loss function in a transformer model might involve the squared difference between the predicted output and the true output, and the gradient of this loss with respect to the model's parameters guides the weight updates.

### Jacobian and Hessian Matrices

The Jacobian matrix generalizes the gradient to vector-valued functions. For a vector-valued function $f : \mathbb{R}^{m \times n} \to \mathbb{R}^p$, the Jacobian matrix $J_f(X)$ is defined as

$$J_f(X) = \frac{\partial f(X)}{\partial X} \in \mathbb{R}^{p \times (mn)},$$

where each row of $J_f(X)$ corresponds to the gradient of one component of the vector-valued function $f$ with respect to the elements of $X$.

In the context of self-attention, consider the softmax function used to compute attention weights. Let $S(X)$ be the similarity matrix resulting from the dot product of queries and keys, normalized by the softmax function:

$$\alpha_{ij} = \frac{\exp(S_{ij})}{\sum_{k=1}^{n} \exp(S_{ik})}.$$

The Jacobian of the softmax function softmax$(S)$ is important for understanding how small changes in the similarity scores $S$ affect the attention weights $\alpha$. The Jacobian matrix $J_{\text{softmax}}(S)$ is given by

$$(J_{\text{softmax}}(S))_{ij,kl} = \alpha_{ij} \left( \delta_{jk} - \alpha_{ik} \right),$$

where $\delta_{jk}$ is the Kronecker delta. This Jacobian matrix plays a crucial role in the backpropagation algorithm, determining how the gradients flow through the softmax layer.

The Hessian matrix further generalizes the concept of the second derivative to functions of matrices. For a scalar-valued function $f : \mathbb{R}^{m \times n} \to \mathbb{R}$, the Hessian matrix $H_f(X)$ is a block matrix where each block corresponds to the second-order partial derivatives with respect to the elements of $X$:

$$H_f(X) = \frac{\partial^2 f(X)}{\partial X \partial X^\top} \in \mathbb{R}^{(mn) \times (mn)}.$$

The Hessian matrix provides information about the curvature of the function $f$, which is essential for second-order optimization methods like Newton's method. In self-attention mechanisms, the Hessian can be used to analyze the stability and convergence of the training process, as well as to detect saddle points and local minima.

For example, if the loss function $L$ in a transformer model is highly curved (i.e., has a large Hessian), the gradient descent steps might need to be adjusted to ensure stable convergence. The Hessian matrix can also help in identifying directions in the parameter space that require more precise adjustments, leading to more efficient training.

## 1.6.2  Optimization and Gradient Flow

Optimization is a critical component of training neural networks, including transformers with self-attention mechanisms. The process involves adjusting model parameters to minimize a loss function, which measures the difference between the model's predictions and the actual data. Gradient flow, convergence analysis, and learning rate schedules are central to understanding how effectively and efficiently a model learns during training.

### Analysis of Convergence

Convergence analysis in the context of gradient-based optimization refers to the study of how and when an iterative optimization algorithm approaches a local or global

minimum of the loss function. For a transformer model utilizing self-attention mechanisms, the goal is to ensure that the parameters converge to values that minimize the loss function, leading to optimal model performance.

Let $L(\theta)$ be the loss function, where $\theta \in \mathbb{R}^d$ represents the vector of model parameters. The gradient descent algorithm updates the parameters iteratively as follows:

$$\theta_{t+1} = \theta_t - \eta \nabla L(\theta_t),$$

where $\eta > 0$ is the learning rate, and $\nabla L(\theta_t)$ is the gradient of the loss function with respect to $\theta$ at iteration $t$.

For convergence of gradient descent to a local minimum, certain conditions on the loss function $L(\theta)$ and the learning rate $\eta$ must be satisfied:

1. Smoothness: The loss function $L(\theta)$ is typically assumed to be smooth, meaning that it has Lipschitz-continuous gradients. Formally, there exists a constant $L > 0$ such that

$$\|\nabla L(\theta_1) - \nabla L(\theta_2)\| \leq L \|\theta_1 - \theta_2\| \quad \forall \theta_1, \theta_2 \in \mathbb{R}^d.$$

This smoothness condition ensures that the gradient does not change too rapidly, which is crucial for the stability of the gradient descent algorithm.

2. Convexity: If the loss function $L(\theta)$ is convex, meaning that for all $\theta_1, \theta_2 \in \mathbb{R}^d$ and $\lambda \in [0, L]$:

$$1(\lambda \theta_1 + (1 - \lambda)\theta_2) \leq \lambda L(\theta_1) + (1 - \lambda)L(\theta_2),$$

then gradient descent is guaranteed to converge to a global minimum. However, in the case of deep neural networks, the loss function is typically non-convex, and the convergence analysis focuses on local minima or saddle points.

3. Learning Rate: The learning rate $\eta$ plays a critical role in convergence. If $\eta$ is too large, the algorithm may overshoot the minimum and diverge; if it is too small, the convergence will be slow. For convergence, $\eta$ must satisfy the condition:

$$0 < \eta < \frac{2}{L},$$

where $L$ is the Lipschitz constant of the gradient. Under this condition, gradient descent can be shown to converge to a stationary point, where $\nabla L(\theta) = 0$.

The rate at which gradient descent converges depends on the properties of the loss function. For a convex and smooth loss function, the convergence rate is linear, meaning that the distance to the minimum decreases geometrically with each iteration:

$$L(\theta_t) - L(\theta^*) \leq (1 - \eta\mu)^t (L(\theta_0) - L(\theta^*)),$$

where $\theta^*$ is the optimal parameter vector and $\mu > 0$ is a parameter related to the strong convexity of $L(\theta)$. In practice, due to the non-convexity of neural network

loss functions, the convergence rate can be slower, and additional techniques such as momentum or adaptive learning rates may be used to accelerate convergence.

**Learning Rate Schedules**

The learning rate $\eta$ is a hyperparameter that significantly impacts the efficiency and success of the optimization process. A fixed learning rate may not be optimal throughout the entire training process, leading to the use of learning rate schedules that adjust $\eta$ dynamically during training.

Types of Learning Rate Schedules:

1. In step decay, the learning rate is reduced by a constant factor $\gamma$ after a fixed number of epochs. Formally, if $\eta_0$ is the initial learning rate, the learning rate at epoch $t$ is given by

$$\eta_t = \eta_0 \cdot \gamma^{\lfloor \frac{t}{T} \rfloor},$$

where $T$ is the number of epochs after which the learning rate is decayed and $\lfloor \cdot \rfloor$ denotes the floor function. This schedule helps in making large steps initially to escape from local minima or saddle points, followed by smaller steps as the algorithm approaches a minimum.

2. Exponential decay continuously decreases the learning rate according to an exponential function:

$$\eta_t = \eta_0 \cdot \exp(-\lambda t),$$

where $\lambda > 0$ is the decay rate. This schedule provides a smooth reduction in the learning rate, which can be advantageous when fine-tuning the model near a minimum.

3. Polynomial decay reduces the learning rate according to a polynomial function of the epoch number:

$$\eta_t = \eta_0 \cdot \left(1 - \frac{t}{T_{\max}}\right)^p,$$

where $T_{\max}$ is the total number of training epochs and $p$ is the power of the polynomial. This schedule allows for a gradual reduction in the learning rate, which can be particularly useful in the later stages of training.

4. Cosine annealing gradually reduces the learning rate following the cosine function:

$$\eta_t = \eta_{\min} + \frac{1}{2}(\eta_{\max} - \eta_{\min})\left(1 + \cos\left(\frac{t\pi}{T_{\max}}\right)\right),$$

where $\eta_{\min}$ and $\eta_{\max}$ are the minimum and maximum learning rates, respectively. Cosine annealing allows for a non-monotonic reduction in the learning rate, potentially allowing the optimization process to escape from shallow local minima.

5. Warm restarts involve periodically resetting the learning rate to a higher value after certain intervals, allowing the optimization process to explore new regions of

the parameter space. This approach can be combined with any of the above schedules, leading to a cyclical learning rate schedule that alternates between exploration and exploitation phases.

Implications: The choice of learning rate schedule can have significant implications for the convergence and final performance of the model. A well-chosen schedule can accelerate convergence, prevent the optimization from getting stuck in poor local minima, and lead to better generalization. Conversely, an inappropriate schedule may cause the optimization to converge too slowly or even diverge. The mathematical analysis of learning rate schedules often involves examining the behavior of the gradient flow under different schedules, using tools from differential equations and dynamical systems. For instance, the convergence of gradient descent under a decaying learning rate schedule can be analyzed using Lyapunov functions, which provide a measure of the stability of the optimization process.

## 1.7　Positional Encodings: A Mathematical Perspective

Positional encodings are essential for certain tasks in transformer models because they provide a way to incorporate the order of elements in a sequence, which is crucial for tasks like language modeling and translation. Unlike recurrent neural networks, transformers process input sequences in parallel, and thus require a method to inject sequential information into the model. This is achieved through positional encodings, which are often formulated using sinusoidal functions, providing a natural connection to Fourier analysis.

### 1.7.1　Fourier Analysis of Positional Encodings

Fourier analysis is a powerful mathematical tool that decomposes functions into their constituent frequencies. It is particularly useful for analyzing periodic functions, which are central to the construction of positional encodings in transformers. The use of sinusoidal functions in positional encodings can be understood through the lens of Fourier series and transforms, which allow us to express a function as a sum of sines and cosines, capturing both local and global features of the sequence.

**Fourier Series and Transforms**

The Fourier series is a way to represent a periodic function $f(x)$ defined on the interval $[0, 2\pi]$ (or equivalently on any interval $[a, a + 2\pi]$) as an infinite sum of sines and cosines. Specifically, if $f(x)$ is a periodic function with period $2\pi$, its Fourier series is given by

$$f(x) = \frac{a_0}{2} + \sum_{n=1}^{\infty} (a_n \cos(nx) + b_n \sin(nx)),$$

where the coefficients $a_n$ and $b_n$ are determined by

$$a_n = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos(nx)\, dx, \quad b_n = \frac{1}{\pi} \int_0^{2\pi} f(x) \sin(nx)\, dx.$$

The Fourier series expresses the function as a superposition of harmonics, where each term corresponds to a specific frequency component of the original function. For functions that are not inherently periodic, the Fourier transform generalizes this idea to express the function in terms of continuous frequencies. The Fourier transform of a function $f(x)$ is defined as

$$\hat{f}(\xi) = \int_{-\infty}^{\infty} f(x)e^{-2\pi i \xi x}\, dx,$$

where $\hat{f}(\xi)$ represents the amplitude of the frequency component $\xi$. The inverse Fourier transform allows us to reconstruct the original function from its frequency components:

$$f(x) = \int_{-\infty}^{\infty} \hat{f}(\xi)e^{2\pi i \xi x}\, d\xi.$$

In the context of positional encodings, the use of Fourier series or transforms is motivated by the desire to represent the position of each element in a sequence as a combination of sinusoidal functions. These functions naturally encode positional information in a way that is smooth and periodic, capturing the relative positions of elements in the sequence.

**Periodic Functions and Signal Processing**

The sinusoidal functions used in positional encodings are periodic, meaning they repeat their values at regular intervals. This periodicity is particularly useful in representing the positions within a sequence, as it allows the model to capture both local and global positional relationships through different frequency components. In signal processing, periodic functions like sines and cosines are fundamental because they serve as the basic building blocks for more complex signals. By decomposing a signal into its constituent frequencies using Fourier analysis, we can analyze and manipulate the signal in an intuitive way.

In transformers, the positional encoding for a given position $p$ and dimension $i$ is typically defined as

$$PE(p, 2i) = \sin\left(\frac{p}{10000^{2i/d}}\right), \quad PE(p, 2i+1) = \cos\left(\frac{p}{10000^{2i/d}}\right),$$

where $d$ is the dimensionality of the model. These encodings are designed so that each dimension of the positional encoding corresponds to a different frequency component, with higher dimensions encoding higher frequencies.

The periodic nature of these functions ensures that positional encodings are smooth and continuous, allowing the model to generalize well to unseen positions. Moreover, the combination of sines and cosines at different frequencies enables the model to capture a wide range of positional relationships, from short-range dependencies to long-range interactions.

Example: Consider a simple sequence of positions $p = 0, 1, 2, \ldots, N - 1$. The positional encoding for the position $p$ can be represented as a vector:

$$
\text{PE}(p) = \left[ \sin\left( \frac{p}{10000^{0/d}} \right), \cos\left( \frac{p}{10000^{0/d}} \right), \sin\left( \frac{p}{10000^{2/d}} \right), \cos\left( \frac{p}{10000^{2/d}} \right), \ldots \right].
$$

This encoding scheme ensures that each position is mapped to a unique point in the high-dimensional space, with the difference between positions being captured by the differences in their sinusoidal components. The periodicity of these functions means that the encodings will eventually repeat, which can be advantageous for tasks that involve cyclic or repetitive patterns.

The use of Fourier analysis in positional encodings highlights the deep connection between signal processing and neural networks, where the goal is to encode information in a way that preserves important features while allowing for efficient computation and generalization. By leveraging the mathematical properties of Fourier series and transforms, transformers can effectively represent the order of elements in a sequence, enabling them to perform well on a wide range of tasks that require understanding of sequential data.

### 1.7.2  Lie Groups and Lie Algebras

Lie groups and Lie algebras [10, 19, 21, 35, 59] form a deep mathematical framework for understanding continuous symmetries, which are prevalent in many areas of mathematics and physics. Their application to transformers, particularly in understanding symmetries in data and model architecture, provides a foundation for analyzing and designing these models. Lie groups are mathematical structures that combine the properties of groups (algebraic operations) with the properties of smooth manifolds (geometric structures). They are used to describe continuous symmetries, such as rotations, translations, and scalings, which are crucial in many areas of mathematics and physics. The corresponding Lie algebras provide a linearized approximation of these symmetries, making them powerful tools for analysis [24, 31].

**Basic Definitions and Properties**

A Lie group is a group $G$ that is also a smooth manifold, meaning that the group operations (multiplication and inversion) are smooth (differentiable) maps. Formally, a Lie group $G$ is a set with two operations:

1. A binary operation (group multiplication) $\cdot : G \times G \to G$, which is smooth.
2. An inversion operation $\text{inv} : G \to G, g \mapsto g^{-1}$, which is also smooth.

   Example: The set of $n \times n$ invertible matrices $\text{GL}(n, \mathbb{R})$, known as the general linear group, is a Lie group. The group operations are matrix multiplication and matrix inversion, both of which are smooth maps on the manifold of $n \times n$ matrices.

   The Lie algebra associated with a Lie group $G$ is a vector space that captures the infinitesimal symmetries of the group. It can be thought of as the tangent space to the Lie group at the identity element. The Lie algebra $\mathfrak{g}$ of a Lie group $G$ is defined as the set of all tangent vectors at the identity element of $G$, with a bilinear operation called the Lie bracket:

$$\mathfrak{g} = T_e G, \quad [X, Y] = \frac{\partial}{\partial t}\bigg|_{t=0} \left(\exp(tX) \cdot \exp(tY) \cdot \exp(-tX) \cdot \exp(-tY)\right),$$

where $\exp : \mathfrak{g} \to G$ is the exponential map that connects the Lie algebra to the Lie group.

Properties:

1. Closure: For any two elements $X, Y \in \mathfrak{g}$, their Lie bracket $[X, Y]$ is also in $\mathfrak{g}$.
2. Bilinearity: The Lie bracket is bilinear, meaning that for any scalars $a, b \in \mathbb{R}$ and elements $X, Y, Z \in \mathfrak{g}$:

$$[aX + bY, Z] = a[X, Z] + b[Y, Z], \quad [Z, aX + bY] = a[Z, X] + b[Z, Y].$$

3. Jacobi Identity: The Lie bracket satisfies the Jacobi identity:

$$[X, [Y, Z]] + [Y, [Z, X]] + [Z, [X, Y]] = 0.$$

4. Antisymmetry: The Lie bracket is antisymmetric:

$$[X, Y] = -[Y, X].$$

These properties define the structure of the Lie algebra, making it a fundamental object in the study of continuous symmetries.

**Representations of Lie Groups**

A representation of a Lie group $G$ on a vector space $V$ is a homomorphism $\rho : G \to \mathrm{GL}(V)$ from $G$ to the general linear group $\mathrm{GL}(V)$ of all invertible linear transformations on $V$. This representation allows the elements of the Lie group to be expressed as matrices, enabling the study of group actions in a linear algebraic framework. The corresponding representation of a Lie algebra $\mathfrak{g}$ is a homomorphism $\phi : \mathfrak{g} \to \mathrm{End}(V)$, where $\mathrm{End}(V)$ is the space of all linear endomorphisms (linear maps from $V$ to itself). The representation $\phi$ satisfies

$$\phi([X, Y]) = \phi(X)\phi(Y) - \phi(Y)\phi(X)$$

for all $X, Y \in \mathfrak{g}$.

Example: Consider the Lie group SO(3), the group of rotations in three-dimensional space. The corresponding Lie algebra $\mathfrak{so}(3)$ consists of all skew-symmetric $3 \times 3$ matrices. A representation of SO(3) on $\mathbb{R}^3$ is given by the action of rotation matrices on vectors in $\mathbb{R}^3$. The corresponding representation of $\mathfrak{so}(3)$ is given by the action of the skew-symmetric matrices on vectors in $\mathbb{R}^3$.

The importance of representations in the context of Lie groups and Lie algebras lies in their ability to linearize the action of groups, making it easier to study their properties and symmetries. Representations allow us to express abstract group elements as concrete matrices, enabling the use of linear algebra tools to analyze group actions.

**Applications to Transformers**

In the context of transformers, Lie groups and their representations can be used to model and exploit symmetries in the data and the model architecture. For example:

1. Symmetries in Data: Many types of data, such as images, signals, and sequences, exhibit symmetries that can be modeled using Lie groups. For instance, rotational symmetry in images can be described by the Lie group SO(2), and translational symmetry can be described by the Euclidean group. By understanding and incorporating these symmetries into the model, transformers can be designed to be more robust and efficient. In natural language processing, certain grammatical structures exhibit symmetries that can be modeled using Lie groups. For example, the symmetry between active and passive voice in a sentence can be represented as a group action that transforms one structure into the other.

2. Symmetries in Model Architecture: The attention mechanism in transformers can be analyzed through the lens of Lie groups, particularly in understanding how different transformations of the input data (such as rotations, translations, or permutations) affect the output of the model. By designing the attention mechanism to be equivariant or invariant under certain group actions, the model can be made more robust to these transformations. The concept of Lie algebraic layers in neural networks, where the transformations applied in each layer correspond to elements of

a Lie algebra, provides a framework for designing architectures that respect certain symmetries and have desirable properties, such as equivariance or invariance.

Example: Consider a transformer model designed for image processing. If the input images exhibit rotational symmetry, the model can be made equivariant to rotations by incorporating a representation of the rotation group SO(2) into the attention mechanism. This can be achieved by designing the attention weights to be invariant under rotations or by using Lie algebraic layers that respect the rotational symmetry of the data.

### *1.7.3   Harmonic Analysis on Groups*

Harmonic analysis on groups extends the ideas of Fourier analysis to more general settings, where the underlying space is not the real line or Euclidean space but rather a group, often a Lie group. This extension is crucial for understanding how functions can be decomposed into basic, symmetric components, particularly in the context of rotational symmetries and other transformations. The tools of harmonic analysis, such as spherical harmonics and Wigner-D functions, provide a deep mathematical framework for studying symmetries in data, which can be exploited in transformer architectures to enhance model expressivity and robustness. Harmonic analysis on groups involves representing functions defined on a group as a sum or integral of basic functions that respect the group's structure. This approach generalizes classical Fourier analysis to functions defined on more complex domains, such as spheres or rotation groups. The study of such representations is fundamental in physics, chemistry, and signal processing, and has important implications in machine learning, particularly in understanding how symmetries can be incorporated into models like transformers.

#### **Spherical Harmonics**

Spherical harmonics [27, 44] are special functions defined on the surface of a sphere that form an orthogonal basis for the space of square-integrable functions on the sphere $S^2$. They arise naturally in the solution of partial differential equations, such as Laplace's equation, in spherical coordinates. Mathematically, spherical harmonics are the eigenfunctions of the Laplace operator on the sphere [1, 3].

Let $Y_\ell^m(\theta, \phi)$ denote the spherical harmonic of degree $\ell$ and order $m$, where $\theta \in [0, \pi]$ is the polar angle, and $\phi \in [0, 2\pi]$ is the azimuthal angle. The spherical harmonics $Y_\ell^m$ are defined as

$$Y_\ell^m(\theta, \phi) = \sqrt{\frac{(2\ell + 1)(\ell - m)!}{4\pi(\ell + m)!}} P_\ell^m(\cos\theta) e^{im\phi},$$

where $P_\ell^m(x)$ are the associated Legendre polynomials, defined as

$$P_\ell^m(x) = (-1)^m (1 - x^2)^{m/2} \frac{d^m}{dx^m} P_\ell(x)$$

with $P_\ell(x)$ being the Legendre polynomials of degree $\ell$. The indices $\ell$ and $m$ satisfy $\ell \geq 0$ and $-\ell \leq m \leq \ell$.

Properties of Spherical Harmonics:

1. Orthogonality: Spherical harmonics satisfy an orthogonality condition over the sphere:

$$\int_0^{2\pi} \int_0^\pi Y_\ell^m(\theta, \phi) \overline{Y_{\ell'}^{m'}(\theta, \phi)} \sin\theta \, d\theta \, d\phi = \delta_{\ell\ell'} \delta_{mm'},$$

where $\delta_{\ell\ell'}$ and $\delta_{mm'}$ are Kronecker deltas, and the bar denotes complex conjugation.

2. Completeness: Any square-integrable function $f(\theta, \phi)$ on the sphere can be expanded as a series of spherical harmonics:

$$f(\theta, \phi) = \sum_{\ell=0}^\infty \sum_{m=-\ell}^\ell c_\ell^m Y_\ell^m(\theta, \phi),$$

where the coefficients $c_\ell^m$ are given by the inner product:

$$c_\ell^m = \int_0^{2\pi} \int_0^\pi f(\theta, \phi) \overline{Y_\ell^m(\theta, \phi)} \sin\theta \, d\theta \, d\phi.$$

Spherical harmonics are particularly useful in scenarios where the data exhibits rotational symmetry, such as in 3D computer vision or molecular modeling. In transformer models, incorporating spherical harmonics can help the model process and recognize patterns that are invariant under rotations, leading to better generalization and robustness. For example, when dealing with 3D data, a transformer model can use spherical harmonics to encode the rotational properties of the input, allowing the attention mechanism to focus on features that are important regardless of the object's orientation. This approach leverages the inherent symmetry of the data, making the model more efficient and effective.

## Wigner-D Functions

Wigner-D functions [60, 61] generalize the concept of spherical harmonics to more complex group actions, particularly in the context of the rotation group SO(3). They are used to represent the rotations of quantum states and are essential in the study of angular momentum in quantum mechanics. In harmonic analysis, Wigner-D functions provide a way to express the action of the rotation group on functions defined on the sphere [12, 17].

Let $D^j_{mm'}(\alpha, \beta, \gamma)$ denote the Wigner-D function for a rotation parameterized by Euler angles $(\alpha, \beta, \gamma)$, where $j$ is the total angular momentum quantum number, and $m, m'$ are magnetic quantum numbers. The Wigner-D function is defined as

$$D^j_{mm'}(\alpha, \beta, \gamma) = e^{-im\alpha} d^j_{mm'}(\beta) e^{-im'\gamma},$$

where $d^j_{mm'}(\beta)$ are the Wigner small d-matrices, which depend only on the angle $\beta$. These small d-matrices are defined by

$$d^j_{mm'}(\beta) = \sum_k \frac{(-1)^k \sqrt{(j+m)!(j-m)!(j+m')!(j-m')!}}{(j+m-k)!(j-m'-k)!k!(k+m'-m)!} \left(\cos\frac{\beta}{2}\right)^{2j+m-m'-2k} \left(\sin\frac{\beta}{2}\right)^{2k+m'-m}.$$

The Wigner-D functions satisfy important orthogonality relations:

$$\int_0^{2\pi} \int_0^\pi \int_0^{2\pi} D^j_{mm'}(\alpha, \beta, \gamma) \overline{D^{j'}_{m''m'''}(\alpha, \beta, \gamma)} \, d\alpha \, d\beta \, d\gamma = \frac{8\pi^2}{2j+1} \delta_{jj'} \delta_{mm''} \delta_{m'm'''}.$$

Wigner-D functions are essential for representing and analyzing data that is subject to rotations, particularly in 3D spaces. In transformer models, these functions can be used to design attention mechanisms that are invariant to rotations, which is critical in fields like robotics, molecular dynamics, and computer graphics. For instance, in a transformer model designed to process 3D point clouds, Wigner-D functions can be used to ensure that the attention mechanism correctly accounts for the rotational symmetries of the input data. By incorporating Wigner-D functions into the model, one can achieve rotational invariance, meaning that the model's output remains consistent regardless of how the input is rotated.

The use of Wigner-D functions and spherical harmonics in transformers extends the model's capability to understand and exploit the symmetries inherent in the data. These tools provide a mathematical framework for analyzing how the model processes data with continuous symmetries, leading to more robust and intelligent architectures.

## 1.8  Geometric Structures in Transformers

The geometric structures underlying transformer models, particularly in the context of embedding spaces and manifolds, play a crucial role in how these models represent and process data. Understanding these structures from a mathematical perspective allows us to analyze and improve the way transformers learn and utilize these embeddings, especially in high-dimensional spaces where traditional Euclidean intuition may fail. The themes of geometry, symmetry, and intelligence are central to this exploration.

### 1.8.1  Embedding Spaces and Manifolds

In machine learning, embeddings are representations of objects (such as words, images, or nodes in a graph) as vectors in a high-dimensional space. These embedding spaces are often structured in a way that reflects the relationships between the objects they represent. When the data has an underlying geometric or topological structure, it is useful to think of the embedding space as a manifold—a smooth, curved space that locally resembles Euclidean space but may have a more complex global structure.

**Manifold Learning**

Manifold learning [8, 49, 56] is a branch of machine learning and data science that focuses on identifying and exploiting the low-dimensional manifold structure within high-dimensional data. The idea is that although the data may lie in a high-dimensional space, it is often constrained to a lower-dimensional manifold within that space [16, 53].

A manifold $M$ of dimension $d$ is a topological space that locally resembles $\mathbb{R}^d$. Formally, for every point $p \in M$, there exists a neighborhood $U \subseteq M$ and a homeomorphism (a continuous bijection with a continuous inverse) $\varphi : U \to \mathbb{R}^d$. This map $\varphi$ is called a chart, and the collection of all such charts $\{(U_i, \varphi_i)\}$ that cover $M$ forms an atlas for the manifold.

Example: Consider the unit circle $S^1$ as a manifold embedded in $\mathbb{R}^2$. Locally, the circle resembles the real line $\mathbb{R}^1$, and it can be covered by two charts that map parts of the circle to open intervals on the real line.

Manifold learning techniques aim to discover the underlying manifold structure from high-dimensional data. Two popular methods include

1. Isomap: Isomap seeks to preserve the geodesic distances between data points. Geodesic distance is the shortest path between two points on a manifold, taking into account the curvature of the space. Isomap first constructs a neighborhood graph of the data, approximates the geodesic distances, and then applies classical multidimensional scaling (MDS) to find a low-dimensional embedding that preserves these distances.

Given a dataset $X = \{x_1, x_2, \ldots, x_n\}$, construct a graph $G = (V, E)$ where each vertex $v_i$ corresponds to a data point $x_i$ and edges $(v_i, v_j)$ are weighted by the Euclidean distance $\|x_i - x_j\|$ if $x_j$ is among the $k$-nearest neighbors of $x_i$. Compute the shortest path distances $d_G(x_i, x_j)$ on this graph, which approximate the geodesic distances on the manifold. Apply MDS to find a low-dimensional embedding $Y$ that minimizes

$$\sum_{i,j} \left( d_G(x_i, x_j) - \|y_i - y_j\| \right)^2,$$

where $y_i$ and $y_j$ are the corresponding points in the low-dimensional embedding.

2. Locally Linear Embedding (LLE): LLE preserves the local linear structure of the data by assuming that each data point and its neighbors lie on or near a locally linear patch of the manifold. The algorithm seeks a low-dimensional embedding that preserves these local relationships.

For each data point $x_i$, find weights $w_{ij}$ that best reconstruct $x_i$ as a linear combination of its neighbors:

$$\min_w \sum_i \left\| x_i - \sum_j w_{ij} x_j \right\|^2, \quad \text{subject to} \sum_j w_{ij} = 1.$$

Then, find a low-dimensional embedding $Y = \{y_1, y_2, \ldots, y_n\}$ that preserves these reconstruction weights:

$$\min_Y \sum_i \left\| y_i - \sum_j w_{ij} y_j \right\|^2.$$

These manifold learning techniques are powerful because they reveal the intrinsic geometry of the data, which can be crucial for tasks like dimensionality reduction, visualization, and unsupervised learning.

In the context of transformers, manifold learning can be applied to the design of embedding spaces, particularly when the data exhibits a non-trivial geometric structure. For instance, word embeddings in natural language processing often lie on a low-dimensional manifold within the high-dimensional space in which they are embedded. By understanding this manifold structure, one can design more efficient and effective embedding layers that capture the essential relationships between words or other data points. Moreover, attention mechanisms in transformers can be interpreted as a form of manifold learning, where the goal is to learn a mapping from the input sequence to a manifold that captures the relevant relationships for the task at hand. This perspective can lead to new insights into how to design attention mechanisms that better exploit the underlying geometry of the data.

**Geometry of High-Dimensional Embeddings**

High-dimensional embeddings are central to transformer models, as they allow the model to represent complex data in a form that can be processed by the self-attention mechanism. However, working in high-dimensional spaces introduces unique challenges, particularly related to the geometry of these spaces.

Curse of Dimensionality: As the dimensionality of the embedding space increases, many aspects of geometry behave counterintuitively. For example, the volume of a high-dimensional sphere becomes concentrated near its surface, and the distance between randomly chosen points tends to become uniform. This phenomenon is

known as the curse of dimensionality. Mathematically, consider a $d$-dimensional Euclidean space $\mathbb{R}^d$. The volume of a ball of radius $r$ in this space is given by

$$V_d(r) = \frac{\pi^{d/2} r^d}{\Gamma\left(\frac{d}{2} + 1\right)},$$

where $\Gamma$ is the Gamma function. As $d$ increases, the volume $V_d(r)$ grows rapidly, but the volume within any small distance $\epsilon$ from the surface of the ball also grows, meaning that most of the volume is near the boundary. This has implications for nearest-neighbor search, clustering, and other tasks that rely on geometric relationships in high-dimensional spaces.

The geometry of high-dimensional embeddings in transformers is crucial for understanding how the model processes and transforms data. For example, the self-attention mechanism in transformers can be seen as a process of mapping high-dimensional embeddings to a lower-dimensional space (or manifold) where the relevant relationships between data points are more apparent. In this context, the choice of embedding space and the way embeddings are learned and manipulated can significantly impact the model's performance. For instance, the attention mechanism may need to account for the concentration of measure to avoid issues where all attention scores become similar, leading to a loss of discriminative power. One way to address these challenges is to design embedding spaces with explicit geometric structures, such as hyperbolic spaces, which have been shown to better capture hierarchical relationships in data. Hyperbolic embeddings can lead to more effective attention mechanisms, particularly in tasks that involve tree-like structures or other forms of hierarchical data.

## *1.8.2 Symmetries and Transformations*

Symmetry plays a fundamental role in mathematics and physics, providing a powerful framework for understanding the invariance and transformation properties of systems. In the context of transformers and machine learning, symmetries can be exploited to design models that are robust, efficient, and capable of generalizing well to unseen data. This section explores the mathematical foundation of symmetries, focusing on group actions on manifolds and their applications in transformer architectures.

**Group Actions on Manifolds**

A group action is a formal way of describing how the elements of a group $G$ systematically "act" on the elements of a set $M$, which in this context is often a manifold [9, 11, 36, 43]. Group actions provide a way to model symmetries and transformations of geometric objects. Formally, a group action of a group $G$ on a manifold $M$ is a map $\phi : G \times M \to M$ that satisfies the following properties:

1. Identity: The identity element $e$ of $G$ acts as the identity transformation on $M$, i.e., for all $x \in M$:

$$\phi(e, x) = x.$$

2. Compatibility: For all $g, h \in G$ and $x \in M$:

$$\phi(g, \phi(h, x)) = \phi(gh, x).$$

When a group $G$ acts on a manifold $M$, it induces a structure on $M$ that reflects the symmetries of $G$. The orbits of this action, which are the sets of points in $M$ that can be reached from one another by the action of elements in $G$, provide insight into the geometry of $M$.

Example: Consider the rotation group SO(2), which consists of all rotations in the plane. This group acts on the 2D Euclidean space $\mathbb{R}^2$ by rotating vectors around the origin. For a point $x \in \mathbb{R}^2$, the action of a rotation $R_\theta$ by an angle $\theta$ is given by

$$\phi(R_\theta, x) = R_\theta x = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}.$$

This action preserves the distance from the origin, reflecting the rotational symmetry of the Euclidean space.

Orbits and Stabilizers: For a point $x \in M$, the orbit of $x$ under the action of $G$ is the set:

$$\text{Orb}(x) = \{\phi(g, x) \mid g \in G\}.$$

The stabilizer (or isotropy group) of $x$ is the set of elements in $G$ that leave $x$ fixed:

$$\text{Stab}(x) = \{g \in G \mid \phi(g, x) = x\}.$$

The orbit-stabilizer theorem states that there is a bijection between the orbit of $x$ and the coset space $G/\text{Stab}(x)$, which provides a powerful way to understand the structure of the manifold $M$ in terms of the group $G$.

**Theorem 1.4** (Orbit-Stabilizer Theorem) *Let $G$ be a group acting on a set $M$, and let $x \in M$. Then there is a bijection between the orbit of $x$ and the coset space $G/\text{Stab}(x)$, given by*

$$\phi : G/\text{Stab}(x) \to \text{Orb}(x), \quad g\,\text{Stab}(x) \mapsto \phi(g, x),$$

*where Stab(x) is the stabilizer of x.*

Group actions are fundamental in differential geometry, where they are used to study the symmetry properties of manifolds. For example, the action of the rotation group SO(3) on the 2-sphere $S^2$ describes the rotational symmetries of the sphere. These symmetries can be leveraged in machine learning models, particularly in tasks involving data with inherent geometric structures, such as 3D point clouds or molecular data.

## Applications in Transformer Architectures

In transformer architectures, understanding and incorporating symmetries through group actions can lead to more robust and efficient models. This is especially relevant when dealing with data that possesses intrinsic symmetries, such as images, sequences, or graphs.

1. Equivariance and Invariance: Equivariance and invariance are crucial concepts in the design of neural networks, including transformers. A function $f : M \to N$ is said to be equivariant with respect to a group action if, for all $g \in G$ and $x \in M$,

$$f(\phi(g, x)) = \psi(g, f(x)),$$

where $\psi : G \times N \to N$ is the corresponding action on the target space $N$. If $\psi(g, f(x)) = f(x)$ for all $g \in G$, then $f$ is invariant under the group action.

Example: In the context of image processing, if $G$ is the group of translations, a convolutional layer in a neural network is designed to be translation-equivariant, meaning that translating the input image results in a corresponding translation of the output feature map. This property is critical for tasks where the position of features should not affect their detection.

In transformers, the self-attention mechanism can be designed to be equivariant to certain transformations of the input. For example, if the input data has a known symmetry, such as rotational or permutation symmetry, the attention mechanism can be modified to respect this symmetry. Let $X = \{x_1, x_2, \ldots, x_n\}$ be the input sequence, and let $G$ be a group acting on this sequence, such as a permutation group that shuffles the elements of $X$. The self-attention mechanism can be made equivariant to the group action if, for a group element $g \in G$,

$$\text{Attention}(\phi(g, Q), \phi(g, K), \phi(g, V)) = \phi(g, \text{Attention}(Q, K, V)),$$

where $Q, K, V$ are the query, key, and value matrices, and $\phi(g, \cdot)$ denotes the action of $g$ on these matrices. This equivariance ensures that the model's output is consistent with the symmetries of the input, which can improve generalization and robustness, especially in tasks where such symmetries are present.

2. Symmetry-Aware Transformers: In tasks involving structured data, such as graphs or 3D objects, symmetry-aware transformers can be designed by incorporating group representations and equivariant layers into the model. For instance, in a graph transformer, the attention mechanism can be designed to be equivariant to node permutations, ensuring that the model's output is invariant to the ordering of the nodes. To implement a symmetry-aware transformer, one can use tools from representation theory, where the input features are mapped into a space where the group $G$ acts linearly. The attention weights and output features are then computed in a way that respects the group action, ensuring equivariance or invariance as required.

Example: For a transformer designed to process 3D molecular structures, the model can incorporate the rotational symmetries of the molecule by representing the atomic positions using spherical harmonics or Wigner-D functions (as discussed in the previous section). The attention mechanism can then be made equivariant to rotations, allowing the model to correctly account for the orientation of the molecule.

3. Enhancing Model Robustness: By explicitly incorporating symmetries into transformer architectures, one can enhance the model's robustness to transformations of the input data. This approach is particularly useful in domains where the data exhibits natural symmetries, such as robotics, physics, and computer vision.

### 1.8.3   Implications for Model Design

Incorporating geometric principles into the design of transformer architectures can significantly enhance their performance and robustness, especially in tasks where the data exhibits inherent symmetries and geometric structures. By embedding geometric priors into model architectures and ensuring invariances to transformations like rotations and translations, we can create models that are not only more aligned with the underlying data but also more capable of generalizing to new, unseen scenarios.

**Geometric Priors in Model Architectures**

Geometric priors are assumptions about the underlying structure or symmetry of the data that can be embedded into the architecture of a model. These priors leverage known properties of the data, such as its invariance under certain transformations, to constrain the model's hypothesis space, leading to more efficient learning and improved generalization.

Consider a dataset $X$ where each data point $x \in X$ lies on or near a manifold $M \subset \mathbb{R}^d$. A geometric prior can be encoded in the model by designing the model's architecture to respect the structure of $M$. For example, if $M$ has a known symmetry group $G$, such as the rotation group SO(3), the model can be constructed to be equivariant to the actions of $G$.

Example: Let $G = \text{SO}(3)$ act on $\mathbb{R}^3$, representing rotations in three-dimensional space. A model designed to process 3D point clouds can incorporate a geometric prior

by ensuring that its operations are equivariant to $G$. Specifically, if $f : \mathbb{R}^3 \to \mathbb{R}^n$ is a function representing a layer in the model, we require that for any rotation $R \in G$ and any point $x \in \mathbb{R}^3$:

$$f(Rx) = Rf(x).$$

This ensures that the model's output respects the rotational symmetry of the input data.

Geometric Deep Learning: Geometric priors are central to the field of geometric deep learning, where the goal is to design neural networks that operate on non-Euclidean domains such as graphs, manifolds, and other geometric structures. In transformers, these priors can guide the design of attention mechanisms, embedding layers, and other components to ensure that they respect the underlying geometry of the data.

One approach to incorporating geometric priors is to use convolutional filters that are designed to be equivariant to specific transformations. For example, in a spherical convolutional neural network, the filters are defined on the sphere $S^2$ and are equivariant to rotations. Mathematically, if $\sigma : S^2 \to \mathbb{R}^n$ represents a filter on the sphere, the convolution operation $*$ with an input function $f : S^2 \to \mathbb{R}^m$ can be written as

$$(f * \sigma)(x) = \int_{S^2} f(y)\sigma(R^{-1}y)\,d\mu(y),$$

where $R \in \mathrm{SO}(3)$ is a rotation, and $d\mu(y)$ is the surface measure on the sphere.

By designing models with such equivariant operations, we encode the geometric prior directly into the architecture, ensuring that the model respects the symmetries of the data.

**Rotational and Translational Invariances**

Rotational and translational invariances are specific types of geometric invariances that are particularly important in tasks involving spatial data, such as image recognition, 3D modeling, and physical simulations. Ensuring that a model is invariant to these transformations allows it to recognize patterns regardless of their orientation or position, leading to more robust and generalizable performance.

Rotational Invariance: A function $f : \mathbb{R}^d \to \mathbb{R}^n$ is said to be rotationally invariant if, for any rotation $R \in \mathrm{SO}(d)$ and any point $x \in \mathbb{R}^d$,

$$f(Rx) = f(x).$$

Rotational invariance ensures that the output of the function remains unchanged under rotations of the input. In the context of transformers, rotational invariance can be particularly useful when dealing with data that has no preferred orientation, such as molecular structures or astronomical images.

Example: Consider a transformer designed to classify 3D molecular structures. To ensure rotational invariance, the attention mechanism can be designed to compute attention scores based on distances between atoms rather than their absolute positions. If $x_i, x_j \in \mathbb{R}^3$ are the positions of two atoms, the attention score could be based on the distance $\|x_i - x_j\|$, which is invariant under rotations.

Translational Invariance: A function $f : \mathbb{R}^d \to \mathbb{R}^n$ is translationally invariant if, for any translation vector $t \in \mathbb{R}^d$ and any point $x \in \mathbb{R}^d$,

$$f(x + t) = f(x).$$

Translational invariance ensures that the function's output does not change when the input is shifted in space. This property is crucial in tasks like image recognition, where objects can appear at different locations within an image.

Example: In image processing, convolutional layers inherently provide translational invariance, as the same filter is applied across the entire image. In transformers, translational invariance can be achieved by using relative positional encodings, where the attention mechanism focuses on the relative positions of elements in the sequence rather than their absolute positions.

To implement rotational and translational invariances in transformers, one approach is to design the self-attention mechanism and positional encodings to be invariant to these transformations. For rotational invariance, the attention mechanism can be modified to use features like distances, angles, or other rotationally invariant quantities. For translational invariance, relative positional encodings can be used to ensure that the attention scores depend only on the relative positions of the elements in the sequence.

**Definition 1.1** (*Equivariance and Invariance in Attention*) Let $G$ be a group acting on the input space $X$, and let $\phi : G \times X \to X$ be the group action. An attention block is said to be equivariant to the group $G$ if the self-attention mechanism satisfies

$$\text{Attention}(\phi(g, Q), \phi(g, K), \phi(g, V)) = \phi(g, \text{Attention}(Q, K, V)),$$

for all $g \in G$. If $\phi(g, y) = y$ for all $g \in G$ and $y \in \mathbb{R}^d$, the model is invariant to the group $G$.

## 1.9  Function Approximation Theory

Approximation theory [13, 46, 48] is a foundational branch of mathematical analysis that deals with how functions can be approximated by simpler functions, such as polynomials, trigonometric functions, or other basis functions. This theory is essential in many areas of numerical analysis, signal processing, and machine learning, where the goal is to approximate complex functions with models that are computationally feasible. In the context of transformers and neural networks, approximation

theory provides the mathematical underpinning for understanding how well a model can represent the underlying function or data distribution.

### 1.9.1  Introduction to Approximation Theory

Approximation theory is concerned with the approximation of functions by simpler or more convenient functions, often within a specified normed space. The most common types of approximations involve polynomial and trigonometric functions, which are used due to their well-understood properties and the ability to provide good approximations under certain conditions.

Let $f : [a, b] \to \mathbb{R}$ be a continuous function. The goal of approximation theory is to find a sequence of simpler functions $\{f_n\}$ that converge to $f$ in some sense, typically measured by a norm such as the $L^p$ norm or the maximum norm (sup norm).

#### Polynomial and Trigonometric Approximations

Polynomial Approximation: One of the most fundamental results in approximation theory is the Weierstrass approximation theorem, which states that any continuous function defined on a closed interval can be uniformly approximated by polynomials. Formally, if $f$ is a continuous function on $[a, b]$, then for every $\epsilon > 0$, there exists a polynomial $P_n(x)$ such that

$$\|f - P_n\|_\infty = \sup_{x \in [a,b]} |f(x) - P_n(x)| < \epsilon.$$

This result implies that polynomials are dense in the space of continuous functions with respect to the uniform norm.

Example: Consider the function $f(x) = \sin(x)$ on the interval $[0, \pi]$. Using Taylor series expansion around $x = 0$, we can approximate $\sin(x)$ by a polynomial:

$$\sin(x) \approx x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots$$

As more terms are included in the series, the polynomial approximation becomes increasingly accurate.

Trigonometric Approximation: For periodic functions, trigonometric polynomials provide a natural and powerful tool for approximation. A trigonometric polynomial of degree $n$ is a function of the form:

$$T_n(x) = a_0 + \sum_{k=1}^{n} (a_k \cos(kx) + b_k \sin(kx)),$$

where $a_k$ and $b_k$ are coefficients determined by the Fourier series of the function. The Fourier series of a function $f(x)$ with period $2\pi$ is given by

$$f(x) = \frac{a_0}{2} + \sum_{k=1}^{\infty} (a_k \cos(kx) + b_k \sin(kx)),$$

where the coefficients $a_k$ and $b_k$ are computed as

$$a_k = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos(kx)\, dx, \quad b_k = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin(kx)\, dx.$$

Example: The function $f(x) = |x|$ on $[-\pi, \pi]$ can be approximated by its Fourier series. Since $f(x)$ is an even function, all the sine coefficients $b_k$ are zero, and the series involves only cosine terms.

The convergence of these series can be studied using various norms, with the $L^2$ norm being particularly important due to its connection with Parseval's theorem, which states that the sum of the squares of the Fourier coefficients equals the $L^2$ norm of the function:

$$\|f\|_2^2 = \sum_{k=0}^{\infty} (a_k^2 + b_k^2).$$

**Jackson's Theorems**

Jackson's theorems [32] provide quantitative results on the rate of approximation of continuous functions by polynomials or trigonometric polynomials. These theorems establish a relationship between the smoothness of a function and the rate at which its best approximation by polynomials converges to the function itself.

**Theorem 1.5** (Jackson's Polynomial Approximation) *Let $f : [-1, 1] \to \mathbb{R}$ be a continuous function with a modulus of continuity $\omega(f, \delta)$, defined as*

$$\omega(f, \delta) = \sup_{|x-y| \leq \delta} |f(x) - f(y)|.$$

*Then there exists a polynomial $P_n(x)$ of degree n such that*

$$\|f - P_n\|_\infty \leq C \frac{\omega(f, 1/n)}{n},$$

*where C is a constant independent of n and $f$.*

This result shows that the rate of convergence of the polynomial approximation depends on the smoothness of the function $f$.

Example: If $f(x) = |x|$ on $[-1, 1]$, then the modulus of continuity $\omega(f, \delta) = \delta$. According to Jackson's theorem, the error in approximating $f(x)$ by a polynomial of degree $n$ is on the order of $1/n$.

**Theorem 1.6** (Jackson's Trigonometric Approximation) *A similar result holds for the approximation of periodic functions by trigonometric polynomials. If $f$ is a periodic function with modulus of continuity $\omega(f, \delta)$, then there exists a trigonometric polynomial $T_n(x)$ of degree $n$ such that*

$$\|f - T_n\|_\infty \leq C \frac{\omega(f, 1/n)}{n}.$$

This result is particularly important in the analysis of Fourier series and the approximation of periodic functions in signal processing and harmonic analysis.

Implications for Machine Learning: In machine learning, approximation theory provides a theoretical foundation for understanding how well a model can approximate a target function. For example, in neural networks, the universal approximation theorem states that a sufficiently large neural network can approximate any continuous function on a compact domain. This result is analogous to the Weierstrass approximation theorem but in the context of neural networks. However, the practical performance of neural networks depends on more than just the existence of an approximation; it also depends on the rate of convergence and the smoothness of the function being approximated. Jackson's theorems provide insight into how the smoothness of the target function affects the quality of approximation, which in turn influences how deep or complex a neural network needs to be to achieve a certain level of accuracy.

### *1.9.2  Universal Approximation Theorems*

Universal approximation theorems [6, 15, 30, 41] form the theoretical backbone of many neural network architectures, including transformers. These theorems provide guarantees that neural networks can approximate a wide range of functions to arbitrary accuracy, under certain conditions. Understanding these theorems is essential for grasping the potential and limitations of neural networks as function approximators, as well as their implications for complex architectures like transformers.

**Theorem 1.7** (Universal Approximation) *Let $\sigma : \mathbb{R} \to \mathbb{R}$ be a non-constant, bounded, and continuous activation function. Then, for any continuous function $f : [a, b] \to \mathbb{R}$ and any $\epsilon > 0$, there exists a neural network with one hidden layer that approximates $f$ to within $\epsilon$ in the $L^\infty$ norm:*

$$\| f(x) - \sum_{i=1}^{N} c_i \sigma (a_i x + b_i) \|_\infty < \epsilon$$

*for some coefficients* $c_i, a_i, b_i \in \mathbb{R}$.

**Neural Networks as Universal Approximators**

The universal approximation theorem, in its most basic form, asserts that a feed-forward neural network with a single hidden layer can approximate any continuous function on a compact domain, given a sufficient number of hidden units. This result is both powerful and foundational, as it assures us that even simple neural networks are, in principle, capable of representing highly complex functions.

Let $\sigma : \mathbb{R} \to \mathbb{R}$ be a non-constant, bounded, and continuous activation function. Consider a function $f : [a, b] \subset \mathbb{R}^n \to \mathbb{R}$ that is continuous. The universal approximation theorem states that for any $\epsilon > 0$, there exists a neural network $\hat{f}(x)$ of the form:

$$\hat{f}(x) = \sum_{i=1}^{N} c_i \sigma (\langle w_i, x \rangle + b_i),$$

where $c_i \in \mathbb{R}$, $w_i \in \mathbb{R}^n$, and $b_i \in \mathbb{R}$ are the parameters of the network, such that

$$\sup_{x \in [a,b]} |f(x) - \hat{f}(x)| < \epsilon.$$

This result implies that the set of functions representable by such a neural network is dense in the space of continuous functions on $[a, b]$ with respect to the sup norm $\| \cdot \|_\infty$. The theorem does not specify how large the network must be to achieve a given level of accuracy, nor does it provide insights into the efficiency or training complexity of such networks.

Proof Sketch: The proof involves constructing a sequence of neural network approximations to a given continuous function $f(x)$ and showing that this sequence converges uniformly to $f(x)$. The construction leverages the fact that the activation function $\sigma(x)$ can be used to approximate simple step functions, which in turn can be used to approximate more general functions through a combination of translations and scalings.

Example: Consider the function $f(x) = |x|$ on the interval $[-1, 1]$. A neural network with a ReLU activation function $\sigma(x) = \max(0, x)$ can approximate $f(x)$ as follows:

$$f(x) = x \cdot \mathbf{1}_{x \geq 0} - x \cdot \mathbf{1}_{x < 0}.$$

This can be rewritten as

$$f(x) = \sigma(x) + \sigma(-x)$$

demonstrating that even simple piecewise linear functions like the absolute value function can be exactly represented by a neural network with ReLU activation.

Extensions: The universal approximation theorem has been extended in various ways, including to networks with more than one hidden layer (deep networks) and to different types of activation functions, such as sigmoid or hyperbolic tangent functions. These extensions show that deep neural networks, with sufficient depth and width, are also universal approximators.

**Implications for Transformers**

Transformers, as deep neural network architectures, inherit the universal approximation capabilities discussed in the universal approximation theorem. However, their specific architecture, which relies on self-attention mechanisms and layer normalization, introduces additional layers of complexity and expressiveness.

Self-Attention Mechanism as a Function Approximator: The self-attention mechanism in transformers can be viewed as a special case of a neural network layer, where the attention weights dynamically adjust based on the input sequence. This mechanism allows the model to focus on different parts of the input sequence, effectively approximating functions that capture dependencies and interactions between different elements of the sequence.

Mathematically, let $X = \{x_1, x_2, \ldots, x_n\}$ be an input sequence, and let $W_Q, W_K, W_V$ be the learned weight matrices for the query, key, and value projections. The self-attention mechanism computes the output as

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right) V,$$

where $Q = XW_Q$, $K = XW_K$, and $V = XW_V$. The softmax operation ensures that the attention weights sum to 1, providing a weighted sum of the value vectors. This mechanism is highly expressive, as it can capture complex dependencies within the sequence, effectively learning to approximate functions that map input sequences to output sequences.

Transformers as Universal Approximators: Transformers, like other deep neural networks, are capable of universal approximation. The specific architecture of transformers, with multiple layers of self-attention followed by feedforward neural networks, allows them to approximate a wide range of functions. The stacking of layers in transformers provides additional depth, which is crucial for capturing hierarchical structures in data.

Implications for Model Design: The universal approximation property of transformers suggests that, in principle, they can approximate any function on a compact domain, given sufficient layers and attention heads. However, in practice, the ability to achieve this approximation depends on factors such as

1. Model Capacity: The number of layers, attention heads, and hidden units determines the capacity of the transformer to approximate complex functions. Increasing

the model capacity not only improves its expressiveness but also increases the risk of overfitting and computational cost.

2. Training Efficiency: The training process, including the choice of optimization algorithms and learning rate schedules, affects how well the transformer can approximate the target function. Convergence guarantees, as discussed in optimization theory, play a critical role in determining whether the model reaches a good approximation.

3. Generalization: While transformers can approximate functions on the training data, their ability to generalize to unseen data is crucial for real-world applications. This generalization depends on the regularization techniques used, such as dropout or weight decay, and the availability of sufficient and diverse training data.

**Theorem 1.8** (Universal Approximation for Transformers) *Let $f : [a, b]^n \to \mathbb{R}^m$ be a continuous function on a compact domain. For any $\epsilon > 0$, there exists a transformer model $T_\theta$ with a sufficient number of layers and attention heads, and parameters $\theta$, such that*

$$\sup_{x \in [a,b]^n} \| f(x) - T_\theta(x) \| < \epsilon.$$

This theorem highlights that transformers, given appropriate architecture and training, can approximate complex functions to arbitrary accuracy, aligning with the broader results of universal approximation in neural networks.

### *1.9.3   Expressivity in Transformer Models*

Expressivity in transformer models refers to the ability of these models to represent a wide range of functions, capturing complex relationships within data. This expressivity is influenced by various factors, including the depth and width of the model, as well as the role of non-linearities in the architecture. A deep understanding of these factors is essential for designing transformers that are not only powerful but also efficient in their computational requirements.

**Depth Versus Width in Model Design**

In the design of neural networks, including transformers, two critical architectural aspects are depth (the number of layers) and width (the number of units or neurons in each layer). These two dimensions influence the model's capacity to approximate functions and capture complex patterns in the data.

Depth in a transformer model refers to the number of layers in the network. Each layer typically consists of a self-attention mechanism followed by a feedforward neural network. Depth allows the model to build hierarchical representations of the

data, with deeper layers capturing increasingly abstract features. Mathematically, let $f : \mathbb{R}^n \to \mathbb{R}^m$ be a function that the transformer aims to approximate. The output of a transformer with $L$ layers can be represented as

$$f(x) \approx T_L(T_{L-1}(\ldots T_1(x) \ldots)),$$

where each $T_i(x)$ represents the transformation applied by the $i$th layer, including self-attention and feedforward operations.

Theoretical results suggest that increasing the depth of a network enhances its expressivity, often exponentially. A seminal result by [55] shows that deep networks can represent highly oscillatory functions that shallow networks cannot approximate efficiently, even with an exponentially larger number of units.

**Theorem 1.9** *Let $f_d$ be a function representable by a deep network with d layers, each with a fixed number of units. There exists a family of functions $f_d(x)$ such that $\sup_x |f_d(x)| = 1$, but any network with $d - 1$ layers requires exponentially more units to approximate $f_d$ within a given error $\epsilon$.*

This result indicates that depth can provide an exponential increase in expressivity, allowing deep networks to approximate functions with complex structures that shallow networks cannot.

Width, on the other hand, refers to the number of units in each layer of the network. While increasing width can also enhance the model's capacity, its impact is different from that of depth. Wide networks can represent more features at each level of abstraction but may struggle to capture deep, hierarchical structures.

A fundamental result in approximation theory is that sufficiently wide networks with a single hidden layer can approximate any continuous function on a compact domain (the universal approximation theorem). However, this comes at the cost of potentially requiring an exponentially large number of units, especially if the target function has high complexity.

Expressivity and the Trade-off: The expressivity of a transformer model is therefore a balance between depth and width. Deeper models can capture more complex, hierarchical relationships, while wider models can represent more features simultaneously. The choice between depth and width is often dictated by the specific task, the nature of the data, and computational constraints.

In practical terms, deep transformers with moderate width are often favored in tasks requiring the modeling of long-range dependencies and intricate patterns, as the layers can progressively refine the representations. However, the increased depth comes with challenges, such as the vanishing gradient problem, which can impede training.

Example: Consider the task of machine translation, where the goal is to map a sequence of words in one language to a sequence in another language. A deep transformer can effectively capture the hierarchical structure of the input sentence (e.g., syntactic dependencies) and produce a context-aware representation that facilitates

accurate translation. If the transformer were shallow but wide, it might struggle to capture these dependencies effectively, leading to inferior translation quality.

**Role of Non-linearities**

Non-linearities in neural networks, including transformers, are crucial for enabling these models to approximate complex functions. Without non-linear activation functions, a neural network composed of linear layers would be equivalent to a single linear transformation, regardless of the depth, and thus would be incapable of representing non-linear relationships in the data.

A non-linear activation function $\sigma : \mathbb{R} \to \mathbb{R}$ is applied element-wise to the output of each layer in a neural network. Common activation functions include the rectified linear unit (ReLU), sigmoid, and hyperbolic tangent (tanh). The non-linearity introduced by these functions allows the network to approximate non-linear functions. Mathematically, consider a feedforward neural network layer given by

$$z^{(l+1)} = \sigma(W^{(l)} z^{(l)} + b^{(l)}),$$

where $W^{(l)}$ is the weight matrix, $z^{(l)}$ is the input to the layer, $b^{(l)}$ is the bias term, and $\sigma$ is the activation function. The non-linearity $\sigma$ is what allows the network to break free from the limitations of linear transformations.

Impact of Non-Linearities on Expressivity: The inclusion of non-linearities greatly increases the expressivity of a neural network. A network with non-linear activation functions can approximate a much broader class of functions compared to a purely linear network. This is critical in transformers, where non-linearities enable the self-attention mechanism and the feedforward layers to capture complex patterns and relationships in the data.

In a transformer model, the feedforward network following the self-attention mechanism typically includes a non-linear activation function like ReLU. This non-linearity allows the model to refine the representation of the input sequence in a non-linear way, enhancing its ability to model intricate dependencies and interactions between elements of the sequence.

Deep Versus Shallow Networks and Non-Linearities: Deep networks with multiple layers of non-linear transformations can approximate functions that require multiple levels of abstraction, something that shallow networks with the same number of non-linear units cannot do efficiently. This layered structure, combined with non-linearities, allows transformers to progressively build complex representations from simpler ones, capturing the hierarchical nature of many tasks, such as language processing and image recognition.

**Theorem 1.10** (Depth Efficiency of Deep Networks) *Let $\sigma : \mathbb{R} \to \mathbb{R}$ be a fixed non-linear activation function, and consider the function class $\mathcal{F}_{L,d}$ representing the set of functions $f : \mathbb{R}^n \to \mathbb{R}$ that can be realized by a feedforward neural network with L layers, each layer containing at most d units, and each unit applying the*

*activation function $\sigma$. For any positive integer $L_1$, there exists a function $f \in \mathcal{F}_{L_2,d_2}$
for some $L_2 > L_1$ and $d_2 \geq d_1$, such that for any function $\tilde{f} \in \mathcal{F}_{L_1,d_1}$, where $d_1$ is a
polynomially bounded function of $d_2$, the approximation error satisfies*

$$\inf_{\tilde{f} \in \mathcal{F}_{L_1,d_1}} \| f - \tilde{f} \|_\infty \geq \epsilon$$

*for some $\epsilon > 0$, unless $d_1$ grows exponentially in $L_2 - L_1$. In other words, there
exist functions that can be represented by deep networks with $L_2$ layers and $d_2$ units
per layer that require exponentially more units to be approximated by shallower
networks with $L_1$ layers, even when using the same non-linear activation $\sigma$.*

This theorem underscores the importance of depth and non-linearities in neural
networks, as they enable the representation of functions with complex hierarchical
structures that would be intractable for shallow networks.

Implications for Transformers: In transformers, the role of non-linearities extends
beyond individual layers. The architecture as a whole benefits from the non-linearities
in the self-attention mechanism, where the softmax function introduces a non-
linearity that is critical for differentiating between important and less important
elements of the sequence. Additionally, the non-linear transformations in the feed-
forward layers help to refine these attention-based representations, further enhancing
the model's expressivity.

## 1.10   Optimization Techniques

Optimization techniques are central to training machine learning models, including
transformers. These techniques involve adjusting the model parameters to minimize
a loss function, which measures the discrepancy between the model's predictions and
the actual outcomes. This section delves into the mathematical foundations of various
optimization algorithms, highlighting their convergence properties, computational
efficiency, and applicability in different scenarios.

### *1.10.1   Gradient Descent and Variants*

Gradient descent is the foundational algorithm for optimization in machine learning.
It is an iterative method used to minimize a differentiable function $f(\theta)$, where $\theta$
represents the model parameters.

The basic idea of gradient descent is to update the parameters $\theta$ in the direction of
the negative gradient of the loss function, which points toward the steepest descent.
Mathematically, the update rule for gradient descent is given by

$$\theta_{t+1} = \theta_t - \eta \nabla f(\theta_t),$$

where $\eta > 0$ is the learning rate and $\nabla f(\theta_t)$ is the gradient of the loss function $f$ with respect to $\theta_t$ at iteration $t$.

The convergence of gradient descent depends on the properties of the loss function $f(\theta)$. If $f$ is convex and has Lipschitz-continuous gradients (i.e., there exists a constant $L > 0$ such that $\|\nabla f(\theta_1) - \nabla f(\theta_2)\| \leq L\|\theta_1 - \theta_2\|$ for all $\theta_1, \theta_2$), then gradient descent converges to a global minimum.

**Theorem 1.11** (Convergence of Gradient Descent) *Let $f : \mathbb{R}^n \to \mathbb{R}$ be a convex function with Lipschitz-continuous gradients. Then, for a sufficiently small learning rate $\eta$, gradient descent converges to a global minimum $\theta^*$:*

$$\|\theta_t - \theta^*\| \leq \frac{1}{t}\left(\frac{2L(f(\theta_0) - f(\theta^*))}{\mu}\right),$$

*where $\mu$ is the strong convexity constant of $f$.*

This result guarantees that the distance between the current parameter vector and the optimal parameter vector decreases over time, leading to convergence.

Variants of Gradient Descent:

1. Batch Gradient Descent: In batch gradient descent, the gradient is computed using the entire dataset. This approach provides a stable estimate of the gradient but can be computationally expensive, especially for large datasets.

2. Mini-batch Gradient Descent: Mini-batch gradient descent computes the gradient using a subset of the data, called a mini-batch. This approach balances the computational efficiency of stochastic gradient descent with the stability of batch gradient descent.

Example: Consider a simple linear regression model with parameters $\theta = (\theta_0, \theta_1)$ and a mean squared error loss function:

$$f(\theta) = \frac{1}{2m}\sum_{i=1}^{m}(y_i - (\theta_0 + \theta_1 x_i))^2.$$

The gradient descent update rules for $\theta_0$ and $\theta_1$ are

$$\theta_0 \leftarrow \theta_0 - \eta\frac{1}{m}\sum_{i=1}^{m}(\theta_0 + \theta_1 x_i - y_i),$$

$$\theta_1 \leftarrow \theta_1 - \eta\frac{1}{m}\sum_{i=1}^{m}(\theta_0 + \theta_1 x_i - y_i)x_i.$$

These updates iteratively adjust the parameters to minimize the loss function.

**Stochastic Gradient Descent**

Stochastic Gradient Descent (SGD) is a variant of gradient descent where the gradient is computed using a single data point or a small random subset of the data at each iteration. This introduces stochasticity into the optimization process, which can help the algorithm escape local minima. In SGD, the parameter update rule is given by

$$\theta_{t+1} = \theta_t - \eta \nabla f_i(\theta_t),$$

where $f_i(\theta_t)$ is the loss function evaluated on the $i$th data point.

SGD does not guarantee convergence to the global minimum, but it does converge to a region close to the minimum under certain conditions, especially when using decreasing learning rates. The stochastic nature of SGD allows it to explore the loss surface more effectively than batch gradient descent.

**Theorem 1.12** (Convergence of SGD) *Let $f(\theta)$ be a convex function, and let the learning rate $\eta_t$ satisfy the conditions:*

$$\sum_{t=1}^{\infty} \eta_t = \infty \quad \text{and} \quad \sum_{t=1}^{\infty} \eta_t^2 < \infty.$$

*Then, SGD converges almost surely to a global minimum $\theta^*$:*

$$\lim_{t \to \infty} \mathbb{E}[f(\theta_t)] = f(\theta^*).$$

This result shows that, with an appropriate learning rate schedule, SGD can converge to the optimal solution in expectation.

In the context of deep learning, SGD is often preferred for training large models because it is computationally efficient and can handle large datasets. The noise introduced by the stochastic updates can also help the model generalize better by avoiding overfitting to the training data.

**Adam Optimizer**

The Adam optimizer [34] is an extension of SGD that incorporates adaptive learning rates for each parameter. It combines the benefits of two other extensions: RMSProp and momentum. Adam computes individual adaptive learning rates for different parameters using estimates of the first and second moments of the gradients:

1. First Moment Estimate (Mean):

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_t,$$

where $g_t = \nabla f(\theta_t)$ is the gradient at time step $t$ and $\beta_1$ is the decay rate for the moving average of the gradients.

2. Second Moment Estimate (Variance):

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)g_t^2,$$

where $\beta_2$ is the decay rate for the moving average of the squared gradients.

3. Bias-Corrected Estimates: The estimates $m_t$ and $v_t$ are biased toward zero, especially in the early iterations. The bias-corrected estimates are

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}.$$

4. Parameter Update: The parameters are updated as

$$\theta_{t+1} = \theta_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon},$$

where $\eta$ is the learning rate and $\epsilon$ is a small constant to prevent division by zero.

Advantages of Adam: Adam adjusts the learning rates based on the magnitudes of the gradients, allowing it to perform well on problems with sparse gradients or varying gradient magnitudes. It incorporates momentum, which helps accelerate convergence and smooths the optimization trajectory.

Convergence Analysis: Adam has been shown to perform well in practice, though its theoretical convergence properties are still an area of active research. Under certain conditions, Adam converges to a stationary point, but the choice of hyperparameters $\beta_1, \beta_2, \eta,$, and $\epsilon$ plays a critical role in its performance.

Adam is widely used in deep learning for training complex models like transformers. Its ability to handle noisy gradients and adjust learning rates dynamically makes it particularly effective for optimizing deep neural networks.

**Learning Rate Schedules**

The learning rate is a critical hyperparameter in gradient-based optimization algorithms. The choice of learning rate can significantly impact the convergence speed and the quality of the solution. Learning rate schedules are strategies for adjusting the learning rate during training.

Types of Learning Rate Schedules:

1. Step Decay: The learning rate is reduced by a factor $\gamma$ after a fixed number of epochs $T$:

$$\eta_t = \eta_0 \cdot \gamma^{\lfloor t/T \rfloor},$$

where $\eta_0$ is the initial learning rate.

2. Exponential Decay: The learning rate decreases exponentially with time:

$$\eta_t = \eta_0 \cdot e^{-\lambda t},$$

where $\lambda$ controls the rate of decay.

3. Polynomial Decay: The learning rate decreases according to a polynomial function of the iteration:

$$\eta_t = \eta_0 \left( 1 - \frac{t}{T_{\max}} \right)^p ,$$

where $T_{\max}$ is the maximum number of iterations and $p$ is the power of the polynomial.

4. Cosine Annealing: The learning rate follows a cosine function:

$$\eta_t = \eta_{\min} + \frac{1}{2}(\eta_{\max} - \eta_{\min}) \left( 1 + \cos \left( \frac{\pi t}{T_{\max}} \right) \right) ,$$

where $\eta_{\min}$ and $\eta_{\max}$ are the minimum and maximum learning rates, respectively.

The learning rate schedule can have significant implications for the convergence behavior of the optimization algorithm. A well-chosen schedule can accelerate convergence, prevent the algorithm from getting stuck in poor local minima, and lead to better generalization.

In training transformers, a learning rate schedule is often combined with the Adam optimizer. The initial learning rate is typically set high to allow rapid progress, then gradually reduced to fine-tune the model as it approaches convergence.

**Remark (Convergence with Learning Rate Schedules):** Let $\eta_t$ be a learning rate schedule that satisfies the conditions:

$$\sum_{t=1}^{\infty} \eta_t = \infty \quad \text{and} \quad \sum_{t=1}^{\infty} \eta_t^2 < \infty.$$

Then, gradient-based optimization algorithms converge to a stationary point $\theta^*$ of the loss function $f(\theta)$.

## 1.10.2   Saddle Points and Local Minima

In non-convex optimization, which is common in training deep neural networks, the loss landscape is often highly complex, containing numerous saddle points, local minima, and potentially flat or sharp minima. Understanding the behavior of optimization algorithms in such landscapes is crucial for designing effective training strategies.

### Saddle Points and Local Minima

A saddle point in a function $f : \mathbb{R}^n \to \mathbb{R}$ is a critical point $\theta^*$ where the gradient $\nabla f(\theta^*) = 0$, but unlike a local minimum, the Hessian matrix $H(\theta^*)$ has both

positive and negative eigenvalues. This means that $f$ decreases in some directions and increases in others near $\theta^*$. Mathematically, if $\theta^*$ is a saddle point, there exist directions $v_1$ and $v_2$ such that

$$v_1^\top H(\theta^*)v_1 > 0 \quad \text{and} \quad v_2^\top H(\theta^*)v_2 < 0.$$

Saddle points are particularly problematic in optimization because standard gradient descent algorithms may get stuck or exhibit slow convergence near these points, especially in high-dimensional settings where saddle points are more prevalent.

Example: Consider the function $f(x, y) = x^2 - y^2$. The point $(0, 0)$ is a saddle point, as $\nabla f(0, 0) = 0$, but the Hessian matrix

$$H(0, 0) = \begin{pmatrix} 2 & 0 \\ 0 & -2 \end{pmatrix}$$

has eigenvalues $2$ and $-2$, indicating directions of both ascent and descent.

A local minimum $\theta^*$ of a function $f$ is a point where $f(\theta^*) \leq f(\theta)$ for all $\theta$ in some neighborhood of $\theta^*$. At a local minimum, the Hessian matrix $H(\theta^*)$ is positive semi-definite, meaning all its eigenvalues are non-negative:

$$v^\top H(\theta^*)v \geq 0 \quad \text{for all } v \in \mathbb{R}^n.$$

In non-convex optimization, the landscape often contains many local minima, and some of these minima may correspond to poor solutions, particularly when they are sharp.

**Remark (Second-Order Condition for Local Minima):** Let $f : \mathbb{R}^n \to \mathbb{R}$ be a twice-differentiable function. A point $\theta^*$ is a local minimum if:

1. $\nabla f(\theta^*) = 0$ (first-order condition).
2. $H(\theta^*)$ is positive semi-definite (second-order condition).

If $H(\theta^*)$ is positive definite, $\theta^*$ is a strict local minimum.

## Flat Versus Sharp Minima

Flat minima refer to regions of the loss landscape where the loss function is relatively constant over a large area, meaning the eigenvalues of the Hessian at these points are small. Optimization algorithms are generally more robust when they converge to flat minima, as these minima tend to generalize better to unseen data. Mathematically, a minimum $\theta^*$ is flat if the second derivative of the loss function is close to zero in most directions:

$$v^\top H(\theta^*)v \approx 0 \quad \text{for most directions } v.$$

Sharp minima, on the other hand, are characterized by a steep increase in the loss function when moving away from the minimum, indicating that the Hessian at these

points has large eigenvalues. Sharp minima often lead to poor generalization because the model is overly sensitive to small changes in the parameters.

The nature of the minima found by an optimization algorithm can significantly affect the model's performance. For instance, sharp minima are associated with overfitting, where the model captures noise in the training data rather than the underlying pattern. Flat minima, being more stable, typically correspond to models that generalize better.

In deep learning, regularization techniques such as dropout, weight decay, and data augmentation are used to encourage the optimization process to find flat minima, leading to models that are less likely to overfit and more likely to perform well on unseen data.

In high-dimensional optimization, models that converge to flat minima often exhibit better generalization properties. Formally, let $\theta^*$ and $\theta^{\#}$ be flat and sharp minima, respectively. If $H(\theta^*)$ has smaller eigenvalues than $H(\theta^{\#})$, the model's expected generalization error is lower at $\theta^*$ than at $\theta^{\#}$.

### 1.10.3  Convergence Analysis

Convergence analysis of optimization algorithms in non-convex settings is more challenging than in convex settings due to the presence of multiple local minima, saddle points, and other critical points [20, 33, 40, 47]. This section discusses the behavior of gradient-based methods in non-convex landscapes and the role of hyperparameters such as learning rates and momentum in ensuring convergence.

**Convergence in Non-convex Settings**

In non-convex optimization, the loss function $f(\theta)$ may have multiple local minima and saddle points, making it difficult to guarantee convergence to a global minimum. However, under certain conditions, gradient-based methods can still achieve convergence to a critical point that is a local minimum or a saddle point.

SGD is widely used in non-convex optimization due to its ability to escape saddle points, thanks to its inherent noise. The noise introduced by stochastic updates allows SGD to avoid getting stuck at saddle points and move toward regions of the loss landscape that contain local minima.

**Theorem 1.13** (Convergence of SGD in Non-Convex Settings) *Let $f(\theta)$ be a non-convex function with Lipschitz-continuous gradients. Assume that the learning rate $\eta_t$ satisfies $\sum_{t=1}^{\infty} \eta_t = \infty$ and $\sum_{t=1}^{\infty} \eta_t^2 < \infty$. Then SGD converges to a critical point $\theta^*$, which could be a local minimum or a saddle point:*

$$\lim_{t \to \infty} \mathbb{E}[\|\nabla f(\theta_t)\|^2] = 0.$$

This theorem indicates that, with an appropriate learning rate schedule, SGD converges to a point where the gradient is zero, though it does not guarantee that this point is a global minimum.

Escaping Saddle Points: Recent research has shown that SGD and other gradient-based methods can escape saddle points efficiently, particularly when the learning rate is adaptive or when momentum is used. The noise in SGD can push the algorithm away from saddle points, while momentum can help maintain the optimization trajectory, preventing it from stagnating near saddle points.

In training deep neural networks, it is common to observe that SGD does not converge to the lowest possible loss but instead hovers around a local minimum. This behavior is often beneficial, as it helps the model generalize better by avoiding sharp minima.

**Role of Learning Rates and Momentum**

The learning rate $\eta$ is a crucial hyperparameter that determines the step size of the optimization algorithm. If the learning rate is too large, the algorithm may oscillate around minima or diverge. If it is too small, the algorithm may converge very slowly or get stuck in a suboptimal region.

As discussed earlier, learning rate schedules involve adjusting the learning rate during training to balance exploration and exploitation. A common practice is to start with a high learning rate to make rapid progress and then gradually reduce it to fine-tune the solution.

Momentum is a technique used to accelerate gradient descent by adding a fraction of the previous update to the current update. This method helps to smooth the optimization trajectory and can speed up convergence, particularly in regions of the loss landscape that are shallow or contain small gradients.

The update rule with momentum is given by

$$v_t = \beta v_{t-1} + \nabla f(\theta_t),$$

$$\theta_{t+1} = \theta_t - \eta v_t,$$

where $v_t$ is the velocity, $\beta$ is the momentum coefficient (typically close to 1), and $\eta$ is the learning rate.

**Remark (Convergence with Momentum):** Under certain conditions, momentum-based gradient descent can achieve faster convergence than standard gradient descent. Specifically, let $\theta_t$ be the parameter at iteration $t$, and assume $f(\theta)$ is strongly convex with Lipschitz-continuous gradients. Then, gradient descent with momentum converges to the global minimum $\theta^*$ at a rate of

$$\|\theta_t - \theta^*\| \leq \mathcal{O}(\beta^t),$$

where $\beta$ is the momentum coefficient. This result shows that momentum can accelerate convergence by smoothing the trajectory and reducing oscillations.

In deep learning, the combination of momentum with an adaptive learning rate optimizer like Adam can lead to rapid and stable convergence. The momentum helps navigate through flat regions and escape saddle points, while the adaptive learning rate ensures that the step size is appropriate for the local geometry of the loss landscape.

## 1.11    Measure Theory and Information Theory

Measure theory and information theory provide the mathematical foundation for understanding and quantifying uncertainty, randomness, and information in various contexts [2, 14, 23, 37]. These concepts are essential in many areas of machine learning and data science, including probability theory, entropy, and inference methods. In this section, we explore these topics, focusing on their mathematical underpinnings and their relevance to the design and analysis of machine learning models, including transformers.

### 1.11.1    Basic Probability Concepts

Probability theory is a branch of mathematics that deals with the analysis of random phenomena. The mathematical framework for probability theory is built on measure theory, where probabilities are assigned to events in a measurable space.

A probability space $(\Omega, \mathcal{F}, \mathbb{P})$ consists of three components:

1. Sample Space ($\Omega$): The set of all possible outcomes of a random experiment.
2. Sigma-Algebra ($\mathcal{F}$): A collection of subsets of $\Omega$, called events, that is closed under complement and countable unions. This ensures that the probability of any event can be well defined.
3. Probability Measure ($\mathbb{P}$): A function $\mathbb{P} : \mathcal{F} \to [0, 1]$ that assigns a probability to each event in $\mathcal{F}$. The measure $\mathbb{P}$ satisfies the following axioms:

   (a) $\mathbb{P}(\Omega) = 1$.
   (b) For any countable sequence of disjoint events $\{A_i\} \subset \mathcal{F}$,

$$\mathbb{P}\left( \bigcup_{i=1}^{\infty} A_i \right) = \sum_{i=1}^{\infty} \mathbb{P}(A_i).$$

A random variable is a measurable function $X : \Omega \to \mathbb{R}$ that assigns a real number to each outcome in the sample space. The distribution of a random variable

is described by its probability distribution function $F_X(x)$ or probability density function $f_X(x)$ in the continuous case:

$$F_X(x) = \mathbb{P}(X \le x), \quad f_X(x) = \frac{dF_X(x)}{dx}.$$

The expectation $\mathbb{E}[X]$, variance $\text{Var}(X)$, and higher moments provide important characteristics of the distribution of $X$.

Example: Consider a random variable $X$ representing the outcome of a fair die roll. The sample space is $\Omega = \{1, 2, 3, 4, 5, 6\}$, and the probability measure assigns $\mathbb{P}(\{i\}) = \frac{1}{6}$ for each $i \in \Omega$. The expectation of $X$ is given by

$$\mathbb{E}[X] = \sum_{i=1}^{6} i \cdot \mathbb{P}(X = i) = \frac{1}{6} \cdot (1 + 2 + 3 + 4 + 5 + 6) = 3.5.$$

Two fundamental results in probability theory are the Law of Large Numbers (LLN) and the Central Limit Theorem (CLT):

**Theorem 1.14** (Law of Large Numbers) *If $X_1, X_2, \ldots$ are independent and identically distributed (i.i.d.) random variables with finite expectation $\mathbb{E}[X_i] = \mu$, then*

$$\frac{1}{n} \sum_{i=1}^{n} X_i \to \mu \quad \text{as } n \to \infty \text{ (almost surely).}$$

**Theorem 1.15** (Central Limit Theorem) *If $X_1, X_2, \ldots$ are i.i.d. random variables with finite mean $\mu$ and variance $\sigma^2$, then the normalized sum:*

$$\frac{1}{\sqrt{n}} \left( \sum_{i=1}^{n} X_i - n\mu \right) \to \mathcal{N}(0, \sigma^2) \quad \text{as } n \to \infty,$$

*where $\mathcal{N}(0, \sigma^2)$ is a normal distribution with mean 0 and variance $\sigma^2$.*

These results are critical in understanding the behavior of random variables and form the basis for many statistical methods.

**Entropy and Information Gain**

Entropy is a measure of uncertainty or randomness in a probability distribution. It is a fundamental concept in information theory, introduced by Claude Shannon, and is used to quantify the amount of information contained in a random variable.

Shannon Entropy: The entropy $H(X)$ of a discrete random variable $X$ with probability mass function $p(x) = \mathbb{P}(X = x)$ is defined as

$$H(X) = -\sum_{x \in \mathcal{X}} p(x) \log p(x),$$

where the logarithm is typically taken base 2 and the entropy is measured in bits. For a continuous random variable with probability density function $f(x)$, the differential entropy is defined as

$$h(X) = -\int_{\mathcal{X}} f(x) \log f(x) \, dx.$$

Entropy represents the average amount of information (in bits) produced by a random process. A higher entropy value indicates greater uncertainty or variability in the outcomes.

Example: For a fair coin toss, the entropy is

$$H(X) = -\left( \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right) = 1 \text{ bit.}$$

Information Gain: Information gain measures the reduction in entropy when moving from a prior probability distribution to a posterior distribution after observing some evidence. It is commonly used in decision tree algorithms in machine learning.

Let $X$ be a random variable representing the class labels, and $Y$ be a feature. The information gain $IG(X, Y)$ is defined as

$$IG(X, Y) = H(X) - H(X|Y),$$

where $H(X|Y)$ is the conditional entropy of $X$ given $Y$:

$$H(X|Y) = \sum_{y \in \mathcal{Y}} p(y) H(X|Y = y).$$

In a decision tree, the feature that provides the highest information gain is chosen for splitting the data, as it reduces the uncertainty (entropy) the most.

**Bayesian Inference**

Bayesian inference is a method of statistical inference in which Bayes' theorem is used to update the probability of a hypothesis as more evidence or information becomes available.

Bayes' theorem relates the conditional probability of an event $A$ given $B$ to the conditional probability of $B$ given $A$, the prior probability of $A$, and the prior probability of $B$:

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A)\mathbb{P}(A)}{\mathbb{P}(B)}.$$

In Bayesian inference, the prior probability $\mathbb{P}(A)$ represents our initial belief about the hypothesis $A$. The likelihood $\mathbb{P}(B|A)$ represents the probability of observing the evidence $B$ given that $A$ is true. The posterior probability $\mathbb{P}(A|B)$ is the updated probability of $A$ after observing $B$.

Example: Suppose a medical test for a disease has a 95% accuracy and 1% of the population has the disease. If a person tests positive, the posterior probability that the person has the disease can be calculated using Bayes' theorem.

**Bayesian Networks** Bayesian networks are graphical models that represent the probabilistic relationships among a set of variables. Each node represents a random variable, and the edges represent conditional dependencies. Bayesian networks are widely used in various applications, including decision-making, diagnostics, and machine learning.

Example: In a Bayesian network for a weather prediction system, nodes could represent variables like temperature, humidity, and wind speed, with edges representing the dependencies between them.

### KL Divergence and Cross-Entropy

Kullback–Leibler (KL) divergence is a measure of how one probability distribution diverges from a second, reference probability distribution. For two discrete probability distributions $P$ and $Q$ defined on the same probability space, the KL divergence $D_{KL}(P||Q)$ is given by

$$D_{KL}(P||Q) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)}.$$

For continuous distributions, the sum is replaced by an integral:

$$D_{KL}(P||Q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \, dx.$$

KL divergence is non-negative and is zero if and only if $P = Q$ almost everywhere. It is not symmetric, meaning $D_{KL}(P||Q) \neq D_{KL}(Q||P)$.

In machine learning, KL divergence is used in variational inference to measure the difference between the approximate posterior distribution and the true posterior distribution.

**Cross-Entropy** It is a related concept that measures the difference between two probability distributions for a given random variable or set of events. For two distributions $P$ and $Q$, the cross-entropy $H(P, Q)$ is defined as

$$H(P, Q) = - \sum_{x \in \mathcal{X}} P(x) \log Q(x).$$

Cross-entropy can be decomposed into the entropy of $P$ and the KL divergence between $P$ and $Q$:

$$H(P, Q) = H(P) + D_{KL}(P||Q).$$

In machine learning, cross-entropy is often used as a loss function for classification tasks. The goal is to minimize the cross-entropy between the true labels (distribution $P$) and the predicted probabilities (distribution $Q$).

Example: In binary classification, if the true label is $y \in \{0, 1\}$ and the predicted probability of the positive class is $\hat{y}$, the cross-entropy loss is

$$\text{Loss} = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})].$$

This loss function penalizes predictions that deviate from the true labels, driving the model to improve its predictions.

**Theorem 1.16** (Information Inequality) *The information inequality, also known as Gibbs' inequality, states that the KL divergence between two distributions $P$ and $Q$ is always non-negative:*

$$D_{KL}(P||Q) \geq 0,$$

*with equality if and only if $P = Q$ almost everywhere.*

This theorem underpins the use of KL divergence as a measure of how one distribution approximates another.

### 1.11.2 Foundations of Measure Theory

Measure theory is a fundamental branch of mathematics that extends the notion of length, area, and volume to more complex sets, providing a foundation for integration and probability [18, 26, 50, 51]. This chapter delves into the core concepts of measure theory, focusing on the Lebesgue measure and integration, as well as measure-preserving transformations. These concepts are not only central to analysis but also play a crucial role in probability theory, information theory, and various fields where continuous distributions are studied.

**Lebesgue Measure and Integration**

The Lebesgue measure is a mathematical construct that generalizes the notion of length, area, and volume to more complex sets in $\mathbb{R}^n$. Unlike the Riemann integral, which relies on partitioning the domain of a function into intervals, the Lebesgue measure allows for the measurement of more irregular sets by considering their "size" in a more general sense. Formally, a measure $\mu$ on a set $X$ is a function that assigns

a non-negative real number or $+\infty$ to each measurable subset of $X$, satisfying the following properties:

1. Non-negativity: $\mu(A) \geq 0$ for all measurable sets $A \subseteq X$.
2. Null empty set: $\mu(\emptyset) = 0$.
3. Countable additivity ($\sigma$-additivity): For any countable collection $\{A_i\}$ of disjoint measurable sets,

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i).$$

The Lebesgue measure $\lambda^n$ on $\mathbb{R}^n$ is the standard measure that extends the concept of length in $\mathbb{R}^1$, area in $\mathbb{R}^2$, and volume in $\mathbb{R}^3$ to higher dimensions. For a set $A \subseteq \mathbb{R}^n$, the Lebesgue measure $\lambda^n(A)$ is defined in such a way that it coincides with the usual notions of length, area, or volume for simple geometric shapes.

Example: Consider the interval $A = [0, 1] \subset \mathbb{R}$. The Lebesgue measure $\lambda^1(A)$ of this interval is simply its length, which is $1 - 0 = 1$.

**Lebesgue Integration** It is a method of integration that, unlike Riemann integration, focuses on measuring the "size" of the set where the function takes on certain values rather than partitioning the domain into intervals. This approach allows for the integration of a broader class of functions, particularly those with discontinuities or unbounded intervals. Let $(X, \mathcal{F}, \mu)$ be a measure space, where $X$ is a set, $\mathcal{F}$ is a $\sigma$-algebra of subsets of $X$, and $\mu$ is a measure on $\mathcal{F}$. A function $f : X \to \mathbb{R}$ is called measurable if for every Borel set $B \subseteq \mathbb{R}$, the preimage $f^{-1}(B)$ is in $\mathcal{F}$.

The Lebesgue integral of a non-negative measurable function $f : X \to [0, \infty]$ with respect to the measure $\mu$ is defined as

$$\int_X f \, d\mu = \sup\left\{\int_X g \, d\mu \mid g \text{ is a simple function}, 0 \leq g \leq f\right\},$$

where a simple function $g : X \to \mathbb{R}$ is a finite linear combination of indicator functions of measurable sets:

$$g(x) = \sum_{i=1}^{n} a_i \mathbf{1}_{A_i}(x),$$

with $a_i \geq 0$ and $A_i \in \mathcal{F}$.

For general integrable functions $f : X \to \mathbb{R}$, the integral is defined by decomposing $f$ into its positive and negative parts:

$$f = f^+ - f^-, \quad \text{where } f^+ = \max(f, 0) \text{ and } f^- = \max(-f, 0),$$

and then integrating each part separately:

$$\int_X f \, d\mu = \int_X f^+ \, d\mu - \int_X f^- \, d\mu,$$

provided that at least one of these integrals is finite.

**Dominated Convergence Theorem** One of the central results in Lebesgue integration is the dominated convergence theorem, which provides conditions under which the limit of a sequence of integrable functions can be interchanged with the integral.

**Theorem 1.17**  (Dominated Convergence Theorem) *Let* $\{f_n\}$ *be a sequence of measurable functions such that* $f_n \to f$ *almost everywhere, and suppose there exists an integrable function* $g$ *such that* $|f_n(x)| \leq g(x)$ *for all n and almost every x. Then*

$$\lim_{n \to \infty} \int_X f_n \, d\mu = \int_X \lim_{n \to \infty} f_n \, d\mu = \int_X f \, d\mu.$$

This theorem is particularly important in the context of probability theory and stochastic processes, where it allows for the exchange of limits and integrals under certain conditions.

Example: Consider the sequence of functions $f_n(x) = x^n$ on the interval $[0, 1]$. As $n$ increases, $f_n(x)$ converges pointwise to the function $f(x) = 0$ for $x \in [0, 1)$ and $f(1) = 1$. Since $f_n(x) \leq 1$ for all $x \in [0, 1]$, the dominated convergence theorem allows us to conclude that

$$\lim_{n \to \infty} \int_0^1 x^n \, dx = \int_0^1 \lim_{n \to \infty} x^n \, dx = \int_0^1 0 \, dx = 0.$$

**Measure-Preserving Transformations**

Measure-Preserving Transformations: A transformation $T : X \to X$ on a measure space $(X, \mathcal{F}, \mu)$ is said to be measure-preserving if, for every measurable set $A \in \mathcal{F}$, the measure of $A$ is equal to the measure of its image under $T$:

$$\mu(T^{-1}(A)) = \mu(A).$$

Measure-preserving transformations are crucial in the study of dynamical systems, ergodic theory, and probability theory, as they describe systems where the total measure (interpreted as total probability, volume, etc.) is conserved over time.

Example: A simple example of a measure-preserving transformation is the rotation of the unit circle $S^1$. Consider the circle as the interval $[0, 1)$ with endpoints identified, equipped with the Lebesgue measure. The rotation by a fixed angle $\alpha$, given by $T(x) = (x + \alpha) \mod 1$, preserves the Lebesgue measure because the "length" of any interval on the circle is unchanged by the rotation.

In ergodic theory, measure-preserving transformations are studied to understand the long-term behavior of dynamical systems. A measure-preserving transformation $T$ is called ergodic if any $T$-invariant set (a set $A$ such that $T^{-1}(A) = A$) has measure 0 or 1. This means that the system, when observed over a long time, cannot be decomposed into simpler invariant subsystems.

**Theorem 1.18** (Birkhoff's Ergodic Theorem) *Let $T : X \to X$ be a measure-preserving transformation on a probability space $(X, \mathcal{F}, \mu)$, and let $f : X \to \mathbb{R}$ be an integrable function. Then, the time average of $f$ along the orbits of $T$ converges almost everywhere to the space average (expected value):*

$$\lim_{N \to \infty} \frac{1}{N} \sum_{n=0}^{N-1} f(T^n(x)) = \int_X f \, d\mu \quad \text{for } \mu\text{-almost every } x \in X.$$

This theorem provides a link between the temporal behavior of a system and its spatial properties, and is foundational in the study of statistical mechanics and information theory.

Measure-preserving transformations and ergodic theory have applications in various fields, including statistical mechanics, where they are used to justify the assumption that time averages can be replaced by ensemble averages. In information theory, these concepts underpin the analysis of data compression algorithms and the behavior of stochastic processes over time.

Example: In the analysis of random walks, measure-preserving transformations are used to study the recurrence properties of the walk, determining whether it returns to a given state with probability 1 or eventually escapes to infinity.

### 1.11.3   Mutual Information

Mutual information measures the amount of information that one random variable contains about another. It quantifies the reduction in uncertainty about one variable given knowledge of the other. Mathematically, the mutual information $I(X; Y)$ between two random variables $X$ and $Y$ is defined as

$$I(X; Y) = \int_{X \times Y} p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right) dx \, dy,$$

where $p(x, y)$ is the joint probability density function of $X$ and $Y$, and $p(x)$ and $p(y)$ are the marginal probability density functions of $X$ and $Y$, respectively.

Alternatively, mutual information can be expressed in terms of entropy:

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X),$$

where $H(X)$ and $H(Y)$ are the entropies of $X$ and $Y$, and $H(X|Y)$ and $H(Y|X)$ are the conditional entropies.

Properties of mutual information:

1. Non-negativity: $I(X; Y) \geq 0$, with equality if and only if $X$ and $Y$ are independent.
2. Symmetry: $I(X; Y) = I(Y; X)$.
3. Data Processing Inequality: If $X \rightarrow Y \rightarrow Z$ forms a Markov chain, then $I(X; Z) \leq I(X; Y)$.

Example: Consider two binary random variables $X$ and $Y$, each taking values in $\{0, 1\}$. If $X$ and $Y$ are perfectly correlated (i.e., $X = Y$), the mutual information $I(X; Y)$ is maximized, indicating that knowing $X$ completely determines $Y$. Conversely, if $X$ and $Y$ are independent, $I(X; Y) = 0$, reflecting no mutual dependence.

**Theorem 1.19** (Chain Rule for Mutual Information) *Let $X, Y, Z$ be random variables. The mutual information satisfies the following chain rule:*

$$I(X; Y, Z) = I(X; Y) + I(X; Z|Y).$$

This theorem is useful for decomposing the mutual information between a set of variables into more manageable components, particularly in analyzing complex models like transformers.

**Applications to Transformer Models**

In transformer models, mutual information plays a crucial role in understanding how information is distributed and processed across different layers and attention heads. It can be used to measure the amount of information that the model retains about the input as it propagates through the layers, providing insights into the model's ability to capture relevant dependencies.

The Information Bottleneck (IB) principle, introduced by [57], provides a framework for understanding the trade-off between compression and relevance in machine learning models. The IB principle suggests that a good representation $Z$ of the input $X$ should maximize the mutual information $I(Z; Y)$ with the output $Y$, while minimizing the mutual information $I(Z; X)$ with the input $X$:

$$\min I(Z; X) \quad \text{subject to } I(Z; Y) \geq \text{constant.}$$

In transformers, this principle can be applied to analyze how different layers balance the retention of useful information versus the compression of irrelevant details.

The attention mechanism in transformers can be analyzed using mutual information to quantify how much information from the input sequence is captured by each

attention head. By measuring the mutual information between the input tokens and the output of each attention layer, one can identify which attention heads are most effective in capturing long-range dependencies or specific patterns in the data.

**Remark (Mutual Information in Attention Layers):** Let $X = \{x_1, x_2, \ldots, x_n\}$ be an input sequence and $Y = \{y_1, y_2, \ldots, y_m\}$ be the output sequence produced by an attention layer. The mutual information between $X$ and $Y$ is given by

$$I(X; Y) = \sum_{i=1}^{n} \sum_{j=1}^{m} I(x_i; y_j | \text{context}),$$

where context refers to the set of all other tokens considered by the attention mechanism. This decomposition allows for a detailed analysis of how information is distributed across different tokens and attention heads.

Example: In natural language processing tasks, one might measure the mutual information between specific words in the input sentence and the output representation to understand how well the transformer captures syntactic or semantic relationships.

### 1.11.4   Complexity and Generalization

The complexity of a model refers to its ability to fit a wide range of functions or data patterns. A key challenge in model design is to balance complexity with generalization, ensuring that the model performs well not only on the training data but also on unseen data. This section explores the concept of Kolmogorov complexity and how generalization bounds can be established for machine learning models, including transformers.

**Kolmogorov Complexity**

Kolmogorov complexity, also known as algorithmic complexity, measures the complexity of an object (such as a string or a function) by the length of the shortest program that can produce that object on a universal Turing machine. Formally, the Kolmogorov complexity $K(x)$ of a string $x$ is defined as

$$K(x) = \min\{|p| : U(p) = x\},$$

where $U$ is a fixed universal Turing machine, $p$ is a program that generates $x$, and $|p|$ denotes the length of $p$ in bits.

Kolmogorov complexity provides a formal measure of the "compressibility" of an object. If $x$ can be generated by a short program, it has low complexity; if the

shortest program is nearly as long as $x$ itself, the object is incompressible and has high complexity.

Properties of Kolmogorov complexity:

1. Incomputability: Kolmogorov complexity is not computable in general, meaning there is no algorithm that can determine the exact complexity of any given string.
2. Invariance: The complexity $K_U(x)$ depends on the choice of the universal Turing machine $U$, but the difference between complexities for different choices of $U$ is bounded by a constant independent of $x$.
3. Prefix Complexity: A variant of Kolmogorov complexity, called prefix complexity $K_{prefix}(x)$, considers only prefix-free programs (where no program is a prefix of another), and is used in areas such as information theory and coding.

Example: The string $x = 01010101\ldots01$ (with length $n$) has low Kolmogorov complexity because it can be generated by a short program that outputs "01" repeatedly. In contrast, a random string of the same length would likely have high complexity, as it cannot be compressed into a shorter program.

Kolmogorov complexity provides a theoretical foundation for understanding the trade-off between model complexity and generalization. Models that are too complex may fit the training data perfectly but fail to generalize to new data (overfitting), while models that are too simple may underfit the data, missing important patterns.

**Remark (Kolmogorov Complexity and Overfitting):** Let $\mathcal{H}$ be a hypothesis class and $h \in \mathcal{H}$ be a hypothesis selected by a learning algorithm based on the training data $D$. If the Kolmogorov complexity $K(h)$ is high relative to the complexity of the data $K(D)$, the hypothesis $h$ is more likely to overfit the data. Conversely, if $K(h)$ is low, the hypothesis is more likely to generalize well.


**Generalization Bounds**

Generalization refers to the ability of a model to perform well on unseen data, not just on the data it was trained on. Generalization bounds provide theoretical guarantees on how well a model trained on a finite sample will perform on the overall data distribution.

In the Probably Approximately Correct (PAC) learning framework, we seek to understand the conditions under which a learning algorithm can produce a hypothesis that, with high probability, performs close to optimally on the underlying distribution. Let $\mathcal{H}$ be a hypothesis class, $D$ be the training data drawn from a distribution $\mathcal{D}$, and $\epsilon$ be the allowable error. A hypothesis $h$ is said to be $\epsilon$-good if:

$$\mathbb{P}_{x \sim \mathcal{D}}[h(x) \neq f(x)] \leq \epsilon,$$

where $f(x)$ is the true underlying function. The PAC learning framework provides bounds on the sample size $m$ needed to guarantee that the learning algorithm produces an $\epsilon$-good hypothesis with high probability.

The Vapnik–Chervonenkis (VC) dimension is a measure of the capacity or complexity of a hypothesis class $\mathcal{H}$. It is defined as the size of the largest set of points that can be shattered by $\mathcal{H}$, where "shatter" means that for every possible labeling of the points, there exists a hypothesis in $\mathcal{H}$ that correctly classifies them.

**Theorem 1.20** (Generalization Bound via VC Dimension) *Let $\mathcal{H}$ be a hypothesis class with VC dimension $d_{VC}$. Then, for any $\epsilon > 0$ and $\delta > 0$, with probability at least $1 - \delta$, the generalization error of a hypothesis h learned from m samples is bounded by*

$$\mathbb{P}_{x \sim \mathcal{D}}[h(x) \neq f(x)] \leq \hat{\epsilon} + O\left(\sqrt{\frac{d_{VC} \log(m/d_{VC}) + \log(1/\delta)}{m}}\right),$$

*where $\hat{\epsilon}$ is the empirical error on the training set. This bound indicates that the generalization error decreases as the number of samples increases, and it depends on the complexity of the hypothesis class as measured by the VC dimension.*

In transformers, the generalization ability is influenced by both the architecture's capacity (e.g., number of layers, attention heads) and the regularization techniques used during training. The VC dimension or similar complexity measures can be used to understand the trade-offs involved in designing transformer models that generalize well to new data. For a transformer model, one might analyze the VC dimension of the attention mechanism or the overall architecture to predict how well the model is likely to generalize based on the amount of training data available.

## 1.12 Backpropagation and Autodiff

Backpropagation is a cornerstone of modern deep learning, providing an efficient way to compute gradients of loss functions with respect to model parameters. These gradients are crucial for training neural networks using optimization algorithms like gradient descent [7, 22, 39, 52]. In this chapter, we explore the mathematical foundations of backpropagation, focusing on the chain rule in matrix calculus and the propagation of errors through neural networks.

### 1.12.1 Backpropagation

Backpropagation is essentially an application of the chain rule of calculus, extended to functions represented by neural networks. It allows for the computation of gradients by systematically applying the chain rule from the output layer back to the input layer, hence the name "backpropagation."

**Chain Rule in Matrix Calculus**

The chain rule is a fundamental tool in calculus, used to differentiate composite functions. In the context of neural networks, where functions are often represented as compositions of multiple layers, the chain rule is applied repeatedly to compute the derivative of the loss function with respect to each parameter. Let $f : \mathbb{R}^m \to \mathbb{R}^n$ and $g : \mathbb{R}^n \to \mathbb{R}^p$ be differentiable functions. The composite function $h(x) = g(f(x))$ maps $\mathbb{R}^m$ to $\mathbb{R}^p$. The derivative (or Jacobian matrix) of $h(x)$ with respect to $x$ is given by the chain rule:

$$\frac{\partial h}{\partial x} = \frac{\partial g}{\partial f} \cdot \frac{\partial f}{\partial x},$$

where $\frac{\partial g}{\partial f}$ is the Jacobian matrix of $g$ with respect to $f(x)$ and $\frac{\partial f}{\partial x}$ is the Jacobian matrix of $f$ with respect to $x$.

Example: Consider a neural network with a single hidden layer. Let the input $x \in \mathbb{R}^m$ be mapped to a hidden representation $h$ via a weight matrix $W_1$ and an activation function $\sigma$:

$$h = \sigma(W_1 x).$$

The output $y$ is then computed as

$$y = W_2 h = W_2 \sigma(W_1 x).$$

If the loss function $L(y, \hat{y})$ measures the discrepancy between the predicted output $y$ and the true output $\hat{y}$, the derivative of the loss with respect to $W_1$ involves applying the chain rule:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h} \cdot \frac{\partial h}{\partial W_1}.$$

Here, $\frac{\partial L}{\partial y}$ is the gradient of the loss with respect to the output, $\frac{\partial y}{\partial h} = W_2^\top$, and $\frac{\partial h}{\partial W_1}$ involves differentiating through the activation function $\sigma$.

Matrix Chain Rule: When dealing with matrices, the chain rule is extended to account for the dimensions and interactions of matrix operations. Suppose $A = f(X)$ and $B = g(A)$ where $X, A, B$ are matrices and $f, g$ are matrix functions. The derivative of $B$ with respect to $X$ is computed using:

$$\frac{\partial B}{\partial X} = \frac{\partial B}{\partial A} \cdot \frac{\partial A}{\partial X},$$

where $\frac{\partial A}{\partial X}$ and $\frac{\partial B}{\partial A}$ are tensor derivatives that account for the multi-dimensional nature of matrix operations.

**Theorem 1.21** (Chain Rule for Scalar-Valued Functions) *Let $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}$ be differentiable functions. The gradient of the composite function $h(x) = g(f(x))$ with respect to $x$ is given by*

$$\nabla h(x) = \nabla f(x)^\top \cdot \nabla g(f(x)),$$

*where $\nabla f(x)$ is the Jacobian matrix of $f$ and $\nabla g(f(x))$ is the gradient of $g$ with respect to $f(x)$.*

The chain rule is particularly powerful in neural networks because it allows the gradient of the loss function with respect to any parameter to be computed efficiently by propagating gradients backward through the network. This efficiency is crucial given the large number of parameters typically involved in deep learning models.

**Error Propagation in Neural Networks**

In backpropagation, errors are propagated backward through the network to update the weights in a way that minimizes the loss function. This process involves computing the gradient of the loss with respect to each weight in the network, starting from the output layer and moving backward to the input layer.

Forward Pass: During the forward pass, the input data $x$ is passed through the network layer by layer, and the output $y$ is computed. Each layer applies a linear transformation followed by a non-linear activation function:

$$h^{(l)} = \sigma(W^{(l)} h^{(l-1)} + b^{(l)}),$$

where $h^{(l)}$ is the output of the $l$th layer, $W^{(l)}$ and $b^{(l)}$ are the weight matrix and bias vector for that layer, and $\sigma$ is the activation function.

Backward Pass: In the backward pass, the error at the output layer is computed as the gradient of the loss function with respect to the output. This error is then propagated backward through the network to update the weights. Let $L$ be the loss function, and let $\delta^{(l)}$ denote the error term at layer $l$, defined as the gradient of the loss with respect to the linear combination of inputs to that layer:

$$\delta^{(l)} = \frac{\partial L}{\partial z^{(l)}} = \frac{\partial L}{\partial h^{(l)}} \cdot \sigma'(z^{(l)}),$$

where $z^{(l)} = W^{(l)} h^{(l-1)} + b^{(l)}$ is the input to the activation function at layer $l$.

The weight updates are then computed using:

$$\frac{\partial L}{\partial W^{(l)}} = \delta^{(l)} \cdot (h^{(l-1)})^\top,$$

$$\frac{\partial L}{\partial b^{(l)}} = \delta^{(l)}.$$

Recursive Error Propagation: The error $\delta^{(l)}$ at layer $l$ can be expressed recursively in terms of the error at the subsequent layer $\delta^{(l+1)}$:

$$\delta^{(l)} = (W^{(l+1)})^\top \delta^{(l+1)} \cdot \sigma'(z^{(l)}).$$

This recursion allows the error to be efficiently propagated from the output layer back to the input layer.

Remark: For a neural network with $L$ layers, the gradient of the loss function $L$ with respect to the weights $W^{(l)}$ in the $l$th layer is given by

$$\frac{\partial L}{\partial W^{(l)}} = \delta^{(l)} \cdot (h^{(l-1)})^\top,$$

where $\delta^{(l)} = \frac{\partial L}{\partial z^{(l)}}$ is the backpropagated error and $h^{(l-1)}$ is the output of the previous layer.

Example: Consider a simple neural network with a single hidden layer. The network output is $y = \sigma(W_2 \sigma(W_1 x))$, and the loss function is the mean squared error $L = \frac{1}{2}(y - \hat{y})^2$. During backpropagation, the error is computed at the output layer and propagated backward to update the weights $W_1$ and $W_2$.

If $\sigma$ is the ReLU activation function, then the gradient of the loss with respect to $W_2$ is

$$\frac{\partial L}{\partial W_2} = (y - \hat{y}) \cdot \sigma(W_1 x)^\top,$$

and the gradient with respect to $W_1$ is

$$\frac{\partial L}{\partial W_1} = ((y - \hat{y}) \cdot W_2 \cdot \sigma'(W_1 x)) \cdot x^\top.$$

### 1.12.2   Automatic Differentiation

Automatic differentiation (autodiff) is a computational technique that efficiently computes derivatives of functions specified by computer programs. Unlike symbolic differentiation, which involves manipulating mathematical expressions, or numerical differentiation, which approximates derivatives using finite differences, autodiff provides exact derivatives with minimal computational overhead. This is particularly important in deep learning, where models like transformers require the computation of gradients for optimization.

**Forward and Reverse Mode Differentiation**

In forward mode automatic differentiation, the derivative of each operation is computed simultaneously as the operation is performed, propagating derivatives from the inputs to the output. Let $f : \mathbb{R}^n \to \mathbb{R}^m$ be a function composed of intermediate

variables $z_1, z_2, \ldots, z_k$, where each $z_i$ depends on a subset of the previous variables. Forward mode calculates the derivative of the output with respect to each input by applying the chain rule in a straightforward manner.

Mathematically, for each intermediate variable $z_i$, the derivative $\dot{z}_i$ with respect to a specific input $x_j$ is computed as

$$\dot{z}_i = \sum_{k \in \text{parents}(z_i)} \frac{\partial z_i}{\partial z_k} \dot{z}_k,$$

where $\dot{z}_k$ is the derivative of the parent variable $z_k$ with respect to $x_j$.

Example: Consider the function $f(x_1, x_2) = \sin(x_1) + x_1 x_2^2$. To compute the derivative of $f$ with respect to $x_1$ using forward mode, we start by defining the intermediate variables:

$$z_1 = \sin(x_1), \quad z_2 = x_1 x_2^2.$$

The derivative with respect to $x_1$ is

$$\frac{df}{dx_1} = \frac{dz_1}{dx_1} + \frac{dz_2}{dx_1} = \cos(x_1) + x_2^2.$$

Reverse mode automatic differentiation, commonly used in backpropagation, computes the derivative of the output with respect to each input by propagating derivatives backward from the output to the inputs. This mode is particularly efficient when the function has many inputs and a single output, as it computes the gradient with respect to all inputs in a single backward pass. Mathematically, reverse mode computes the derivative $\bar{z}_i$ of the output with respect to each intermediate variable $z_i$ using the chain rule in reverse:

$$\bar{z}_i = \sum_{k \in \text{children}(z_i)} \bar{z}_k \frac{\partial z_k}{\partial z_i},$$

where $\bar{z}_k$ is the derivative of the output with respect to the child variable $z_k$.

Example: For the same function $f(x_1, x_2) = \sin(x_1) + x_1 x_2^2$, reverse mode starts by computing the derivative of $f$ with respect to the output and then propagates it backward to each input:

$$\bar{z}_1 = 1, \quad \bar{z}_2 = 1.$$

Then, the derivatives with respect to the inputs are

$$\frac{df}{dx_1} = \bar{z}_1 \cdot \frac{dz_1}{dx_1} + \bar{z}_2 \cdot \frac{dz_2}{dx_1} = \cos(x_1) + x_2^2,$$

$$\frac{df}{dx_2} = \bar{z}_2 \cdot \frac{dz_2}{dx_2} = 2x_1 x_2.$$

Let $f : \mathbb{R}^n \to \mathbb{R}$ be a differentiable function. The reverse mode of automatic differentiation computes the gradient $\nabla f(x)$ with respect to all inputs $x \in \mathbb{R}^n$ with a computational cost proportional to that of evaluating $f(x)$. Specifically, the time complexity of reverse mode is $\mathcal{O}(n)$, making it highly efficient for scalar-valued functions with many inputs.

In transformer models, which typically involve complex compositions of linear transformations, attention mechanisms, and non-linear activations, reverse mode automatic differentiation is used to compute gradients during backpropagation. This allows for efficient optimization of model parameters using gradient-based methods.

For example, in the multi-head self-attention mechanism, the output is a weighted sum of value vectors, where the weights are computed based on the similarity between queries and keys. To optimize this mechanism, gradients of the loss function with respect to the query, key, and value matrices are required. Reverse mode autodiff efficiently computes these gradients by propagating errors back through the layers and attention heads.

Let $\mathcal{L}$ be the loss function for a transformer model, and let $\Theta = \{W_Q, W_K, W_V, W_O, \ldots\}$ represent the set of trainable parameters. The gradient of $\mathcal{L}$ with respect to each parameter $\theta \in \Theta$ is computed using reverse mode autodiff as

$$\nabla_\theta \mathcal{L} = \frac{\partial \mathcal{L}}{\partial y} \cdot \frac{\partial y}{\partial \theta},$$

where $y$ is the model output. The efficiency of this computation is critical for training large-scale transformer models with millions or billions of parameters.

### 1.12.3  Optimization Challenges in Transformers

Training deep neural networks, including transformers, presents significant optimization challenges. Among these are the issues of vanishing and exploding gradients, which can hinder the training process, especially in very deep networks [5, 28, 45]. Various techniques, such as gradient clipping, have been developed to mitigate these issues and ensure stable training.

**Vanishing and Exploding Gradients**

The vanishing gradient problem occurs when the gradients of the loss function with respect to the model parameters become very small during backpropagation, effectively halting the learning process. This problem is particularly prevalent in deep networks, where gradients must be propagated through many layers.

Consider a deep neural network where the activation function is sigmoid $\sigma(x) = \frac{1}{1+e^{-x}}$. The derivative of the sigmoid function is

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)).$$

Since $\sigma(x)$ is bounded between 0 and 1, its derivative is also bounded, and for most values of $x$, $\sigma'(x)$ is close to zero. During backpropagation, the gradient at each layer is multiplied by the derivative of the activation function from the previous layer. In deep networks, this repeated multiplication can lead to an exponential decay of the gradient, making it nearly zero by the time it reaches the earlier layers:

$$\frac{\partial \mathcal{L}}{\partial W^{(l)}} \approx \prod_{k=l}^{L} \sigma'(z^{(k)}) \cdot \frac{\partial \mathcal{L}}{\partial z^{(L)}},$$

where $L$ is the total number of layers. If $\sigma'(z^{(k)})$ is small, the product can approach zero, leading to vanishing gradients.

Conversely, the exploding gradient problem arises when the gradients grow exponentially during backpropagation, leading to numerical instability. This typically happens when the weights of the network are initialized poorly or when the loss landscape has steep regions. In the case of exploding gradients, the derivative of the loss with respect to the weights can grow exponentially if the derivatives of the activation functions are large or if the weight matrices have large eigenvalues. This causes the gradients to increase exponentially as they are propagated backward:

$$\frac{\partial \mathcal{L}}{\partial W^{(l)}} \approx \prod_{k=l}^{L} \sigma'(z^{(k)}) \cdot \frac{\partial \mathcal{L}}{\partial z^{(L)}}.$$

If the factors $\sigma'(z^{(k)})$ or the weight matrices are large, the product can lead to very large gradients, causing the learning process to diverge.

Let $f(x) = W_L \sigma(W_{L-1} \ldots \sigma(W_1 x) \ldots)$ be a deep neural network. The norm of the gradient of the loss $\mathcal{L}$ with respect to the input $x$ satisfies

$$\|\nabla_x \mathcal{L}\| = \|\prod_{l=1}^{L} W_l^\top \sigma'(z^{(l)}) \cdot \nabla_y \mathcal{L}\|,$$

where $y = f(x)$ and $\sigma'(z^{(l)})$ denotes the derivative of the activation function at layer $l$. If $\|W_l\| > 1$ for most layers, the norm of the gradient can grow exponentially, leading to exploding gradients. Conversely, if $\|W_l\| < 1$, the gradient can vanish.

**Gradient Clipping and Other Techniques**

Gradient clipping is a technique used to prevent exploding gradients by capping the gradients at a predefined threshold. This is particularly useful in very deep networks or recurrent neural networks where the gradients can become excessively large. Let

$g$ be the gradient of the loss with respect to a parameter $\theta$. In gradient clipping, the gradient is rescaled if its norm exceeds a certain threshold $c$:

$$\tilde{g} = \begin{cases} g & \text{if } \|g\| \leq c, \\ c \cdot \frac{g}{\|g\|} & \text{if } \|g\| > c. \end{cases}$$

This ensures that the gradients remain within a manageable range, preventing numerical instability during training.

Example: Consider a transformer model with gradients $\nabla_\theta \mathcal{L}$ for each parameter $\theta$. If the norm of the gradient $\|\nabla_\theta \mathcal{L}\|$ exceeds a threshold $c$, gradient clipping rescales the gradient:

$$\tilde{\nabla_\theta} \mathcal{L} = \frac{c}{\|\nabla_\theta \mathcal{L}\|} \cdot \nabla_\theta \mathcal{L},$$

ensuring that the gradient update remains stable.

Other Techniques:

1. Weight Initialization: Proper weight initialization can mitigate the vanishing and exploding gradient problems. Techniques like Xavier initialization and He initialization are designed to keep the variance of the activations and gradients stable across layers.

Xavier Initialization: For a layer with $n_{\text{in}}$ input units, Xavier initialization sets the weights $W$ to be drawn from a uniform distribution:

$$W \sim \mathcal{U}\left(-\frac{\sqrt{6}}{\sqrt{n_{\text{in}} + n_{\text{out}}}}, \frac{\sqrt{6}}{\sqrt{n_{\text{in}} + n_{\text{out}}}}\right),$$

where $n_{\text{out}}$ is the number of output units. This ensures that the variance of the activations remains constant across layers.

2. Batch Normalization: Batch normalization normalizes the activations within each mini-batch to have zero mean and unit variance. This helps to mitigate the vanishing and exploding gradient problems by ensuring that the input to each layer remains stable during training. Given a mini-batch $\{x_1, x_2, \ldots, x_m\}$, batch normalization transforms each activation $x_i$ as

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}},$$

where $\mu_B$ and $\sigma_B^2$ are the mean and variance of the mini-batch, and $\epsilon$ is a small constant for numerical stability.

## 1.13 Statistical Learning Theory

Statistical learning theory provides a foundation for understanding and analyzing the performance of machine learning models. It addresses questions of generalization, capacity control, and risk minimization, offering theoretical tools to ensure that models not only fit the training data but also perform well on unseen data [42, 58]. In this chapter, we delve into the foundational concepts of statistical learning theory, including risk minimization, VC dimension, and Rademacher complexity, with a focus on their applications to modern machine learning models such as transformers.

### *1.13.1 Foundations of Statistical Learning*

**Risk Minimization**

In statistical learning theory, the goal of a learning algorithm is to find a hypothesis $h$ from a hypothesis class $\mathcal{H}$ that minimizes the expected risk, also known as the true risk or population risk. The true risk $R(h)$ of a hypothesis $h$ is defined as the expected loss over the distribution of the data:

$$R(h) = \mathbb{E}_{(x,y)\sim\mathcal{D}}[\ell(h(x), y)],$$

where $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ is the loss function, $(x, y)$ is a data point drawn from the distribution $\mathcal{D}$, and $h(x)$ is the prediction made by the hypothesis $h$.

Empirical Risk Minimization (ERM): Since the true distribution $\mathcal{D}$ is typically unknown, the true risk cannot be directly minimized. Instead, we minimize the empirical risk $\hat{R}_n(h)$, which is the average loss over the training sample $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$:

$$\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^{n} \ell(h(x_i), y_i).$$

The empirical risk serves as an approximation to the true risk, and the hypothesis that minimizes the empirical risk is chosen as the model. The ERM principle seeks to find:

$$h_{\text{ERM}} = \arg\min_{h\in\mathcal{H}} \hat{R}_n(h).$$

**Remark (Uniform Convergence):** For ERM to be effective, we require that the empirical risk $\hat{R}_n(h)$ converges uniformly to the true risk $R(h)$ as the sample size $n$ increases. Formally, with high probability, the following bound holds for all $h \in \mathcal{H}$:

$$\sup_{h\in\mathcal{H}} |R(h) - \hat{R}_n(h)| \le \epsilon(n),$$

where $\epsilon(n)$ decreases as $n$ increases, reflecting the idea that with more data, the empirical risk becomes a better approximation of the true risk.

Example: Consider a binary classification problem where the loss function is the 0-1 loss:

$$\ell(h(x), y) = \mathbf{1}_{\{h(x) \neq y\}}.$$

The true risk in this case corresponds to the probability that the hypothesis misclassifies a randomly chosen example:

$$R(h) = \mathbb{P}_{(x,y)\sim\mathcal{D}}(h(x) \neq y),$$

while the empirical risk is simply the proportion of errors on the training sample:

$$\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}_{\{h(x_i) \neq y_i\}}.$$

Minimizing the empirical risk in this setting corresponds to finding the hypothesis that makes the fewest errors on the training data.

**Empirical Risk Versus True Risk**

Generalization Error: The difference between the true risk and the empirical risk is known as the generalization error:

$$\text{Generalization Error}(h) = R(h) - \hat{R}_n(h).$$

The generalization error quantifies how well the performance of a hypothesis on the training data translates to its performance on unseen data. A small generalization error indicates that the hypothesis generalizes well, while a large generalization error suggests overfitting.

Generalization Bounds: Statistical learning theory provides bounds on the generalization error, often in terms of the complexity of the hypothesis class $\mathcal{H}$. These bounds ensure that, with high probability, the true risk of the hypothesis chosen by ERM is close to the empirical risk:

$$R(h_{\text{ERM}}) \leq \hat{R}_n(h_{\text{ERM}}) + \mathcal{O}\left(\frac{\text{Complexity}(\mathcal{H})}{\sqrt{n}}\right),$$

where $\text{Complexity}(\mathcal{H})$ is a measure of the capacity of the hypothesis class, such as the VC dimension or Rademacher complexity.

**Theorem 1.22** (Hoeffding's Inequality) *One of the fundamental results that provides a generalization bound is Hoeffding's inequality. For any fixed hypothesis $h \in \mathcal{H}$, Hoeffding's inequality states that*

$$\mathbb{P}\left(|R(h) - \hat{R}_n(h)| > \epsilon\right) \le 2 \exp\left(-2n\epsilon^2\right).$$

This bound shows that as the sample size $n$ increases, the probability that the empirical risk deviates significantly from the true risk decreases exponentially.

Example: Consider a linear classifier with hypothesis space $\mathcal{H} = \{h(x) = \text{sign}(w^\top x) : w \in \mathbb{R}^d\}$. The generalization error depends on the complexity of $\mathcal{H}$, which, in this case, can be controlled by the norm of the weight vector $w$. Regularization techniques, such as $\ell_2$ regularization, are used to control the capacity of the hypothesis class, thereby reducing the generalization error.

## 1.13.2  VC Dimension and Capacity Control

### Definition and Properties of VC Dimension

The VC dimension as discussed earlier is a measure of the capacity or complexity of a hypothesis class $\mathcal{H}$. It is defined as the maximum number of points that can be shattered by $\mathcal{H}$. A set of points is said to be shattered by $\mathcal{H}$ if, for every possible labeling of the points, there exists a hypothesis in $\mathcal{H}$ that correctly classifies them. Formally, the VC dimension of $\mathcal{H}$, denoted $\text{VC}(\mathcal{H})$, is

$$\text{VC}(\mathcal{H}) = \max\{m \in \mathbb{N} : \exists S \subset \mathbb{R}^d, |S| = m, S \text{ is shattered by } \mathcal{H}\}.$$

Properties of VC Dimension:

1. Upper Bound: For any hypothesis class $\mathcal{H}$, the VC dimension provides an upper bound on the capacity of $\mathcal{H}$. If $\text{VC}(\mathcal{H}) = d$, then $\mathcal{H}$ cannot shatter any set of $d + 1$ points.
2. Implications for Generalization: A high VC dimension implies that the hypothesis class is very flexible and capable of fitting complex patterns, which increases the risk of overfitting. Conversely, a low VC dimension indicates that the hypothesis class is more constrained, which can reduce the risk of overfitting but may also limit the ability to capture complex patterns.
3. Relation to Sample Size: The VC dimension provides a guideline for the sample size needed to ensure good generalization. Specifically, to achieve a small generalization error with high probability, the sample size $n$ should be larger than the VC dimension of the hypothesis class.

**Theorem 1.23** (Sauer–Shelah Lemma) *The Sauer–Shelah lemma provides a combinatorial bound on the number of distinct labelings that a hypothesis class $\mathcal{H}$ can produce on a sample of size n, based on its VC dimension:*

$$|\mathcal{H}_S| \le \sum_{i=0}^{d} \binom{n}{i},$$

*where $\mathcal{H}_S$ is the set of labelings of the sample S by $\mathcal{H}$, and $d = VC(\mathcal{H})$.*

This result is crucial for deriving generalization bounds in terms of the VC dimension.

Example: Consider the hypothesis class of linear classifiers in $\mathbb{R}^2$:

$$\mathcal{H} = \{h(x) = \text{sign}(w_1 x_1 + w_2 x_2 + b) : (w_1, w_2) \in \mathbb{R}^2, b \in \mathbb{R}\}.$$

The VC dimension of this class is 3 because it can shatter any set of 3 points in general position (no 3 points are collinear) but cannot shatter 4 points.

**VC Dimension of Transformer Models**

The VC dimension of deep neural networks, including transformers, is typically very high due to the large number of parameters and the flexibility of the network architecture. However, the exact VC dimension of a specific model can be challenging to determine due to the complexity of the network and the interactions between different layers.

For neural networks with a fixed architecture, the VC dimension can be bounded in terms of the number of parameters and the depth of the network. For a network with $W$ parameters and $L$ layers, the VC dimension is often proportional to $W \log W$, reflecting the fact that deep networks have high capacity but also require careful regularization to avoid overfitting.

Example: Consider a transformer model with multiple layers of self-attention and feedforward networks. Each layer introduces a large number of parameters, contributing to the overall capacity of the model. The VC dimension of such a model would be very high, indicating a strong ability to fit complex data, but also a high risk of overfitting without proper regularization and sufficient training data.

**Theorem 1.24** (VC Dimension of Neural Networks) *Let $\mathcal{H}$ be the hypothesis class of functions representable by a neural network with W parameters. Then, the VC dimension $VC(\mathcal{H})$ satisfies*

$$VC(\mathcal{H}) = \mathcal{O}(W \log W).$$

This bound highlights the importance of controlling the number of parameters in the network to manage the capacity and generalization ability of the model.

## *1.13.3  Rademacher Complexity*

Rademacher complexity is a measure of the richness of a hypothesis class $\mathcal{H}$ based on its ability to fit random noise. Unlike VC dimension, which is combinatorial,

Rademacher complexity provides a data-dependent measure of complexity, making it more flexible and often more informative in practical scenarios. Formally, the empirical Rademacher complexity of a hypothesis class $\mathcal{H}$ with respect to a sample $S = \{x_1, \ldots, x_n\}$ is defined as

$$\hat{\mathcal{R}}_n(\mathcal{H}) = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(x_i) \right],$$

where $\sigma = (\sigma_1, \ldots, \sigma_n)$ are i.i.d. Rademacher variables, each taking values $+1$ or $-1$ with equal probability. The Rademacher complexity measures how well the hypothesis class $\mathcal{H}$ can fit a random labeling of the data.

Properties of Rademacher Complexity:

1. Data Dependence: Rademacher complexity depends on the specific sample $S$, allowing it to capture the complexity of $\mathcal{H}$ relative to the given data.

2. Bounding Generalization Error: Rademacher complexity provides a bound on the generalization error of the hypothesis chosen by ERM. Specifically, with high probability, the following bound holds:

$$R(h) \leq \hat{R}_n(h) + 2\hat{\mathcal{R}}_n(\mathcal{H}) + \mathcal{O}\left(\sqrt{\frac{1}{n}}\right).$$

This bound shows that the generalization error is controlled by the empirical risk and the Rademacher complexity of the hypothesis class.

For a hypothesis class $\mathcal{H}$ with functions bounded by $B$, the Rademacher complexity satisfies the following bound:

$$\hat{\mathcal{R}}_n(\mathcal{H}) \leq \frac{B}{\sqrt{n}}.$$

This bound highlights the inverse relationship between the sample size and the Rademacher complexity, reflecting the fact that larger samples provide better generalization guarantees.

Example: Consider the class of linear classifiers $\mathcal{H} = \{h(x) = w^\top x : \|w\| \leq 1\}$. The Rademacher complexity of this class is bounded by

$$\hat{\mathcal{R}}_n(\mathcal{H}) \leq \frac{1}{\sqrt{n}}.$$

This indicates that as the sample size increases, the complexity of the hypothesis class relative to the data decreases, leading to better generalization.

**Applications to Generalization Analysis**

In transformer models, Rademacher complexity can be used to analyze the generalization ability of different layers or components of the model. For instance, one might compute the Rademacher complexity of the self-attention mechanism or the feedforward layers to assess their capacity to fit the training data relative to random noise.

Example: Suppose we have a transformer model with a large number of attention heads and layers. The Rademacher complexity can help determine whether the model's capacity is too high relative to the amount of training data, indicating a potential risk of overfitting. By regularizing the model or adjusting the architecture (e.g., reducing the number of attention heads), we can control the Rademacher complexity and improve generalization.

Let $\mathcal{H}_{\text{transformer}}$ be the hypothesis class corresponding to a transformer model. The generalization error of the model is bounded by

$$R(h_{\text{ERM}}) \leq \hat{R}_n(h_{\text{ERM}}) + 2\hat{\mathcal{R}}_n(\mathcal{H}_{\text{transformer}}) + \mathcal{O}\left(\sqrt{\frac{1}{n}}\right).$$

This bound underscores the importance of controlling the complexity of the transformer model to ensure good generalization.

In practice, Rademacher complexity provides a tool for selecting model architectures and regularization strategies that balance fit and generalization. By evaluating the Rademacher complexity on a validation set, practitioners can make informed decisions about model adjustments, such as pruning layers, reducing the number of parameters, or applying stronger regularization.

## 1.14  Probabilistic Perspectives on Transformers

Transformers, like other deep learning models, are often viewed through the lens of deterministic optimization. However, a probabilistic perspective, particularly one grounded in Bayesian inference, offers a powerful framework for understanding and quantifying uncertainty in these models. This chapter explores the application of Bayesian inference to transformers, focusing on the computation of posterior distributions and the implementation of Bayesian neural networks (BNNs) within transformer architectures. These approaches provide insights into model uncertainty, leading to more robust and interpretable predictions.

## *1.14.1  Bayesian Inference in Transformers*

Bayesian inference is a statistical method that updates the probability estimate for a hypothesis as more evidence or information becomes available. It combines prior knowledge with observed data to produce a posterior distribution, which quantifies uncertainty about model parameters. This framework is particularly useful in deep learning, where models are often highly parameterized, and understanding uncertainty is crucial for tasks such as decision-making and risk assessment.

### Posterior Distributions and Uncertainty Quantification

Posterior Distributions: In a Bayesian framework, the goal is to compute the posterior distribution of the model parameters given the observed data. Let $\theta$ denote the parameters of a transformer model, and let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}$ represent the observed data, where $x_i$ is the input and $y_i$ is the corresponding output. The posterior distribution $p(\theta|\mathcal{D})$ is given by Bayes' theorem:

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)\, p(\theta)}{p(\mathcal{D})},$$

where $p(\mathcal{D}|\theta)$ is the likelihood, representing the probability of the data given the parameters $\theta$, $p(\theta)$ is the prior distribution, representing the initial beliefs about the parameters before observing the data, and $p(\mathcal{D})$ is the marginal likelihood or evidence, which normalizes the posterior distribution. The posterior distribution $p(\theta|\mathcal{D})$ encapsulates all the information about the parameters after observing the data, including any uncertainty.

Example: Consider a simple linear regression model where $y_i = \theta_0 + \theta_1 x_i + \epsilon_i$, with $\epsilon_i$ representing Gaussian noise. The likelihood function for a single observation $(x_i, y_i)$ is

$$p(y_i|x_i, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \theta_0 - \theta_1 x_i)^2}{2\sigma^2}\right).$$

Given $n$ independent observations, the joint likelihood is

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{n} p(y_i|x_i, \theta).$$

If we assume a prior distribution $p(\theta)$ that is Gaussian (e.g., $p(\theta) \sim \mathcal{N}(\mu_0, \Sigma_0)$), the posterior distribution can be computed analytically or approximated using methods such as Markov Chain Monte Carlo (MCMC) or variational inference.

Uncertainty Quantification: The posterior distribution provides a natural way to quantify uncertainty in model predictions. Instead of a single point estimate for $\theta$, the posterior distribution offers a range of plausible values, each associated with a

probability. This allows for uncertainty-aware predictions, where the model output is expressed as a distribution rather than a deterministic value.

Given a new input $x^*$, the predictive distribution for the output $y^*$ is obtained by marginalizing over the posterior distribution of the parameters:

$$p(y^*|x^*, \mathcal{D}) = \int p(y^*|x^*, \theta) p(\theta|\mathcal{D}) d\theta.$$

This predictive distribution captures both the uncertainty in the model parameters and the inherent noise in the data.

Example: In the case of the linear regression model, the predictive distribution for a new input $x^*$ is a Gaussian distribution with mean and variance given by

$$\mathbb{E}[y^*|x^*, \mathcal{D}] = \mathbb{E}[\theta_0|\mathcal{D}] + \mathbb{E}[\theta_1|\mathcal{D}]x^*,$$

$$\mathrm{Var}(y^*|x^*, \mathcal{D}) = \sigma^2 + x^{*\top} \mathrm{Var}(\theta|\mathcal{D})x^*.$$

Let $\theta$ be the parameters of a Bayesian model with posterior distribution $p(\theta|\mathcal{D})$. The posterior predictive distribution for a new data point $x^*$ is given by

$$p(y^*|x^*, \mathcal{D}) = \int p(y^*|x^*, \theta) p(\theta|\mathcal{D}) d\theta.$$

This distribution quantifies the uncertainty in predictions by integrating over the posterior distribution of the parameters.

**Bayesian Neural Networks**

A Bayesian neural network (BNN) extends the standard neural network framework by treating the network weights as random variables with a prior distribution. Instead of learning a single set of weights, BNNs learn a distribution over weights, which allows for uncertainty estimation in the model's predictions. Let $\theta = \{W_l\}_{l=1}^{L}$ represent the set of weights in a neural network with $L$ layers. In a BNN, we specify a prior distribution over the weights $p(\theta)$ and use Bayesian inference to compute the posterior distribution given the data:

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta) p(\theta)}{p(\mathcal{D})}.$$

Training a BNN involves approximating the posterior distribution $p(\theta|\mathcal{D})$, as exact computation is often intractable due to the high dimensionality of the weight space and the complexity of the likelihood function. Common methods for approximating the posterior include

1. Variational Inference (VI): In VI, the true posterior $p(\theta|\mathcal{D})$ is approximated by a simpler distribution $q(\theta|\phi)$, where $\phi$ are the variational parameters. The goal is to minimize the Kullback–Leibler (KL) divergence between $q(\theta|\phi)$ and the true posterior:

$$\phi^* = \arg\min_{\phi} \mathrm{KL}(q(\theta|\phi) \| p(\theta|\mathcal{D})).$$

The variational approximation is often chosen to be a Gaussian distribution with mean and variance as the variational parameters.

2. Monte Carlo Dropout: Monte Carlo dropout is a practical approximation technique where dropout is applied during both training and inference. The posterior distribution is approximated by sampling multiple forward passes with dropout enabled, effectively averaging predictions over different network configurations.

Example: Consider a BNN with a single hidden layer and a ReLU activation function. The weights $W_1$ and $W_2$ are treated as random variables with Gaussian priors:

$$W_1 \sim \mathcal{N}(0, \sigma_1^2 I), \quad W_2 \sim \mathcal{N}(0, \sigma_2^2 I).$$

The likelihood function is defined based on the output of the network and the observed data. Using variational inference, we approximate the posterior distribution over $W_1$ and $W_2$ by optimizing the variational parameters.

Uncertainty in Predictions: The predictive distribution in a BNN is obtained by marginalizing over the posterior distribution of the weights:

$$p(y^*|x^*, \mathcal{D}) = \int p(y^*|x^*, \theta) p(\theta|\mathcal{D}) d\theta.$$

In practice, this integral is approximated by sampling from the posterior distribution of the weights and averaging the predictions:

$$p(y^*|x^*, \mathcal{D}) \approx \frac{1}{M} \sum_{i=1}^{M} p(y^*|x^*, \theta_i),$$

where $\theta_i$ are samples from the approximate posterior distribution.

In transformer models, Bayesian inference can be applied to the weights of the attention layers and feedforward networks. By modeling the weights as random variables and learning their posterior distributions, we can quantify uncertainty in the model's predictions, which is particularly useful in applications such as natural language processing (NLP) and decision-making under uncertainty.

For instance, in a Bayesian transformer, the weights of the self-attention mechanism can be treated as random variables with a prior distribution. The posterior distribution of these weights is learned using variational inference, allowing the model to produce uncertainty-aware predictions for tasks such as machine translation or text classification.

Let $\theta = \{W_l\}_{l=1}^{L}$ be the weights of a BNN. The posterior predictive distribution for a new input $x^*$ is given by

$$p(y^*|x^*, \mathcal{D}) = \int p(y^*|x^*, \theta) p(\theta|\mathcal{D}) d\theta.$$

This distribution captures both the model uncertainty (due to limited data) and the epistemic uncertainty (due to uncertainty in the model parameters).

## 1.14.2  PAC-Bayes Generalization Bounds

PAC-Bayesian (PAC-Bayes) theory provides a framework for deriving generalization bounds that combine elements of both PAC (Probably Approximately Correct) learning theory and Bayesian inference. These bounds offer insights into how well a model trained on a finite dataset is expected to perform on unseen data. The PAC-Bayes approach is particularly relevant for understanding the generalization behavior of complex models like transformers, where the capacity of the model and the amount of training data both play critical roles. PAC-Bayes theory blends the probabilistic interpretation of Bayesian methods with the performance guarantees of PAC learning. The core idea is to derive bounds on the generalization error of a learned hypothesis, not just in the worst-case scenario, but with respect to a prior distribution over the hypothesis space. These bounds are particularly useful for models with high capacity, such as transformers, where traditional generalization bounds may be too loose.

### PAC-Bayes Theorem and Applications

The PAC-Bayes theorem provides a bound on the generalization error of a hypothesis $h$ selected from a hypothesis space $\mathcal{H}$, where the hypothesis is drawn from a posterior distribution $Q$ that depends on the training data. The bound is given in terms of the KL divergence between the posterior distribution $Q$ and a prior distribution $P$ over $\mathcal{H}$, as well as the empirical risk $\hat{R}_n(h)$ and the true risk $R(h)$. Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}$ be the training data, and $\ell(h(x), y)$ be the loss function for a hypothesis $h$. The empirical risk of $h$ is defined as

$$\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^{n} \ell(h(x_i), y_i).$$

The PAC-Bayes theorem states that, with probability at least $1 - \delta$ over the choice of the training data, for any posterior distribution $Q$ over $\mathcal{H}$ and prior distribution $P$ over $\mathcal{H}$, the following bound holds:

$$\mathbb{E}_{h \sim Q}[R(h)] \leq \mathbb{E}_{h \sim Q}[\hat{R}_n(h)] + \sqrt{\frac{\text{KL}(Q\|P) + \log \frac{n}{\delta}}{2n}},$$

where $\text{KL}(Q\|P)$ is the KL divergence between the posterior $Q$ and the prior $P$.

Interpreting the Bound:

1. Empirical Risk $\hat{R}_n(h)$: Represents the average error of the hypothesis on the training data.
2. True Risk $R(h)$: Represents the expected error on new, unseen data.
3. KL Divergence $\text{KL}(Q\|P)$: Measures how much the posterior distribution $Q$ diverges from the prior distribution $P$. A smaller KL divergence indicates that the posterior is close to the prior, leading to tighter generalization bounds.
4. Sample Size $n$: The bound tightens as the sample size increases, reflecting that more data leads to better generalization.

Applications of PAC-Bayes Theorem:

1. Model Selection: PAC-Bayes bounds can be used to select models by comparing the generalization bounds for different posterior distributions. For instance, in a transformer model, one might consider different architectures or regularization strategies and choose the one with the lowest PAC-Bayes bound.

2. Regularization: The KL divergence term in the PAC-Bayes bound encourages the posterior distribution to remain close to the prior, effectively regularizing the model. This is particularly useful in preventing overfitting, as it penalizes models that deviate too much from the prior knowledge encoded in $P$.

Example: Consider a binary classification problem where the hypothesis space $\mathcal{H}$ consists of linear classifiers. The prior distribution $P$ might be a Gaussian distribution centered on a particular weight vector, representing prior knowledge about the likely configuration of the classifier. The posterior $Q$ is then learned from the data, and the PAC-Bayes bound provides a guarantee on how well the classifier is expected to perform on new data.

**Theorem 1.25** (PAC-Bayes Theorem) *Let $\mathcal{H}$ be a hypothesis class, $P$ a prior distribution over $\mathcal{H}$, and $Q$ a posterior distribution over $\mathcal{H}$. With probability at least $1 - \delta$, the generalization error satisfies*

$$\mathbb{E}_{h \sim Q}[R(h)] \leq \mathbb{E}_{h \sim Q}[\hat{R}_n(h)] + \sqrt{\frac{KL(Q\|P) + \log \frac{n}{\delta}}{2n}}.$$

This bound provides a principled way to quantify the generalization error of models, taking into account both the empirical performance and the complexity of the model as measured by the KL divergence.

**Interpreting PAC-Bayes Bounds in Transformers**

Transformers are powerful models with high capacity, often leading to concerns about overfitting, especially when the amount of training data is limited. PAC-Bayes bounds offer a way to analyze and control the generalization behavior of transformers by incorporating prior knowledge and regularizing the learned parameters.

1. Prior Distribution in Transformers: In the context of transformers, the prior distribution $P$ can encode prior beliefs about the model parameters, such as the weights of the attention layers or the feedforward networks. For example, one might use a Gaussian prior centered around zero or around pre-trained weights from a related task.

2. Posterior Distribution: The posterior distribution $Q$ is learned from the data, often through techniques like variational inference or Monte Carlo methods. The PAC-Bayes bound then provides a measure of how much the learned model (represented by $Q$) deviates from the prior knowledge (represented by $P$).

3. Generalization Analysis: The PAC-Bayes bound helps in understanding the trade-off between fitting the training data and maintaining generalization. A small KL divergence between $Q$ and $P$ suggests that the model has not deviated too much from the prior, which can be a sign of good generalization. Conversely, a large KL divergence might indicate overfitting.

Example: Suppose a transformer model is trained on a natural language processing task, such as sentiment analysis. The prior distribution $P$ might reflect a pre-trained transformer on a large language corpus, while the posterior $Q$ is fine-tuned on the sentiment analysis dataset. The PAC-Bayes bound can be used to assess whether the fine-tuned model is likely to generalize well to new text data.

Let $\Theta$ represent the set of parameters in a transformer model, and let $P(\Theta)$ and $Q(\Theta)$ denote the prior and posterior distributions over these parameters. The PAC-Bayes generalization bound for the transformer model is given by

$$\mathbb{E}_{\Theta \sim Q}[R(\Theta)] \leq \mathbb{E}_{\Theta \sim Q}[\hat{R}_n(\Theta)] + \sqrt{\frac{\mathrm{KL}(Q(\Theta) \| P(\Theta)) + \log \frac{n}{\delta}}{2n}}.$$

This bound highlights the importance of selecting appropriate priors and regularization strategies to control the generalization error in transformer models.

The PAC-Bayes framework naturally incorporates regularization by penalizing posterior distributions that deviate significantly from the prior. This can be implemented in transformers through techniques like weight decay, which can be viewed as imposing a Gaussian prior on the model weights. They can guide the tuning of hyperparameters, such as learning rates and regularization strengths. By evaluating how these choices affect the KL divergence and the empirical risk, one can optimize the generalization performance of the transformer model.

## *1.14.3 Generalization Performance in Practice*

The theoretical foundations of generalization provide a robust framework for understanding how machine learning models, including transformers, perform on unseen data. However, practical implementation often involves empirical validation, model selection, and hyperparameter tuning to optimize generalization performance. This section explores empirical studies on generalization and discusses strategies for model selection and hyperparameter tuning in the context of transformers, using mathematical analysis to support these practices. Generalization performance refers to a model's ability to perform well on new, unseen data, beyond the training set. While theoretical bounds such as PAC-Bayes offer insights into generalization, empirical studies provide practical evidence and guidance for enhancing model performance. The interplay between theory and practice is crucial in developing models that are both robust and efficient.

**Empirical Studies on Generalization**

Empirical studies on generalization often involve comparing the performance of models on training data versus validation and test data. These studies help identify factors that contribute to overfitting, underfitting, or optimal generalization. For transformers, key factors include the size of the model, the amount of training data, the choice of regularization techniques, and the complexity of the task.

Let $\mathcal{D}_{\text{train}}$ represent the training dataset, $\mathcal{D}_{\text{val}}$ the validation dataset, and $\mathcal{D}_{\text{test}}$ the test dataset. The generalization error $\epsilon_{\text{gen}}$ is typically defined as the difference between the test error $\hat{R}_{\text{test}}(h)$ and the training error $\hat{R}_{\text{train}}(h)$:

$$\epsilon_{\text{gen}} = \hat{R}_{\text{test}}(h) - \hat{R}_{\text{train}}(h),$$

where $h$ is the hypothesis (or model) selected during training. A small generalization error indicates that the model has successfully learned to generalize from the training data to new data, while a large generalization error suggests overfitting.

Example: Consider a transformer model trained on a dataset of English sentences for a machine translation task. The training error is measured as the average translation error on the training set, while the validation and test errors are measured on held-out datasets. An empirical study might involve training the model with varying amounts of data or different regularization strengths and observing how the generalization error changes.

Overfitting and Underfitting: Overfitting occurs when the model is too complex relative to the amount of training data, leading to low training error but high test error. This is often indicated by a large generalization error. Underfitting occurs when the model is too simple to capture the underlying patterns in the data, leading to high training and test errors.

Regularization: Empirical studies often explore the impact of regularization techniques on generalization. Regularization methods, such as weight decay or dropout, penalize model complexity, effectively reducing the variance and improving generalization. The effect of regularization can be analyzed mathematically by adding a penalty term to the loss function:

$$\mathcal{L}_{\text{reg}}(h) = \mathcal{L}(h) + \lambda \Omega(h),$$

where $\lambda$ is the regularization strength and $\Omega(h)$ is a regularization term (e.g., $\|h\|_2^2$ for weight decay).

**Remark (Regularization and Generalization):** Let $\hat{R}_\lambda(h)$ denote the regularized empirical risk. The generalization bound for the regularized model is

$$R(h) \le \hat{R}_\lambda(h) + \mathcal{O}\left(\sqrt{\frac{\text{Complexity}(\mathcal{H}) + \lambda \Omega(h)}{n}}\right),$$

where $n$ is the sample size and Complexity($\mathcal{H}$) is a measure of the model class complexity. This bound shows that regularization can reduce the generalization error by controlling model complexity.

Empirical studies on transformers have shown that

1. Model Size Versus Data Size: Larger transformers generally perform better when trained on large datasets, but they are prone to overfitting on smaller datasets.
2. Layer Normalization and Dropout: Techniques like layer normalization and dropout are crucial for stabilizing training and improving generalization.
3. Transfer Learning: Fine-tuning pre-trained transformers on specific tasks often leads to better generalization than training from scratch, due to the regularization effect of the pre-training phase.

Example: An empirical study might involve training several transformer models of varying sizes on a dataset like GLUE (General Language Understanding Evaluation) and evaluating their performance on the test set. The study could measure how generalization performance changes with different amounts of pre-training data, regularization techniques, and hyperparameters.

## Model Selection and Hyperparameter Tuning

Model selection involves choosing the best hypothesis from a set of candidates based on their expected generalization performance. In practice, this is often done using cross-validation, where the data is split into several folds, and the model is trained and validated on different subsets. The model with the best average performance across the folds is selected.

Let $\mathcal{H} = \{h_1, h_2, \ldots, h_k\}$ be a set of candidate models, and let $\mathcal{D}_{\text{val}}$ be the validation set. The selected model $h^*$ is the one that minimizes the validation error:

$$h^* = \arg\min_{h \in \mathcal{H}} \hat{R}_{\text{val}}(h).$$

The challenge is to ensure that the validation set is representative of the test set so that the selected model generalizes well.

Hyperparameter tuning is the process of optimizing the parameters that control the learning process, such as the learning rate, regularization strength, and network architecture. This is typically done using grid search, random search, or more sophisticated methods like Bayesian optimization.

Let $\lambda \in \Lambda$ represent a hyperparameter. The goal is to find the optimal $\lambda^*$ that minimizes the validation error:

$$\lambda^* = \arg\min_{\lambda \in \Lambda} \hat{R}_{\text{val}}(h(\lambda)),$$

where $h(\lambda)$ is the model trained with hyperparameter $\lambda$.

In a transformer model, the learning rate $\eta$ and dropout rate $p$ are crucial hyperparameters. An empirical study might involve training several models with different combinations of $\eta$ and $p$, evaluating their performance on a validation set, and selecting the combination that results in the best generalization.

Cross-validation is a common technique for model selection and hyperparameter tuning. In $k$-fold cross-validation, the data is split into $k$ subsets, and the model is trained on $k-1$ subsets while validating on the remaining subset. This process is repeated $k$ times, and the average validation error is used to select the best model or hyperparameters.

**Theorem 1.26** (Generalization in Cross-Validation) *Let $\hat{R}_{CV}(h)$ be the cross-validated error for a model h. The generalization error is bounded by*

$$R(h) \le \hat{R}_{CV}(h) + \mathcal{O}\left(\sqrt{\frac{1}{n}}\right),$$

*where n is the size of the dataset.*

This bound indicates that cross-validation provides a reliable estimate of the generalization error. A practical example of hyperparameter tuning in transformers might involve optimizing the number of attention heads and the size of the feedforward network. The model is trained on different configurations and validated using cross-validation. The configuration with the lowest cross-validated error is chosen for deployment.

# References

1. Andrews, L.C.: Special Functions of Mathematics for Engineers. Oxford University Press (1998)
2. Ash, R.B., Doléans-Dade, C.A.: Probability and Measure Theory. Academic (2000)

3. Atkinson, K.E., Han, W.: Spherical Harmonics and Approximations on the Unit Sphere: An Introduction. Springer Monographs in Mathematics (2012)
4. Axler, S.: Linear Algebra Done Right, 4th edn. Springer (2023). ISBN 978-3-031-41025-3. https://doi.org/10.1007/978-3-031-41026-0
5. Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization (2016). arXiv:1607.06450
6. Barron, A.R.: Universal approximation bounds for superpositions of a sigmoidal function. IEEE Trans. Inf. Theory **39**(3), 930–945 (1993)
7. Baydin, A.G., Pearlmutter, B.A., Radul, A.A., Siskind, J.M.: Automatic differentiation in machine learning: a survey. J. Mach. Learn. Res. **18**(1), 5595–5637 (2017). ISSN 1532-4435
8. Belkin, M., Niyogi, P.: Laplacian eigenmaps for dimensionality reduction and data representation. Neural Comput. **15**(6), 1373–1396 (2003)
9. Bott, R., Tu, L.W.: Differential Forms in Algebraic Topology. Springer (2014)
10. Bourbaki, N.: Lie Groups and Lie Algebras. Springer (1989)
11. Bredon, G.E.: Introduction to Compact Transformation Groups. Academic (1972)
12. Brink, D.M., Satchler, G.R.: Angular Momentum. Oxford University Press (1993)
13. Cheney, W., Light, W.: A Course in Approximation Theory. American Mathematical Society (2009)
14. Cover, T.M., Thomas, J.A.: Elements of Information Theory. Wiley (2006)
15. Cybenko, G.: Approximation by superpositions of a sigmoidal function. Math. Control Signals Syst. **2**(4), 303–314 (1989)
16. Donoho, D.L., Grimes, C.: Hessian eigenmaps: locally linear embedding techniques for high-dimensional data. Proc. Natl. Acad. Sci. **100**(10), 5591–5596 (2003)
17. Edmonds, A.R.: Angular Momentum in Quantum Mechanics. Princeton University Press (1955)
18. Folland, G.B.: Real Analysis: Modern Techniques and Their Applications. Wiley (1999)
19. Fulton, W., Harris, J.: Representation Theory: A First Course. Springer (2013)
20. Ge, R., Jin, C., Zheng, Y.: Escaping from saddle points—online stochastic gradient for tensor decomposition. In: Conference on Learning Theory (COLT), pp. 797–842 (2015)
21. Gilmore, R.: Lie Groups, Lie Algebras, and Some of Their Applications. Dover Publications (2008)
22. Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press (2016)
23. Gray, R.M.: Entropy and Information Theory. Springer (2011)
24. Hall, B.C.: Lie Groups, Lie Algebras, and Representations: An Elementary Introduction. Graduate Texts in Mathematics, 2nd edn. Springer (2016). ISBN 978-3-319-37433-8. https://doi.org/10.1007/978-3-319-13467-3
25. Halmos, P.R.: Finite-Dimensional Vector Spaces. Springer (1974)
26. Halmos, P.R.: Measure Theory. Springer (1974)
27. Hobson, E.W.: The Theory of Spherical and Ellipsoidal Harmonics. Cambridge University Press (1931)
28. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural Comput. **9**(8), 1735–1780 (1997)
29. Hoffman, K., Kunze, R.: Linear Algebra, 2nd edn. Pearson Education India (2015). ISBN 978-9332550070
30. Hornik, K.: Approximation capabilities of multilayer feedforward networks. Neural Netw. **4**(2), 251–257 (1991)
31. Humphreys, J.E.: Introduction to Lie Algebras and Representation Theory. Graduate Texts in Mathematics, 1st edn. Springer (1973). ISBN 978-0-387-90052-0. https://doi.org/10.1007/978-1-4612-6398-2
32. Jackson, D.: The Theory of Approximation. American Mathematical Society (1930)
33. Jin, C., Ge, R., Netrapalli, P., Kakade, S.M., Jordan, M.I.: How to escape saddle points efficiently. In: Proceedings of the 34th International Conference on Machine Learning (ICML), pp. 1724–1732 (2017)
34. Kingma, D.P., Ba, J.: Adam: a method for stochastic optimization (2014). arXiv:1412.6980

35. Knapp, A.W.: Representation Theory of Semisimple Groups: An Overview Based on Examples. Princeton University Press (1986)
36. Kobayashi, S., Nomizu, K.: Foundations of Differential Geometry, vol 1. Interscience Publishers (1963)
37. Kolmogorov, A.N.: Foundations of the Theory of Probability. Chelsea Publishing Company (1950)
38. Lang, S.: Linear Algebra, 3rd edn. Springer (2010). ISBN 978-1-4419-3081-1. https://doi.org/10.1007/978-1-4757-1949-9
39. LeCun, Y.: A theoretical framework for back-propagation. In: Proceedings of the 1988 Connectionist Models Summer School, pp. 21–28 (1988)
40. Lee, J.D., Simchowitz, M., Jordan, M.I., Recht, B.: Gradient descent only converges to minimizers. In: Conference on Learning Theory (COLT), pp. 1246–1257 (2016)
41. Leshno, M., Lin, V.C., Pinkus, A., Schocken, S.: Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. Neural Netw. **6**(6), 861–867 (1993)
42. Mohri, M.: Afshin Rostamizadeh, and Ameet Talwalkar. MIT Press, Foundations of machine learning (2012)
43. Montgomery, D., Samelson, H.: Transformation groups of spheres. Ann. Math. **44**, 454 (1943). https://api.semanticscholar.org/CorpusID:124533884
44. Müller, C.: Spherical Harmonics. Springer (1966)
45. Pascanu, R., Mikolov, T., Bengio, Y.: On the difficulty of training recurrent neural networks. In: Proceedings of the 30th International Conference on Machine Learning (ICML), pp. 1310–1318 (2013)
46. Michael, J.D.: Powell. Cambridge University Press, Approximation Theory and Methods (1981)
47. Reddi, S.J., Hefny, A., Sra, S., Póczós, B., Smola, A.: Stochastic variance reduction for non-convex optimization. In: Proceedings of the 33rd International Conference on International Conference on Machine Learning, vol. 48. ICML'16, pp. 314–323. JMLR.org (2016)
48. Rivlin, T.J.: An Introduction to the Approximation of Functions. Dover Publications (2003)
49. Roweis, S.T., Saul, L.K.: Nonlinear dimensionality reduction by locally linear embedding. Science **290**(5500), 2323–2326 (2000)
50. Royden, H.L.: Real Analysis. Macmillan (1988)
51. Rudin, W.: Real and Complex Analysis. McGraw-Hill (1987)
52. Rumelhart, D.E., Hinton, G.E., Williams, R.J.: Learning representations by back-propagating errors. Nature **323**(6088), 533–536 (1986)
53. Saul, L.K., Roweis, S.T.: Think globally, fit locally: unsupervised learning of low dimensional manifolds. J. Mach. Learn. Res. **4**(null), 119–155 (2003). ISSN 1532-4435. https://doi.org/10.1162/153244304322972667
54. Strang, G.: Introduction to Linear Algebra, 4th edn. Wellesley-Cambridge Press (2009). ISBN 978-0980232714
55. Telgarsky, M.: Benefits of depth in neural networks. In: 29th Annual Conference on Learning Theory. In: Proceedings of Machine Learning Research, vol. 49, pp. 1517–1539. Columbia University, New York, New York, USA. Accessed 23–26 Jun 2016. PMLR. https://proceedings.mlr.press/v49/telgarsky16.html
56. Tenenbaum, J.B., de Silva, V., Langford, J.C.: A global geometric framework for nonlinear dimensionality reduction. Science **290**(5500), 2319–2323 (2000)
57. Tishby, N., Pereira, F.C., Bialek, W.: The information bottleneck method (2000). ArXiv:physics/0004057
58. Vapnik, V.N.: The Nature of Statistical Learning Theory. Springer (1995)
59. Varadarajan, V.S.: Lie Groups, Lie Algebras, and Their Representations. Springer (1984)
60. Varshalovich, D.A., Moskalev, A.N., Khersonskii, V.K.: Quantum Theory of Angular Momentum. World Scientific (1988)
61. Wigner, E.P.: Group Theory and its Application to the Quantum Mechanics of Atomic Spectra. Academic (1959)

# Chapter 2
# Word Embeddings and Positional Encoding

## 2.1 Word Embeddings

Word embeddings are fundamental in NLP, providing a way to represent words as vectors in a continuous vector space. This representation enables the capture of semantic relationships between words, allowing models to perform complex linguistic tasks such as translation, sentiment analysis, and information retrieval. This section delves into the mathematical foundations of word embeddings, focusing on vector space models and the Word2Vec framework, which includes the Skip-gram and Continuous Bag of Words (CBOW) models.

The mathematical foundation of word embeddings lies in the representation of words as vectors in a high-dimensional vector space. The goal is to map words to vectors in such a way that the geometric relationships between the vectors capture meaningful semantic relationships between the words.

### 2.1.1 Vector Space Models

In vector space models, each word $w$ in a vocabulary $\mathcal{V}$ is represented by a vector $\mathbf{v}_w$ in a $d$-dimensional vector space $\mathbb{R}^d$. The embedding space is typically learned from a large corpus of text, with the vectors being adjusted to reflect the contextual usage of words. Given a corpus $C$ consisting of a sequence of words $\{w_1, w_2, \ldots, w_n\}$, the objective is to find a mapping $\phi : \mathcal{V} \to \mathbb{R}^d$ such that semantically similar words have embeddings that are close in the vector space.

A common measure of similarity between two word vectors $\mathbf{v}_w$ and $\mathbf{v}_{w'}$ is the cosine similarity, defined as

$$\cos(\mathbf{v}_w, \mathbf{v}_{w'}) = \frac{\mathbf{v}_w \cdot \mathbf{v}_{w'}}{\|\mathbf{v}_w\| \|\mathbf{v}_{w'}\|},$$

where $\mathbf{v}_w \cdot \mathbf{v}_{w'}$ is the dot product of the two vectors, and $\|\mathbf{v}_w\|$ and $\|\mathbf{v}_{w'}\|$ are their magnitudes. Cosine similarity ranges from –1 to 1, with 1 indicating that the vectors are identical, –1 indicating that they are diametrically opposed, and 0 indicating orthogonality.

Example: Consider the words "king" and "queen." In a well-trained embedding space, the vector $\mathbf{v}_{king}$ should be close to $\mathbf{v}_{queen}$, reflecting the semantic similarity between these words. Moreover, the vector difference $\mathbf{v}_{king} - \mathbf{v}_{queen}$ might be similar to the vector difference $\mathbf{v}_{man} - \mathbf{v}_{woman}$, capturing the analogy "king is to queen as man is to woman."

One approach to constructing word embeddings is to use a word co-occurrence matrix $M$, where each entry $M_{ij}$ represents the number of times word $w_i$ co-occurs with word $w_j$ within a specified context window. The context window is typically defined as a fixed number of words before and after the target word. Given the co-occurrence matrix $M$, a variety of methods can be applied to derive word vectors, such as matrix factorization techniques (e.g., singular value decomposition) or probabilistic models.

**Theorem 2.1** (Low-Rank Approximation) *Let $M$ be the word co-occurrence matrix, and suppose $M$ can be approximated by a low-rank matrix $M_k$ of rank $k$:*

$$M \approx M_k = U_k \Sigma_k V_k^\top,$$

*where $U_k$ and $V_k$ are matrices whose columns are the left and right singular vectors of $M$, and $\Sigma_k$ is a diagonal matrix containing the top $k$ singular values. The rows of $U_k$ and $V_k$ provide $k$-dimensional word embeddings that capture the most significant relationships in the co-occurrence data.*

The spectral properties of the word co-occurrence matrix provide insights into the structure of the embedding space. The singular values in $\Sigma_k$ represent the importance of the corresponding singular vectors in capturing the variance in the data. By selecting the top $k$ singular values, we obtain a reduced-dimensional representation that retains the most significant information.

In practice, one might construct a co-occurrence matrix from a corpus of English text and then apply Singular Value Decomposition (SVD) to obtain a low-rank approximation. The resulting word vectors can then be used in various NLP tasks, such as word similarity, analogy solving, and text classification.

### 2.1.2   Word2Vec: Skip-gram and CBOW

Word2Vec is a popular framework for learning word embeddings, introduced by [3]. It consists of two main models: Skip-gram and Continuous Bag of Words (CBOW). Both models are based on the distributional hypothesis, which posits that words appearing in similar contexts tend to have similar meanings.

The Skip-gram model aims to predict the context words surrounding a target word. Given a word $w_t$ in the corpus, the model predicts the words $w_{t-j}, \ldots, w_{t-1}, w_{t+1}, \ldots, w_{t+j}$ within a context window of size $2j$. Formally, the objective is to maximize the following log-likelihood:

$$\mathcal{L}_{\text{SG}} = \frac{1}{T} \sum_{t=1}^{T} \sum_{-j \leq k \leq j, k \neq 0} \log p(w_{t+k}|w_t),$$

where $p(w_{t+k}|w_t)$ is the probability of observing the context word $w_{t+k}$ given the target word $w_t$. This probability is modeled using the softmax function:

$$p(w_{t+k}|w_t) = \frac{\exp(\mathbf{v}_{w_{t+k}}^{\top} \mathbf{v}_{w_t})}{\sum_{w' \in \mathcal{V}} \exp(\mathbf{v}_{w'}^{\top} \mathbf{v}_{w_t})},$$

where $\mathbf{v}_{w_t}$ and $\mathbf{v}_{w_{t+k}}$ are the embeddings of the target and context words, respectively.

Computing the softmax function for all words in the vocabulary is computationally expensive. To mitigate this, Word2Vec uses negative sampling, where the model only updates a small number of negative samples (words that do not appear in the context) along with the positive context words. The negative sampling objective is

$$\mathcal{L}_{\text{NS}} = \log \sigma(\mathbf{v}_{w_{t+k}}^{\top} \mathbf{v}_{w_t}) + \sum_{i=1}^{K} \mathbb{E}_{w_i \sim P_n(w)}[\log \sigma(-\mathbf{v}_{w_i}^{\top} \mathbf{v}_{w_t})],$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function, $K$ is the number of negative samples, and $P_n(w)$ is the noise distribution.

The CBOW model is the inverse of Skip-gram. Instead of predicting context words from a target word, CBOW predicts the target word from its surrounding context. The objective function for CBOW is

$$\mathcal{L}_{\text{CBOW}} = \frac{1}{T} \sum_{t=1}^{T} \log p(w_t|w_{t-j}, \ldots, w_{t-1}, w_{t+1}, \ldots, w_{t+j}),$$

where the context words are combined (e.g., averaged) to predict the target word $w_t$.

Both Skip-gram and CBOW models implicitly factorize a word-context co-occurrence matrix, with the learned word vectors corresponding to the low-dimensional representations of words that capture the underlying semantic structure. The embeddings are learned by optimizing the respective objective functions using SGD or other optimization methods.

**Remark (Representation Power of Word2Vec):** Let $\mathbf{v}_w$ be the embedding of word $w$ learned by Word2Vec. The learned embeddings capture not only word similarity but also syntactic and semantic relationships between words. Specifically, for

words $w_1$, $w_2$, $w_3$, $w_4$, if the relationship $w_1$ is to $w_2$ as $w_3$ is to $w_4$ holds (e.g., "king" is to "queen" as "man" is to "woman"), then

$$\mathbf{v}_{w_1} - \mathbf{v}_{w_2} \approx \mathbf{v}_{w_3} - \mathbf{v}_{w_4}.$$

This property is known as the vector offset property and is a key feature of the embeddings produced by Word2Vec.

Example: Consider the analogy task: "king" is to "queen" as "man" is to "woman." The Word2Vec model learns embeddings such that

$$\mathbf{v}_{\text{king}} - \mathbf{v}_{\text{queen}} \approx \mathbf{v}_{\text{man}} - \mathbf{v}_{\text{woman}}.$$

This relationship can be tested by finding the word vector that is closest to the result of $\mathbf{v}_{\text{king}} - \mathbf{v}_{\text{man}} + \mathbf{v}_{\text{woman}}$, which should ideally be $\mathbf{v}_{\text{queen}}$.

Next we'll continue exploring the mathematical foundations and different models used to generate word embeddings, with a focus on GloVe, FastText, and contextual embeddings like ELMo and BERT. These models extend the principles introduced in the Word2Vec framework by incorporating global corpus statistics, subword information, and contextual dependencies, providing richer and more nuanced word representations.

### 2.1.3   GloVe: Global Vectors for Word Representation

GloVe is based on the idea that word embeddings can be learned by leveraging global co-occurrence statistics from a corpus. Unlike Word2Vec, which relies on local context windows, GloVe constructs embeddings by directly factorizing a word co-occurrence matrix. The core idea is to model the ratios of co-occurrence probabilities rather than the probabilities themselves. Let $X$ be the word co-occurrence matrix where $X_{ij}$ denotes the number of times word $w_j$ appears in the context of word $w_i$. GloVe seeks to find word vectors $\mathbf{v}_w$ and context vectors $\mathbf{v}_c$ such that their dot product approximates the logarithm of the co-occurrence counts:

$$\mathbf{v}_w^\top \mathbf{v}_c + b_w + b_c = \log(X_{wc}),$$

where $b_w$ and $b_c$ are bias terms associated with the word and context, respectively.

The objective function in GloVe is to minimize the weighted least squares error:

$$J = \sum_{i,j=1}^{|V|} f(X_{ij}) \left( \mathbf{v}_i^\top \mathbf{v}_j + b_i + b_j - \log(X_{ij}) \right)^2,$$

where $f(X_{ij})$ is a weighting function that controls the influence of different co-occurrence pairs. Typically, the weighting function is defined as

$$f(X_{ij}) = \begin{cases} \left(\frac{X_{ij}}{X_{\max}}\right)^{\alpha} & \text{if } X_{ij} < X_{\max}, \\ 1 & \text{if } X_{ij} \geq X_{\max}, \end{cases}$$

where $\alpha$ and $X_{\max}$ are hyperparameters. This function downweights the influence of very frequent word pairs and emphasizes mid-frequency pairs, which are considered more informative.

The logarithmic function in the objective reflects the intuition that the difference in word meanings should be proportional to the logarithm of their co-occurrence probability. The symmetry of the co-occurrence matrix leads to symmetric embeddings, i.e., $\mathbf{v}_w \approx \mathbf{v}_c$, resulting in embeddings that capture global corpus-wide relationships between words.

**Theorem 2.2** (Existence of GloVe Embeddings) *Given a co-occurrence matrix X and appropriate choices of the weighting function $f(X_{ij})$, there exist word vectors $\mathbf{v}_w$ and context vectors $\mathbf{v}_c$ that satisfy the GloVe objective function. The uniqueness of the solution depends on the specific choice of initialization and optimization algorithm.*

Example: For words "ice" and "steam," GloVe might learn embeddings $\mathbf{v}_{\text{ice}}$ and $\mathbf{v}_{\text{steam}}$ such that the difference between these vectors reflects the difference between the co-occurrence contexts of these words (e.g., "cold" is more associated with "ice" than "steam").

Advantages of GloVe:

1. Global Information: GloVe captures both local context and global statistical information, providing more balanced embeddings.
2. Interpretability: The model's reliance on co-occurrence ratios makes the learned embeddings interpretable, often reflecting intuitive analogies.

### 2.1.4 FastText Embeddings

FastText, an extension of Word2Vec developed by Facebook AI, addresses the limitation of treating each word as an atomic unit by incorporating subword information ([1, 2]). Words are represented not only by their vectors but also by vectors of their constituent character n-grams. This allows FastText to generate embeddings for out-of-vocabulary (OOV) words and to capture morphological information.

Let $w$ be a word and $\mathcal{G}(w)$ be the set of character n-grams generated from $w$. In FastText, the word embedding $\mathbf{v}_w$ is computed as the sum of the embeddings of its n-grams:

$$\mathbf{v}_w = \sum_{g \in \mathcal{G}(w)} \mathbf{v}_g,$$

where $\mathbf{v}_g$ is the vector representation of n-gram $g$. The set $\mathcal{G}(w)$ typically includes all n-grams of length 3 to 6, as well as the word itself.

FastText extends the Skip-gram model by predicting the context words not only from the word $w$ itself but also from its subwords. The objective is to maximize the log-likelihood:

$$\mathcal{L}_{\text{FastText}} = \frac{1}{T} \sum_{t=1}^{T} \sum_{-j \le k \le j, k \ne 0} \log p(w_{t+k} | \mathbf{v}_w),$$

where $\mathbf{v}_w$ is the subword-based embedding defined above.

Advantages of FastText:

1. Handling Morphology: By considering subword information, FastText captures morphological variations (e.g., "run," "running," "runner").
2. OOV Words: FastText can generate embeddings for words not seen during training by summing the embeddings of their n-grams.

**Remark (FastText Embedding Construction):** Let $\mathcal{V}_{\text{train}}$ be the training vocabulary, and let $w_{\text{new}}$ be a word not in $\mathcal{V}_{\text{train}}$. The embedding $\mathbf{v}_{w_{\text{new}}}$ can be constructed as

$$\mathbf{v}_{w_{\text{new}}} = \sum_{g \in \mathcal{G}(w_{\text{new}})} \mathbf{v}_g,$$

where $\mathbf{v}_g$ are the pre-trained n-gram embeddings from the model. This ensures that $w_{\text{new}}$ has a meaningful embedding even if it was not present in the training data.

Example: For the word "unhappiness," which may not be in the training vocabulary, FastText would construct its embedding by summing the embeddings of its n-grams, such as "un," "happiness," "ness," etc., capturing both the root word "happy" and the prefix/suffix information.

### 2.1.5  Contextual Word Embeddings (e.g., ELMo, BERT)

Traditional word embeddings like Word2Vec, GloVe, and FastText generate static vectors for each word, meaning the vector for a word like "bank" is the same regardless of whether the context refers to a financial institution or the side of a river. Contextual word embeddings, introduced with models like ELMo and BERT, address this limitation by generating word representations that are context dependent.

ELMo (Embeddings from Language Models) embeddings ([6]) are generated from a deep, bi-directional LSTM language model. For a given word $w_t$ in a sentence, ELMo generates its embedding by considering the entire sentence, both before and after $w_t$:

$$\mathbf{v}_{\text{ELMo}}(w_t) = f(\{\mathbf{h}_t^{(k)}\}_{k=1}^{L}),$$

where $\mathbf{h}_t^{(k)}$ is the hidden state at position $t$ in the $k$-th layer of the LSTM, and $f(\cdot)$ is a function (e.g., a weighted sum) that combines these hidden states into a single embedding.

BERT (Bidirectional Encoder Representations from Transformers) extends the idea of contextual embeddings by using a transformer architecture, which allows the model to capture bidirectional context at every layer. For a given word $w_t$ in a sentence, BERT's embedding is obtained from the final hidden layer of the transformer:

$$\mathbf{v}_{\text{BERT}}(w_t) = \text{Transformer}(w_t, \{w_1, \ldots, w_T\}),$$

where the transformer function processes the entire sentence $\{w_1, \ldots, w_T\}$ and produces context-aware embeddings for each word.

BERT embeddings are inherently contextual, meaning that the same word can have different embeddings depending on its usage in a sentence. This is achieved through the self-attention mechanism, which computes attention scores for each word relative to every other word in the sentence.

Let $w_t$ and $w_t'$ be two occurrences of the same word in different sentences $S$ and $S'$. The embeddings $\mathbf{v}_{\text{BERT}}(w_t)$ and $\mathbf{v}_{\text{BERT}}(w_t')$ satisfy

$$\mathbf{v}_{\text{BERT}}(w_t) \neq \mathbf{v}_{\text{BERT}}(w_t'),$$

if the contexts $S$ and $S'$ differ significantly. This property allows BERT to capture the polysemy and context-specific meanings of words.

Example: For the word "bank" in the sentences "I deposited money in the bank" and "The boat was near the river bank," BERT will generate different embeddings for "bank," reflecting its distinct meanings in these contexts.

Applications of Contextual Embeddings:

1. Named Entity Recognition (NER): Contextual embeddings allow for more accurate NER by considering the context in which a word appears.
2. Question Answering (QA): BERT has been particularly successful in QA tasks, where understanding the context of a question and its potential answers is crucial.

## 2.1.6   Properties of Embedding Spaces

Word embeddings map words to points in a high-dimensional vector space. The relationships between these points reflect linguistic properties, and these relationships can be analyzed using mathematical tools such as cosine similarity and Euclidean distance. Additionally, tasks like word analogy and metrics like embedding quality scores help evaluate the effectiveness of these embeddings.

### Cosine Similarity

Cosine similarity is a measure of similarity between two non-zero vectors in an inner product space. It is defined as the cosine of the angle between the two vectors,

which is equivalent to the dot product of the vectors normalized by their magnitudes. Mathematically, for two vectors $\mathbf{v}_1$ and $\mathbf{v}_2$ in a vector space $\mathbb{R}^d$, the cosine similarity $\cos(\mathbf{v}_1, \mathbf{v}_2)$ is given by

$$\cos(\mathbf{v}_1, \mathbf{v}_2) = \frac{\mathbf{v}_1 \cdot \mathbf{v}_2}{\|\mathbf{v}_1\| \|\mathbf{v}_2\|} = \frac{\sum_{i=1}^{d} v_{1,i} v_{2,i}}{\sqrt{\sum_{i=1}^{d} v_{1,i}^2} \sqrt{\sum_{i=1}^{d} v_{2,i}^2}},$$

where $\mathbf{v}_1 \cdot \mathbf{v}_2$ denotes the dot product, and $\|\mathbf{v}_1\|$ and $\|\mathbf{v}_2\|$ are the Euclidean norms (magnitudes) of $\mathbf{v}_1$ and $\mathbf{v}_2$, respectively.

Properties:

1. Range: Cosine similarity ranges from –1 to 1. A cosine similarity of 1 indicates that the vectors are identical in direction, 0 indicates orthogonality, and –1 indicates that the vectors point in opposite directions.
2. Interpretation: In the context of word embeddings, a high cosine similarity between two word vectors suggests that the words have similar meanings or appear in similar contexts.

Example: Consider the word vectors $\mathbf{v}_{king}$ and $\mathbf{v}_{queen}$. If these vectors are trained effectively, the cosine similarity $\cos(\mathbf{v}_{king}, \mathbf{v}_{queen})$ should be close to 1, reflecting their semantic similarity.

**Theorem 2.3** (Cauchy–Schwarz Inequality) *The Cauchy–Schwarz inequality underpins the cosine similarity measure. It states that for any vectors $\mathbf{v}_1, \mathbf{v}_2$ in $\mathbb{R}^d$,*

$$|\mathbf{v}_1 \cdot \mathbf{v}_2| \leq \|\mathbf{v}_1\| \|\mathbf{v}_2\|,$$

*with equality if and only if $\mathbf{v}_1$ and $\mathbf{v}_2$ are linearly dependent. This ensures that the cosine similarity is well defined and bounded between –1 and 1.*

**Euclidean Distance**

Euclidean distance is a measure of the "straight-line" distance between two points in a Euclidean space. For vectors $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^d$, the Euclidean distance $d(\mathbf{v}_1, \mathbf{v}_2)$ is defined as

$$d(\mathbf{v}_1, \mathbf{v}_2) = \|\mathbf{v}_1 - \mathbf{v}_2\| = \sqrt{\sum_{i=1}^{d} (v_{1,i} - v_{2,i})^2}.$$

In word embeddings, Euclidean distance is often used to measure the dissimilarity between words. Words with similar meanings tend to have smaller Euclidean distances between their embeddings, while words with different meanings have larger distances.

Example: Consider the word vectors $\mathbf{v}_{\text{cat}}$ and $\mathbf{v}_{\text{dog}}$. The Euclidean distance $d(\mathbf{v}_{\text{cat}}, \mathbf{v}_{\text{dog}})$ should be relatively small, reflecting the semantic similarity between these words.

**Theorem 2.4** (Pythagorean Theorem) *In a Euclidean space, if $\mathbf{v}_1$ and $\mathbf{v}_2$ are orthogonal (i.e., $\mathbf{v}_1 \cdot \mathbf{v}_2 = 0$), then the Euclidean distance satisfies*

$$d(\mathbf{v}_1, \mathbf{v}_2)^2 = \|\mathbf{v}_1\|^2 + \|\mathbf{v}_2\|^2.$$

This relationship is fundamental in understanding the geometric structure of the embedding space.

## Word Analogy

A word analogy task tests the ability of word embeddings to capture linguistic relationships. Given a pair of words $(a : b)$, the task is to find a word $d$ such that the relationship $c : d$ is analogous to $a : b$. The underlying assumption is that the vector difference $\mathbf{v}_a - \mathbf{v}_b$ should be approximately equal to $\mathbf{v}_c - \mathbf{v}_d$.

Given vectors $\mathbf{v}_a, \mathbf{v}_b, \mathbf{v}_c$, the analogy task involves solving

$$\mathbf{v}_d = \mathbf{v}_c + (\mathbf{v}_b - \mathbf{v}_a).$$

The word $d$ is then chosen as the word whose vector $\mathbf{v}_d$ is closest to this target vector, typically using cosine similarity or Euclidean distance.

Example: In the famous analogy "king is to queen as man is to woman," the task is to find $\mathbf{v}_{\text{woman}}$ given $\mathbf{v}_{\text{king}} - \mathbf{v}_{\text{queen}} \approx \mathbf{v}_{\text{man}} - \mathbf{v}_{\text{woman}}$. The word vector closest to $\mathbf{v}_{\text{queen}} - \mathbf{v}_{\text{king}} + \mathbf{v}_{\text{man}}$ is expected to be $\mathbf{v}_{\text{woman}}$.

**Remark (Vector Offset Property):** Let $\mathbf{v}_a, \mathbf{v}_b, \mathbf{v}_c, \mathbf{v}_d$ be word vectors such that the relationships $a : b$ and $c : d$ are analogous. Then, under the vector offset property, the relationship satisfies

$$\mathbf{v}_a - \mathbf{v}_b \approx \mathbf{v}_c - \mathbf{v}_d,$$

which implies

$$\mathbf{v}_d \approx \mathbf{v}_c + (\mathbf{v}_b - \mathbf{v}_a).$$

This property is crucial for solving analogy tasks and is a hallmark of well-structured embedding spaces.

## Embedding Quality Metrics

Intrinsic evaluation metrics assess the quality of word embeddings based on their geometric properties, without reference to downstream tasks. Common metrics include

1. Word Similarity: Measures the correlation between the cosine similarity of word pairs and human-judged similarity scores. Spearman's rank correlation is often used to compare the rankings of word pairs.

2. Analogy Accuracy: The percentage of correct answers in analogy tasks, such as "king is to queen as man is to woman."

3. Cluster Tightness: Evaluates how well word embeddings cluster semantically similar words together. Can be quantified using measures like intra-cluster distance and inter-cluster distance.

Extrinsic Evaluation Metrics: Extrinsic metrics evaluate embeddings based on their performance in downstream tasks, such as text classification, sentiment analysis, or named entity recognition (NER). Common metrics include

1. Task-Specific Accuracy: The accuracy of a model that uses word embeddings as features in a supervised task.

2. Transfer Learning Performance: The effectiveness of pre-trained embeddings when fine-tuned on a specific task.

Let $\{(w_i, w_j)\}$ be a set of word pairs with human-judged similarity scores $s_{ij}$. The word similarity score $\rho$ is given by the Spearman rank correlation between $s_{ij}$ and the cosine similarity of the embeddings:

$$\rho = 1 - \frac{6 \sum_{(i,j)} (r_{s_{ij}} - r_{\cos(w_i, w_j)})^2}{n(n^2 - 1)},$$

where $r_{s_{ij}}$ and $r_{\cos(w_i, w_j)}$ are the ranks of the human-judged similarity and the cosine similarity, respectively, and $n$ is the number of word pairs.

Example: Suppose we have a set of word pairs with human-judged similarity scores. We compute the cosine similarity between the corresponding word embeddings and evaluate the correlation using Spearman's rank correlation. A high correlation indicates that the embeddings capture human-perceived similarities effectively.

## 2.1.7  Advanced Embedding Techniques

Traditional word embeddings, such as those produced by Word2Vec or GloVe, treat words as atomic units, which can lead to challenges in handling out-of-vocabulary (OOV) words, morphological variations, and multilingual contexts. Advanced embedding techniques address these issues by integrating finer grained linguistic units and by mapping words across languages in a shared vector space.

**Subword Information**

Word embeddings that rely solely on word-level representations struggle with OOV words and morphological variations. For instance, words like "running," "runner," and "ran" share a common root but may have distinct embeddings in a purely word-based model. By incorporating subword information, embeddings can better capture morphological patterns and generalize to OOV words.

Let $w$ be a word in the vocabulary $\mathcal{V}$, and let $\mathcal{G}(w)$ denote the set of subword units (such as character n-grams) derived from $w$. In a subword-based model, the word embedding $\mathbf{v}_w$ is computed as the sum (or another aggregation function) of the embeddings of its constituent subwords:

$$\mathbf{v}_w = \sum_{g \in \mathcal{G}(w)} \mathbf{v}_g,$$

where $\mathbf{v}_g$ is the vector representation of the subword $g$.

Example: Consider the word "happiness." It can be decomposed into subwords like "happ," "iness," "ness," etc. The embedding for "happiness" is then the sum of the embeddings for these subwords:

$$\mathbf{v}_{\text{happiness}} = \mathbf{v}_{\text{happ}} + \mathbf{v}_{\text{iness}} + \mathbf{v}_{\text{ness}}.$$

This allows the model to capture both the root word "happy" and the suffix "ness," leading to more robust embeddings that can generalize across different forms of the word.

FastText is a prominent example of a model that incorporates subword information. The model learns embeddings not only for entire words but also for character n-grams, allowing it to generate embeddings for OOV words by summing the vectors of their n-grams.

**Theorem 2.5**  (Universal Approximation with Subwords) *Let $\mathcal{V}$ be a finite vocabulary and $\mathcal{G}(w)$ the set of subwords for each word $w \in \mathcal{V}$. The space of all possible word embeddings $\mathbb{R}^d$ can be approximated to arbitrary precision by the sum of subword embeddings $\mathbf{v}_g$ for sufficiently large $d$. Formally, for any word embedding $\mathbf{v}_w \in \mathbb{R}^d$, there exists a set of subword embeddings $\{\mathbf{v}_g\}_{g \in \mathcal{G}(w)}$ such that*

$$\mathbf{v}_w = \sum_{g \in \mathcal{G}(w)} \mathbf{v}_g + \epsilon,$$

*where $\epsilon$ is an arbitrarily small error term.*

This theorem underscores the ability of subword-based models to approximate word embeddings effectively, even for words not seen during training.

**Character-Level Embeddings**

Character-level embeddings provide an even finer granularity than subword models by representing words as sequences of individual characters. This approach is particularly useful for languages with rich morphology, for handling OOV words, and for processing noisy text (e.g., social media data) where spelling variations are common.

Let $w$ be a word composed of characters $\{c_1, c_2, \ldots, c_n\}$. The character-level embedding $\mathbf{v}_w$ is derived by applying a function $f$ to the embeddings of the individual characters $\mathbf{v}_{c_i}$:

$$\mathbf{v}_w = f(\mathbf{v}_{c_1}, \mathbf{v}_{c_2}, \ldots, \mathbf{v}_{c_n}),$$

where $f$ can be a recurrent neural network (RNN), convolutional neural network (CNN), or another appropriate function.

Example: For the word "apple," with characters $\{a, p, p, l, e\}$, the character-level embedding might be computed as

$$\mathbf{v}_{\text{apple}} = \text{RNN}(\mathbf{v}_a, \mathbf{v}_p, \mathbf{v}_p, \mathbf{v}_l, \mathbf{v}_e),$$

where the RNN processes the sequence of character embeddings to produce a single word-level embedding.

Convolutional Neural Networks (CNNs) for Character Embeddings: CNNs are often used to extract features from character sequences. The input sequence is convolved with multiple filters to produce a set of feature maps, which are then aggregated (e.g., via max-pooling) to produce the final word embedding:

$$\mathbf{v}_w = \text{max-pool}(\text{ReLU}(\mathbf{W} * \mathbf{v}_{\text{chars}} + \mathbf{b})),$$

where $\mathbf{W}$ is a filter matrix, $*$ denotes convolution, and $\mathbf{v}_{\text{chars}}$ is the matrix of character embeddings.

**Remark (Character Embedding Generalization):** Character-level models possess a strong generalization capability because they are not tied to a specific vocabulary. Given a finite character set $\mathcal{C}$ and any word $w$ composed of characters from $\mathcal{C}$, the embedding $\mathbf{v}_w$ can be constructed by the character model, ensuring that even unseen words can be embedded meaningfully. Formally, for any word $w = \{c_1, c_2, \ldots, c_n\}$,

$$\mathbf{v}_w = f(\mathbf{v}_{c_1}, \mathbf{v}_{c_2}, \ldots, \mathbf{v}_{c_n}),$$

where $f$ is a function that can generalize across the entire character set $\mathcal{C}$.

Character-level models are particularly effective in languages with complex morphology, such as Finnish or Turkish, where words can have many inflected forms. By processing the character sequence directly, these models can capture the morphological structure without requiring extensive vocabulary coverage.

**Cross-Lingual Embeddings**

Cross-lingual embeddings map words from different languages into a shared vector space, enabling direct comparison and translation between languages. These embeddings are essential for multilingual NLP tasks, such as machine translation, cross-lingual information retrieval, and multilingual sentiment analysis.

Let $\mathcal{V}^{(1)}$ and $\mathcal{V}^{(2)}$ be the vocabularies of two languages, $L_1$ and $L_2$, respectively. Cross-lingual embeddings aim to find mappings $\phi_1 : \mathcal{V}^{(1)} \to \mathbb{R}^d$ and $\phi_2 : \mathcal{V}^{(2)} \to \mathbb{R}^d$ such that words with similar meanings in $L_1$ and $L_2$ have close embeddings in the shared space $\mathbb{R}^d$.

Alignment Techniques:

1. Supervised Alignment: Given a bilingual dictionary $\{(w_i^{(1)}, w_i^{(2)})\}$ with pairs of equivalent words $w_i^{(1)} \in \mathcal{V}^{(1)}$ and $w_i^{(2)} \in \mathcal{V}^{(2)}$, the embeddings are aligned by minimizing the distance between the mapped pairs:

$$\min_{\mathbf{W}} \sum_i \|\mathbf{W}\phi_1(w_i^{(1)}) - \phi_2(w_i^{(2)})\|^2,$$

where $\mathbf{W}$ is a linear transformation matrix that aligns the embeddings of $L_1$ to the space of $L_2$.

2. Unsupervised Alignment: In the absence of a bilingual dictionary, unsupervised techniques align the embeddings by matching the geometric structures of the monolingual spaces. Techniques such as adversarial training and Procrustes alignment are used to find the optimal mapping.

**Theorem 2.6** (Existence of a Shared Embedding Space) *Given two languages $L_1$ and $L_2$, with corresponding embedding spaces $\mathbb{R}^{d_1}$ and $\mathbb{R}^{d_2}$, there exists a shared embedding space $\mathbb{R}^d$ and mappings $\phi_1 : \mathbb{R}^{d_1} \to \mathbb{R}^d$ and $\phi_2 : \mathbb{R}^{d_2} \to \mathbb{R}^d$ such that semantically equivalent words from $L_1$ and $L_2$ are close in the shared space. Formally, for equivalent words $w_1$ and $w_2$,*

$$\|\phi_1(w_1) - \phi_2(w_2)\| \le \epsilon,$$

*where $\epsilon$ is a small error term that depends on the alignment technique.*

Example: In a cross-lingual setting, the word "apple" in English might be mapped close to "manzana" in Spanish. The shared embedding space allows direct comparisons between words from different languages, facilitating tasks like cross-lingual sentiment analysis.

Applications:

1. Machine Translation: Cross-lingual embeddings are used to map source-language words to target-language words, improving translation quality by leveraging shared semantic structures.

2. Cross-lingual Information Retrieval: These embeddings enable the retrieval of documents in one language based on queries in another, expanding the reach of information access across languages.

## 2.2   Positional Encoding

In transformer architectures, positional encoding is a crucial component that allows the model to incorporate information about the order of tokens in a sequence. Since transformers lack the sequential inductive bias of RNNs, positional encoding provides a mechanism for capturing the relative and absolute positions of tokens, enabling the model to understand the structure of the input data. This section explores the need for positional encoding, its mathematical formulation, and various techniques for implementing it.

### 2.2.1   Need for Positional Encoding

Transformers process input sequences in parallel, unlike RNNs or CNNs, which have an inherent notion of sequence order due to their sequential or hierarchical structure. However, language and many other types of data are inherently sequential, where the order of elements carries significant meaning. For instance, in the sentence "The cat sat on the mat," the meaning depends on the specific order of the words. Without positional information, a transformer would treat the sentence as a bag of words, losing the sequential context that is critical for understanding.

Given a sequence of tokens $\{x_1, x_2, \ldots, x_n\}$, a transformer model processes these tokens independently and simultaneously. To enable the model to recognize the position of each token, we need to inject positional information into the model's input. This is accomplished by adding or concatenating a positional encoding vector to the input embeddings.

Mathematically, let $\mathbf{E} = [\mathbf{e}_1, \mathbf{e}_2, \ldots, \mathbf{e}_n]$ represent the input embeddings, where $\mathbf{e}_i$ is the embedding of token $x_i$. The positional encoding adds a positional vector $\mathbf{p}_i$ to each $\mathbf{e}_i$, resulting in a modified input:

$$\mathbf{E}' = [\mathbf{e}_1 + \mathbf{p}_1, \mathbf{e}_2 + \mathbf{p}_2, \ldots, \mathbf{e}_n + \mathbf{p}_n].$$

This allows the transformer to use positional information when processing the sequence.

Example: Consider the sequence "I saw the man with the telescope." The meaning of the sentence depends on whether "with the telescope" modifies "saw" or "man." Without positional encoding, a transformer might not be able to disambiguate these meanings because it lacks information about the relative positions of the tokens.

## *2.2.2   Mathematical Formulation*

Positional encoding can be implemented in several ways, each with its own mathematical formulation. The choice of encoding affects how the model perceives positional relationships within the input sequence. Below, we explore three common approaches: sinusoidal positional encoding, learned positional encoding, and relative positional encoding.

### Sinusoidal Positional Encoding

The sinusoidal positional encoding method, introduced in the original Transformer paper by [7], uses sinusoidal functions to encode positions. For a sequence of length $n$, the positional encoding for position $i$ is a vector $\mathbf{p}_i$ defined as

$$\mathbf{p}_i = \left[ \sin\left( \frac{i}{10000^{2j/d}} \right), \cos\left( \frac{i}{10000^{2j/d}} \right) \right]_{j=0}^{d/2-1},$$

where $i$ is the position index, $j$ is the dimension index, and $d$ is the dimensionality of the positional encoding.

The encoding is constructed so that each dimension of the positional vector corresponds to a sinusoid with a different frequency.

Properties:

1. Periodicity: The sinusoidal functions are periodic, which means that the positional encoding is able to capture relative positions. Specifically, for any two positions $i$ and $k$, the difference $\mathbf{p}_i - \mathbf{p}_k$ encodes the relative distance between $i$ and $k$.

2. Smoothness: The encoding changes smoothly with position, which is beneficial for capturing the continuous nature of sequential data.

3. Interpretability: The use of sine and cosine functions of varying frequencies ensures that each position has a unique encoding, and the relative differences between positions are captured in a manner that is independent of the absolute sequence length.

**Remark (Orthogonality of Encodings):** For any fixed position $i$, the sinusoidal positional encoding $\mathbf{p}_i$ is nearly orthogonal to the encoding of a position $i + k$ for large $k$. Specifically, the dot product $\mathbf{p}_i \cdot \mathbf{p}_{i+k}$ decreases as $k$ increases. This property ensures that the positional encodings are distinct and help the model differentiate between different positions in the sequence.

Example: For a sequence of length 10 and embedding dimension $d = 16$, the positional encoding for the first position might be

$$\mathbf{p}_1 = [\sin(1), \cos(1), \sin(1/10000^{1/8}), \cos(1/10000^{1/8}), \dots]$$

The encoding for the second position would have slightly different values due to the changing argument in the sine and cosine functions, ensuring that each position has a unique encoding.

**Learned Positional Encoding**

In learned positional encoding, the positional vectors $\mathbf{p}_i$ are treated as parameters that are learned during the training process, similar to how the word embeddings $\mathbf{e}_i$ are learned. Formally, let $\mathbf{P} \in \mathbb{R}^{n \times d}$ be the matrix of positional encodings, where each row $\mathbf{p}_i$ corresponds to the positional encoding for position $i$:

$$\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \ldots, \mathbf{p}_n]^\top.$$

During training, the matrix $\mathbf{P}$ is optimized to minimize the model's loss function, along with other model parameters.

Properties:

1. Flexibility: Unlike sinusoidal encoding, learned positional encoding does not impose any predefined structure on the encodings. The model can learn the optimal positional encodings based on the task and dataset.

2. Task-Specific Adaptation: Since the positional encodings are learned, they can adapt to the specific characteristics of the training data, potentially capturing more complex positional relationships.

3. Scalability: While learned encodings are flexible, they require the model to learn additional parameters, which can increase the computational and memory requirements, especially for long sequences.

Example: In a machine translation task, the model might learn that certain positions in the input sequence are more important for aligning with positions in the output sequence. The learned positional encodings could reflect this by assigning higher weights or more distinct vectors to those positions.

**Remark (Expressiveness of Learned Positional Encoding):** Given a sufficiently large model and enough training data, learned positional encoding can approximate any smooth function of position. Specifically, for any smooth function $f : \mathbb{R} \to \mathbb{R}^d$, there exists a set of learned positional encodings $\mathbf{P}$ such that

$$\mathbf{p}_i \approx f(i),$$

for all positions $i$.

This highlights the expressiveness and adaptability of learned positional encoding, which can theoretically match or exceed the performance of fixed encodings.

**Relative Positional Encoding**

Relative positional encoding encodes the relative position $k = j - i$ between two tokens at positions $i$ and $j$ rather than their absolute positions. This approach is particularly useful in tasks where the model needs to focus on the relative order of tokens rather than their exact positions. Let $\mathbf{r}_k$ represent the encoding of the relative

position $k$. The model then incorporates these relative encodings into the attention mechanism, modifying the attention score $\alpha_{ij}$ between tokens $x_i$ and $x_j$:

$$\alpha_{ij} = \mathbf{q}_i^\top (\mathbf{k}_j + \mathbf{r}_{j-i}),$$

where $\mathbf{q}_i$ and $\mathbf{k}_j$ are the query and key vectors associated with tokens $x_i$ and $x_j$, respectively.

Properties:

1. Translation Invariance: Relative positional encoding naturally handles shifts in the input sequence, making the model invariant to translations of the input.

2. Efficiency in Long Sequences: By focusing on relative rather than absolute positions, the model can efficiently manage long sequences, as the number of unique relative positions is limited by the sequence length.

3. Simplicity: Relative positional encoding can be simpler to implement in certain models, as it does not require a separate positional encoding matrix for each position in the sequence.

**Theorem 2.7** (Equivariance of Relative Positional Encoding) *Let $\mathbf{r}_k$ be the relative positional encoding for position difference $k$. Then, for any shift $s$, the relative positional encodings satisfy*

$$\mathbf{r}_{k+s} = \mathbf{r}_k + \mathbf{r}_s,$$

*ensuring that the relative positional information is preserved under translations of the input sequence.*

Example: In a document where the order of sentences matters more than their absolute positions, relative positional encoding allows the model to focus on the relationships between sentences rather than their positions within the document. For instance, in tasks like summarization or dialog generation, understanding relative order can be crucial.

## 2.2.3 Properties and Analysis

The properties of positional encodings significantly influence how a model processes and understands sequential data. Through frequency analysis, we can explore how different frequencies in sinusoidal encodings affect the model's ability to capture patterns. Additionally, by examining the impact on model performance and interpretability, we gain insights into how positional encodings contribute to the overall effectiveness and transparency of transformer models.

**Frequency Analysis**

Sinusoidal positional encodings rely on encoding positions using sine and cosine functions with varying frequencies. The choice of these frequencies plays a critical role in how well the model can differentiate between positions, especially when sequences are long. The encoding for a position $i$ and dimension $j$ in a $d$-dimensional space is given by

$$\mathbf{p}_i^{(2j)} = \sin\left(\frac{i}{10000^{2j/d}}\right), \quad \mathbf{p}_i^{(2j+1)} = \cos\left(\frac{i}{10000^{2j/d}}\right),$$

where $j$ ranges from 0 to $\frac{d}{2} - 1$.

The parameter $10000^{2j/d}$ controls the frequency of the sinusoidal functions. Lower dimensions (small $j$) correspond to higher frequencies, capturing fine-grained positional differences, while higher dimensions (large $j$) correspond to lower frequencies, capturing broader positional relationships.

We can interpret sinusoidal positional encoding as a Fourier series representation, where the position $i$ is decomposed into a sum of sinusoids with different frequencies. The Fourier series is particularly effective for representing periodic signals, which aligns with the periodic nature of language and other sequential data.

**Remark (Fourier Representation of Positional Encodings):** Let $\mathbf{p}_i$ be the positional encoding for position $i$. Then, the encoding can be viewed as a truncated Fourier series:

$$\mathbf{p}_i = \sum_{k=0}^{d/2-1} \left[ a_k \cos\left(\frac{2\pi ki}{L}\right) + b_k \sin\left(\frac{2\pi ki}{L}\right) \right],$$

where $a_k$ and $b_k$ are coefficients that depend on the dimension $j$, and $L$ is the length of the sequence. This series captures both high-frequency and low-frequency components of the sequence positions.

Implications:

1. Resolution of Positional Differences: Lower dimensions, with higher frequencies, allow the model to resolve small positional differences, making them suitable for capturing fine-grained word order. Higher dimensions capture broader patterns and can identify when tokens are in roughly the same position within the sequence.

2. Periodicity: The use of sine and cosine functions inherently introduces periodicity, meaning that the model can capture recurring patterns in the data, such as those found in cyclical sequences (e.g., periodic time series data).

Example: For a sequence length $L = 100$ and embedding dimension $d = 16$, the lower frequency components (corresponding to small $j$) would allow the model to differentiate between positions 1 and 2, while the higher frequency components would differentiate between positions 50 and 100.

**Impact on Model Performance**

The effectiveness of positional encoding directly influences the model's performance on tasks that require understanding the order and structure of input data. For example, in tasks such as machine translation or text summarization, capturing the correct sequence of tokens is essential for generating accurate outputs.

Empirically, transformers using sinusoidal positional encodings have demonstrated strong performance across a range of NLP tasks. However, the choice of positional encoding type—sinusoidal versus learned—can have different impacts depending on the nature of the task and the characteristics of the data.

The impact of positional encoding on model performance can be assessed using standard NLP metrics such as BLEU score for translation, ROUGE score for summarization, and accuracy for classification tasks. These metrics reflect how well the model preserves and interprets the sequential information in the data.

**Remark (Generalization Bounds with Positional Encoding):** Let $\mathcal{H}_{\text{PE}}$ denote the hypothesis class of transformers with positional encoding. The generalization error $R(h)$ of a model $h \in \mathcal{H}_{\text{PE}}$ is bounded by

$$R(h) \leq \hat{R}_n(h) + \mathcal{O}\left(\sqrt{\frac{\text{Complexity}(\mathcal{H}_{\text{PE}})}{n}}\right),$$

where $\hat{R}_n(h)$ is the empirical risk, and $\text{Complexity}(\mathcal{H}_{\text{PE}})$ is a measure of the model's complexity, including the influence of positional encoding. This bound shows that positional encoding contributes to the model's ability to generalize by influencing the complexity of the hypothesis space.

Example: In a translation task, a model with well-designed positional encodings might exhibit a lower BLEU score when these encodings are removed or poorly tuned, indicating the critical role of positional information in maintaining translation quality.

**Interpretability of Positional Encodings**

The interpretability of positional encodings stems from their mathematical structure. Sinusoidal encodings, due to their periodic nature, are particularly interpretable, as each dimension of the encoding has a clear geometric meaning related to frequency components of the position.

Positional encodings can be visualized as heatmaps, where each row corresponds to a position $i$ and each column to a dimension $j$. These visualizations often reveal regular patterns that reflect how the model perceives positional information.

The positional encodings $\mathbf{p}_i$ for different positions $i$ and $j$ are nearly orthogonal for large positional differences, which ensures that the model can distinguish between distant positions effectively. Formally,

$$\mathbf{p}_i \cdot \mathbf{p}_j \approx 0 \quad \text{for } |i - j| \gg 0.$$

This orthogonality is a key factor in making positional encodings interpretable, as it ensures that the model can separate and process different positions distinctly.

Implications for model interpretability:

1. Traceability: The clear mathematical structure of sinusoidal encodings allows researchers to trace how positional information flows through the model, aiding in the interpretation of model outputs.

2. Insights into model behavior: By analyzing the patterns in the positional encodings, one can gain insights into how the model attends to different parts of the input sequence. This can be particularly useful in tasks like attention visualization, where understanding the positional contributions to the attention mechanism is crucial.

Example: In tasks like document classification, where the order of sentences can influence the final prediction, interpretability of positional encodings can help identify whether the model is correctly attending to the beginning, middle, or end of the document, providing transparency in the decision-making process.

## 2.2.4   Applications and Variations

Positional encoding is adapted to different types of data beyond textual sequences. Each adaptation retains the core principles of encoding positional or structural information while modifying the encoding mechanism to suit the specific data type. Below, we delve into how positional encoding is applied in vision transformers, temporal sequences, and hierarchical data structures.

### Positional Encoding in Vision Transformers

Vision Transformers (ViTs) apply transformer architectures to image data, which differs significantly from text due to its two-dimensional nature. Unlike sequences of words, images are composed of pixels arranged in a grid. To effectively use transformers in vision tasks, it is necessary to encode the spatial positions of image patches, preserving the spatial relationships inherent in images.

In vision transformers, an image is divided into patches, and each patch is treated as a token in a sequence. Let $\mathcal{I}$ be an image of size $H \times W$ (height $H$ and width $W$), divided into $N$ patches, where each patch is of size $P \times P$. The sequence length is $N = \frac{HW}{P^2}$, and each patch is flattened into a vector $\mathbf{p}_i \in \mathbb{R}^{P^2 \cdot C}$, where $C$ is the number of channels (e.g., RGB).

To incorporate spatial information, positional encodings $\mathbf{e}_i^{\text{pos}}$ are added to the patch embeddings $\mathbf{p}_i$:

$$\mathbf{z}_i = \mathbf{p}_i + \mathbf{e}_i^{\text{pos}},$$

where $\mathbf{z}_i$ is the input to the transformer. The positional encodings $\mathbf{e}_i^{\text{pos}}$ must capture both the row and column positions of the patches in the original image.

2D Positional Encoding: The positional encoding can be extended to two dimensions, where separate encoding vectors are used for the row and column indices:

$$\mathbf{e}_i^{\text{pos}} = \mathbf{e}_{\text{row}(i)} + \mathbf{e}_{\text{col}(i)},$$

where $\text{row}(i)$ and $\text{col}(i)$ are the row and column indices of the $i$-th patch. Each encoding vector $\mathbf{e}_{\text{row}}$ and $\mathbf{e}_{\text{col}}$ can be sinusoidal, learned, or based on other techniques suitable for encoding two-dimensional positions.

Let $\mathbf{e}_{\text{row}(i)}$ and $\mathbf{e}_{\text{col}(i)}$ be the positional encodings for the row and column indices of the $i$-th patch in a grid. The resulting 2D positional encoding $\mathbf{e}_i^{\text{pos}}$ uniquely identifies each position within the grid, ensuring that

$$\mathbf{e}_i^{\text{pos}} \neq \mathbf{e}_j^{\text{pos}} \quad \text{for } i \neq j.$$

This ensures that each patch retains its unique position within the image, enabling the model to learn spatial relationships effectively.

Example: For a $224 \times 224$ image with $16 \times 16$ patches, there are $14 \times 14 = 196$ patches. Each patch position is encoded using the 2D positional encoding, allowing the vision transformer to process the image as a sequence of patches while preserving the spatial structure.

### Temporal Positional Encoding in Time Series

Time series data, unlike static images or text, has a continuous and often irregular temporal dimension. For tasks such as forecasting, anomaly detection, or temporal classification, it is crucial to encode temporal information that reflects the progression and spacing of events in time.

Given a time series $\{x_1, x_2, \ldots, x_T\}$, where $x_t$ represents the value at time $t$, temporal positional encoding must reflect both the sequence order and the potentially varying time intervals between observations.

Sinusoidal Temporal Encoding: A straightforward extension of sinusoidal positional encoding can be applied to time series by treating each time step as a position:

$$\mathbf{p}_t^{(2j)} = \sin\left(\frac{t}{10000^{2j/d}}\right), \quad \mathbf{p}_t^{(2j+1)} = \cos\left(\frac{t}{10000^{2j/d}}\right),$$

where $t$ is the time index and $j$ indexes the dimension. This encoding assumes uniform time intervals.

Continuous Temporal Encoding: For time series with non-uniform intervals, the encoding can be adjusted by directly encoding the actual time $t_i$ of each observation:

$$\mathbf{p}_{t_i}^{(2j)} = \sin\left(\frac{t_i}{\tau^{2j/d}}\right), \quad \mathbf{p}_{t_i}^{(2j+1)} = \cos\left(\frac{t_i}{\tau^{2j/d}}\right),$$

where $\tau$ is a scaling factor chosen based on the range of time values $t_i$. This allows the encoding to capture non-uniform temporal intervals.

For a time series with non-uniform intervals, let $t_i$ be the time of the $i$-th observation. The continuous temporal encoding $\mathbf{p}_{t_i}$ provides a representation that preserves the temporal order and spacing of observations. For any two observations at times $t_i$ and $t_j$:

$$\|\mathbf{p}_{t_i} - \mathbf{p}_{t_j}\| \propto |t_i - t_j|,$$

ensuring that the temporal distance between observations is reflected in their positional encodings.

Example: Consider a time series of daily stock prices with missing weekends. A continuous temporal encoding would assign different positional vectors to Friday and Monday, reflecting the 2-day gap, unlike a simple index-based encoding.

### Hierarchical Positional Encoding

In many datasets, information is organized hierarchically, such as paragraphs within a document, sections within a book, or folders within a filesystem. Hierarchical positional encoding provides a way to encode multi-level structures, enabling transformers to process data at multiple levels of granularity.

Let the data be organized into a hierarchy with $L$ levels, where each level $\ell$ contains $n_\ell$ units. For example, in a document, $L$ could represent sections, paragraphs, and sentences. Each unit at level $\ell$ is encoded with a positional vector $\mathbf{e}_{\ell,i}$, where $i$ indexes the units at that level.

Hierarchical Encoding Structure: The overall positional encoding for a unit at the deepest level (e.g., a word) is the sum of the positional encodings from all levels:

$$\mathbf{p}_{\text{word}} = \sum_{\ell=1}^{L} \mathbf{e}_{\ell,i_\ell},$$

where $i_\ell$ indexes the position within the $\ell$-th level.

Properties:

1. Multi-Level Information: Hierarchical encoding captures information at multiple levels, allowing the model to distinguish not only between positions within a sequence but also between different hierarchical contexts.

2. Flexibility: The encoding scheme can be adapted to different hierarchical structures, whether they are balanced (e.g., binary trees) or unbalanced (e.g., document structures).

**Remark (Uniqueness of Hierarchical Positional Encoding):** For a hierarchical structure with $L$ levels, let $\mathbf{e}_{\ell,i_\ell}$ be the positional encoding at level $\ell$. The sum of these encodings uniquely identifies each position within the hierarchy, ensuring that

$$\mathbf{p}_i \neq \mathbf{p}_j \quad \text{for any } i \neq j \text{ at the same or different levels.}$$

This guarantees that each unit in the hierarchy is uniquely represented, preserving the structure's integrity.

Example: In a book, sentences within the same paragraph might have similar lower-level positional encodings, but their higher-level encodings will differ based on their section or chapter, allowing the model to understand both the local and global contexts.

## 2.3 Integration of Word Embeddings and Positional Encoding

The integration of word embeddings and positional encoding is a critical aspect of transformer architectures. This integration enables transformers to effectively process sequences by embedding both the content and the positional information of tokens, allowing the model to capture both the semantic and structural aspects of the input data. In this section, we explore how word embeddings and positional encoding are combined within the transformer framework, focusing on the embedding layer and the positional encoding layer.

### 2.3.1 Combining Embeddings and Positional Encoding in Transformers

In transformer models, each input token is represented by a combination of its word embedding and its positional encoding. This combined representation allows the model to retain information about both the identity and the order of tokens within the sequence. The mathematical formulation and integration of these components are essential to the success of the transformer architecture.

**Embedding Layer in Transformers**

The embedding layer in a transformer maps each token in the input sequence to a high-dimensional vector space, where similar tokens have similar representations. Given a vocabulary $\mathcal{V}$ with $|\mathcal{V}|$ tokens, the embedding layer is defined by an embedding matrix $\mathbf{E} \in \mathbb{R}^{|\mathcal{V}| \times d}$, where $d$ is the dimensionality of the embedding space. For a sequence of tokens $\{x_1, x_2, \ldots, x_n\}$, the embedding layer produces a sequence of embeddings $\{\mathbf{e}_1, \mathbf{e}_2, \ldots, \mathbf{e}_n\}$, where

$$\mathbf{e}_i = \mathbf{E}[x_i],$$

and $\mathbf{E}[x_i]$ denotes the embedding vector corresponding to token $x_i$.

Properties of the Embedding Layer:

1. Semantic Representation: The embedding layer captures semantic similarities between tokens, such that tokens with similar meanings are mapped to nearby points in the embedding space.

2. Learnability: The embedding matrix $\mathbf{E}$ is learned during the training process, allowing the model to adapt the embeddings to the specific task and dataset.

3. Dimensionality: The choice of the embedding dimensionality $d$ affects the expressiveness and capacity of the model. Higher dimensions can capture more nuanced semantic relationships but also increase the model's complexity.

**Remark (Universal Approximation with Embeddings):** Given a sufficiently large embedding dimension $d$, the embedding layer can approximate any continuous mapping from tokens to a high-dimensional space. Formally, for any continuous function $f : \mathcal{V} \to \mathbb{R}^d$, there exists an embedding matrix $\mathbf{E}$ such that

$$\mathbf{E}[x_i] \approx f(x_i),$$

for all tokens $x_i \in \mathcal{V}$. This theorem highlights the ability of the embedding layer to capture complex semantic relationships within the vocabulary.

Example: Consider the words "cat" and "dog." In a well-trained embedding space, $\mathbf{E}[cat]$ and $\mathbf{E}[dog]$ will be close to each other, reflecting the semantic similarity between these tokens.

**Positional Encoding Layer in Transformers**

The positional encoding layer introduces positional information into the sequence of embeddings produced by the embedding layer. For a sequence of length $n$, the positional encoding layer generates a sequence of positional vectors $\{\mathbf{p}_1, \mathbf{p}_2, \ldots, \mathbf{p}_n\}$, where each positional vector $\mathbf{p}_i$ corresponds to the position of the $i$-th token. The final input to the transformer model is obtained by adding the positional encoding to the embedding:

$$\mathbf{z}_i = \mathbf{e}_i + \mathbf{p}_i,$$

where $\mathbf{z}_i$ is the combined representation for the $i$-th token, capturing both its semantic meaning (through $\mathbf{e}_i$) and its position within the sequence (through $\mathbf{p}_i$).

Properties of the Positional Encoding Layer:

1. Order Sensitivity: The addition of positional encodings ensures that the model is sensitive to the order of tokens, allowing it to differentiate between sequences with the same tokens in different orders.

2. Non-Interference: The positional encoding is typically designed so that it does not dominate the embedding. For example, in sinusoidal encoding, the encoding values are scaled to be of similar magnitude to the embeddings, ensuring that positional and semantic information are balanced.

3. Fixed versus Learned: The positional encoding can be either fixed (e.g., sinusoidal) or learned. Fixed encodings are predefined and do not change during training, while learned encodings are optimized along with other model parameters.

**Remark (Injectivity of Combined Representation):** Let $\mathbf{z}_i = \mathbf{e}_i + \mathbf{p}_i$ be the combined representation of a token at position $i$. The injectivity of the combined representation function ensures that no two different positions $i$ and $j$ produce the same combined vector:

$$\mathbf{z}_i \neq \mathbf{z}_j \quad \text{for } i \neq j.$$

This guarantees that the model can uniquely identify each token based on both its content and position.

Example: For the sentence "The quick brown fox," the embeddings $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4\}$ might represent the words "The," "quick," "brown," and "fox." The positional encodings $\{\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3, \mathbf{p}_4\}$ ensure that the sequence "The quick brown fox" is distinguished from "fox brown quick The."

The integration of word embeddings and positional encodings is crucial for the transformer's ability to process sequential data. Without positional encoding, the model would lose the ability to understand the order of tokens, reducing its effectiveness in tasks that rely on sequence structure, such as translation, text summarization, and language modeling.

Example: Consider the embeddings $\mathbf{e}_i$ for a sequence of four tokens, along with their corresponding positional encodings $\mathbf{p}_i$. The final input sequence to the transformer could be represented as

$$\mathbf{z}_1 = \mathbf{e}_1 + \mathbf{p}_1,$$
$$\mathbf{z}_2 = \mathbf{e}_2 + \mathbf{p}_2,$$
$$\mathbf{z}_3 = \mathbf{e}_3 + \mathbf{p}_3,$$
$$\mathbf{z}_4 = \mathbf{e}_4 + \mathbf{p}_4.$$

Each $\mathbf{z}_i$ is then processed by the transformer layers, which can now leverage both the content and positional information to make predictions.

## 2.3.2   Impact on Attention Mechanisms

The attention mechanism in transformers allows the model to weigh the importance of different tokens in a sequence relative to each other. Without positional encoding, the model would treat all tokens as independent and identically distributed, losing the crucial information about their order in the sequence. By integrating positional information, the attention mechanism can differentiate between tokens based on their positions, thereby improving the model's ability to capture the underlying structure of the data.

### Enhancing Self-Attention with Positional Information

Self-attention is a mechanism that computes a weighted sum of all tokens in a sequence, where the weights (attention scores) are determined by the relevance of each token to the others. For a sequence of tokens $\{x_1, x_2, \ldots, x_n\}$ with corresponding combined embeddings $\{\mathbf{z}_1, \mathbf{z}_2, \ldots, \mathbf{z}_n\}$, the self-attention mechanism computes the output $\mathbf{y}_i$ for each token $x_i$ as follows:

$$\mathbf{y}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{z}_j,$$

where $\alpha_{ij}$ is the attention score between tokens $x_i$ and $x_j$. The attention score is computed using the scaled dot-product attention:

$$\alpha_{ij} = \frac{\exp\left(\mathbf{q}_i^\top \mathbf{k}_j\right)}{\sum_{k=1}^{n} \exp\left(\mathbf{q}_i^\top \mathbf{k}_k\right)},$$

where $\mathbf{q}_i = \mathbf{W}_q \mathbf{z}_i$ and $\mathbf{k}_j = \mathbf{W}_k \mathbf{z}_j$ are the query and key vectors for tokens $x_i$ and $x_j$, respectively, and $\mathbf{W}_q$ and $\mathbf{W}_k$ are learned weight matrices.

Positional encoding $\mathbf{p}_i$ influences the attention mechanism by modifying the embeddings $\mathbf{z}_i = \mathbf{e}_i + \mathbf{p}_i$, where $\mathbf{e}_i$ is the word embedding and $\mathbf{p}_i$ is the positional encoding. This addition ensures that the attention mechanism considers both the content of the tokens and their positions in the sequence.

The positional encoding affects the dot product $\mathbf{q}_i^\top \mathbf{k}_j$ by introducing terms that depend on the positions $i$ and $j$. This alters the attention scores $\alpha_{ij}$, making them sensitive to the relative positions of the tokens. For example, if $x_j$ is farther from $x_i$ in the sequence, the positional term may decrease the attention score, reducing the influence of distant tokens on the current token's representation.

**Remark (Positional Sensitivity in Self-Attention):**  Let $\mathbf{p}_i$ and $\mathbf{p}_j$ be the positional encodings for tokens $x_i$ and $x_j$. The attention score $\alpha_{ij}$ in self-attention is influenced by the difference $\mathbf{p}_i - \mathbf{p}_j$. Formally, the attention score can be expressed as

$$\alpha_{ij} = \frac{\exp\left((\mathbf{e}_i + \mathbf{p}_i)^\top \mathbf{W}_q^\top \mathbf{W}_k (\mathbf{e}_j + \mathbf{p}_j)\right)}{\sum_{k=1}^n \exp\left((\mathbf{e}_i + \mathbf{p}_i)^\top \mathbf{W}_q^\top \mathbf{W}_k (\mathbf{e}_k + \mathbf{p}_k)\right)},$$

which shows that the attention mechanism considers not only the semantic content $\mathbf{e}_i^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{e}_j$ but also the positional difference $\mathbf{p}_i^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{p}_j$.

Example: Consider a sentence "The cat sat on the mat." Without positional encoding, the model might struggle to differentiate between similar sentences like "The mat sat on the cat." By incorporating positional encoding, the attention mechanism can assign higher weights to tokens that follow the expected word order, reducing the likelihood of such confusion.

**Role in Multi-head Attention**

Multi-head attention extends the self-attention mechanism by allowing the model to focus on different parts of the sequence simultaneously. Instead of computing a single set of attention scores, multi-head attention computes multiple sets of scores, each corresponding to a different "head" in the model. Given $h$ attention heads, the multi-head attention output for a token $x_i$ is computed as

$$\mathbf{y}_i^{\text{multi}} = \mathbf{W}_o \left[ \mathbf{y}_i^{(1)} \| \mathbf{y}_i^{(2)} \| \cdots \| \mathbf{y}_i^{(h)} \right],$$

where $\mathbf{y}_i^{(k)}$ is the output of the $k$-th attention head, $\|$ denotes concatenation, and $\mathbf{W}_o$ is a learned output matrix.

Each attention head can capture different aspects of the token relationships, and positional encoding ensures that each head remains sensitive to the sequence order. For example, one head might focus on short-range dependencies, while another might capture long-range dependencies, with positional encoding providing the necessary context.

The incorporation of positional encoding into multi-head attention affects the calculation of attention scores for each head:

$$\alpha_{ij}^{(k)} = \frac{\exp\left((\mathbf{W}_q^{(k)} \mathbf{z}_i)^\top \left(\mathbf{W}_k^{(k)} \mathbf{z}_j\right)\right)}{\sum_{l=1}^n \exp\left((\mathbf{W}_q^{(k)} \mathbf{z}_i)^\top \left(\mathbf{W}_k^{(k)} \mathbf{z}_l\right)\right)},$$

where $\mathbf{W}_q^{(k)}$ and $\mathbf{W}_k^{(k)}$ are the query and key weight matrices for the $k$-th head. The positional encodings $\mathbf{p}_i$ and $\mathbf{p}_j$ influence these scores, allowing each head to focus on different positional aspects of the sequence.

Let $\alpha_{ij}^{(k)}$ be the attention score for the $k$-th head in a multi-head attention mechanism. The positional encoding $\mathbf{p}_i$ ensures that each head can focus on different

positional relationships, leading to diverse attention patterns. Formally, for different heads $k$ and $l$:

$$\alpha_{ij}^{(k)} \neq \alpha_{ij}^{(l)} \quad \text{if } \mathbf{W}_q^{(k)} \neq \mathbf{W}_q^{(l)} \text{ or } \mathbf{W}_k^{(k)} \neq \mathbf{W}_k^{(l)}.$$

This diversity allows the model to capture a wide range of positional and semantic relationships within the sequence.

Example: In a translation task, one attention head might focus on the subject of a sentence, while another might focus on the verb. Positional encoding ensures that each head understands the roles of tokens based on their positions, allowing the model to generate more accurate translations.

## 2.4   Mathematical Analysis and Performance Metrics

The evaluation of embedding quality is critical to understanding how well word embeddings and positional encodings capture the underlying structure of the data. By employing mathematical analysis and well-defined performance metrics, we can assess the effectiveness of these embeddings in various tasks. This section explores both intrinsic and extrinsic evaluation methods, providing a comprehensive framework for evaluating embedding quality.

### 2.4.1   Evaluation of Embedding Quality

Word embeddings and positional encodings are fundamental to the performance of models in NLP and other domains. To assess the quality of these embeddings, we employ two main categories of evaluation: intrinsic methods, which assess the embeddings independently of specific tasks, and extrinsic methods, which evaluate embeddings based on their performance in downstream tasks.

#### Intrinsic Evaluation Methods

Intrinsic evaluation methods focus on analyzing the properties of embeddings without reference to any specific application. These methods are based on the assumption that good embeddings should capture meaningful relationships between words or other tokens in a way that reflects semantic or syntactic properties.

One of the most common intrinsic evaluation methods involves measuring the similarity between word pairs using cosine similarity or Euclidean distance. Given two word vectors $\mathbf{v}_1$ and $\mathbf{v}_2$, the cosine similarity $\cos(\mathbf{v}_1, \mathbf{v}_2)$ is defined as

$$\cos(\mathbf{v}_1, \mathbf{v}_2) = \frac{\mathbf{v}_1 \cdot \mathbf{v}_2}{\|\mathbf{v}_1\| \|\mathbf{v}_2\|},$$

where $\|\mathbf{v}_1\|$ and $\|\mathbf{v}_2\|$ are the Euclidean norms of $\mathbf{v}_1$ and $\mathbf{v}_2$, respectively. Higher cosine similarity indicates greater semantic similarity.

For analogy tasks, the evaluation is based on how well the embeddings solve word analogy problems of the form "A is to B as C is to D." The task is to find a word $D$ such that

$$\mathbf{v}_B - \mathbf{v}_A \approx \mathbf{v}_D - \mathbf{v}_C.$$

The quality of the embeddings is assessed by the proportion of correct answers.

For any word vectors $\mathbf{v}_A, \mathbf{v}_B, \mathbf{v}_C$, and $\mathbf{v}_D$ that satisfy a semantic analogy, the vector offset property holds

$$\mathbf{v}_B - \mathbf{v}_A \approx \mathbf{v}_D - \mathbf{v}_C.$$

This property allows us to assess the semantic consistency of embeddings through analogy tasks.

Another intrinsic method involves clustering word embeddings and measuring the coherence of the clusters. Words that are semantically similar should be grouped together in the embedding space. Clustering metrics such as intra-cluster variance and inter-cluster separation are used to quantify the quality of these clusters.

Given a clustering $\mathcal{C} = \{C_1, C_2, \ldots, C_k\}$ of word vectors, the intra-cluster variance $V_{\text{intra}}$ is defined as

$$V_{\text{intra}} = \frac{1}{k} \sum_{i=1}^{k} \frac{1}{|C_i|} \sum_{\mathbf{v}_j \in C_i} \|\mathbf{v}_j - \mathbf{c}_i\|^2,$$

where $\mathbf{c}_i$ is the centroid of cluster $C_i$.

Example: If we cluster word embeddings for "dog," "cat," "wolf," and "car," we expect "dog," "cat," and "wolf" to form one cluster, while "car" forms another. Low intra-cluster variance and high inter-cluster separation indicate that the embeddings correctly capture the semantic distinctions.

**Extrinsic Evaluation Methods**

Extrinsic evaluation methods assess the quality of embeddings based on their performance in downstream tasks. The assumption is that better embeddings will lead to improved performance in tasks such as text classification, named entity recognition (NER), machine translation, and question answering.

For each downstream task, the embeddings are used as input features for a model, and the performance of the model is measured using task-specific metrics. For example, in text classification, accuracy, precision, recall, and F1-score are commonly used metrics. For NER, the model's ability to correctly identify and classify named entities is evaluated.

Let $T$ be a task, and let $\mathcal{M}_T(\mathbf{E})$ denote the performance metric for task $T$ when using embeddings $\mathbf{E}$. The goal of extrinsic evaluation is to maximize $\mathcal{M}_T(\mathbf{E})$ across a set of tasks:

$$\mathbf{E}^* = \arg\max_{\mathbf{E}} \sum_T \mathcal{M}_T(\mathbf{E}).$$

Transfer learning is another important aspect of extrinsic evaluation. In this scenario, embeddings are pre-trained on a large corpus and then fine-tuned on a specific downstream task. The quality of the embeddings is assessed by how well they transfer to the new task.

Let $\mathcal{L}_{\text{pre}}$ and $\mathcal{L}_{\text{finetune}}$ be the loss functions during pre-training and fine-tuning, respectively. Pre-trained embeddings $\mathbf{E}$ generalize well to a downstream task if the following holds:

$$\mathcal{L}_{\text{finetune}}(\mathbf{E}) \leq \mathcal{L}_{\text{finetune}}(\mathbf{E}_{\text{rand}}) + \epsilon,$$

where $\mathbf{E}_{\text{rand}}$ are randomly initialized embeddings and $\epsilon$ is a small positive constant. This indicates that pre-trained embeddings lead to better generalization compared to random initialization.

Example: In a sentiment analysis task, pre-trained embeddings might be fine-tuned on a labeled dataset to classify reviews as positive or negative. The performance of the model, as measured by accuracy or F1-score, reflects the quality of the embeddings.

Intrinsic and extrinsic methods serve complementary roles in the evaluation of embeddings. Intrinsic methods offer insight into the geometric properties of the embedding space, while extrinsic methods provide evidence of how well these embeddings perform in practical applications. Together, these methods form a comprehensive evaluation framework that ensures embeddings are both theoretically sound and practically effective.

### 2.4.2   Impact of Embeddings on Downstream Tasks

Embeddings serve as the foundational input to models in many NLP tasks. The quality of these embeddings directly influences the performance of models in tasks such as sentiment analysis, machine translation, and named entity recognition. The mathematical relationship between the quality of embeddings and downstream task performance can be understood through the concept of representation power.

Let $\mathbf{E} \in \mathbb{R}^{n \times d}$ represent the embedding matrix, where $n$ is the vocabulary size and $d$ is the embedding dimension. For a given downstream task $T$, the performance $\mathcal{M}_T(\mathbf{E})$ is a function of the embeddings:

$$\mathcal{M}_T(\mathbf{E}) = \frac{1}{|D_T|} \sum_{(x,y) \in D_T} \ell(f(x; \mathbf{E}), y),$$

where $D_T$ is the dataset for task $T$, $f(x; \mathbf{E})$ is the model's prediction based on input $x$ using embeddings $\mathbf{E}$, $y$ is the true label, and $\ell$ is the loss function (e.g., cross-entropy for classification tasks).

The effectiveness of embeddings $\mathbf{E}$ can be assessed by their ability to linearly separate classes or capture relevant features for the task. This is often analyzed using the notion of linear separability. Let $\mathbf{w} \in \mathbb{R}^d$ be a weight vector for a linear classifier. The classification decision is given by

$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{E}[x]).$$

The embeddings $\mathbf{E}$ are considered effective if the classes $y$ and $\hat{y}$ are well separated in the embedding space. The quality of this separation can be quantified using a margin $\gamma$:

$$\gamma = \min_{(x,y) \in D_T} \frac{y \cdot \mathbf{w}^\top \mathbf{E}[x]}{\|\mathbf{w}\|}.$$

A larger margin $\gamma$ indicates better separation and, consequently, higher task performance.

Example: In sentiment analysis, where the task is to classify text as positive or negative, the embeddings must capture sentiment-related features such that positive and negative texts are easily separable by a linear classifier. If the embeddings are well structured, the performance metric, such as accuracy or F1-score, will be high.

**Theorem 2.8** (Performance Bound with Embeddings) *Let $\mathbf{E}^*$ be the optimal embedding matrix for a task $T$. The performance $\mathcal{M}_T(\mathbf{E})$ of any embedding $\mathbf{E}$ is bounded by*

$$\mathcal{M}_T(\mathbf{E}) \leq \mathcal{M}_T(\mathbf{E}^*) + \epsilon,$$

*where $\epsilon$ depends on the quality of $\mathbf{E}$ relative to $\mathbf{E}^*$. This bound emphasizes that embeddings closer to the optimal $\mathbf{E}^*$ yield better task performance.*

### 2.4.3 Performance Metrics for Positional Encoding

The effectiveness of positional encoding can be assessed through task-specific performance metrics and analyses of generalization and robustness.

**Task-Specific Performance Metrics**

Positional encoding is evaluated based on how well it enables a model to understand and process sequences. The effectiveness of positional encoding can be measured by comparing model performance on tasks that rely heavily on sequential information.

For a given task $T$ with performance metric $\mathcal{M}_T$, the evaluation focuses on how the presence of positional encoding **P** improves performance:

$$\Delta \mathcal{M}_T = \mathcal{M}_T(\mathbf{E} + \mathbf{P}) - \mathcal{M}_T(\mathbf{E}),$$

where $\mathcal{M}_T(\mathbf{E} + \mathbf{P})$ represents the performance with positional encoding and $\mathcal{M}_T(\mathbf{E})$ without it.

Tasks such as machine translation, text summarization, and language modeling are highly dependent on the correct interpretation of sequence order. In these tasks, positional encoding helps the model maintain context and coherence in the generated outputs. The evaluation can be performed using metrics like BLEU for translation, ROUGE for summarization, or perplexity for language modeling.

Example: In a machine translation task, the BLEU score measures the overlap between machine-generated translations and reference translations. Positional encoding is expected to improve the BLEU score by ensuring that the model correctly captures the order of words and phrases, which is crucial for producing coherent translations.

Let $T$ be a sequence-dependent task with performance metric $\mathcal{M}_T$. The inclusion of positional encoding **P** improves the performance bound:

$$\mathcal{M}_T(\mathbf{E} + \mathbf{P}) \geq \mathcal{M}_T(\mathbf{E}) + \gamma,$$

where $\gamma$ is a positive constant representing the benefit of positional encoding in maintaining sequence information.

**Generalization and Robustness Analysis**

The generalization capability of a model refers to its performance on unseen data, while robustness indicates its ability to handle perturbations or variations in the input. Positional encoding should enhance both aspects by providing the model with a more structured understanding of sequences.

Let $\mathcal{L}_{\text{train}}(\mathbf{E}, \mathbf{P})$ and $\mathcal{L}_{\text{test}}(\mathbf{E}, \mathbf{P})$ represent the training and test loss functions, respectively. The generalization gap $\Delta \mathcal{L}$ is defined as

$$\Delta \mathcal{L} = \mathcal{L}_{\text{test}}(\mathbf{E}, \mathbf{P}) - \mathcal{L}_{\text{train}}(\mathbf{E}, \mathbf{P}).$$

A smaller generalization gap indicates better generalization. To analyze robustness, we introduce perturbations $\delta$ to the input sequence $x$, and measure the change in model performance:

$$\mathcal{R}(\mathbf{E}, \mathbf{P}) = \max_\delta \mathcal{M}_T(f(x + \delta; \mathbf{E} + \mathbf{P}), y) - \mathcal{M}_T(f(x; \mathbf{E} + \mathbf{P}), y).$$

A lower value of $\mathcal{R}(\mathbf{E}, \mathbf{P})$ indicates higher robustness. For any sequence-dependent task $T$, the inclusion of positional encoding $\mathbf{P}$ reduces the generalization gap:

$$\Delta\mathcal{L}(\mathbf{E} + \mathbf{P}) \leq \Delta\mathcal{L}(\mathbf{E}),$$

indicating that positional encoding helps the model generalize better by providing a clearer structure to the sequence data.

Example: In a text summarization task, a model trained with positional encoding is likely to produce more coherent summaries on unseen documents, demonstrating better generalization. Additionally, the model's robustness can be tested by introducing noise or slight reordering of input sentences, and observing how well the model maintains the integrity of the summary.

## 2.5  Data Handling and Preprocessing

Data handling and preprocessing are critical steps in preparing data for machine learning models, particularly for NLP tasks. Proper preprocessing ensures that the data is in a suitable format for training and that the model can learn effectively from the input data. This section focuses on two key aspects of data preprocessing: normalization and standardization of data, and text data tokenization, with a detailed exploration of Byte-Pair Encoding (BPE) and WordPiece Tokenization. In the context of NLP, data preprocessing involves transforming raw text data into a structured format that models can effectively process. This typically includes steps such as normalization, standardization, and tokenization. Each of these steps plays a crucial role in ensuring that the data is consistent, reducing variability that is not relevant to the task, and breaking down the text into manageable pieces that the model can interpret.

### 2.5.1  Data Normalization and Standardization

Data normalization and standardization are techniques used to adjust the scale of features in a dataset, ensuring that they contribute equally to the model's learning process. This is particularly important in NLP, where different features (e.g., word frequencies, embedding values) may have different scales.

Normalization typically rescales the data to a fixed range, usually [0, 1] or [−1, 1]. For a feature $x$ in the dataset, normalization is mathematically defined as

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)},$$

where $\min(x)$ and $\max(x)$ are the minimum and maximum values of $x$ in the dataset. Normalization ensures that all features are on the same scale, which can prevent features with larger ranges from dominating the learning process.

Standardization transforms the data to have a mean of 0 and a standard deviation of 1. For a feature $x$, the standardized value $x'$ is calculated as

$$x' = \frac{x - \mu_x}{\sigma_x},$$

where $\mu_x$ is the mean of $x$, and $\sigma_x$ is the standard deviation of $x$. Standardization is particularly useful when the data follows a normal distribution, as it maintains the relative relationships between data points while adjusting for scale.

Let $\mathbf{X}$ be a matrix of features where each column represents a feature vector $x$. The standardized data matrix $\mathbf{X}'$ is given by

$$\mathbf{X}' = \frac{\mathbf{X} - \mu_{\mathbf{X}}}{\sigma_{\mathbf{X}}},$$

where $\mu_{\mathbf{X}}$ and $\sigma_{\mathbf{X}}$ are the mean and standard deviation vectors of $\mathbf{X}$. The covariance matrix of the standardized data $\mathbf{X}'$ is the correlation matrix, which has ones on the diagonal and values between –1 and 1 off-diagonal, reflecting the linear relationships between features.

Example: In an NLP task, consider a feature representing word frequency in a document. If the word frequencies vary greatly, normalization ensures that all words contribute proportionally to the model, rather than allowing frequent words to dominate due to their higher counts.

### 2.5.2　Text Data Tokenization

Tokenization is the process of converting raw text into smaller units, called tokens, which can be words, subwords, or characters. In NLP, tokenization is a crucial step that allows models to handle text data by breaking it down into these manageable units. There are various methods for tokenization, each with its advantages and limitations. In this section, we will explore two widely used tokenization methods: Byte-Pair Encoding (BPE) and WordPiece Tokenization.

**Byte-Pair Encoding (BPE)**

Byte-Pair Encoding (BPE) is a subword tokenization method that iteratively merges the most frequent pairs of characters or character sequences in the text. The goal is to create a vocabulary that includes both characters and common subwords, reducing the number of unknown tokens while keeping the vocabulary size manageable.

Algorithm:

1. Initialization: Start with a vocabulary that includes all individual characters in the text.

2. Pair Merging: Identify the most frequent pair of adjacent tokens (characters or subwords) in the text and merge them into a new token.

3. Vocabulary Update: Add the new token to the vocabulary and replace all instances of the merged pair in the text with the new token.

4. Iteration: Repeat the pair merging and vocabulary update steps until the desired vocabulary size is reached.

Let $\mathcal{V}$ be the initial vocabulary of characters, and $T = \{t_1, t_2, \ldots, t_n\}$ represent the text sequence. At each iteration, identify the most frequent pair $(t_i, t_{i+1})$ and merge it into a new token $t_{i,i+1}$:

$$\mathcal{V} \leftarrow \mathcal{V} \cup \{t_{i,i+1}\},$$

$$T \leftarrow T - \{t_i, t_{i+1}\} + \{t_{i,i+1}\}.$$

The process continues until the vocabulary $\mathcal{V}$ reaches a predefined size.

BPE effectively compresses the text by reducing the average length of tokens. Let $L(T)$ be the average token length in text $T$ before applying BPE, and $L'(T)$ be the average length after applying BPE. Then

$$L'(T) \leq L(T),$$

indicating that BPE reduces the average token length, leading to more compact text representation while maintaining the ability to reconstruct the original text.

Example: Consider the text "banana bandana." Initially, the vocabulary consists of individual characters b, a, n, d. After the first iteration, the most frequent pair "an" is merged into a new token "an." The updated text and vocabulary reflect this change, gradually building more complex subwords.

## WordPiece Tokenization

WordPiece Tokenization ([8]), similar to BPE, is a subword tokenization method that builds a vocabulary of subwords to handle rare or unseen words. It was originally developed for models like BERT and has become a standard in many transformer-based architectures.

Algorithm:

1. Initialization: Start with a vocabulary that includes all individual characters and some common words.

2. Subword Splitting: Split each word in the text into the longest possible subwords that are in the current vocabulary.

3. Vocabulary Update: For each word, identify potential subword splits and calculate their likelihood based on the current vocabulary. Add the most likely subwords to the vocabulary.

4. Iteration: Continue splitting words and updating the vocabulary until the vocabulary reaches the desired size.

Let $w$ be a word to tokenize, and $\mathcal{V}$ be the current vocabulary. WordPiece aims to split $w$ into subwords $\{s_1, s_2, \ldots, s_k\}$ such that

$$w = s_1 \circ s_2 \circ \cdots \circ s_k,$$

where $\circ$ denotes concatenation. The likelihood of a subword split $\{s_1, s_2, \ldots, s_k\}$ is calculated as

$$P(w) = \prod_{i=1}^{k} P(s_i),$$

where $P(s_i)$ is the probability of subword $s_i$ in the current vocabulary. The goal is to maximize this likelihood while keeping the vocabulary size manageable.

Let $S(w) = \{s_1, s_2, \ldots, s_k\}$ be the optimal subword split for word $w$ that maximizes the likelihood $P(w)$. Then

$$S(w) = \arg \max_{S} \prod_{i=1}^{k} P(s_i),$$

where $S$ is the set of all possible subword splits of $w$. This theorem ensures that the selected subword split is the most likely given the current vocabulary, contributing to more accurate tokenization.

Example: For the word "unhappiness," WordPiece might tokenize it into subwords "un," "##happi," and "##ness," where "##" indicates that the subword is not a standalone word. This tokenization captures the morphemes and meaningful subword units, aiding in better understanding by the model.

### 2.5.3  Handling Missing Data

In any data-driven modeling task, missing data and the need for data augmentation are common challenges. Proper handling of missing data is crucial for maintaining the integrity and performance of models, while data augmentation techniques are essential for enhancing the robustness and generalization of models. This section explores imputation techniques for handling missing data and various data augmentation methods for both text and image data. Missing data is a pervasive issue in real-world datasets, and how it is handled can significantly affect model performance. The presence of missing data can introduce bias, reduce the representativeness of the

sample, and lead to inaccurate or invalid conclusions. To address this, several impu-
tation techniques are employed to estimate the missing values based on the available
data.

**Imputation Techniques**

Imputation techniques aim to replace missing data points with plausible values to
create a complete dataset that can be used for analysis. The choice of imputation
method depends on the nature of the data and the assumptions that can be reasonably
made about the missingness mechanism.

Mean/Median Imputation: One of the simplest imputation techniques is to replace
missing values with the mean or median of the observed values in the same feature.
For a feature $X$ with missing values, the imputed value $\hat{x}_{\text{mean}}$ is given by

$$\hat{x}_{\text{mean}} = \frac{1}{n} \sum_{i=1}^{n} x_i,$$

where $x_i$ are the observed values and $n$ is the number of observed values. Median
imputation replaces missing values with the median $\hat{x}_{\text{median}}$ of the observed values.

Let $X$ be a random variable with missing values imputed using the mean $\hat{x}_{\text{mean}}$.
The bias introduced by mean imputation in the estimation of the population mean
$\mu_X$ is given by

$$\text{Bias}(\hat{\mu}_X) = \frac{n_{\text{miss}}}{n} \left( \mu_X - \hat{x}_{\text{mean}} \right),$$

where $n_{\text{miss}}$ is the number of missing values and $n$ is the total number of observations.
This bias tends to zero as the proportion of missing data decreases.

K-Nearest Neighbors (K-NN) Imputation: K-NN imputation estimates the miss-
ing values based on the values of the k-nearest neighbors in the feature space. For a
missing value $x_i$ in feature $X$, the imputed value $\hat{x}_i$ is given by

$$\hat{x}_i = \frac{1}{k} \sum_{j \in \mathcal{N}(i)} x_j,$$

where $\mathcal{N}(i)$ represents the indices of the k-nearest neighbors to $x_i$ in the feature
space. This method assumes that similar instances have similar values, making it
suitable for datasets where the features are correlated.

Multiple Imputation: Multiple imputation involves creating several different plau-
sible imputed datasets and combining the results obtained from each to account for
the uncertainty of the imputed values. Let $D_1, D_2, \ldots, D_m$ be the imputed datasets.
The final estimate $\hat{\theta}$ of a parameter $\theta$ is obtained by averaging the estimates $\hat{\theta}_j$ from
each imputed dataset:

$$\hat{\theta} = \frac{1}{m} \sum_{j=1}^{m} \hat{\theta}_j,$$

where $\hat{\theta}_j$ is the estimate from the j-th imputed dataset.

Under the assumption of missing at random (MAR), multiple imputation provides consistent estimates of population parameters. Formally, if the data is MAR and the imputation model is correctly specified, then

$$\lim_{n \to \infty} \hat{\theta} = \theta,$$

where $\theta$ is the true population parameter.

Example: In a medical dataset where some patients' blood pressure readings are missing, K-NN imputation might be used to estimate the missing values based on the readings of similar patients, while multiple imputation could provide a range of possible values to reflect the uncertainty in the imputation.

### 2.5.4   Data Augmentation

Data augmentation is a technique used to increase the diversity and size of a dataset by generating new data points from the existing data. This is particularly important in machine learning, where large and varied datasets are needed to train models that generalize well to new, unseen data. Different techniques are used for augmenting text and image data, each designed to preserve the underlying structure of the data while introducing variation.

**Techniques for Text Data**

Text data augmentation involves generating new sentences or phrases that retain the original meaning while introducing variability in the form of synonyms, paraphrasing, or sentence structure modifications.

Synonym Replacement: Synonym replacement involves substituting words in a sentence with their synonyms. Given a sentence $S = \{w_1, w_2, \ldots, w_n\}$, a new sentence $S'$ is generated by replacing a word $w_i$ with a synonym $w_i'$ drawn from a thesaurus or word embedding space:

$$S' = \{w_1, \ldots, w_{i-1}, w_i', w_{i+1}, \ldots, w_n\}.$$

This method assumes that the meaning of the sentence remains unchanged when synonyms are substituted.

Back-Translation: Back-translation generates augmented data by translating a sentence into another language and then translating it back into the original language. Let $T_1$ and $T_2$ be translation functions for two languages. A sentence $S$ is first translated into a different language $L_2$ and then back to the original language $L_1$:

$$S' = T_1(T_2(S)).$$

The back-translated sentence $S'$ often differs in structure and word choice from the original sentence $S$, providing a diverse augmentation while preserving the original meaning.

Let $S$ be an original sentence and $S'$ its back-translated version. The lexical diversity $D(S, S')$ introduced by back-translation can be quantified as

$$D(S, S') = \frac{\sum_{i=1}^{n} \mathbf{1}(w_i \neq w_i')}{n},$$

where $\mathbf{1}(w_i \neq w_i')$ is an indicator function that equals 1 if the word $w_i$ in $S$ differs from the corresponding word $w_i'$ in $S'$, and $n$ is the length of the sentence. This diversity metric captures the variation introduced by back-translation.

Example: For the sentence "The cat sat on the mat," back-translation might generate "The feline rested on the rug." Although the structure and word choice differ, the meaning remains consistent, providing a useful augmentation.

**Techniques for Image Data**

Image data augmentation involves applying transformations to images to generate new variations while preserving the essential features that define the objects in the images. Common techniques include rotation, scaling, flipping, and adding noise.

Rotation and Scaling: Let $I(x, y)$ represent the intensity of an image at pixel coordinates $(x, y)$. Rotation by an angle $\theta$ and scaling by a factor $\alpha$ transform the image coordinates according to

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \alpha \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}.$$

The transformed image $I'(x', y')$ retains the content of the original image but with altered orientation and size.

Flipping: Flipping an image horizontally or vertically changes the pixel coordinates as follows:

Horizontal flip: $x' = -x$.

Vertical flip: $y' = -y$.

These transformations preserve the structure of the image while introducing variability in the dataset.

Adding Noise: Noise can be added to images by perturbing the pixel values. Let $I(x, y)$ be the original pixel value at $(x, y)$. Additive Gaussian noise is defined as

$$I'(x, y) = I(x, y) + \epsilon(x, y),$$

where $\epsilon(x, y) \sim \mathcal{N}(0, \sigma^2)$ is Gaussian noise with mean 0 and variance $\sigma^2$. The noise introduces randomness while maintaining the overall content of the image.

**Remark (Invariance under Image Augmentation):** For an image classification model, let $f(I)$ be the predicted class label for an image $I$. An augmentation transformation $T$ should satisfy the invariance property:

$$f(I) = f(T(I)),$$

indicating that the model's prediction remains consistent across different augmentations of the same image.

Example: Consider an image of a cat. Augmentations such as rotating the image by 15 degrees, flipping it horizontally, or adding slight Gaussian noise should still allow the model to correctly classify the image as "cat."

# References

1. Bojanowski, P., Grave, E., Joulin, A., Mikolov, T.: Enriching word vectors with subword information. Trans. Assoc. Comput. Linguist. **5**, 135–146 (2017)
2. Armand, J., Edouard, G., Piotr, B., Tomas, M.: Bag of tricks for efficient text classification (2016). arXiv preprint arXiv:1607.01759
3. Tomas, M., Kai, C., Corrado, G.S., Dean, J.: Efficient estimation of word representations in vector space. In: Proceedings of the International Conference on Learning Representations (ICLR) (2013a)
4. Tomas, M., Ilya, S., Kai, C., Greg, C., Jeffrey, D.: Distributed representations of words and phrases and their compositionality. In: Proceedings of the 27th International Conference on Neural Information Processing Systems–Volume 2, NIPS'13, pp. 3111–3119. Curran Associates Inc, Red Hook, NY, USA (2013b)
5. Jeffrey, P., Richard, S., Manning, C.D.: Glove: Global vectors for word representation. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 1532–1543 (2014)
6. Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L.: Deep contextualized word representations. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, vol. 1, pp. 2227–2237 (2018)
7. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in Neural Information Processing Systems, vol. 30 (2017)
8. Wu, Y., Schuster, M., Chen, Z., Le, Q.V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., Klingner, J., Shah, A., Johnson, M., Liu, X., Kaiser, L., Gouws, S., Kato, Y., Kudo, T., Kazawa, H., Stevens, K., Kurian, G., Patil, N., Wang, W., Young, C., Smith, J.R., Riesa, J., Rudnick, A., Vinyals, O., Corrado, G.S., Hughes, M., Dean, J.: Google's neural machine translation system: Bridging the gap between human and machine translation (2016). ArXiv arXiv:1609.08144

# Chapter 3
# Attention Mechanisms: Theory and Variations

## 3.1  Attention Mechanisms

Attention mechanisms are a cornerstone of modern deep learning architectures, particularly in natural language processing and computer vision. They allow models to focus on specific parts of the input data, enabling more effective and context-aware processing. This section explores the theoretical underpinnings of different types of attention mechanisms, focusing on scalar, vector, and matrix attention. Each type of attention mechanism is examined in detail, with mathematical formulations and examples to illustrate their principles and applications.

### 3.1.1  Scalar, Vector, and Matrix Attention

Attention mechanisms can be categorized based on the dimensionality of the attention scores they produce. Scalar, vector, and matrix attention differ in how they represent the relevance of different parts of the input, and each has its unique applications and advantages.

### 3.1.2  Scalar Attention

Scalar attention is the most straightforward form of attention, where a single scalar weight is assigned to each element of the input sequence. This weight represents the relevance of the corresponding element in the sequence to the task at hand.

Given an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, scalar attention computes a set of attention scores $\alpha = \{\alpha_1, \alpha_2, \ldots, \alpha_n\}$, where each $\alpha_i$ is a scalar value. These scores

are typically derived from a similarity measure between a query vector $\mathbf{q}$ and each element $x_i$ in the sequence:

$$\alpha_i = \text{softmax}(e_i), \quad \text{where} \quad e_i = \mathbf{q}^\top \mathbf{k}_i,$$

where $\mathbf{k}_i$ is the key vector associated with $x_i$, and $\text{softmax}(e_i)$ ensures that the attention scores sum to one:

$$\alpha_i = \frac{\exp(e_i)}{\sum_{j=1}^n \exp(e_j)}.$$

The final output of scalar attention is a weighted sum of the input elements, where each element $x_i$ is scaled by its corresponding attention score $\alpha_i$:

$$\mathbf{z} = \sum_{i=1}^n \alpha_i x_i.$$

This output $\mathbf{z}$ is a context vector that aggregates the information from the input sequence, focusing more on the relevant elements as determined by the attention scores.

Let $\mathbf{z}$ be the output of scalar attention for the input sequence $\mathbf{x}$. Then, $\mathbf{z}$ is a convex combination of the input elements:

$$\mathbf{z} = \sum_{i=1}^n \alpha_i x_i, \quad \text{where} \quad \sum_{i=1}^n \alpha_i = 1 \quad \text{and} \quad \alpha_i \geq 0.$$

This convex combination ensures that the output $\mathbf{z}$ lies within the convex hull of the input sequence $\mathbf{x}$, meaning it is a weighted average of the input elements.

Example: In a machine translation task, scalar attention might be used to determine which words in the source sentence should be most emphasized when generating each word in the target sentence. The attention scores $\alpha_i$ reflect the importance of each source word for producing the current target word.

### 3.1.3   Vector Attention

Vector attention extends scalar attention by assigning a vector of attention scores to each element in the input sequence. This approach allows the model to capture more complex relationships between the elements by considering multiple aspects of relevance simultaneously.

Given an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, vector attention computes a set of attention vectors $\alpha = \{\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \ldots, \boldsymbol{\alpha}_n\}$, where each $\boldsymbol{\alpha}_i$ is a vector in $\mathbb{R}^d$. These vectors are typically derived from a similarity measure between a query vector $\mathbf{q}$ and each

element $x_i$ in the sequence, followed by a softmax operation applied element-wise:

$$\boldsymbol{\alpha}_i = \text{softmax}(\mathbf{W}\mathbf{k}_i),$$

where $\mathbf{k}_i$ is the key vector associated with $x_i$, and $\mathbf{W}$ is a learned weight matrix that projects the key vectors into the same space as the query vector.

The final output of vector attention is a weighted sum of the input elements, where each element $x_i$ is scaled by its corresponding attention vector $\boldsymbol{\alpha}_i$:

$$\mathbf{z} = \sum_{i=1}^{n} \boldsymbol{\alpha}_i \odot x_i,$$

where $\odot$ denotes element-wise multiplication. This allows the model to focus on different aspects of the input elements simultaneously.

Let $\mathbf{z}$ be the output of vector attention for the input sequence $\mathbf{x}$. Then, $\mathbf{z}$ is a linear combination of the input elements, with each dimension of the output being influenced by a different combination of input elements:

$$\mathbf{z}_j = \sum_{i=1}^{n} \alpha_{ij} x_{ij},$$

where $\alpha_{ij}$ is the $j$-th component of the attention vector $\boldsymbol{\alpha}_i$. This multi-dimensional attention allows the model to capture more nuanced relationships between the input elements.

Example: In image captioning, vector attention might be used to focus on different regions of an image when generating each word of the caption. The attention vectors $\boldsymbol{\alpha}_i$ could represent different visual features such as color, shape, and texture, enabling the model to generate more detailed and accurate descriptions.

### 3.1.4  Matrix Attention

Matrix attention further generalizes attention mechanisms by assigning a matrix of attention scores to each pair of elements in the input sequence. This allows the model to capture interactions between pairs of elements, rather than just between individual elements and a query.

Given an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, matrix attention computes a set of attention matrices $\mathbf{A} = \{\mathbf{A}_{ij}\}$, where each $\mathbf{A}_{ij} \in \mathbb{R}^{d \times d}$ represents the interaction between elements $x_i$ and $x_j$. These matrices are derived from similarity measures between pairs of elements, possibly involving queries and keys:

$$\mathbf{A}_{ij} = \text{softmax}(\mathbf{q}_i^{\top} \mathbf{W} \mathbf{k}_j),$$

where $\mathbf{q}_i$ and $\mathbf{k}_j$ are the query and key vectors for elements $x_i$ and $x_j$, respectively, and $\mathbf{W}$ is a learned weight matrix.

The final output of matrix attention is a transformed version of the input sequence, where each element is influenced by the interactions captured in the attention matrices:

$$\mathbf{z}_i = \sum_{j=1}^{n} \mathbf{A}_{ij} x_j.$$

This output $\mathbf{z}_i$ reflects the combined influence of all other elements in the sequence on the element $x_i$, with each influence weighted by the corresponding attention matrix $\mathbf{A}_{ij}$.

Let $\mathbf{z}_i$ be the output of matrix attention for the input sequence $\mathbf{x}$. The output can be expressed as a bilinear form involving the input elements and the attention matrices:

$$\mathbf{z}_i = \sum_{j=1}^{n} x_j^{\top} \mathbf{A}_{ij} x_j.$$

This bilinear form captures the second-order interactions between input elements, allowing the model to encode more complex relationships.

In graph neural networks, matrix attention can be used to model the relationships between nodes in a graph. Each node's representation is updated based on its connections to other nodes, with the strength and type of connection captured by the attention matrices $\mathbf{A}_{ij}$.

## 3.2 Self-Attention

Self-attention is a foundational mechanism in modern deep learning architectures, particularly in the context of NLP and computer vision. It allows a model to weigh the relevance of different parts of the same input sequence when encoding that sequence into a representation. This section explores the mathematical foundation of self-attention, including its specific variants like dot-product attention, scaled dot-product attention, and masked self-attention.

Self-attention mechanisms are designed to allow each element in a sequence to attend to every other element in that sequence, effectively enabling the model to capture dependencies regardless of their distance in the sequence.

### 3.2.1 Mathematical Definition

Given an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$ where each $x_i \in \mathbb{R}^d$ represents a d-dimensional vector, the goal of self-attention is to produce a new sequence of the

same length, where each element $z_i$ is a weighted sum of all elements in the input sequence, with the weights determined by the relevance of the elements to each other.

The self-attention mechanism can be mathematically defined as

$$\mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{v}_j,$$

where $\alpha_{ij}$ are the attention weights and $\mathbf{v}_j$ are value vectors derived from the input $x_j$. The weights $\alpha_{ij}$ are calculated based on the compatibility between a query vector $\mathbf{q}_i$ associated with $x_i$ and a key vector $\mathbf{k}_j$ associated with $x_j$:

$$\alpha_{ij} = \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{k=1}^{n} \exp(\mathbf{q}_i^\top \mathbf{k}_k)},$$

where $\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i$, $\mathbf{k}_j = \mathbf{W}_k \mathbf{x}_j$, and $\mathbf{v}_j = \mathbf{W}_v \mathbf{x}_j$ are linear transformations of the input vectors using learned weight matrices $\mathbf{W}_q$, $\mathbf{W}_k$, and $\mathbf{W}_v$.

The attention weights $\alpha_{ij}$ are non-negative and sum to one for each $i$:

$$\sum_{j=1}^{n} \alpha_{ij} = 1, \quad \alpha_{ij} \geq 0.$$

This ensures that the output $\mathbf{z}_i$ is a convex combination of the value vectors $\mathbf{v}_j$, meaning the model can focus on the most relevant parts of the input sequence.

Example: In a sentence, self-attention allows each word to consider every other word when forming its contextualized representation. For instance, in the sentence "The cat sat on the mat," the representation of "sat" can attend to both "cat" and "mat" to better capture the meaning.

### 3.2.2 Dot-Product Attention

Dot-product attention is a specific form of self-attention where the similarity between query and key vectors is computed using the dot product. For an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, the attention score $e_{ij}$ between elements $x_i$ and $x_j$ is given by

$$e_{ij} = \mathbf{q}_i^\top \mathbf{k}_j.$$

The attention weights $\alpha_{ij}$ are then obtained by applying the softmax function to the scores:

$$\alpha_{ij} = \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{k=1}^{n} \exp(\mathbf{q}_i^\top \mathbf{k}_k)}.$$

The final output for each element is a weighted sum of the value vectors:

$$\mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{v}_j.$$

Let $d$ be the dimensionality of the input vectors. The dot-product attention mechanism computes the attention scores in $O(nd^2)$ time for an input sequence of length $n$. This efficiency is one reason dot-product attention is widely used in practice.

Example: In machine translation, dot-product attention might be used to align words in the source and target languages, determining which source words should be emphasized when generating each word in the translation.

### 3.2.3 Scaled Dot-Product Attention

Scaled dot-product attention is a variation of dot-product attention that addresses the issue of large dot-product values when the dimensionality of the input vectors is high. By scaling the dot products, the softmax function produces more balanced gradients, leading to more stable training.

Given an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, the attention score $e_{ij}$ in scaled dot-product attention is computed as

$$e_{ij} = \frac{\mathbf{q}_i^\top \mathbf{k}_j}{\sqrt{d}},$$

where $d$ is the dimensionality of the key vectors. The attention weights $\alpha_{ij}$ are then computed as

$$\alpha_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i^\top \mathbf{k}_j}{\sqrt{d}}\right)}{\sum_{k=1}^{n} \exp\left(\frac{\mathbf{q}_i^\top \mathbf{k}_k}{\sqrt{d}}\right)}.$$

The final output is computed as

$$\mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{v}_j.$$

Scaling by $\frac{1}{\sqrt{d}}$ ensures that the dot products $e_{ij}$ have a standard deviation close to 1, preventing the gradients from becoming too small or too large during training. This leads to more stable optimization and faster convergence.

### *3.2.4 Masked Self-Attention*

Masked self-attention is a variant of self-attention where certain elements of the input sequence are masked (i.e., ignored) during the computation of attention scores. This is particularly useful in tasks like autoregressive language modeling, where future tokens should not be considered when predicting the current token.

Given an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, the attention score $e_{ij}$ is computed as

$$e_{ij} = \begin{cases} \mathbf{q}_i^\top \mathbf{k}_j & \text{if } j \leq i, \\ -\infty & \text{if } j > i, \end{cases}$$

where $j > i$ indicates that $x_j$ is a future token relative to $x_i$. The attention weights $\alpha_{ij}$ are then computed using the softmax function, where the softmax is only applied to the non-masked scores:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{i} \exp(e_{ik})}.$$

The final output is computed as

$$\mathbf{z}_i = \sum_{j=1}^{i} \alpha_{ij} \mathbf{v}_j.$$

Masked self-attention enforces a causal structure in the output sequence, ensuring that each element $z_i$ only depends on the previous elements in the sequence:

$$\mathbf{z}_i = f(x_1, x_2, \ldots, x_i),$$

where $f$ is the function implemented by the attention mechanism. This causality is essential for tasks like language modeling, where predictions should not be influenced by future tokens.

In autoregressive language models like GPT, masked self-attention ensures that when predicting the next word in a sentence, the model only considers the words that have already been generated, preventing information leakage from future words.

## 3.3 Multi-head Attention

Multi-head attention is an extension of the basic self-attention mechanism, designed to improve the model's ability to focus on different parts of the input simultaneously. By using multiple attention heads, this mechanism allows the model to capture diverse aspects of the input data, leading to richer and more nuanced representations. In this section, we explore the mathematical principles behind multi-head attention,

including parallel attention heads, concatenation and projection, and the role of head diversity and specialization.

Multi-head attention involves running multiple self-attention mechanisms, or "heads," in parallel, each focusing on different parts or aspects of the input. The outputs of these heads are then concatenated and linearly projected to produce the final output.

### 3.3.1  Parallel Attention Heads

Given an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, multi-head attention computes multiple sets of attention scores and value vectors, with each set corresponding to a different attention head. For the $h$-th head, the query, key, and value vectors are computed as

$$\mathbf{q}_i^{(h)} = \mathbf{W}_q^{(h)}\mathbf{x}_i, \quad \mathbf{k}_j^{(h)} = \mathbf{W}_k^{(h)}\mathbf{x}_j, \quad \mathbf{v}_j^{(h)} = \mathbf{W}_v^{(h)}\mathbf{x}_j,$$

where $\mathbf{W}_q^{(h)}, \mathbf{W}_k^{(h)}, \mathbf{W}_v^{(h)}$ are the learned weight matrices specific to the $h$-th head. The attention score $\alpha_{ij}^{(h)}$ and the output $\mathbf{z}_i^{(h)}$ for each head are computed as:

$$\alpha_{ij}^{(h)} = \frac{\exp\left(\frac{\mathbf{q}_i^{(h)\top}\mathbf{k}_j^{(h)}}{\sqrt{d_k}}\right)}{\sum_{k=1}^{n} \exp\left(\frac{\mathbf{q}_i^{(h)\top}\mathbf{k}_k^{(h)}}{\sqrt{d_k}}\right)},$$

$$\mathbf{z}_i^{(h)} = \sum_{j=1}^{n} \alpha_{ij}^{(h)}\mathbf{v}_j^{(h)},$$

where $d_k$ is the dimensionality of the key vectors.

For any two attention heads $h$ and $h'$, the corresponding attention scores $\alpha_{ij}^{(h)}$ and $\alpha_{ij}^{(h')}$ are computed independently, leading to diverse attention patterns:

$$\alpha_{ij}^{(h)} \text{ and } \alpha_{ij}^{(h')} \text{ are independent if } \mathbf{W}_q^{(h)} \neq \mathbf{W}_q^{(h')} \text{ or } \mathbf{W}_k^{(h)} \neq \mathbf{W}_k^{(h')}.$$

This independence allows the model to focus on different aspects of the input simultaneously, enhancing its ability to capture complex patterns.

Example: In a machine translation model, one attention head might focus on word order, another on syntactic dependencies, and another on semantic content. Each head attends to different elements of the input sentence, providing a richer representation for the translation task.

### *3.3.2 Concatenation and Projection*

After computing the outputs $\mathbf{z}_i^{(h)}$ from all $H$ heads, these outputs are concatenated to form a single vector. This concatenated vector is then projected back into the original dimensional space using a learned weight matrix.

Formally, the concatenated vector $\mathbf{z}_i^{\text{concat}}$ is given by

$$\mathbf{z}_i^{\text{concat}} = \mathbf{z}_i^{(1)} \parallel \mathbf{z}_i^{(2)} \parallel \ldots \parallel \mathbf{z}_i^{(H)},$$

where $\parallel$ denotes vector concatenation. The final output $\mathbf{z}_i$ is obtained by applying a linear transformation to the concatenated vector:

$$\mathbf{z}_i = \mathbf{W}_o \mathbf{z}_i^{\text{concat}},$$

where $\mathbf{W}_o$ is the learned projection matrix.

Let $d_h$ be the dimensionality of each head's output vector $\mathbf{z}_i^{(h)}$, and let $d_o$ be the dimensionality of the final output vector $\mathbf{z}_i$. Then

$$\dim(\mathbf{z}_i^{\text{concat}}) = H \times d_h, \quad \dim(\mathbf{z}_i) = d_o.$$

The projection matrix $\mathbf{W}_o$ transforms the concatenated vector from $H \times d_h$ dimensions back to $d_o$ dimensions, ensuring that the final output has the same dimensionality as the input sequence elements.

In a transformer model, after each attention head has processed its part of the input sequence, the results are concatenated and projected to form a unified representation that integrates the diverse information captured by each head. This projection step ensures that the output remains compatible with subsequent layers in the model.

### *3.3.3 Head Diversity and Specialization*

One of the key advantages of multi-head attention is the ability of different heads to specialize in capturing different types of relationships within the input data. This specialization emerges naturally from the independent training of each head, which allows them to focus on different patterns and features.

Let $H$ be the number of attention heads, and let $\mathbf{z}_i^{(h)}$ represent the output of the $h$-th head. The diversity of the heads can be quantified by the variance in their outputs:

$$\text{Diversity}(\mathbf{z}_i^{(1)}, \mathbf{z}_i^{(2)}, \ldots, \mathbf{z}_i^{(H)}) = \frac{1}{H} \sum_{h=1}^{H} \left\| \mathbf{z}_i^{(h)} - \overline{\mathbf{z}_i} \right\|^2,$$

where $\overline{\mathbf{z}}_i$ is the mean output vector across all heads:

$$\overline{\mathbf{z}}_i = \frac{1}{H} \sum_{h=1}^{H} \mathbf{z}_i^{(h)}.$$

A higher diversity score indicates that the heads are capturing different aspects of the input, leading to a more comprehensive overall representation. If the attention heads in multi-head attention are trained independently, then each head tends to specialize in a different feature or aspect of the input sequence. Formally, the expected specialization $\mathbb{E}[\text{Spec}(h)]$ of the $h$-th head is maximized when the outputs $\mathbf{z}_i^{(h)}$ are orthogonal to each other:

$$\mathbb{E}[\text{Spec}(h)] = \max \quad \text{subject to} \quad \mathbf{z}_i^{(h)} \perp \mathbf{z}_i^{(h')}, \ \forall h \neq h'.$$

This orthogonality ensures that each head contributes unique information to the final representation.

Example: In a sentiment analysis task, one attention head might specialize in detecting negations, another in identifying strong emotional words, and another in understanding the overall sentence structure. The diversity and specialization of these heads enable the model to capture a wide range of linguistic cues, leading to more accurate sentiment predictions.

## 3.4  Cross-Attention

Cross-attention is a specialized form of attention that allows a model to focus on and integrate information from two different sequences or sets of data. It is a crucial component in encoder–decoder architectures, where the model must align and combine information from an input sequence (e.g., a source sentence) with an output sequence (e.g., a target sentence). This section provides a mathematical exploration of cross-attention, including its formulation and applications in encoder–decoder models.

Cross-attention differs from self-attention in that it operates on two distinct sets of input sequences. It computes the relevance of elements from one sequence with respect to elements of another sequence, allowing the model to transfer and align information between the two sequences.

Given two sequences, an *encoder sequence* $\mathbf{x}^{(e)} = \{x_1^{(e)}, x_2^{(e)}, \ldots, x_m^{(e)}\}$ and a *decoder sequence* $\mathbf{x}^{(d)} = \{x_1^{(d)}, x_2^{(d)}, \ldots, x_n^{(d)}\}$, cross-attention computes a weighted sum of the encoder sequence elements for each element in the decoder sequence. This enables the decoder to focus on relevant parts of the encoder's output when generating its own sequence.

The cross-attention mechanism computes the attention score $\alpha_{ij}$ between the $i$-th element of the decoder sequence and the $j$-th element of the encoder sequence as

$$\alpha_{ij} = \frac{\exp\left(\mathbf{q}_i^{(d)\top}\mathbf{k}_j^{(e)}\right)}{\sum_{k=1}^{m}\exp\left(\mathbf{q}_i^{(d)\top}\mathbf{k}_k^{(e)}\right)},$$

where $\mathbf{q}_i^{(d)} = \mathbf{W}_q\mathbf{x}_i^{(d)}$ is the query vector derived from the $i$-th decoder element, and $\mathbf{k}_j^{(e)} = \mathbf{W}_k\mathbf{x}_j^{(e)}$ is the key vector derived from the $j$-th encoder element, with $\mathbf{W}_q$ and $\mathbf{W}_k$ being learned weight matrices.

The output of the cross-attention mechanism for the $i$-th decoder element is then computed as

$$\mathbf{z}_i^{(d)} = \sum_{j=1}^{m}\alpha_{ij}\mathbf{v}_j^{(e)},$$

where $\mathbf{v}_j^{(e)} = \mathbf{W}_v\mathbf{x}_j^{(e)}$ is the value vector associated with the $j$-th encoder element and $\mathbf{W}_v$ is another learned weight matrix.

The output $\mathbf{z}_i^{(d)}$ of cross-attention for the $i$-th decoder element is a convex combination of the value vectors $\mathbf{v}_j^{(e)}$ from the encoder sequence:

$$\mathbf{z}_i^{(d)} = \sum_{j=1}^{m}\alpha_{ij}\mathbf{v}_j^{(e)}, \quad \text{where} \quad \sum_{j=1}^{m}\alpha_{ij} = 1 \quad \text{and} \quad \alpha_{ij} \geq 0.$$

This convex combination ensures that the decoder's output at each step is a weighted average of the encoder's output, with the weights reflecting the relevance of each encoder element to the current decoding step.

Example: In a machine translation model, cross-attention allows the decoder to focus on relevant words or phrases from the source sentence (encoded by the encoder) while generating the target sentence. For instance, when translating "The cat sat on the mat" to another language, the decoder attends to the words "cat" and "sat" to generate the corresponding translation for the verb phrase in the target language.

### 3.4.1 Applications in Encoder–Decoder Models

Encoder–decoder models, particularly in tasks like machine translation and sequence-to-sequence learning, rely heavily on cross-attention to transfer information from the input (encoder) to the output (decoder). In these models, the encoder first processes the input sequence into a set of context-aware representations. The decoder then uses cross-attention to dynamically attend to these representations while generating the output sequence.

Formally, let $\mathbf{H}^{(e)} = \{\mathbf{h}_1^{(e)}, \mathbf{h}_2^{(e)}, \ldots, \mathbf{h}_m^{(e)}\}$ be the set of hidden states produced by the encoder for the input sequence. During each decoding step, the decoder state $\mathbf{s}_i^{(d)}$ is updated by attending to the encoder states:

$$\mathbf{s}_i^{(d)} = \text{DecoderRNN}\left(\mathbf{x}_i^{(d)}, \mathbf{z}_i^{(d)}\right),$$

where $\mathbf{z}_i^{(d)}$ is the cross-attention output for the $i$-th decoder step, and DecoderRNN represents the recurrent neural network (RNN) or any other processing unit used in the decoder.

Let $\mathbf{H}^{(e)}$ be the encoder hidden states and $\mathbf{S}^{(d)}$ be the decoder hidden states. Cross-attention ensures that each decoder state $\mathbf{s}_i^{(d)}$ depends on a weighted combination of the encoder states:

$$\mathbf{s}_i^{(d)} = f\left(\mathbf{x}_i^{(d)}, \sum_{j=1}^{m} \alpha_{ij}\mathbf{h}_j^{(e)}\right),$$

where $f$ is a function determined by the decoder architecture. This dependence guarantees that the decoder generates outputs that are contextually aligned with the input sequence.

Example: In an image captioning task, the encoder might process an image into a set of feature vectors representing different regions of the image. The decoder then uses cross-attention to focus on specific regions of the image while generating each word of the caption. For example, when generating the word "dog," the decoder attends to the region of the image where the dog is located, ensuring that the generated caption accurately describes the image.

## 3.5 Efficiency and Variations

Attention mechanisms, while powerful, can be computationally expensive, particularly when applied to long sequences. This section explores various efficient attention mechanisms designed to reduce the computational complexity of traditional self-attention while maintaining or even enhancing performance. We will examine several approaches, including Sparse Attention, Linformer, Longformer, Reformer, and Performer, each of which introduces unique strategies to optimize attention calculations.

### 3.5.1 Efficient Attention Mechanisms

The self-attention mechanism in its basic form has a computational complexity of $O(n^2d)$, where $n$ is the sequence length and $d$ is the dimensionality of the embeddings. This quadratic dependence on the sequence length becomes a bottleneck for long sequences, necessitating the development of more efficient alternatives. The mechanisms discussed here aim to reduce this complexity, often by approximating the full attention matrix or by limiting the number of interactions considered.

### 3.5.2 Sparse Attention

Sparse attention mechanisms reduce the computational burden by restricting the attention computation to a subset of the elements in the input sequence. Instead of computing attention weights for every pair of elements in the sequence, sparse attention only computes weights for a predefined set of pairs, effectively reducing the number of interactions from $O(n^2)$ to $O(n \cdot k)$, where $k$ is the number of non-zero attention weights per query.

Formally, let $\mathcal{S}_i \subset \{1, 2, \ldots, n\}$ be the set of indices for which the attention is computed for the $i$-th element. The attention score $\alpha_{ij}$ is computed as

$$\alpha_{ij} = \begin{cases} \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{k \in \mathcal{S}_i} \exp(\mathbf{q}_i^\top \mathbf{k}_k)}, & \text{if } j \in \mathcal{S}_i, \\ 0, & \text{if } j \notin \mathcal{S}_i. \end{cases}$$

The output for the $i$-th element is then computed as

$$\mathbf{z}_i = \sum_{j \in \mathcal{S}_i} \alpha_{ij} \mathbf{v}_j.$$

If the size of the sparse set $\mathcal{S}_i$ is $k$, then the computational complexity of sparse attention is $O(n \cdot k \cdot d)$, where $k \ll n$. This represents a significant reduction in complexity, particularly for large sequences.

Example: In natural language processing tasks, sparse attention can be implemented by allowing each word to attend only to nearby words within a fixed window or to a predefined set of important words (e.g., the beginning of the sentence).

### 3.5.3 Linformer

Linformer ([4]) is an efficient attention mechanism that approximates the full self-attention by projecting the sequence length $n$ down to a fixed lower-dimensional space $k$, where $k \ll n$. This reduces the quadratic complexity of self-attention to linear.

Given an input sequence $\mathbf{x} \in \mathbb{R}^{n \times d}$, Linformer introduces a projection matrix $\mathbf{P} \in \mathbb{R}^{n \times k}$ that projects the sequence dimension down to $k$:

$$\mathbf{PK}, \quad \mathbf{PV},$$

where $\mathbf{K}$ and $\mathbf{V}$ are the key and value matrices, respectively. The attention score computation is then based on these projected matrices:

$$\mathbf{A} = \text{softmax}\left(\mathbf{Q}(\mathbf{PK})^\top\right),$$

where $\mathbf{Q}$ is the query matrix. The output is

$$\mathbf{Z} = \mathbf{A}(\mathbf{PV}),$$

with $\mathbf{Z} \in \mathbb{R}^{n \times d}$.

The computational complexity of Linformer is $O(ndk)$, where $k$ is the reduced dimension. Since $k$ is typically much smaller than $n$, this results in significant computational savings.

Example: In sequence modeling tasks like document classification, Linformer can be applied to long documents where reducing the sequence length during attention computation leads to faster processing without sacrificing much accuracy.

### 3.5.4   Longformer

Longformer ([1]) introduces an efficient attention mechanism by combining sparse attention patterns with global attention. The model is designed to handle very long sequences by mixing local attention (where each token attends to a small neighborhood) and global attention (where specific tokens attend to all other tokens).

For a sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, Longformer defines a local attention window $w$ such that each token $x_i$ attends to the tokens within a window $[i - w, i + w]$. Global attention is applied to a small subset of important tokens (e.g., CLS tokens in classification tasks), which attend to all tokens in the sequence.

The attention score for token $x_i$ in the local window is given by

$$\alpha_{ij} = \frac{\exp\left(\mathbf{q}_i^\top \mathbf{k}_j\right)}{\sum_{k=i-w}^{i+w} \exp\left(\mathbf{q}_i^\top \mathbf{k}_k\right)},$$

and for global attention:

$$\alpha_{ij} = \frac{\exp\left(\mathbf{q}_i^\top \mathbf{k}_j\right)}{\sum_{k=1}^{n} \exp\left(\mathbf{q}_i^\top \mathbf{k}_k\right)} \quad \text{for global tokens.}$$

The computational complexity of Longformer is $O(nwd)$ for local attention, where $w$ is the window size, and $O(nd)$ for global attention applied to a small number of tokens, typically resulting in a much lower overall complexity compared to full attention.

In tasks like document classification or language modeling on long texts, Longformer can efficiently handle sequences of thousands of tokens, applying detailed attention locally while still capturing global context through selectively applied global attention.

### *3.5.5  Reformer*

Reformer ([3]) is an efficient attention mechanism that reduces the quadratic complexity of self-attention using two key innovations: locality-sensitive hashing (LSH) and reversible residual layers.

LSH reduces the number of interactions by hashing the queries and keys into buckets such that only queries and keys in the same bucket are considered for attention computation. Formally, the queries $\mathbf{Q}$ and keys $\mathbf{K}$ are hashed using an LSH function $h$:

$$h(\mathbf{Q}), h(\mathbf{K}) \in \{1, 2, \ldots, B\},$$

where $B$ is the number of buckets. The attention is then computed only within each bucket, significantly reducing the number of computations.

The reversible residual layers allow the model to save memory during training by recomputing activations from the output rather than storing them, further enhancing efficiency.

The computational complexity of Reformer is $O(n \log n \cdot d)$, where the $\log n$ factor comes from the LSH-based bucket sorting. This is a substantial reduction from the quadratic complexity of traditional self-attention.

Reformer is particularly useful in tasks like image processing or long text modeling, where both the sequence length and the need for efficiency are high. The LSH-based approach ensures that only relevant interactions are computed, leading to significant speedups.

### *3.5.6  Performer*

Performer ([2]) introduces a method called FAVOR+ (Fast Attention Via Orthogonal Random Features), which approximates the softmax attention mechanism with a linear attention mechanism that scales linearly with the sequence length.

For a sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, Performer approximates the attention score $\mathbf{A}$ using orthogonal random features $\phi(\mathbf{q}_i)$ and $\phi(\mathbf{k}_j)$, where $\phi$ is a feature mapping:

$$\mathbf{A}_{ij} = \phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_j).$$

The output is computed as

$$\mathbf{Z} = \phi(\mathbf{Q})^\top \mathbf{V} \phi(\mathbf{K})^\top,$$

where $\mathbf{Q}$, $\mathbf{K}$, and $\mathbf{V}$ are the query, key, and value matrices, respectively.

The computational complexity of Performer is $O(nd^2)$, where $d$ is the dimensionality of the embeddings, making it linear in sequence length and highly efficient for long sequences.

Performer is ideal for tasks requiring the processing of very long sequences, such as DNA sequence analysis or long document processing, where maintaining efficiency without sacrificing performance is crucial.

## 3.6  Variations in Attention

Attention mechanisms have evolved significantly since their inception, with various modifications aimed at improving their efficiency, expressiveness, and adaptability. This section explores some of the most impactful variations in attention, including global and local attention, hierarchical attention, adaptive attention span, and dynamic convolutions. Each variation is designed to address specific challenges or limitations of traditional self-attention, and we will examine them with a focus on their mathematical foundations.

### *3.6.1  Global and Local Attention*

Global and local attention refer to two distinct approaches in which attention can be applied across sequences. Global attention allows each element of the sequence to attend to every other element, while local attention restricts the attention scope to a smaller, localized window around each element.

For an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, global attention computes the attention score $\alpha_{ij}$ between every pair of elements $x_i$ and $x_j$:

$$\alpha_{ij} = \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{k=1}^{n} \exp(\mathbf{q}_i^\top \mathbf{k}_k)},$$

where $\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i$ and $\mathbf{k}_j = \mathbf{W}_k \mathbf{x}_j$ are the query and key vectors, respectively.

In contrast, local attention restricts the computation to a fixed window size $w$, so that each element $x_i$ only attends to elements within a window $[i - w, i + w]$:

$$\alpha_{ij} = \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{k=\max(1,i-w)}^{\min(n,i+w)} \exp(\mathbf{q}_i^\top \mathbf{k}_k)} \quad \text{for} \quad j \in [i - w, i + w].$$

The computational complexity of global attention is $O(n^2 d)$, whereas local attention reduces this to $O(nwd)$, where $w$ is the window size. This reduction is particularly beneficial when $n$ is large, and $w$ is much smaller than $n$.

In text processing, global attention might be used in tasks where understanding the entire context is crucial, such as in document classification, while local attention is more suitable for tasks like language modeling, where dependencies are typically local.

### 3.6.2  Hierarchical Attention

Hierarchical attention introduces a multi-level approach to attention, where attention mechanisms are applied at different levels of abstraction within the data. This is particularly useful in tasks that involve multi-scale data, such as document-level text processing or video analysis.

Given a sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$ with a hierarchical structure, let $\mathbf{z}^{(l)}$ represent the output of the attention mechanism at level $l$. The hierarchical attention mechanism first applies attention at the lowest level (e.g., word level):

$$\mathbf{z}_i^{(1)} = \sum_{j=1}^{n} \alpha_{ij}^{(1)} \mathbf{v}_j^{(1)},$$

where $\alpha_{ij}^{(1)}$ are the attention weights computed at level 1. The output $\mathbf{z}^{(1)}$ then forms the input to the next level of attention (e.g., sentence level):

$$\mathbf{z}^{(2)} = \sum_{j=1}^{m} \alpha_{ij}^{(2)} \mathbf{z}_j^{(1)},$$

and so on, until the highest level of abstraction is reached.

Let $L$ be the number of levels in the hierarchy. The expressive power of hierarchical attention is given by the combination of attentions at each level:

$$\mathbf{z}^{(L)} = f \left( \sum_{j_1=1}^{n_1} \alpha_{i_1 j_1}^{(1)} \cdots \sum_{j_L=1}^{n_L} \alpha_{i_L j_L}^{(L)} \mathbf{v}_{j_L}^{(L)} \right),$$

where $f$ is the function implemented by the hierarchical structure. This layered approach allows for a richer representation of complex, multi-scale data.

Example: In document classification, hierarchical attention might first attend to words within sentences, then to sentences within paragraphs, and finally to paragraphs within the document, allowing the model to capture nuanced relationships at multiple levels.

### 3.6.3  Adaptive Attention Span

Adaptive attention span introduces a mechanism where the attention span, or the window size over which attention is computed, is dynamically adjusted based on the needs of the task. This allows the model to allocate more computational resources to important parts of the sequence while reducing the attention span for less critical parts.

Given a sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, the adaptive attention span mechanism computes a dynamic window size $w_i$ for each element $x_i$, typically based on a learned gating function:

$$w_i = \text{softmax}(\mathbf{W}_w \mathbf{x}_i),$$

where $\mathbf{W}_w$ is a learned weight matrix. The attention score $\alpha_{ij}$ is then computed over the dynamically determined window $[i - w_i, i + w_i]$:

$$\alpha_{ij} = \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{k=i-w_i}^{i+w_i} \exp(\mathbf{q}_i^\top \mathbf{k}_k)}.$$

Let $w_{\text{max}}$ be the maximum possible attention span. The computational complexity of adaptive attention span is $O(n w_{\text{max}} d)$, but, in practice, this complexity is reduced due to the adaptive nature of $w_i$, which can be much smaller than $w_{\text{max}}$ for many elements.

In language modeling, adaptive attention span allows the model to focus on longer dependencies for complex sentences while reducing the attention span for simpler sentences or less critical words, improving both efficiency and performance.

### 3.6.4   Dynamic Convolutions

Dynamic convolutions combine the strengths of CNNs with the flexibility of attention mechanisms. In dynamic convolutions, the convolutional filters are dynamically generated based on the input sequence, allowing the model to adapt its filters to different parts of the input.

For an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, a dynamic convolution mechanism generates a set of filters $\mathbf{F}_i$ for each element $x_i$:

$$\mathbf{F}_i = \mathbf{W}_F \mathbf{x}_i,$$

where $\mathbf{W}_F$ is a learned weight matrix. The convolution operation is then applied using these dynamically generated filters:

$$\mathbf{z}_i = \sum_{j=-w}^{w} \mathbf{F}_i \cdot \mathbf{x}_{i+j},$$

where $w$ is the window size and $\cdot$ denotes the convolution operation.

Dynamic convolutions increase the expressive power of the model by allowing the convolutional filters to adapt to the input. Formally, the output $\mathbf{z}_i$ is a function of both the input and the dynamically generated filters:

$$\mathbf{z}_i = f\left(\mathbf{W}_F \mathbf{x}_i, \mathbf{x}_{i-w:i+w}\right),$$

where $f$ is the convolution operation. This allows the model to capture more complex patterns and dependencies within the input.

In speech recognition, dynamic convolutions can adapt to different phonemes or speech patterns, allowing the model to process varying speech inputs more effectively by adjusting the convolutional filters dynamically.

## 3.7  Mathematical Properties

Understanding the mathematical properties of attention mechanisms is crucial for analyzing their efficiency and scalability in practical applications. This section delves into the complexity analysis of attention mechanisms, focusing on time complexity and space complexity. A mathematical treatment of these aspects provides insights into the trade-offs involved in using different attention mechanisms, particularly in terms of their computational requirements.

### *3.7.1  Complexity Analysis*

The complexity analysis of attention mechanisms involves evaluating the computational resources required to execute the attention operations. This includes both the time complexity, which measures how the computation time scales with the size of the input, and the space complexity, which assesses the amount of memory required.

**Time Complexity**

Time complexity measures the amount of time an algorithm takes to complete as a function of the size of the input. For attention mechanisms, the primary factors influencing time complexity are the length of the input sequence $n$ and the dimensionality of the feature vectors $d$.

In traditional self-attention, for an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, the attention scores are computed between every pair of elements, resulting in $n^2$ comparisons. For each comparison, a dot product operation is performed, which has a time complexity of $O(d)$. Therefore, the overall time complexity of the self-attention mechanism is

$$\text{Time Complexity} = O(n^2 \cdot d).$$

Let $T(n, d)$ denote the time complexity of the self-attention mechanism. Then

$$T(n, d) = O(n^2 \cdot d).$$

This quadratic time complexity in terms of the sequence length $n$ becomes a significant bottleneck for large sequences, motivating the development of more efficient attention mechanisms.

Examples:

1. Global Attention: As discussed earlier, global attention computes attention scores between all pairs of elements in the sequence, leading to a time complexity of $O(n^2 \cdot d)$.

2. Local Attention: By restricting attention to a local window of size $w$, the time complexity is reduced to $O(n \cdot w \cdot d)$, where $w$ is typically much smaller than $n$.

3. Linformer: By projecting the sequence length $n$ down to a fixed lower-dimensional space $k$, Linformer reduces the time complexity to $O(n \cdot k \cdot d)$, where $k$ is a small constant relative to $n$.

## Space Complexity

Space complexity measures the amount of memory required to execute an algorithm as a function of the input size. In the context of attention mechanisms, the key factors influencing space complexity are again the sequence length $n$ and the dimensionality of the feature vectors $d$.

For self-attention, the primary space requirement comes from storing the attention scores, which are typically represented as an $n \times n$ matrix. Each element of this matrix requires $O(1)$ space, and thus the overall space complexity for storing the attention scores is

$$\text{Space Complexity} = O(n^2).$$

In addition to the attention scores, the query, key, and value matrices each require $O(n \cdot d)$ space. Therefore, the total space complexity of the self-attention mechanism is

$$\text{Total Space Complexity} = O(n^2 + n \cdot d).$$

Let $S(n, d)$ denote the space complexity of the self-attention mechanism. Then

$$S(n, d) = O(n^2 + n \cdot d).$$

This quadratic space complexity in terms of the sequence length $n$ is another factor that limits the scalability of self-attention, particularly when dealing with very large sequences.

Examples:

1. Global Attention: Similar to time complexity, global attention also requires storing an $n \times n$ matrix of attention scores, leading to a space complexity of $O(n^2)$.

2. Sparse Attention: By restricting attention to a subset of elements, sparse attention mechanisms reduce the space complexity to $O(n \cdot k)$, where $k$ is the number of non-zero attention weights per query.

3. Linformer: Linformer further reduces space complexity by projecting the sequence into a lower-dimensional space, resulting in a space complexity of $O(n \cdot k)$ for storing the attention scores.

## 3.7.2 Gradient Analysis

Understanding the flow of gradients in attention mechanisms is critical for ensuring effective training, particularly in deep neural networks. This analysis involves examining how gradients propagate through the network, how they can be affected by the structure of the attention mechanism, and how techniques like gradient clipping and normalization can be employed to address potential issues such as vanishing or exploding gradients.

### Gradient Flow in Attention Mechanisms

In attention mechanisms, the gradient flow is influenced by the softmax function used to compute attention weights and the dot products between query and key vectors. Given the attention score $\alpha_{ij}$ between elements $x_i$ and $x_j$:

$$\alpha_{ij} = \frac{\exp\left(\mathbf{q}_i^\top \mathbf{k}_j\right)}{\sum_{k=1}^n \exp\left(\mathbf{q}_i^\top \mathbf{k}_k\right)},$$

the gradient of $\alpha_{ij}$ with respect to the query vector $\mathbf{q}_i$ can be computed using the chain rule:

$$\frac{\partial \alpha_{ij}}{\partial \mathbf{q}_i} = \alpha_{ij}\left(\mathbf{k}_j - \sum_{k=1}^n \alpha_{ik}\mathbf{k}_k\right).$$

This expression shows that the gradient depends not only on the specific key vector $\mathbf{k}_j$ but also on the weighted sum of all key vectors, highlighting the interconnected nature of the gradient flow in attention mechanisms.

Let $\mathbf{g}_{\mathbf{q}_i}$ denote the gradient of the attention mechanism with respect to the query vector $\mathbf{q}_i$. The magnitude of the gradient is given by

$$\|\mathbf{g}_{\mathbf{q}_i}\| = \|\alpha_{ij}(\mathbf{k}_j - \bar{\mathbf{k}}_i)\|,$$

where $\bar{\mathbf{k}}_i = \sum_{k=1}^n \alpha_{ik}\mathbf{k}_k$ is the weighted average of the key vectors. The magnitude of the gradient can become small if the key vectors are similar, leading to potential issues with gradient flow.

In practice, when key vectors are very similar across different elements of the sequence, the gradient may diminish, leading to slow learning. This phenomenon

is particularly relevant in deep attention-based models where multiple layers of attention mechanisms are stacked.

**Gradient Clipping and Normalization**

Gradient clipping and normalization are techniques used to control the gradient magnitude during training, particularly in attention mechanisms where the gradients can become very large or very small due to the dot products involved.

Gradient clipping involves capping the gradients at a maximum value to prevent them from becoming too large, which can destabilize the training process. Formally, let $\mathbf{g}$ be the gradient vector, and $\lambda$ be the clipping threshold. The clipped gradient $\mathbf{g}'$ is given by

$$\mathbf{g}' = \mathbf{g} \cdot \min\left(1, \frac{\lambda}{\|\mathbf{g}\|}\right).$$

This ensures that the magnitude of the gradient $\|\mathbf{g}'\|$ does not exceed $\lambda$, helping to prevent exploding gradients.

Let $\mathbf{g}$ be the original gradient, and $\mathbf{g}'$ the clipped gradient. The relative change in gradient direction is minimized when $\|\mathbf{g}\| \leq \lambda$, ensuring that clipping primarily affects the magnitude rather than the direction of the gradient:

$$\frac{\mathbf{g} \cdot \mathbf{g}'}{\|\mathbf{g}\|\|\mathbf{g}'\|} \approx 1 \quad \text{if} \quad \|\mathbf{g}\| \leq \lambda.$$

Gradient normalization involves scaling the gradients to have a consistent magnitude, typically by normalizing them to unit norm:

$$\mathbf{g}' = \frac{\mathbf{g}}{\|\mathbf{g}\|}.$$

This technique helps ensure stable learning rates and can prevent both vanishing and exploding gradients by keeping the gradient magnitudes within a manageable range.

In deep networks with many layers of attention, gradient clipping and normalization can be essential for maintaining stable training. Without these techniques, gradients can easily become unstable, particularly in the early stages of training.

## 3.7.3 Convergence Properties

Convergence properties of attention mechanisms refer to the behavior of the optimization process as it approaches a minimum in the loss landscape. Convergence is

influenced by the structure of the attention mechanism, the choice of optimization algorithm, and the learning rate schedule.

For an attention mechanism with a loss function $L(\theta)$, where $\theta$ represents the parameters of the model, convergence is achieved when the gradient of the loss with respect to the parameters approaches zero:

$$\lim_{t \to \infty} \|\nabla_\theta L(\theta^{(t)})\| = 0,$$

where $t$ denotes the training iteration.

For a convex loss function $L(\theta)$ with Lipschitz-continuous gradient $\nabla_\theta L(\theta)$, the convergence rate of gradient descent is given by

$$L(\theta^{(t)}) - L(\theta^*) \leq \frac{1}{2\eta t} \|\theta^{(0)} - \theta^*\|^2,$$

where $\theta^*$ is the optimal parameter value, $\eta$ is the learning rate, and $\theta^{(0)}$ is the initial parameter value.

In attention-based models, the convergence rate can be affected by the choice of learning rate schedule, with adaptive learning rates (e.g., Adam) often leading to faster convergence compared to fixed learning rates. Additionally, attention mechanisms with sparse gradients may require more careful tuning to achieve convergence.

### 3.7.4  Stability and Robustness

Stability and robustness refer to the ability of attention mechanisms to maintain performance in the presence of perturbations, such as noise in the input data or variations in model parameters. Stability is particularly important in ensuring that small changes in the input do not lead to disproportionately large changes in the output.

For an attention mechanism $f(\mathbf{x}; \theta)$ with input $\mathbf{x}$ and parameters $\theta$, robustness can be analyzed through the Lipschitz constant $K$, which bounds the change in output relative to the change in input:

$$\|f(\mathbf{x}_1; \theta) - f(\mathbf{x}_2; \theta)\| \leq K \|\mathbf{x}_1 - \mathbf{x}_2\|,$$

where $K$ is the smallest constant for which the inequality holds.

Let $f(\mathbf{x}; \theta)$ be an attention mechanism. If $f$ is Lipschitz continuous with constant $K$, then the attention mechanism is robust to input perturbations, and the change in output is linearly bounded by the change in input:

$$\|f(\mathbf{x}_1; \theta) - f(\mathbf{x}_2; \theta)\| \leq K \|\mathbf{x}_1 - \mathbf{x}_2\|.$$

This property is critical in ensuring that the model's predictions are stable and reliable.

In tasks like image recognition or natural language processing, where input data can be noisy or imprecise, the robustness of attention mechanisms ensures that the model's output remains consistent, even in the face of such challenges.

## 3.8   Theoretical Challenges and Future Directions

As attention mechanisms continue to evolve and play a central role in modern deep learning architectures, several theoretical challenges remain. Addressing these challenges is crucial for the scalability, efficiency, and broader applicability of attention-based models. This section focuses on two key challenges: scalability to long sequences and memory efficiency. We will explore these challenges, using mathematical analysis to highlight the limitations and potential solutions.

### 3.8.1   Scalability Issues

One of the most significant challenges in attention mechanisms is scalability, particularly when dealing with long sequences. Traditional self-attention mechanisms exhibit quadratic complexity with respect to the sequence length, which makes them computationally infeasible for very long sequences. This section delves into the mathematical aspects of scaling attention mechanisms and explores potential avenues for improving scalability.

**Scaling to Long Sequences**

The challenge of scaling attention mechanisms to long sequences arises from the quadratic time and space complexity inherent in traditional self-attention. For an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, the attention mechanism computes the attention scores between every pair of elements, leading to a time complexity of $O(n^2 \cdot d)$ and a space complexity of $O(n^2)$, where $n$ is the sequence length and $d$ is the dimensionality of the embeddings.

Given the quadratic complexity, scaling self-attention to sequences with length $n \gg 1$ becomes computationally expensive. For example, if $n$ increases by a factor of 10, the computational cost increases by a factor of 100, making the mechanism impractical for sequences encountered in tasks like document processing, video analysis, or genomic data analysis.

Approach 1: Sparse attention mechanisms address this issue by reducing the number of interactions considered in the attention computation. Instead of computing attention scores for every pair of elements, sparse attention focuses on a subset $\mathcal{S}_i$ of relevant elements for each query $x_i$:

$$\alpha_{ij} = \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{k \in \mathcal{S}_i} \exp(\mathbf{q}_i^\top \mathbf{k}_k)}, \quad \text{for} \quad j \in \mathcal{S}_i.$$

By choosing $|\mathcal{S}_i| = k$, where $k \ll n$, the time complexity is reduced to $O(n \cdot k \cdot d)$.

Approach 2: Another approach to scaling is through low-rank factorization of the attention matrix. The attention matrix $\mathbf{A}$ is decomposed into lower-dimensional matrices $\mathbf{A} \approx \mathbf{BC}$, where $\mathbf{B} \in \mathbb{R}^{n \times r}$ and $\mathbf{C} \in \mathbb{R}^{r \times n}$, and $r$ is the rank of the decomposition. This reduces the time complexity to $O(n \cdot r \cdot d)$, where $r$ is typically much smaller than $n$.

Let $r$ be the rank of the attention matrix $\mathbf{A}$ after factorization. The time complexity of computing the factorized attention is

$$T(n, r, d) = O(n \cdot r \cdot d),$$

which represents a significant reduction from the original $O(n^2 \cdot d)$ complexity.

In natural language processing, when processing a book-length document, traditional self-attention would be computationally prohibitive. Sparse attention or low-rank factorization can be applied to handle the sequence length more efficiently, making it feasible to process such long sequences.

**Memory Efficiency**

Memory efficiency is another crucial consideration, particularly for training and deploying large models on hardware with limited memory resources. The memory complexity of traditional self-attention is $O(n^2)$, which can quickly become prohibitive as the sequence length $n$ increases.

Memory limitations become particularly pronounced when deploying models on edge devices or when training very large models on GPUs with finite memory capacity. The quadratic memory complexity of self-attention mechanisms often necessitates compromises, such as reducing the batch size or sequence length, which can negatively impact model performance.

Approach 1: Several memory-efficient attention mechanisms have been proposed to address this issue. One approach is to use memory-efficient sparse attention, where the attention computation is restricted to non-zero entries, significantly reducing the memory footprint. For example, if $k$ non-zero entries are considered per query, the memory complexity is reduced to $O(n \cdot k)$.

Approach 2: Another approach to improving memory efficiency is chunking the input sequence into smaller segments and applying attention within each segment independently. This reduces the memory requirement to $O((n/w) \cdot w^2)$, where $w$ is the chunk size. Recurrent mechanisms can be used to propagate information between chunks, maintaining long-range dependencies.

Let $w$ be the chunk size and $n$ the sequence length. The memory complexity of chunked attention is

$$S(n, w) = O\left(\frac{n}{w} \cdot w^2\right) = O(n \cdot w),$$

where $w$ can be chosen to balance between memory usage and attention range.

Example: In speech processing, where long audio sequences need to be processed, chunking the sequence and applying attention within each chunk allow the model to fit within the memory constraints of standard GPUs, while still capturing essential dependencies.

### 3.8.2   Interpretability and Explainability

Interpretability and explainability are crucial for understanding how attention mechanisms make decisions, particularly in high-stakes applications such as medical diagnosis, autonomous driving, and legal decision-making. Despite the effectiveness of attention mechanisms, their often "black-box" nature poses significant challenges to gaining insights into their decision-making processes.

**Visualizing Attention Weights**

Attention weights are central to the functioning of attention mechanisms, as they determine the relative importance of different elements in the input sequence. Given an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$, the attention weights $\alpha_{ij}$ computed between elements $x_i$ and $x_j$ are

$$\alpha_{ij} = \frac{\exp\left(\mathbf{q}_i^\top \mathbf{k}_j\right)}{\sum_{k=1}^{n} \exp\left(\mathbf{q}_i^\top \mathbf{k}_k\right)},$$

where $\mathbf{q}_i$ and $\mathbf{k}_j$ are the query and key vectors, respectively.

Visualizing these weights can provide insights into which parts of the input sequence the model considers important. For example, in a machine translation task, attention weights can reveal how the model aligns words in the source and target languages.

Let $\mathbf{A}$ be the attention matrix with elements $\alpha_{ij}$. The interpretability of attention weights can be enhanced by visualizing the matrix $\mathbf{A}$ as a heatmap, where each entry $\alpha_{ij}$ represents the strength of attention from $x_i$ to $x_j$:

$$\text{Interpretation}(\mathbf{A}) = \text{Heatmap}(\mathbf{A}),$$

where the heatmap visually encodes the attention weights, providing an intuitive representation of the attention mechanism's focus.

In natural language processing tasks, attention weight visualizations can help linguists understand how the model is processing language by showing which words the model is focusing on at each step of the sequence. This can lead to a better understanding of the model's behavior and highlight potential areas for improvement.

**Understanding Head Roles**

In multi-head attention mechanisms, each attention head computes its own set of attention weights, which are then combined to form the final output. Understanding the role of each head—whether it focuses on syntactic relationships, semantic content, or other features—can provide deeper insights into how the model processes information.

Given $H$ attention heads, the attention output for head $h$ is computed as

$$\mathbf{z}_i^{(h)} = \sum_{j=1}^{n} \alpha_{ij}^{(h)} \mathbf{v}_j^{(h)},$$

where $\alpha_{ij}^{(h)}$ are the attention weights for head $h$, and $\mathbf{v}_j^{(h)}$ are the corresponding value vectors. The final output is a combination of all heads:

$$\mathbf{z}_i = \text{Concat}(\mathbf{z}_i^{(1)}, \mathbf{z}_i^{(2)}, \ldots, \mathbf{z}_i^{(H)})\mathbf{W}_o,$$

where $\mathbf{W}_o$ is a learned weight matrix.

Let $\mathbf{z}_i^{(h)}$ represent the output of the $h$-th head. The diversity and specialization of attention heads can be quantified by the variance in their outputs:

$$\text{Diversity}(\mathbf{z}_i^{(1)}, \ldots, \mathbf{z}_i^{(H)}) = \frac{1}{H} \sum_{h=1}^{H} \left\| \mathbf{z}_i^{(h)} - \overline{\mathbf{z}_i} \right\|^2,$$

where $\overline{\mathbf{z}_i}$ is the mean output across all heads. A higher diversity score indicates that the heads are capturing different aspects of the input, leading to a richer overall representation.

Example: In a transformer model used for text classification, different heads might specialize in detecting negations, understanding subject–verb relationships, or capturing overall sentiment. Analyzing the output of each head can reveal these specializations and improve the interpretability of the model.

### 3.8.3    Future Directions in Attention Research

As attention mechanisms continue to be a focal point in deep learning research, several future directions are emerging. These include the development of novel attention mechanisms and the integration of attention with other model architectures to create hybrid models. This section explores these future directions with a focus on their theoretical foundations.

**Novel Attention Mechanisms**

The quest for novel attention mechanisms is driven by the need to address the limitations of current models, such as computational inefficiency, lack of interpretability, and challenges in capturing long-range dependencies. One promising direction is the development of adaptive attention mechanisms that can dynamically adjust their behavior based on the input data.

For instance, an adaptive attention mechanism might adjust the size of the attention window or the number of attention heads based on the complexity of the input sequence. Mathematically, this could involve a gating function $g(\mathbf{x})$ that determines the attention parameters $\theta$ for each input sequence:

$$\theta = g(\mathbf{x}) \quad \text{and} \quad \mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij}(\theta)\mathbf{v}_j(\theta),$$

where $\alpha_{ij}(\theta)$ and $\mathbf{v}_j(\theta)$ are functions of the dynamically adjusted parameters $\theta$.

Let $g(\mathbf{x})$ be a gating function that optimally adjusts the attention parameters based on the input sequence. The expected computational complexity $\mathbb{E}[T(\mathbf{x})]$ of an adaptive attention mechanism can be lower than that of a fixed attention mechanism, particularly for sequences with varying complexity:

$$\mathbb{E}[T(\mathbf{x})] < T_{\text{fixed}}(\mathbf{x}),$$

where $T_{\text{fixed}}(\mathbf{x})$ is the complexity of the corresponding fixed attention mechanism.

In visual tasks, an adaptive attention mechanism could dynamically allocate more attention resources to complex regions of an image, such as areas with high texture or detail, while reducing attention to simpler regions like backgrounds.

**Hybrid Models**

Hybrid models combine attention mechanisms with other types of neural architectures, such as CNNs or RNNs, to leverage the strengths of each approach. For example, a hybrid model might use a CNN to extract local features from an image and an attention mechanism to capture global dependencies.

Let $\mathbf{f}_{CNN}$ represent the feature extraction function of a CNN and $\mathbf{f}_{Attn}$ the function of an attention mechanism. The output $\mathbf{z}$ of a hybrid model can be represented as

$$\mathbf{z} = \mathbf{f}_{Attn}(\mathbf{f}_{CNN}(\mathbf{x})),$$

where $\mathbf{x}$ is the input data. The attention mechanism $\mathbf{f}_{Attn}$ processes the feature maps produced by the CNN to capture long-range dependencies and global context.

The expressivity of the hybrid model $\mathbf{f}_{Hybrid}$ is enhanced by the combination of local feature extraction and global context capture:

$$\mathbf{f}_{Hybrid}(\mathbf{x}) = \mathbf{f}_{Attn}(\mathbf{f}_{CNN}(\mathbf{x})) \quad \text{is more expressive than} \quad \mathbf{f}_{CNN}(\mathbf{x}) \text{ or } \mathbf{f}_{Attn}(\mathbf{x}) \text{ alone.}$$

This combination allows the model to perform better on tasks that require both local and global understanding, such as object detection or language translation.

Example: In video analysis, a hybrid model might use a CNN to process individual frames and an attention mechanism to capture temporal dependencies across frames, leading to more accurate action recognition.

## References

1. Beltagy, I., Peters, M.E., Cohan, A.: Longformer: The Long-document Transformer (2020). arXiv preprint arXiv:2004.05150
2. Choromanski, K., Likhosherstov, V., Dohan, D., Song, X., Gane, A., Sarlos, T., Hawkins, P., Davis, J., Mohiuddin, A., Kaiser, L., Belanger, D., Colwell, L., Weller, A.: Rethinking attention with performers (2020). arXiv preprint arXiv:2009.14794
3. Kitaev, N., Kaiser, L., Levskaya, A.: Reformer: The efficient transformer (2020). arXiv preprint arXiv:2001.04451
4. Wang, S., Li, B.Z., Khabsa, M., Fang, H., Ma, H.: Linformer: Self-attention with linear complexity (2020). arXiv preprint arXiv:2006.04768

# Chapter 4
# Transformer Architecture: Encoder and Decoder

## 4.1 Encoder Structure

The transformer architecture has revolutionized natural language processing by leveraging self-attention mechanisms to capture dependencies in sequential data without relying on recurrent or convolutional layers. The encoder is a critical component of the transformer, responsible for processing the input sequence and producing representations that the decoder or downstream tasks can utilize. This section explores the structure of the encoder, focusing on the composition of its layers, particularly the multi-head attention mechanism and feedforward network, as well as the role of activation functions.

The encoder in a transformer model consists of multiple layers, each composed of a multi-head attention mechanism followed by a feedforward network. The output of each layer is passed to the next, with residual connections and normalization steps ensuring stability and efficient training. The encoder's layered structure allows it to capture complex dependencies and hierarchical representations of the input data (see Fig. 4.1).

### 4.1.1 Layer Composition

Each layer in the encoder is a composite of several sub-layers, designed to process different aspects of the input. The primary sub-layers are the multi-head attention mechanism, which allows the model to focus on different parts of the input sequence simultaneously, and the feedforward network, which applies non-linear transformations to the attention outputs.

**Fig. 4.1** Diagram of a Transformer Encoder, showcasing the interplay of multi-head self-attention and feedforward layers, where linear transformations and positional encodings enable the modeling of contextual relationships within vector spaces

Input Tensor

Input Embedding

Positional Encoding

Multi-Head Attention

Add & Norm

Feed Forward Network

Add & Norm

Output Tensor

## Multi-Head Attention Layer

The multi-head attention mechanism is a cornerstone of the transformer's ability to process sequences in parallel. For an input sequence $\mathbf{X} \in \mathbb{R}^{n \times d}$, where $n$ is the sequence length and $d$ is the dimensionality of the embeddings, the multi-head attention mechanism computes attention across multiple "heads," each focusing on different aspects of the sequence.

Let $h$ denote the number of attention heads. The input $\mathbf{X}$ is first linearly transformed into query, key, and value matrices $\mathbf{Q}$, $\mathbf{K}$, and $\mathbf{V}$, respectively, for each head:

$$\mathbf{Q}_i = \mathbf{X}\mathbf{W}_i^Q, \quad \mathbf{K}_i = \mathbf{X}\mathbf{W}_i^K, \quad \mathbf{V}_i = \mathbf{X}\mathbf{W}_i^V,$$

where $\mathbf{W}_i^Q$, $\mathbf{W}_i^K$, and $\mathbf{W}_i^V$ are learned weight matrices for the $i$-th head. The attention score between queries and keys is computed as

$$\text{Attention}(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i) = \text{softmax}\left(\frac{\mathbf{Q}_i \mathbf{K}_i^\top}{\sqrt{d_k}}\right) \mathbf{V}_i,$$

where $d_k$ is the dimensionality of the keys (typically $d_k = d/h$). This operation results in an output matrix for each head, which is then concatenated and linearly transformed:

$$\mathbf{Z} = \text{Concat}(\mathbf{Z}_1, \mathbf{Z}_2, \ldots, \mathbf{Z}_h)\mathbf{W}^O,$$

where $\mathbf{Z}_i$ is the output of the $i$-th head, and $\mathbf{W}^O$ is a learned weight matrix that projects the concatenated outputs back to the original embedding space.

The multi-head attention mechanism can express a wider range of dependencies compared to a single attention head. Formally, let $\mathcal{A}$ be the set of attention functions expressible by a single head, and $\mathcal{M}$ be the set expressible by a multi-head mechanism. Then

$$\mathcal{A} \subset \mathcal{M},$$

implying that multi-head attention can capture more complex patterns and relationships in the data.

Example: In machine translation, multi-head attention allows the model to focus on different parts of the source sentence when generating each word in the target sentence. One head might focus on subject–verb agreement, while another might capture noun–adjective pairings.

**Mathematical Formulation**

The self-attention mechanism, as implemented in the transformer encoder, is mathematically described by the following steps:

1. Linear Transformations: Each input vector $\mathbf{x}_i$ from the sequence is linearly transformed into a query $\mathbf{q}_i$, a key $\mathbf{k}_i$, and a value $\mathbf{v}_i$ using the learned matrices $\mathbf{W}^Q$, $\mathbf{W}^K$, and $\mathbf{W}^V$, respectively:

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q, \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K, \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V.$$

2. Attention Weight Calculation: The attention weight between two positions $i$ and $j$ in the sequence is computed using the scaled dot product:

$$\alpha_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}}\right)}{\sum_{j'=1}^{n} \exp\left(\frac{\mathbf{q}_i \mathbf{k}_{j'}^\top}{\sqrt{d_k}}\right)},$$

where $\alpha_{ij}$ represents the attention score between $\mathbf{q}_i$ and $\mathbf{k}_j$, normalized across all keys.

3. Weighted Sum of Values: The output for each position $i$ is a weighted sum of the values, weighted by the attention scores:

$$\mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{v}_j.$$

4. Multi-Head Attention Output: The outputs from multiple heads are concatenated and projected back to the original dimension using a learned matrix $\mathbf{W}^O$:

$$\mathbf{Z} = \text{Concat}(\mathbf{z}_1, \mathbf{z}_2, \ldots, \mathbf{z}_h)\mathbf{W}^O.$$

Example: In a sentence such as "The cat sat on the mat," the attention mechanism allows the model to focus on the relevant parts of the sentence for each word being processed. For example, when processing "sat," the model might pay more attention to "cat" than to other words.

## Attention Weight Calculation

The calculation of attention weights is a critical part of the transformer's ability to learn dependencies between different positions in the sequence. The attention weight $\alpha_{ij}$ is computed as

$$\alpha_{ij} = \frac{\exp\left(\mathbf{q}_i \cdot \mathbf{k}_j / \sqrt{d_k}\right)}{\sum_{k=1}^{n} \exp\left(\mathbf{q}_i \cdot \mathbf{k}_k / \sqrt{d_k}\right)},$$

where $\mathbf{q}_i$ and $\mathbf{k}_j$ are the query and key vectors, respectively, for positions $i$ and $j$, and $d_k$ is the dimensionality of the key vectors. The use of the softmax function ensures that the attention weights sum to one, making them interpretable as probabilities.

The stability of the attention mechanism is ensured by the scaling factor $\sqrt{d_k}$. Without this factor, the dot products $\mathbf{q}_i \cdot \mathbf{k}_j$ could grow large in magnitude, leading to small gradients during training. By scaling the dot products, the gradients remain in a stable range, facilitating efficient learning.

Example: If $d_k = 64$, and the dot product $\mathbf{q}_i \cdot \mathbf{k}_j$ is 10, the scaled dot product becomes $10/\sqrt{64} = 1.25$. The softmax function applied to this value will produce a moderate attention weight, avoiding extreme values that could destabilize training.

## Head Concatenation and Projection

After calculating the attention scores for each head, the outputs from all heads are concatenated to form a single vector. This concatenated vector is then projected back into the original embedding space using a learned projection matrix $\mathbf{W}^O$:

$$\mathbf{Z} = \text{Concat}(\mathbf{z}_1, \mathbf{z}_2, \ldots, \mathbf{z}_h)\mathbf{W}^O,$$

where $\mathbf{z}_i$ represents the output from the $i$-th head. The projection matrix $\mathbf{W}^O$ ensures that the multi-head attention output can be directly combined with the outputs of other layers in the transformer.

The projection step ensures that the multi-head attention mechanism retains the ability to express complex relationships while reducing the dimensionality of the concatenated output. Formally, let $\mathbf{Z} \in \mathbb{R}^{n \times (h \cdot d_v)}$ be the concatenated output from $h$ heads. The projection $\mathbf{W}^O$ maps $\mathbf{Z}$ back to $\mathbb{R}^{n \times d}$, preserving the richness of the multi-head attention while matching the model's input dimension.

Example: In machine translation, after processing different linguistic aspects through multiple heads, the projection step combines these aspects into a unified representation that can be used by subsequent layers to generate the final translation.

## Feedforward Layer

Following the multi-head attention mechanism, each layer in the transformer encoder includes a feedforward network (FFN) that applies non-linear transformations to the attention outputs. The FFN consists of two linear transformations with a non-linear activation function in between:

$$\text{FFN}(\mathbf{z}) = \max(0, \mathbf{z}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2,$$

where $\mathbf{W}_1$ and $\mathbf{W}_2$ are weight matrices, $\mathbf{b}_1$ and $\mathbf{b}_2$ are bias vectors, and $\max(0, x)$ denotes the ReLU activation function.

The two-layer FFN with non-linear activation is a universal approximator, meaning it can approximate any continuous function on a compact domain, given sufficient dimensionality. Formally, for any continuous function $f : \mathbb{R}^d \to \mathbb{R}$ and any $\epsilon > 0$, there exist weight matrices $\mathbf{W}_1$, $\mathbf{W}_2$ and bias vectors $\mathbf{b}_1$, $\mathbf{b}_2$ such that

$$\|\text{FFN}(\mathbf{z}) - f(\mathbf{z})\| < \epsilon \quad \text{for all } \mathbf{z} \in \mathbb{R}^d.$$

In a transformer model, the FFN might be responsible for transforming the output of the attention mechanism into a form that captures complex non-linear interactions between the tokens, enabling the model to understand more abstract aspects of the input data.

## Layer Composition and Operation

To recap, each layer in the transformer encoder consists of the following components arranged in sequence:

1. Multi-Head Attention: Computes attention scores and outputs weighted sums of values.

2. Add and Norm: A residual connection adds the input of the layer to the output of the attention mechanism, followed by layer normalization.

3. Feedforward Network: Applies non-linear transformations to the normalized output of the attention mechanism.

4. Add and Norm: A second residual connection adds the input of the FFN to its output, followed by another layer normalization.

Let $\mathbf{X}_{in}$ be the input to a layer. The layer output $\mathbf{X}_{out}$ is given by

$$\mathbf{X}_{out} = \text{LayerNorm}(\text{FFN}(\text{LayerNorm}(\text{MultiHeadAttention}(\mathbf{X}_{in}) + \mathbf{X}_{in})) + \text{LayerNorm}(\mathbf{X}_{in})).$$

The use of residual connections and layer normalization ensures that the layer operations remain stable during training. Residual connections help mitigate the vanishing gradient problem, while layer normalization maintains a consistent range of activations, facilitating efficient gradient-based optimization.

In practice, this layered structure allows the transformer encoder to build hierarchical representations of the input data, where each layer captures increasingly abstract features. The stability provided by residual connections and normalization ensures that these features are learned effectively during training.

**Activation Functions (ReLU, GELU)**

Activation functions introduce non-linearity into the model, enabling it to learn complex mappings from inputs to outputs. The two most commonly used activation functions in transformers are the Rectified Linear Unit (ReLU) and the Gaussian Error Linear Unit (GELU).

The ReLU activation function is defined as

$$\text{ReLU}(x) = \max(0, x).$$

ReLU is computationally efficient and helps mitigate the vanishing gradient problem by allowing gradients to flow through the network for positive inputs.

The GELU activation function is defined as

$$\text{GELU}(x) = x \cdot \Phi(x),$$

where $\Phi(x)$ is the cumulative distribution function of the standard normal distribution. GELU smoothly approximates the ReLU function and has been shown to improve performance in transformer models.

The GELU function is differentiable and smooth, meaning its derivative does not have discontinuities. This smoothness can lead to better convergence properties during training, especially in deep models.

In transformer models, GELU is often preferred over ReLU because its smoothness can lead to more stable training and slightly better performance on complex tasks like language modeling.

### 4.1.2  Normalization and Residual Connections

The stability and trainability of deep neural networks, such as transformers, rely heavily on normalization techniques and the strategic use of residual connections. These components work together to maintain the flow of information and gradients through the network, enabling the training of very deep models without degradation in performance.

**Layer Normalization**

Layer normalization is a technique that normalizes the activations within each layer of the network, ensuring that they have a consistent distribution. Unlike batch normalization, which normalizes across a batch of data, layer normalization operates on each individual sample independently.

Given an input $\mathbf{x} = [x_1, x_2, \ldots, x_d]$ to a layer, where $d$ is the dimensionality of the input, layer normalization transforms $\mathbf{x}$ as follows:

$$\mu_{\mathbf{x}} = \frac{1}{d} \sum_{i=1}^{d} x_i, \quad \sigma_{\mathbf{x}}^2 = \frac{1}{d} \sum_{i=1}^{d} (x_i - \mu_{\mathbf{x}})^2,$$

$$\hat{x}_i = \frac{x_i - \mu_{\mathbf{x}}}{\sqrt{\sigma_{\mathbf{x}}^2 + \epsilon}},$$

$$\text{LayerNorm}(\mathbf{x})_i = \gamma \hat{x}_i + \beta,$$

where $\mu_{\mathbf{x}}$ is the mean, $\sigma_{\mathbf{x}}^2$ is the variance of the input, $\gamma$ and $\beta$ are learnable parameters that scale and shift the normalized output, and $\epsilon$ is a small constant added for numerical stability.

Layer normalization is invariant to affine transformations of the input. That is, for any input $\mathbf{x}$ and affine transformation $\mathbf{Ax} + \mathbf{b}$, the layer normalization output satisfies

$$\text{LayerNorm}(\mathbf{Ax} + \mathbf{b}) = \text{LayerNorm}(\mathbf{x}).$$

This invariance ensures that the normalization is robust to changes in the scale and bias of the input, contributing to more stable training.

In a transformer model, layer normalization is applied after the multi-head attention mechanism and after the feedforward network. This ensures that the outputs of these components have a stable distribution, which is crucial for the model's ability to learn effectively.

## Impact on Training Stability

Layer normalization improves training stability by reducing the internal covariate shift, which refers to changes in the distribution of layer inputs during training. By normalizing the activations, the learning process becomes more predictable, allowing for faster convergence and reducing the sensitivity to hyperparameters such as the learning rate.

Let $\mathbf{z}^{(l)}$ be the activations at layer $l$ of a deep network with layer normalization. The gradient of the loss function $\mathcal{L}$ with respect to the parameters of layer $l$ is given by

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(l)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{(l)}} \cdot \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{W}^{(l)}},$$

where $\mathbf{W}^{(l)}$ are the weights of layer $l$. Layer normalization ensures that $\frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{W}^{(l)}}$ remains well behaved, preventing gradients from exploding or vanishing and thus facilitating convergence.

In deep transformers, the use of layer normalization in every layer helps maintain the consistency of the activations throughout the network, enabling the model to learn complex patterns over many layers without suffering from the vanishing or exploding gradient problems.

## Residual Connections

Residual connections, also known as skip connections, are added between layers in deep neural networks to allow gradients to flow more easily through the network. These connections are particularly important in very deep models, where the risk of gradient vanishing increases with the depth of the network.

In a residual block, the input $\mathbf{x}$ to a layer is added directly to the output $\mathbf{F}(\mathbf{x})$ of that layer:

$$\mathbf{y} = \mathbf{x} + \mathbf{F}(\mathbf{x}),$$

where $\mathbf{F}(\mathbf{x})$ represents the transformation applied by the layer (e.g., the combination of multi-head attention and feedforward network in a transformer).

Let $\mathcal{L}$ be the loss function, and consider a deep network with residual connections. The gradient of the loss with respect to the input $\mathbf{x}$ of a residual block is given by

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}} \cdot \left(1 + \frac{\partial \mathbf{F}(\mathbf{x})}{\partial \mathbf{x}}\right).$$

The identity component in the gradient ensures that gradients can flow directly through the network, bypassing the transformation $\mathbf{F}(\mathbf{x})$ if necessary, which mitigates the vanishing gradient problem.

In a transformer model, residual connections are used around both the multi-head attention mechanism and the feedforward network. These connections allow the model to learn identity mappings more easily, which is important for preserving the original input information while adding new features.

**Skip Connections**

Skip connections are a special case of residual connections where the output of a layer is added to the input of a non-adjacent layer. These connections are particularly useful in networks with hierarchical structures, allowing higher layers to directly access information from lower layers.

If $\mathbf{x}_1$ is the input to an earlier layer and $\mathbf{F}_2(\mathbf{x}_2)$ is the output of a later layer, a skip connection adds $\mathbf{x}_1$ to $\mathbf{F}_2(\mathbf{x}_2)$:

$$\mathbf{y} = \mathbf{F}_2(\mathbf{x}_2) + \mathbf{x}_1.$$

This connection enables the model to combine low-level features with high-level features, improving the network's ability to learn complex patterns.

Skip connections facilitate hierarchical learning by allowing gradients to flow directly between non-adjacent layers. Formally, let $\mathcal{L}$ be the loss function, and let $\mathbf{x}_1$ and $\mathbf{x}_2$ be the inputs to non-adjacent layers with a skip connection. The gradient with respect to $\mathbf{x}_1$ is given by

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}_1} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}} + \frac{\partial \mathcal{L}}{\partial \mathbf{F}_2(\mathbf{x}_2)} \cdot \frac{\partial \mathbf{F}_2(\mathbf{x}_2)}{\partial \mathbf{x}_1},$$

where the direct gradient flow through $\mathbf{y}$ bypasses intermediate layers, ensuring that the model can effectively learn both low-level and high-level features.

In a deep transformer network, skip connections might be used to allow the model to combine low-level syntactic features with high-level semantic features, improving the overall performance on tasks like machine translation or text summarization.

**Effect on Gradient Flow**

Residual and skip connections significantly improve the flow of gradients through the network during backpropagation. This improved gradient flow prevents the gradients

from becoming too small as they are propagated through many layers, a problem known as the vanishing gradient problem.

Consider a deep network without residual connections. The gradient of the loss $\mathcal{L}$ with respect to the parameters of the first layer is given by

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} = \frac{\partial \mathcal{L}}{\partial \mathbf{W}_L} \cdot \prod_{i=1}^{L-1} \frac{\partial \mathbf{z}_{i+1}}{\partial \mathbf{z}_i},$$

where $\mathbf{W}_i$ are the weights of the $i$-th layer and $\mathbf{z}_i$ are the activations at layer $i$. In deep networks, the product of derivatives can become very small, leading to vanishing gradients.

With residual connections, the gradient is modified as follows:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} = \frac{\partial \mathcal{L}}{\partial \mathbf{W}_L} \cdot \left(1 + \prod_{i=1}^{L-1} \frac{\partial \mathbf{z}_{i+1}}{\partial \mathbf{z}_i}\right),$$

where the identity component "1" prevents the gradient from vanishing entirely, ensuring that it remains large enough to update the early layers effectively.

In a transformer with many layers, the use of residual connections ensures that the gradients can flow back through the entire network, allowing the model to learn effectively from the very first layer.

**Normalization Techniques Comparison**

Normalization techniques such as batch normalization and layer normalization serve similar purposes but operate differently and are suited to different types of networks. In the context of transformers, layer normalization is typically preferred due to the nature of the data and the architecture.

Batch normalization normalizes the activations across a mini-batch of data. For an input $\mathbf{x}$ within a batch, the normalized output is

$$\mu_{\text{batch}} = \frac{1}{m} \sum_{i=1}^{m} x_i, \quad \sigma_{\text{batch}}^2 = \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\text{batch}})^2,$$

$$\hat{x}_i = \frac{x_i - \mu_{\text{batch}}}{\sqrt{\sigma_{\text{batch}}^2 + \epsilon}},$$

$$\text{BatchNorm}(\mathbf{x})_i = \gamma \hat{x}_i + \beta,$$

where $m$ is the size of the mini-batch.

Layer normalization, as previously defined, normalizes across the features within each individual data point, rather than across the batch.

Batch normalization is effective in convolutional networks and other architectures where batch-level statistics are meaningful. However, in sequence models like transformers, where each input might have different statistical properties, layer normalization is more effective because it normalizes each data point individually, maintaining consistency across varying input sequences.

Both batch normalization and layer normalization reduce the variance of activations, which stabilizes the learning process. However, the variance reduction occurs at different levels (batch versus feature), influencing the choice of normalization technique based on the architecture.

In transformers, where inputs can vary greatly in length and content, layer normalization ensures that each token is treated consistently, regardless of the batch composition, leading to more stable and efficient training.

## 4.2 Decoder Structure

The decoder in the transformer architecture is designed to generate output sequences based on the encoded representations provided by the encoder. It is especially crucial in tasks such as machine translation, where the decoder must generate a sequence of words in the target language that corresponds to the sequence in the source language. The decoder is composed of multiple layers, each containing a combination of masked multi-head attention, cross-attention, and feedforward networks. This section explores the composition of these layers, focusing on the masked multi-head attention mechanism, cross-attention, and their interaction with the encoder outputs (see Fig. 4.2).

### 4.2.1 Layer Composition

Each layer of the transformer decoder consists of three main sub-layers: masked multi-head attention, cross-attention, and a feedforward network. The masked multi-head attention sub-layer allows the decoder to focus on previous positions in the output sequence, ensuring that future tokens are not considered during generation. The cross-attention sub-layer allows the decoder to attend to the encoder's outputs, integrating information from the source sequence. Finally, the feedforward network applies non-linear transformations to the attention outputs, adding depth and complexity to the model's representations.

**Fig. 4.2** Diagram of a
transformer decoder,
featuring masked multi-head
self-attention and
cross-attention layers,
designed to process
sequential data while
integrating context from the
encoder outputs

Output Sentence

**Output Embedding**

**Positional Encoding**

**Masked Multi-Head Attention**

**Add & Norm**

Refined Embeddings
from Encoder

N x

**Multi-Head Attention**

**Add & Norm**

**Feed Forward Network**

**Add & Norm**

**Linear**

**Softmax**

Output Probabilities

## Masked Multi-Head Attention

The masked multi-head attention mechanism in the decoder is similar to the multi-head attention mechanism in the encoder, with the key difference being the application of a mask to prevent the model from attending to future tokens. This ensures that the decoder generates the output sequence one token at a time.

Let $\mathbf{Y} \in \mathbb{R}^{m \times d}$ represent the input to the decoder at a particular layer, where $m$ is the length of the sequence generated so far and $d$ is the dimensionality of the embeddings. The masked multi-head attention mechanism computes attention in the following steps:

1. Linear Transformations: The input $\mathbf{Y}$ is linearly transformed into query, key, and value matrices for each attention head:

$$\mathbf{Q}_i = \mathbf{Y}\mathbf{W}_i^Q, \quad \mathbf{K}_i = \mathbf{Y}\mathbf{W}_i^K, \quad \mathbf{V}_i = \mathbf{Y}\mathbf{W}_i^V,$$

where $\mathbf{W}_i^Q$, $\mathbf{W}_i^K$, and $\mathbf{W}_i^V$ are learned weight matrices for the $i$-th head.

2. Masking Mechanism: A mask $\mathbf{M}$ is applied to the attention scores to prevent the model from attending to future tokens. The attention score between tokens $i$ and $j$ is computed as

$$\alpha_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}} + M_{ij}\right)}{\sum_{j'=1}^{m} \exp\left(\frac{\mathbf{q}_i \mathbf{k}_{j'}^\top}{\sqrt{d_k}} + M_{ij'}\right)},$$

where $M_{ij}$ is set to a large negative value (e.g., $-\infty$) if $j > i$, effectively masking future positions.

3. Weighted Sum of Values: The masked attention score $\alpha_{ij}$ is then used to compute the weighted sum of the values:

$$\mathbf{z}_i = \sum_{j=1}^{m} \alpha_{ij} \mathbf{v}_j.$$

4. Multi-Head Attention Output: The outputs from all heads are concatenated and linearly transformed:
$$\mathbf{Z} = \text{Concat}(\mathbf{z}_1, \mathbf{z}_2, \ldots, \mathbf{z}_h)\mathbf{W}^O,$$

where $\mathbf{W}^O$ is a learned projection matrix.

The masked multi-head attention mechanism ensures that each position $i$ in the output sequence depends only on the tokens at positions $\leq i$. Formally, for any position $i$, the output $\mathbf{z}_i$ is independent of the inputs $\mathbf{y}_j$ for all $j > i$:

$$P(\mathbf{z}_i \mid \mathbf{y}_{>i}) = P(\mathbf{z}_i).$$

This property ensures that the decoder generates sequences in a causal manner, respecting the order of the sequence.

In machine translation, when generating a sentence in the target language, the masked multi-head attention ensures that the model only considers the words that have been generated so far, avoiding any information about the future tokens that have not yet been generated.

**Masking Mechanism and Its Necessity**

The masking mechanism is essential for ensuring the correct autoregressive behavior in sequence generation tasks. Without masking, the model could access future tokens during training, leading to data leakage and incorrect learning of dependencies.

The mask matrix $\mathbf{M}$ is defined as

$$
M_{ij} = \begin{cases} 0 & \text{if } i \geq j, \\ -\infty & \text{if } i < j, \end{cases}
$$

where $i$ and $j$ index the positions in the sequence. This mask is added to the attention logits before applying the softmax function, ensuring that the attention scores for future positions are effectively zero.

For an autoregressive model, the prediction $\hat{y}_i$ at position $i$ should only depend on the inputs $y_1, y_2, \ldots, y_{i-1}$. The masking mechanism guarantees this by enforcing:

$$
P(\hat{y}_i \mid y_1, y_2, \ldots, y_n) = P(\hat{y}_i \mid y_1, y_2, \ldots, y_{i-1}),
$$

where $n$ is the total sequence length. This ensures that future tokens do not influence the prediction at position $i$.

When generating a sentence, the word at position $i + 1$ should not influence the generation of the word at position $i$. Masking enforces this by ensuring that the model does not consider future positions when computing the attention scores.

**Attention Computation**

Attention computation in the masked multi-head attention follows the same principles as in the encoder, with the key difference being the application of the mask. The attention scores $\alpha_{ij}$ are computed using the masked softmax, ensuring that only past and present tokens are considered for each position in the sequence.

For each attention head, the attention score is calculated as

$$
\alpha_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}} + M_{ij}\right)}{\sum_{j'=1}^{m} \exp\left(\frac{\mathbf{q}_i \mathbf{k}_{j'}^\top}{\sqrt{d_k}} + M_{ij'}\right)},
$$

where $M_{ij}$ ensures that the model does not attend to future tokens.

The masking mechanism introduces sparsity into the attention matrix, where many attention scores are effectively zero due to the mask. Formally, for a sequence of length $m$, the attention matrix $\mathbf{A}$ satisfies

$$\text{Sparsity}(\mathbf{A}) = \frac{1}{m^2} \sum_{i=1}^{m} \sum_{j=i+1}^{m} \mathbf{1}(\alpha_{ij} = 0) \approx \frac{m(m-1)}{2m^2} = \frac{1}{2} - \frac{1}{2m}.$$

As $m$ increases, the attention matrix becomes increasingly sparse, with half of the potential connections being masked.

Example: In a sequence-to-sequence model for text generation, this sparsity ensures that the model focuses on the relevant past context, avoiding the risk of data leakage from future tokens.

**Cross-Attention**

Cross-attention, also known as encoder–decoder attention, allows the decoder to attend to the encoder's output. This mechanism enables the decoder to incorporate information from the input sequence, effectively guiding the generation of the output sequence based on the encoded representations.

Given the encoder output $\mathbf{H} \in \mathbb{R}^{n \times d}$, where $n$ is the length of the input sequence, the cross-attention mechanism in the decoder computes attention between the decoder's queries $\mathbf{Q}_i$ and the encoder's keys and values $\mathbf{K}_j$ and $\mathbf{V}_j$:

$$\alpha_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}}\right)}{\sum_{j'=1}^{n} \exp\left(\frac{\mathbf{q}_i \mathbf{k}_{j'}^\top}{\sqrt{d_k}}\right)},$$

$$\mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{v}_j,$$

where $\mathbf{Q}_i = \mathbf{Y}\mathbf{W}_i^Q$, $\mathbf{K}_j = \mathbf{H}\mathbf{W}_i^K$, and $\mathbf{V}_j = \mathbf{H}\mathbf{W}_i^V$.

Cross-attention facilitates the flow of information from the encoder to the decoder, ensuring that the output sequence is generated in alignment with the input sequence. Formally, the output $\mathbf{z}_i$ at position $i$ in the decoder is a weighted combination of the encoder outputs $\mathbf{H}$, allowing the decoder to incorporate information from the entire input sequence.

Example: In machine translation, the cross-attention mechanism allows the decoder to focus on relevant parts of the source sentence when generating each word in the target sentence. For example, when generating a verb, the decoder might focus on the corresponding noun in the source sentence to ensure grammatical correctness.

**Interaction with Encoder Outputs**

The interaction between the decoder and encoder outputs is mediated through the cross-attention mechanism, where the encoder's representations guide the decoder in generating the output sequence. This interaction ensures that the output sequence is semantically aligned with the input sequence.

The decoder's cross-attention sub-layer takes the encoder output $\mathbf{H}$ and the decoder input $\mathbf{Y}$, computing the attention scores and resulting representations as

$$\mathbf{Z}_{\text{cross}} = \text{CrossAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}),$$

where $\mathbf{Q}$ is derived from $\mathbf{Y}$, and $\mathbf{K}$ and $\mathbf{V}$ are derived from $\mathbf{H}$.

Cross-attention ensures that the semantic content of the output sequence aligns with that of the input sequence. For any position $i$ in the output sequence, the cross-attention output $\mathbf{z}_i$ depends on the relevant encoder outputs, ensuring that the generated sequence is contextually appropriate.

Example: In machine translation, cross-attention helps the model maintain alignment between the source and target languages, ensuring that words and phrases in the output sequence correspond accurately to their counterparts in the input sequence.

**Weight Sharing with Encoder**

In some transformer architectures, weights may be shared between the encoder and decoder, particularly in the attention mechanisms and feedforward networks. Weight sharing reduces the number of parameters in the model, leading to more efficient training and inference.

Weight sharing between the encoder and decoder involves using the same weight matrices for corresponding layers. For example, if the encoder uses weight matrix $\mathbf{W}^Q$ for queries, the decoder might use the same matrix $\mathbf{W}^Q$ for its query transformations:

$$\mathbf{Q}_{\text{encoder}} = \mathbf{X}\mathbf{W}^Q, \quad \mathbf{Q}_{\text{decoder}} = \mathbf{Y}\mathbf{W}^Q.$$

Let $\mathcal{P}_{\text{shared}}$ and $\mathcal{P}_{\text{separate}}$ denote the number of parameters in models with shared and separate weights, respectively. Then

$$\mathcal{P}_{\text{shared}} = \frac{1}{2}\mathcal{P}_{\text{separate}}.$$

Weight sharing reduces the parameter count by half, leading to more efficient use of computational resources.

Example: In a transformer model designed for low-resource settings, weight sharing between the encoder and decoder helps reduce the model size, making it more feasible to train on limited data or deploy on devices with constrained memory.

## *4.2.2   Output Processing*

After the decoder generates the final representations through masked multi-head attention and cross-attention mechanisms, these representations must be processed to produce the final output tokens. This involves a series of mathematical operations that project the decoder's output into a space corresponding to the vocabulary and convert these projections into probabilities that can be used to select the most likely next token in the sequence.

### Linear Transformation

The decoder outputs a sequence of vectors $\mathbf{Z} \in \mathbb{R}^{m \times d}$, where $m$ is the sequence length and $d$ is the dimensionality of the model. To map these vectors into a space that corresponds to the vocabulary size, a linear transformation is applied. Let $\mathbf{W}_{\text{out}} \in \mathbb{R}^{d \times V}$ be the weight matrix associated with this transformation, where $V$ is the size of the output vocabulary. The linear transformation is defined as

$$\mathbf{O} = \mathbf{Z}\mathbf{W}_{\text{out}},$$

where $\mathbf{O} \in \mathbb{R}^{m \times V}$ is the output matrix that contains the unnormalized logits for each position in the sequence.

A linear transformation followed by a non-linear activation function can approximate any continuous function to arbitrary accuracy, given sufficient dimensionality of the transformation matrix. This theorem, based on the universal approximation theorem, ensures that the linear transformation in the transformer is capable of expressing the necessary mappings from the decoder outputs to the output vocabulary space.

Example: In a machine translation task, the linear transformation maps the rich contextual embeddings produced by the decoder into a space where each dimension corresponds to a word in the target-language vocabulary. This mapping is crucial for generating the final translated sentence.

### Projection to Output Vocabulary

The output matrix $\mathbf{O}$ obtained from the linear transformation contains logits for each position in the sequence across the entire vocabulary. Each row $\mathbf{o}_i \in \mathbb{R}^V$ of $\mathbf{O}$ represents the unnormalized scores for each vocabulary word at position $i$ in the sequence. The next step is to project these logits into a probability distribution using the softmax function.

The softmax function is applied to each row $\mathbf{o}_i$ to produce a probability distribution over the vocabulary:

$$P(y_i = w \mid \mathbf{o}_i) = \frac{\exp(o_{iw})}{\sum_{v=1}^{V} \exp(o_{iv})},$$

where $P(y_i = w \mid \mathbf{o}_i)$ is the probability of selecting word $w$ from the vocabulary as the output at position $i$ and $o_{iw}$ is the logit corresponding to word $w$.

The softmax function maps the logits $\mathbf{o}_i$ into a probability distribution over the vocabulary. Formally, for any logits $\mathbf{o}_i \in \mathbb{R}^V$,

$$\sum_{w=1}^{V} P(y_i = w \mid \mathbf{o}_i) = 1,$$

and $P(y_i = w \mid \mathbf{o}_i) \geq 0$ for all $w$. This ensures that the output is a valid probability distribution.

Example: If the logits for a particular position in the sequence strongly favor the word "cat," the softmax function will assign a high probability to "cat" and lower probabilities to other words like "dog" or "fish."

## Softmax Activation

The softmax activation is essential for converting the logits into a probability distribution, which can then be used to make predictions about the next token in the sequence.

The softmax function has several important properties that make it suitable for this task:

1. Differentiability: The softmax function is smooth and differentiable, which is crucial for gradient-based optimization methods.

2. Exponential Sensitivity: The use of the exponential function ensures that large differences in logits lead to large differences in probabilities, making the model sensitive to differences in the unnormalized scores.

Let $\mathbf{p}_i = \text{softmax}(\mathbf{o}_i)$. The gradient of the softmax function with respect to the logits is given by

$$\frac{\partial p_{iw}}{\partial o_{iv}} = p_{iw}(\delta_{wv} - p_{iv}),$$

where $\delta_{wv}$ is the Kronecker delta. This gradient is used during backpropagation to update the weights in the linear transformation matrix.

Example: In a transformer model, during training, the softmax activation helps adjust the logits so that the probability of the correct word in the sequence is maximized. This is done by minimizing the cross-entropy loss between the predicted probabilities and the true one-hot encoded vectors of the target words.

**Probability Distribution Over Vocabulary**

The output of the softmax function is a probability distribution over the vocabulary, from which the model selects the next word in the sequence. This probability distribution reflects the model's confidence in each possible word given the context provided by the previous words and the encoder's output.

For each position $i$ in the sequence, the probability distribution over the vocabulary is given by

$$\mathbf{p}_i = \text{softmax}(\mathbf{o}_i),$$

where $\mathbf{p}_i$ is a vector of probabilities for each word in the vocabulary. The model can then select the word with the highest probability or sample from the distribution.

The entropy $H$ of the probability distribution $\mathbf{p}_i$ is given by

$$H(\mathbf{p}_i) = -\sum_{w=1}^{V} P(y_i = w \mid \mathbf{o}_i) \log P(y_i = w \mid \mathbf{o}_i).$$

High entropy indicates a more uniform distribution (more uncertainty in the prediction), while low entropy indicates a peaked distribution (higher confidence in a specific word).

Example: In a language generation task, if the model is very confident about the next word, the entropy of the distribution will be low, indicating that most of the probability mass is concentrated on a single word.

**Temperature Scaling**

Temperature scaling is a technique used to adjust the sharpness of the probability distribution produced by the softmax function. By controlling the temperature, one can make the model's predictions more or less confident, depending on the task's requirements.

Given the logits $\mathbf{o}_i$, temperature scaling modifies the softmax function as follows:

$$P(y_i = w \mid \mathbf{o}_i, T) = \frac{\exp(o_{iw}/T)}{\sum_{v=1}^{V} \exp(o_{iv}/T)},$$

where $T$ is the temperature parameter. When $T = 1$, the distribution is the same as the standard softmax. When $T > 1$, the distribution becomes softer (more uniform), and when $T < 1$, the distribution becomes sharper (more peaked).

Let $H(T)$ be the entropy of the softmax distribution with temperature $T$. Then

$$\frac{dH(T)}{dT} > 0,$$

indicating that increasing the temperature $T$ increases the entropy of the distribution, making the model's predictions more uncertain.

Example: In text generation, a higher temperature might be used to encourage more diversity in the generated text, leading to more creative outputs, while a lower temperature might be used in tasks requiring more deterministic and accurate predictions.

## 4.3   Mathematical Analysis

The interaction between the encoder and decoder in a transformer model is fundamental to its success in sequence-to-sequence tasks such as machine translation, summarization, and more. These interactions are mediated by the attention mechanisms, which allow the decoder to effectively utilize the information encoded by the encoder. This section analyzes the mathematical foundations of these interactions, focusing on the flow of information between the encoder and decoder, the role of the attention mechanism, and the impact of layer depth on these interactions.

### 4.3.1   Encoder–Decoder Interactions

The relationship between the encoder and decoder in a transformer model is characterized by a bidirectional flow of information, primarily mediated through the cross-attention mechanism. Understanding this interaction is crucial for appreciating how transformers can generate coherent and contextually appropriate outputs based on the input sequence.

**Information Flow Between Encoder and Decoder**

The information flow between the encoder and decoder is governed by the cross-attention mechanism in the decoder. For a given position $i$ in the decoder sequence, the decoder's query $\mathbf{q}_i$ interacts with the keys $\mathbf{K}$ and values $\mathbf{V}$ produced by the encoder:

$$\mathbf{K} = \mathbf{H}\mathbf{W}_K, \quad \mathbf{V} = \mathbf{H}\mathbf{W}_V,$$

where $\mathbf{H}$ represents the output of the encoder. The attention score $\alpha_{ij}$ between the decoder's query $\mathbf{q}_i$ and the encoder's key $\mathbf{k}_j$ is given by

$$\alpha_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}}\right)}{\sum_{j'=1}^{n} \exp\left(\frac{\mathbf{q}_i \mathbf{k}_{j'}^\top}{\sqrt{d_k}}\right)},$$

where $d_k$ is the dimensionality of the keys. The information from the encoder is then aggregated as a weighted sum of the encoder's values:

$$\mathbf{z}_i = \sum_{j=1}^{n} \alpha_{ij} \mathbf{v}_j.$$

This aggregation allows the decoder to utilize information from the entire input sequence while generating each token in the output sequence.

The cross-attention mechanism preserves the essential information from the encoder, allowing it to influence the decoder's output at each position. Formally, the aggregated output $\mathbf{z}_i$ captures the relevant information from the encoder's output, ensuring that the decoder's generation process is contextually grounded in the input sequence.

Example: In machine translation, when the decoder is generating the word "jumps" in the sentence "The cat jumps over the dog," it might focus on the corresponding word "saut" in the French input, ensuring that the translation remains accurate and contextually appropriate.

**Attention Mechanism in Encoder–Decoder Interactions**

The attention mechanism plays a central role in how the encoder and decoder interact. It allows the decoder to selectively focus on different parts of the input sequence based on the current context in the output sequence.

The attention mechanism in the encoder–decoder interaction can be described as a mapping from the decoder's query $\mathbf{q}_i$ to a weighted combination of the encoder's values $\mathbf{v}_j$, where the weights are determined by the similarity between the query and the encoder's keys:

$$\mathbf{z}_i = \sum_{j=1}^{n} \frac{\exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}}\right)}{\sum_{j'=1}^{n} \exp\left(\frac{\mathbf{q}_i \mathbf{k}_{j'}^\top}{\sqrt{d_k}}\right)} \mathbf{v}_j.$$

The similarity measure, typically the dot product, captures the alignment between the current decoder state and the relevant parts of the encoder's output.

The attention mechanism can be viewed as a content-based addressing system, where the query $\mathbf{q}_i$ retrieves the most relevant content from the encoder's output. This ensures that the decoder's focus is dynamically adjusted based on the evolving context during sequence generation.

Example: In a summarization task, the attention mechanism allows the decoder to focus on the key points in the input text, ensuring that the summary captures the most important information.

**Impact of Layer Depth on Interactions**

The depth of the encoder and decoder layers has a significant impact on the quality and richness of the interactions between the encoder and decoder. Deeper layers allow the model to capture more complex patterns and hierarchical structures in the data.

Let $L_e$ and $L_d$ denote the number of layers in the encoder and decoder, respectively. The output of the $l$-th layer in the encoder is given by

$$\mathbf{H}^{(l)} = \text{LayerNorm}\left(\mathbf{H}^{(l-1)} + \text{FFN}\left(\text{MultiHeadAttention}\left(\mathbf{H}^{(l-1)}\right)\right)\right),$$

where $\mathbf{H}^{(0)} = \mathbf{X}$ is the input sequence and FFN denotes the feedforward network. The depth $L_e$ determines the complexity of the representations learned by the encoder.

In the decoder, the output at layer $l$ is similarly defined as

$$\mathbf{Y}^{(l)} = \text{LayerNorm}\left(\mathbf{Y}^{(l-1)} + \text{FFN}\left(\text{MultiHeadAttention}\left(\mathbf{Y}^{(l-1)}, \mathbf{H}^{(L_e)}\right)\right)\right),$$

where the cross-attention mechanism allows the decoder to incorporate the encoded information from the deepest layer of the encoder.

The expressive power of the transformer model increases with the depth of the encoder and decoder layers. Formally, for sufficiently large $L_e$ and $L_d$, the model can approximate any continuous sequence-to-sequence mapping with arbitrary precision, assuming sufficient model capacity (number of attention heads, dimensionality, etc.).

Example: In a complex translation task, a deeper encoder may capture nuanced linguistic features from the source text, which are then leveraged by a deep decoder to generate more accurate and fluent translations.

## *4.3.2 Training Dynamics*

Understanding the training dynamics of transformers requires a deep dive into how loss functions are formulated, how optimization algorithms are applied to minimize these losses, and how regularization techniques are used to prevent overfitting. These components are essential to achieving stable training and convergence in large-scale models.

**Loss Functions and Optimization**

The loss function quantifies the difference between the model's predictions and the true labels, guiding the optimization process during training. In sequence-to-sequence tasks, cross-entropy loss is commonly used, particularly when the goal is to predict a probability distribution over a discrete set of classes, such as words in a vocabulary.

Let $\mathbf{y} = (y_1, y_2, \ldots, y_m)$ be the true sequence of labels, and $\hat{\mathbf{y}} = (\hat{y}_1, \hat{y}_2, \ldots, \hat{y}_m)$ be the predicted sequence generated by the model. The cross-entropy loss $\mathcal{L}$ is defined as

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{i=1}^{m} \log P(y_i \mid \hat{\mathbf{y}}_{<i}),$$

where $P(y_i \mid \hat{\mathbf{y}}_{<i})$ is the probability assigned by the model to the correct token $y_i$ given the previous predicted tokens $\hat{\mathbf{y}}_{<i}$.

The minimization of the cross-entropy loss is equivalent to maximizing the likelihood of the correct sequence under the model's distribution. Formally, minimizing $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$ corresponds to maximizing:

$$\prod_{i=1}^{m} P(y_i \mid \hat{\mathbf{y}}_{<i}).$$

This ensures that the model learns to assign high probabilities to the correct sequence during training.

Example: In machine translation, the cross-entropy loss measures how well the predicted translation aligns with the true translation, guiding the model to improve its predictions over successive iterations.

## Cross-Entropy Loss

The cross-entropy loss is particularly well suited for tasks where the output is a probability distribution over a discrete set of classes, such as in language modeling or translation.

Given a predicted probability distribution $\hat{p}_i = P(\hat{y}_i \mid \hat{\mathbf{y}}_{<i})$ and the true one-hot encoded distribution $p_i$, the cross-entropy loss for a single token is

$$\mathcal{L}_i = -\sum_{w=1}^{V} p_i(w) \log \hat{p}_i(w).$$

For a sequence, the total cross-entropy loss is the sum over all tokens:

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{i=1}^{m} \sum_{w=1}^{V} p_i(w) \log \hat{p}_i(w),$$

where $V$ is the size of the vocabulary.

The cross-entropy loss function is convex with respect to the predicted probabilities $\hat{p}_i$. This convexity ensures that optimization algorithms like gradient descent can efficiently find a global minimum in the loss landscape.

Example: In a transformer-based language model, the cross-entropy loss is minimized by adjusting the model parameters so that the predicted probability distributions over the vocabulary closely match the true distributions, leading to accurate text generation.

**Teacher Forcing in Training**

Teacher forcing is a technique used during training where the true previous token is used as input to the decoder, rather than the model's own prediction. This helps the model converge faster by providing the correct context at each step of the sequence generation.

During training with teacher forcing, the input to the decoder at each time step $i$ is the true token $y_{i-1}$ from the training data, rather than the predicted token $\hat{y}_{i-1}$:

$$\hat{y}_i = f(y_{i-1}, \mathbf{h}_i),$$

where $\mathbf{h}_i$ is the hidden state at time step $i$.

Teacher forcing reduces the exposure bias that can occur when the model is only trained on its own predictions. By providing the correct context, teacher forcing ensures that the model learns to predict the next token given the true previous tokens, leading to faster convergence during training.

Example: In a sequence generation task, such as summarization, teacher forcing helps the model learn the correct structure of summaries by ensuring that each part of the output sequence is conditioned on the correct input.

**Gradient Descent and Backpropagation**

Gradient descent is the primary optimization method used to minimize the loss function in neural networks, including transformers. Backpropagation is the algorithm used to compute the gradients of the loss function with respect to the model parameters.

The gradient of the loss function $\mathcal{L}$ with respect to a parameter $\theta$ is given by

$$\nabla_\theta \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \theta}.$$

The parameter update in gradient descent is then

$$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L},$$

where $\eta$ is the learning rate.

Under certain conditions, such as a sufficiently small learning rate and a convex loss function, gradient descent is guaranteed to converge to a local or global minimum of the loss function.

In a transformer model, gradient descent is used to update the parameters of the attention mechanisms, feedforward networks, and embeddings, gradually reducing the cross-entropy loss and improving the model's performance on the training data.

**Training Stability and Convergence**

Training stability refers to the ability of the model to maintain consistent learning progress without encountering issues such as exploding or vanishing gradients. Convergence refers to the model's ability to reach a minimum in the loss function.

Training stability can be affected by the choice of learning rate, initialization of model parameters, and the architecture of the model. For instance, residual connections and layer normalization in transformers help stabilize training by ensuring that gradients do not vanish or explode.

Gradient clipping is a technique used to maintain training stability by limiting the magnitude of the gradients:

$$\nabla_\theta \mathcal{L} \leftarrow \frac{\nabla_\theta \mathcal{L}}{\max\left(1, \frac{\|\nabla_\theta \mathcal{L}\|}{\text{threshold}}\right)}.$$

This ensures that the updates to the parameters remain within a reasonable range, preventing instability.

In training deep transformer models, gradient clipping is often employed to prevent large updates that could destabilize the training process, especially in early stages when the model parameters are far from their optimal values.

**Regularization Techniques**

Regularization techniques are used to prevent overfitting, ensuring that the model generalizes well to unseen data. In transformers, common regularization methods include dropout and label smoothing.

**Dropout**

Dropout ([2]) is a regularization technique where, during training, a random subset of activations is set to zero, preventing the model from becoming too reliant on any single neuron. Mathematically, for a layer with activations **h**, dropout is applied as

$$\mathbf{h}_{\text{dropout}} = \mathbf{h} \odot \mathbf{r},$$

where **r** is a binary mask vector with each element drawn from a Bernoulli distribution with parameter $p$, and $\odot$ denotes element-wise multiplication.

Dropout acts as a form of ensemble learning, where multiple sub-networks are trained simultaneously. This prevents overfitting and improves the generalization of the model.

In transformer models, dropout is often applied after the attention and feedforward layers to prevent overfitting, particularly in large models trained on limited data.

**Label Smoothing**

Label smoothing ([3]) is a technique that softens the target labels during training, replacing the one-hot encoded labels with a distribution that assigns a small probability to all classes. For a target label $y_i$, label smoothing modifies the target distribution as

$$p_i'(w) = (1 - \epsilon) \cdot p_i(w) + \frac{\epsilon}{V},$$

where $\epsilon$ is the smoothing parameter and $V$ is the vocabulary size.

Label smoothing reduces the confidence of the model in its predictions, which acts as a regularizer and prevents the model from becoming overly confident on the training data. This can improve the generalization performance of the model on unseen data.

In translation tasks, label smoothing helps prevent the model from becoming too certain about specific translations, which can lead to better performance on rare words and out-of-vocabulary words.

### *4.3.3   Hyperparameter Tuning*

Hyperparameters are parameters set before training begins, and they control various aspects of the learning process. Unlike model parameters, which are learned during training, hyperparameters must be carefully selected through experimentation, as they can greatly influence the effectiveness of the training process and the resulting model's performance.

**Impact of Hyperparameters on Performance**

The choice of hyperparameters such as learning rate, batch size, and the number of training epochs can drastically affect the model's ability to converge to a good solution and generalize well to unseen data. Understanding the mathematical relationship between these hyperparameters and the training process is essential for effective model tuning.

Let $\mathcal{L}(\theta; \lambda)$ represent the loss function for a model with parameters $\theta$ and hyperparameters $\lambda$. The goal of hyperparameter tuning is to find the optimal set of hyperparameters $\lambda^*$ that minimizes the expected loss on a validation set $\mathcal{D}_{val}$:

$$\lambda^* = \operatorname{argmin}_\lambda \mathbb{E}_{\mathcal{D}_{val}}[\mathcal{L}(\theta; \lambda)].$$

This optimization problem is often solved through methods such as grid search, random search, or more sophisticated techniques like Bayesian optimization.

The sensitivity of the loss function to changes in hyperparameters can be quantified by the gradient of the loss with respect to the hyperparameters:

$$\frac{\partial \mathcal{L}(\theta; \lambda)}{\partial \lambda}.$$

Large gradients indicate high sensitivity, meaning small changes in the hyperparameter values can lead to significant changes in model performance.

In a transformer model, a small learning rate might lead to slow convergence, while a large learning rate might cause the model to overshoot minima in the loss landscape, leading to instability. The optimal learning rate balances these effects, ensuring stable and efficient convergence.

### Learning Rate Schedules

The learning rate is one of the most crucial hyperparameters, controlling the step size during gradient descent. Learning rate schedules adjust the learning rate during training to improve convergence and prevent the model from getting stuck in suboptimal minima.

A learning rate schedule $\eta(t)$ is a function that defines how the learning rate $\eta$ changes over time $t$ (or training epochs). Common learning rate schedules include

1. Step Decay: The learning rate is reduced by a factor $\gamma$ after every $s$ epochs:

$$\eta(t) = \eta_0 \cdot \gamma^{\lfloor t/s \rfloor},$$

where $\eta_0$ is the initial learning rate.

2. Exponential Decay: The learning rate decays exponentially over time:

$$\eta(t) = \eta_0 \cdot \exp(-\lambda t),$$

where $\lambda$ is the decay rate.

3. Cosine Annealing: The learning rate follows a cosine function, reducing gradually to zero:

$$\eta(t) = \eta_0 \cdot \frac{1}{2}\left(1 + \cos\left(\frac{\pi t}{T}\right)\right),$$

where $T$ is the total number of training epochs.

Under certain conditions, such as a sufficiently small initial learning rate and appropriate decay, learning rate schedules can ensure that gradient descent converges to a local or global minimum of the loss function. The convergence is often faster and more stable with well-chosen learning rate schedules.

In transformer models, a common practice is to use a learning rate warm-up followed by a decay schedule. This approach starts with a very low learning rate, gradually increasing it over the first few epochs (warm-up), and then decays the learning rate as training progresses, allowing the model to settle into a good minimum.

## Batch Size and Training Duration

The batch size determines the number of training examples used to compute the gradient in each iteration. It influences both the stability of the training process and the time required to complete an epoch. Training duration, often measured in terms of the number of epochs, dictates how long the model is trained and can affect its ability to generalize.

Let $B$ be the batch size, $N$ be the total number of training examples, and $E$ be the number of epochs. The number of iterations per epoch is $I = \lceil N/B \rceil$, and the total number of updates to the model parameters during training is $U = I \times E$.

The variance of the gradient estimate decreases with increasing batch size. Specifically, if $\nabla \mathcal{L}(\theta; \mathcal{B})$ is the gradient computed on a batch $\mathcal{B}$, then

$$\text{Var}(\nabla \mathcal{L}(\theta; \mathcal{B})) \propto \frac{1}{B}.$$

This means that larger batch sizes lead to more stable gradient estimates, but they also require more computational resources.

In transformer training, larger batch sizes can help stabilize training, especially when using a high learning rate. However, very large batch sizes may require adjustments to the learning rate to avoid diminishing returns on training speed and model performance.

Training duration $E$ affects the model's ability to generalize. Too few epochs may result in underfitting, where the model fails to capture the underlying patterns in the data. Too many epochs may lead to overfitting, where the model becomes too specialized to the training data and performs poorly on unseen data. The optimal training duration balances these effects, typically determined through validation performance.

During transformer training, monitoring validation loss and other metrics helps determine when to stop training to avoid overfitting. Early stopping is a common technique that halts training once the validation performance ceases to improve.

## 4.4   Advanced Topics

As the original transformer architecture has become a foundational model in deep learning, various enhancements and variants have been proposed to address its limitations and extend its capabilities. These transformer variants, such as Transformer-XL, Memory-Enhanced Transformers, and Universal Transformers, introduce new mechanisms that enable better handling of long-term dependencies, memory retention, and dynamic computation. This section explores these topics, providing a mathematical foundation for each variant and discussing their implications for model performance and flexibility.

### *4.4.1   Transformer Variants*

Transformer variants build on the original architecture by introducing novel components or modifications aimed at addressing specific challenges, such as the fixed-length context window, computational efficiency, and the ability to handle more complex tasks. Each variant introduces new ideas that push the boundaries of what transformers can achieve.

**Transformer-XL**

Transformer-XL ([1]) is designed to overcome the fixed-length context limitation of the original transformer by introducing a mechanism to capture long-range dependencies. It does so by segmenting the input sequence and introducing a recurrence mechanism that enables the model to retain information across segments.

In Transformer-XL, the input sequence is divided into segments $\{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T\}$, where each segment $\mathbf{x}_t$ is processed by the transformer model. The key innovation is the introduction of a hidden state recurrence mechanism that allows information to flow from one segment to the next:

$$\mathbf{h}_t^l = \text{Attention}\left(\mathbf{q}_t^l, [\mathbf{k}_t^l; \mathbf{k}_{t-1}^l], [\mathbf{v}_t^l; \mathbf{v}_{t-1}^l]\right),$$

where $\mathbf{h}_t^l$ is the hidden state at layer $l$ for segment $t$, and $[\mathbf{k}_t^l; \mathbf{k}_{t-1}^l]$ denotes the concatenation of keys from the current and previous segments. This recurrence mechanism enables the model to maintain a memory of past segments, effectively increasing the context length without requiring larger computational resources.

Transformer-XL's recurrence mechanism allows the model to capture dependencies beyond the typical context window. Formally, let $\mathcal{L}_{\text{XL}}(\mathbf{h}_t^l, \mathbf{h}_{t-1}^l)$ denote the loss function that includes the recurrent hidden states. The memory retention property ensures that

$$\frac{\partial \mathcal{L}_{\text{XL}}}{\partial \mathbf{h}_t^l} \approx \frac{\partial \mathcal{L}_{\text{XL}}}{\partial \mathbf{h}_{t-1}^l},$$

indicating that the gradient information is preserved across segments, allowing the model to effectively propagate information over long sequences.

In language modeling, Transformer-XL can model dependencies over longer text sequences, such as paragraphs or entire documents, without being constrained by the fixed context window of traditional transformers. This results in more coherent text generation and better understanding of long-range dependencies.

**Memory-Enhanced Transformers**

Memory-enhanced transformers extend the transformer architecture by incorporating external memory modules, allowing the model to store and retrieve information across long sequences or even entire datasets. These models are particularly useful in tasks requiring long-term memory, such as question answering or summarization.

Memory-enhanced transformers introduce an external memory matrix $\mathbf{M}$ that is dynamically updated during the learning process. The memory matrix $\mathbf{M} \in \mathbb{R}^{N \times d}$ can store $N$ memory slots, each of dimension $d$. The attention mechanism is modified to incorporate this external memory:

$$\mathbf{h}_t^l = \text{Attention}\left(\mathbf{q}_t^l, [\mathbf{k}_t^l; \mathbf{M}], [\mathbf{v}_t^l; \mathbf{M}]\right),$$

where $\mathbf{M}$ serves as an additional source of keys and values, enabling the model to retrieve relevant information from memory.

The inclusion of external memory enhances the model's ability to store and retrieve information across different contexts. Let $\mathcal{L}_{\text{MEM}}(\mathbf{h}_t^l, \mathbf{M})$ represent the loss function with the memory module. The memory augmentation ensures that

$$\frac{\partial \mathcal{L}_{\text{MEM}}}{\partial \mathbf{M}} \neq 0,$$

indicating that the memory matrix actively contributes to the learning process by influencing the gradients, thereby allowing the model to refine its memory over time.

Example: In tasks like open-domain question answering, memory-enhanced transformers can store relevant facts and retrieve them when needed, leading to more accurate and contextually aware answers. The external memory allows the model to "remember" information from previous questions, improving performance over time.

**Universal Transformers**

Universal transformers introduce a dynamic computational mechanism that iteratively refines representations at each position in the sequence. This is achieved by applying the same set of transformer layers multiple times, akin to a recurrent neural network, allowing the model to adaptively learn complex dependencies.

In a universal transformer, the input sequence is processed by a stack of transformer layers that are repeatedly applied to refine the hidden states. Let $\mathbf{h}_t^{(0)} = \mathbf{E}_t$ be the initial embedding of the token at position $t$. The hidden state is updated iteratively:

$$\mathbf{h}_t^{(k+1)} = \text{TransformerLayer}\left(\mathbf{h}_t^{(k)}, \mathbf{H}^{(k)}\right),$$

where $k$ denotes the iteration index and $\mathbf{H}^{(k)}$ represents the set of hidden states for all positions at iteration $k$. The iterative process continues for a fixed number of iterations $K$ or until convergence.

The iterative refinement process in universal transformers allows the model to approximate any continuous sequence transformation with arbitrary precision, given sufficient iterations. Formally, let $\mathcal{T}$ be the set of sequence transformations expressible by a universal transformer with $K$ iterations. Then

$$\mathcal{T}(K) \subseteq \mathcal{T}(K + 1),$$

indicating that increasing the number of iterations $K$ enhances the model's expressive power.

In tasks requiring complex reasoning, such as multi-hop question answering or logic-based tasks, universal transformers can iteratively refine their understanding of the input, leading to more accurate and nuanced predictions. The iterative nature allows the model to process information in a manner that mimics human-like reasoning.

As transformer models have grown in size and complexity, addressing issues of scalability and efficiency has become crucial. This section delves into the mathematical underpinnings of various techniques designed to enhance the scalability and efficiency of transformers, including efficient attention mechanisms, model parallelism, and strategies for training large-scale transformers. These approaches are essential for deploying transformers in large-scale applications while maintaining computational feasibility.

### *4.4.2   Scalability and Efficiency Improvements*

The scalability of transformers is fundamentally linked to the efficiency of their core components, particularly the attention mechanism, and the ability to parallelize computations across large models. Efficient attention mechanisms reduce the computational and memory overhead, while model parallelism and optimized training strategies enable the handling of large-scale models that are essential for state-of-the-art performance.

**Efficient Attention Mechanisms**

The attention mechanism in transformers, while powerful, is computationally expensive, particularly in terms of memory and time complexity. The standard attention mechanism has a quadratic complexity in relation to the sequence length, $O(n^2)$, where $n$ is the sequence length. Efficient attention mechanisms have been proposed to address this issue, reducing the complexity and making transformers more scalable.

The standard self-attention mechanism computes attention scores between all pairs of tokens in the sequence. For an input sequence $\mathbf{X} \in \mathbb{R}^{n \times d}$, the attention mechanism is computed as

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V},$$

where $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$, $\mathbf{K} = \mathbf{X}\mathbf{W}_K$, and $\mathbf{V} = \mathbf{X}\mathbf{W}_V$ are the query, key, and value matrices, respectively.

The computational complexity of the standard attention mechanism is $O(n^2 d)$, where $n$ is the sequence length and $d$ is the dimensionality of the model. This quadratic complexity poses challenges for processing long sequences.

Several efficient attention mechanisms have been proposed to reduce this complexity:

1. Sparse attention restricts the attention computation to a subset of tokens, reducing the complexity to $O(n \cdot k \cdot d)$, where $k$ is the number of tokens attended to by each token. This sparsity can be achieved through fixed patterns or learned patterns.

Let $\mathcal{S}_i \subseteq \{1, 2, \ldots, n\}$ represent the indices of tokens that token $i$ attends to. The sparse attention is computed as

$$\text{SparseAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})_i = \sum_{j \in \mathcal{S}_i} \alpha_{ij} \mathbf{v}_j,$$

where $\alpha_{ij}$ are the attention weights restricted to the sparse subset.

2. Low-rank factorization approximates the attention matrix by decomposing it into the product of lower-rank matrices, reducing the complexity to $O(n \cdot r \cdot d)$, where $r$ is the rank of the approximation.

The attention matrix $\mathbf{A}$ can be factorized as

$$\mathbf{A} \approx \mathbf{BC},$$

where $\mathbf{B} \in \mathbb{R}^{n \times r}$ and $\mathbf{C} \in \mathbb{R}^{r \times n}$. The low-rank approximation reduces the computation required for the attention operation.

3. Linearized attention mechanisms approximate the softmax function with a linear function, reducing the complexity to $O(n \cdot d)$.

The linearized attention can be expressed as

$$\text{LinearAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \phi(\mathbf{Q}) \left( \phi(\mathbf{K})^\top \mathbf{V} \right),$$

where $\phi$ is a feature map that approximates the softmax function.

Linearized attention mechanisms reduce the complexity of the attention operation to linear time $O(n \cdot d)$, making them suitable for very long sequences.

In large-scale NLP tasks, efficient attention mechanisms allow transformers to process entire documents or long texts without running into memory or computational constraints, enabling more effective and scalable models.

**Model Parallelism**

Model parallelism is a strategy used to distribute the computation of large models across multiple devices or processors. This approach is essential for training very large transformer models, where the model parameters and activations may exceed the memory capacity of a single device.

Let $\mathbf{W} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}}$ be a large weight matrix in the model. In model parallelism, this matrix is partitioned across multiple devices:

$$\mathbf{W} = [\mathbf{W}_1, \mathbf{W}_2, \ldots, \mathbf{W}_p],$$

where $p$ is the number of devices and $\mathbf{W}_i$ is the portion of the weight matrix on device $i$. The input $\mathbf{x}$ is similarly partitioned, and the output is computed as

$$\mathbf{y} = \sum_{i=1}^{p} \mathbf{x}_i \mathbf{W}_i.$$

Model parallelism can lead to significant speedup in training, particularly when the model size is much larger than the capacity of a single device. The overall speedup $S$ is given by

$$S \approx \frac{1}{\max_i \left( \text{ComputationTime}_i + \text{CommunicationTime}_i \right)},$$

where ComputationTime$_i$ and CommunicationTime$_i$ are the computation and communication times on device $i$.

In training models like GPT-3, which have hundreds of billions of parameters, model parallelism is essential to distribute the computations across thousands of GPUs, making the training process feasible.

**Training Large-Scale Transformers**

Training large-scale transformers requires not only efficient algorithms and hardware but also careful management of hyperparameters, data pipelines, and optimization strategies. Techniques like gradient accumulation, mixed-precision training, and distributed data parallelism are commonly used.

Consider a large-scale transformer with $N$ layers, each with $d$ dimensionality, trained on a dataset of size $M$. The training process involves minimizing a loss function $\mathcal{L}(\theta)$ with respect to the model parameters $\theta$:

$$\theta^* = \text{argmin}_\theta \mathcal{L}(\theta).$$

Due to the large size of the model, the following techniques are often employed:

1. Gradients are accumulated over multiple mini-batches before performing a parameter update, reducing memory usage and enabling training with smaller batch sizes. Let $\mathcal{B}_i$ be the $i$-th mini-batch and $g_i = \nabla_\theta \mathcal{L}(\mathcal{B}_i)$ be the gradient. The accumulated gradient over $k$ mini-batches is

$$g_{\text{acc}} = \sum_{i=1}^{k} g_i.$$

2. Mixed-precision training involves using lower precision arithmetic (e.g., FP16) for computations while maintaining high precision (e.g., FP32) for certain critical operations, reducing memory usage and increasing computational efficiency. Mixed-precision training reduces memory bandwidth requirements and can lead to faster training times, with the theoretical speedup being proportional to the reduction in memory transfer time:

$$S_{\text{mixed}} = \frac{\text{MemoryTransferTime}_{\text{FP32}}}{\text{MemoryTransferTime}_{\text{FP16}}},$$

while maintaining comparable model accuracy.

3. Data parallelism distributes the training data across multiple devices, with each device independently computing gradients on its subset of data. The gradients are then averaged across devices. Let $\mathcal{D}_i$ be the subset of data on device $i$. The gradient on device $i$ is $g_i = \nabla_\theta \mathcal{L}(\mathcal{D}_i)$. The global gradient is computed as

$$g_{\text{global}} = \frac{1}{p} \sum_{i=1}^{p} g_i,$$

where $p$ is the number of devices.

When training models like BERT or GPT on large datasets, mixed-precision training and distributed data parallelism are often used in combination to maximize efficiency and scalability, enabling the training of models with billions of parameters in a reasonable timeframe.

# References

1. Dai, Z., Yang, Z., Yang, Y., Carbonell, J.G., Le, Q.V., Salakhutdinov, R.: Transformer-xl: Attentive language models beyond a fixed-length context. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 2978–2988 (2019)
2. Nitish, S., Geoffrey, H., Alex, K., Ilya, S., Ruslan, S.: Dropout: A simple way to prevent neural networks from overfitting. J. Mach. Learn. Res. **15**(1), 1929–1958 (2014)
3. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2818–2826 (2016)

# Chapter 5
# Transformers in Natural Language Processing

## 5.1 Language Modeling

Language modeling is one of the foundational tasks in NLP and is crucial for various downstream tasks such as machine translation, summarization, and dialog systems. Transformers, with their self-attention mechanisms and ability to model complex dependencies, have become the standard architecture for state-of-the-art language models. This section explores the mathematical foundations of language modeling using transformers, focusing on sequence-to-sequence models, the encoder–decoder framework, conditional probability and likelihood, and pre-training objectives.

### 5.1.1 Mathematical Formulation

Language modeling involves predicting the probability distribution of the next word in a sequence, given the previous words. This can be formalized mathematically using the principles of probability and information theory, combined with the structural components of transformers.

**Sequence-to-Sequence Models**

Sequence-to-sequence (Seq2Seq) models are a class of models designed to transform one sequence into another, making them ideal for tasks like translation, where the input is a sequence in one language, and the output is a sequence in another language. In transformers, the Seq2Seq model is typically implemented using an encoder–decoder architecture.

Let $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ be the input sequence and $\mathbf{y} = (y_1, y_2, \ldots, y_m)$ be the output sequence. The goal of a Seq2Seq model is to learn a mapping from $\mathbf{x}$ to $\mathbf{y}$ by maximizing the conditional probability:

$$P(\mathbf{y} \mid \mathbf{x}; \theta) = \prod_{i=1}^{m} P(y_i \mid y_1, \ldots, y_{i-1}, \mathbf{x}; \theta),$$

where $\theta$ represents the parameters of the model.

The training of a Seq2Seq model can be viewed as a maximum likelihood estimation (MLE) problem, where the objective is to maximize the likelihood of the observed sequence pairs $(\mathbf{x}, \mathbf{y})$:

$$\theta^* = \arg\max_{\theta} \sum_{(\mathbf{x},\mathbf{y})\in\mathcal{D}} \log P(\mathbf{y} \mid \mathbf{x}; \theta),$$

where $\mathcal{D}$ is the training dataset.

In machine translation, a Seq2Seq transformer model is trained to maximize the likelihood of the correct translation $\mathbf{y}$ given the source sentence $\mathbf{x}$, effectively learning to translate between languages.

### Encoder–Decoder Framework

The encoder–decoder framework is a fundamental structure in Seq2Seq models, particularly in transformers. The encoder processes the input sequence and generates a set of contextualized representations, which the decoder then uses to generate the output sequence (see Fig. 5.1).

Let $\mathbf{H} = (\mathbf{h}_1, \mathbf{h}_2, \ldots, \mathbf{h}_n)$ represent the output of the encoder, where $\mathbf{h}_i$ is the hidden state corresponding to the input token $x_i$. The decoder generates the output sequence $\mathbf{y}$ by computing

$$\mathbf{y}_i = \text{Decoder}(\mathbf{y}_{<i}, \mathbf{H}),$$

where $\mathbf{y}_{<i} = (y_1, y_2, \ldots, y_{i-1})$ represents the previously generated tokens.

The encoder–decoder framework ensures that the information from the entire input sequence $\mathbf{x}$ is preserved and utilized in generating each token of the output sequence $\mathbf{y}$. Formally, the hidden states $\mathbf{H}$ should capture sufficient information about $\mathbf{x}$ such that

$$\text{Mutual Information}(H_i; \mathbf{x}) \geq \text{Mutual Information}(y_i; \mathbf{x}).$$

Example: In a translation task, the encoder processes the entire source sentence to produce contextual embeddings, which the decoder uses to generate the target sentence, ensuring that the output is contextually aligned with the input.

**Fig. 5.1**  Diagram of the transformer encoder–decoder architecture, illustrating the flow of information from the encoder's context-aware representations to the decoder's sequence generation process through cross-attention mechanisms

**Conditional Probability and Likelihood**

In the context of language modeling, conditional probability plays a crucial role in defining how likely a particular sequence of words is, given the preceding words. The likelihood of a sequence is computed as the product of these conditional probabilities.

Given a sequence $\mathbf{y} = (y_1, y_2, \ldots, y_m)$, the conditional probability of the sequence, given the previous tokens, is defined as

$$P(\mathbf{y} \mid \mathbf{x}; \theta) = \prod_{i=1}^{m} P(y_i \mid y_1, \ldots, y_{i-1}, \mathbf{x}; \theta).$$

The log-likelihood $\mathcal{L}(\theta)$ of the sequence under the model is

$$\mathcal{L}(\theta) = \sum_{i=1}^{m} \log P(y_i \mid y_1, \ldots, y_{i-1}, \mathbf{x}; \theta).$$

The chain rule of probability allows the decomposition of the joint probability of a sequence into a product of conditional probabilities:

$$P(\mathbf{y} \mid \mathbf{x}) = P(y_1 \mid \mathbf{x}) \cdot P(y_2 \mid y_1, \mathbf{x}) \cdot \ldots \cdot P(y_m \mid y_1, y_2, \ldots, y_{m-1}, \mathbf{x}),$$

which is the foundation for autoregressive models like transformers.

Example: In language modeling, the transformer uses the chain rule to predict the next word in a sequence, given the previous words. This approach enables the model to generate coherent text by leveraging the conditional dependencies between words.

**Pre-training Objectives**

Pre-training objectives are critical in modern language models, as they enable the model to learn rich representations from large amounts of text data before being fine-tuned on specific tasks. Common pre-training objectives include masked language modeling and autoregressive language modeling.

1. Masked Language Modeling (MLM): In MLM, a subset of tokens in the input sequence is masked, and the model is trained to predict these masked tokens. Let $\mathbf{y} = (y_1, y_2, \ldots, y_m)$ be the sequence and let $\mathbf{y}_{\text{mask}}$ be the masked sequence. The objective is to maximize

$$\mathcal{L}_{\text{MLM}}(\theta) = \mathbb{E}_{\mathbf{y}, \mathbf{y}_{\text{mask}}} \left[ \sum_{i \in \text{mask}} \log P(y_i \mid \mathbf{y}_{\text{mask}}; \theta) \right],$$

where $i$ indexes the masked positions.

2. Autoregressive Language Modeling (ALM): In ALM, the model is trained to predict the next token in the sequence, given the previous tokens. The objective is to maximize

$$\mathcal{L}_{\text{ALM}}(\theta) = \mathbb{E}_{\mathbf{y}} \left[ \sum_{i=1}^{m} \log P(y_i \mid y_1, \ldots, y_{i-1}; \theta) \right].$$

Pre-training with objectives like MLM and ALM enables the model to learn generalizable features that can be transferred to a wide range of downstream tasks. Formally, the pre-training phase aims to learn a representation $\mathbf{H}$ such that

$$\mathbf{H} = \arg\max_{\mathbf{H}} \mathcal{L}_{\text{pretrain}}(\theta),$$

where $\mathcal{L}_{\text{pretrain}}$ is the loss function for the pre-training objective. The learned representations $\mathbf{H}$ are then fine-tuned on task-specific data.

Models like BERT ([2]) and GPT ([3]) use MLM and ALM, respectively, as their pre-training objectives. BERT, trained with MLM, excels at understanding the context within a sentence, while GPT, trained with ALM, is particularly effective at generating coherent text.

**Masked Language Modeling**

MLM is a pre-training objective where a subset of tokens in a sequence is randomly masked, and the model is tasked with predicting these masked tokens based on the context provided by the unmasked tokens. Let $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ be the input sequence and let $\mathbf{m} = (m_1, m_2, \ldots, m_n)$ be a binary mask, where $m_i = 1$ indicates that the token $x_i$ is masked.

The objective of MLM is to maximize the log-likelihood of the masked tokens given their surrounding context:

$$\mathcal{L}_{\text{MLM}}(\theta) = \mathbb{E}_{\mathbf{x},\mathbf{m}} \left[ \sum_{i=1}^{n} m_i \log P(x_i \mid \mathbf{x}_{\backslash i}; \theta) \right],$$

where $\mathbf{x}_{\backslash i}$ denotes the sequence $\mathbf{x}$ with the $i$-th token removed and $\theta$ represents the model parameters.

MLM forces the model to learn rich contextual representations by making predictions based on incomplete information. Formally, the training process encourages the model to maximize the mutual information between the masked token $x_i$ and the surrounding context $\mathbf{x}_{\backslash i}$:

$$\text{Mutual Information}(x_i; \mathbf{x}_{\backslash i}) \geq \text{Mutual Information}(x_i; \mathbf{x}).$$

This ensures that the model captures dependencies that go beyond simple adjacent tokens, learning deeper contextual relationships.

Example: In BERT, MLM is used as a core pre-training objective. Tokens are randomly masked, and the model learns to predict these tokens based on the context provided by the other tokens in the sentence, enabling the model to develop a deep understanding of language structure.

**Next Sentence Prediction (NSP)**

NSP is a pre-training task designed to help the model understand the relationship between pairs of sentences. Given two sentences $\mathbf{s}_1$ and $\mathbf{s}_2$, the model is tasked with predicting whether $\mathbf{s}_2$ logically follows $\mathbf{s}_1$.

The objective is to maximize the log-likelihood of correctly classifying sentence pairs:

$$\mathcal{L}_{\text{NSP}}(\theta) = \mathbb{E}_{(\mathbf{s}_1, \mathbf{s}_2), y} \left[ y \log P(\text{True} \mid \mathbf{s}_1, \mathbf{s}_2; \theta) + (1 - y) \log P(\text{False} \mid \mathbf{s}_1, \mathbf{s}_2; \theta) \right],$$

where $y = 1$ if $\mathbf{s}_2$ follows $\mathbf{s}_1$, and $y = 0$ otherwise.

NSP encourages the model to learn semantic and syntactic dependencies between sentences, which are crucial for tasks like question answering and text summarization. Formally, NSP training maximizes the conditional dependency between consecutive sentences:

$$\text{Mutual Information}(\mathbf{s}_2; \mathbf{s}_1) \geq \text{Mutual Information}(\mathbf{s}_2; \mathbf{s}_1 \mid \mathbf{s}_3),$$

where $\mathbf{s}_3$ is an unrelated sentence.

In BERT, NSP is combined with MLM during pre-training. The model is presented with pairs of sentences and learns to predict whether the second sentence is a continuation of the first, enhancing its ability to understand text coherence and discourse.

**Causal Language Modeling (CLM)**

Causal language modeling, also known as autoregressive language modeling, is a pre-training objective where the model predicts each token in a sequence based on the preceding tokens. This setup aligns with the natural reading order of text, where each word is conditioned on all previous words.

The objective is to maximize the likelihood of the sequence under the model:

$$\mathcal{L}_{\text{CLM}}(\theta) = \mathbb{E}_{\mathbf{x}} \left[ \sum_{i=1}^{n} \log P(x_i \mid x_1, x_2, \ldots, x_{i-1}; \theta) \right].$$

CLM models the sequential dependency between tokens, with each token being conditioned on all preceding tokens. This autoregressive process can be expressed using the chain rule of probability:

$$P(\mathbf{x}; \theta) = \prod_{i=1}^{n} P(x_i \mid x_1, x_2, \ldots, x_{i-1}; \theta).$$

The model is trained to maximize this joint probability, leading to a natural understanding of language in a left-to-right manner.

Example GPT (Generative Pre-trained Transformer) models are trained using CLM. The model generates text by predicting the next word in a sequence, conditioned on all previous words, making it particularly effective at text generation tasks.

**Permutation Language Modeling (PLM)**

PLM is a pre-training objective introduced by the XLNet model. PLM involves predicting tokens in a sequence based on a random permutation of the sequence order, rather than strictly left to right or right to left. This approach enables the model to capture bidirectional context while maintaining the autoregressive property.

Let $\pi$ be a permutation of the indices $\{1, 2, \ldots, n\}$. The objective is to maximize the likelihood of the sequence under the permutation:

$$\mathcal{L}_{\text{PLM}}(\theta) = \mathbb{E}_{\mathbf{x}, \pi} \left[ \sum_{i=1}^{n} \log P(x_{\pi(i)} \mid x_{\pi(1)}, x_{\pi(2)}, \ldots, x_{\pi(i-1)}; \theta) \right].$$

PLM ensures that the model learns to represent sequences in a permutation-invariant manner, capturing dependencies that are not tied to a fixed sequential order. The model is trained to maximize the expected likelihood over all possible permutations:

$$\mathbb{E}_{\pi} \left[ P(\mathbf{x}; \theta) \right] = \sum_{\pi} P(\pi) \prod_{i=1}^{n} P(x_{\pi(i)} \mid x_{\pi(1)}, \ldots, x_{\pi(i-1)}; \theta),$$

where $P(\pi)$ is the probability of the permutation $\pi$.

Example: XLNet ([4]) uses PLM to combine the strengths of both autoregressive and bidirectional models. By predicting tokens based on all possible permutations of a sequence, XLNet captures richer context compared to traditional left-to-right or right-to-left models.

## 5.1.2  Applications in NLP

Transformers have revolutionized various NLP tasks, with applications ranging from
machine translation to text generation. This section explores the mathematical foun-
dations of these applications, focusing on how transformer models are formulated
for translation tasks, the evaluation metrics used to assess their performance, and the
principles underlying text generation.

### Machine Translation

Machine translation (MT) is one of the most prominent applications of transformer
models. The goal of MT is to translate a sentence or document from one language to
another while preserving the meaning and fluency of the original text. Transformers,
with their encoder–decoder architecture, have become the backbone of state-of-the-
art translation systems.

Let $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ be the input sentence in the source language, and
$\mathbf{y} = (y_1, y_2, \ldots, y_m)$ be the output sentence in the target language. The goal of
the translation model is to learn the conditional probability distribution $P(\mathbf{y} \mid \mathbf{x}; \theta)$,
where $\theta$ represents the parameters of the transformer model.

The translation task can be formulated as a sequence-to-sequence problem, where
the model is trained to maximize the log-likelihood of the target sentence given the
source sentence:

$$\mathcal{L}_{\text{MT}}(\theta) = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log P(\mathbf{y} \mid \mathbf{x}; \theta),$$

where $\mathcal{D}$ is the training dataset consisting of pairs of sentences in the source and
target languages.

The consistency of a Seq2Seq model in machine translation requires that the
encoder fully captures the semantic content of the source sentence, and the decoder
uses this information to generate a syntactically and semantically correct target
sentence. Formally, for a well-trained model:

$$\text{Mutual Information}(H_i; \mathbf{x}) \geq \text{Mutual Information}(y_i; \mathbf{x}),$$

where $H_i$ are the encoder hidden states and $y_i$ are the decoder outputs.

Example: In a transformer-based translation system like Google Translate, the
encoder processes the entire input sentence to generate a set of context-aware embed-
dings, which the decoder then uses to produce the translation, one word at a time,
conditioned on the previously generated words.

**Mathematical Formulation of Translation Models**

In machine translation, the transformer model operates within the encoder–decoder framework, where the encoder generates a representation of the source sentence, and the decoder generates the translation by attending to the encoder's representation.

The encoder takes the source sentence $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ and produces a sequence of hidden states $\mathbf{H} = (H_1, H_2, \ldots, H_n)$, where each hidden state $H_i$ is a contextualized representation of the word $x_i$.

The decoder then generates the target sentence $\mathbf{y} = (y_1, y_2, \ldots, y_m)$ by attending to the hidden states $\mathbf{H}$ and predicting each word $y_i$ sequentially:

$$P(y_i \mid y_1, \ldots, y_{i-1}, \mathbf{H}; \theta) = \text{softmax}\left(\mathbf{W}_o \mathbf{h}_i^{(L)}\right),$$

where $\mathbf{h}_i^{(L)}$ is the final hidden state in the decoder and $\mathbf{W}_o$ is the output projection matrix.

The attention mechanism in the transformer allows the decoder to focus on relevant parts of the source sentence when generating each word in the target sentence. Mathematically, the attention weights $\alpha_{ij}$ between the decoder's query $q_i$ and the encoder's key $k_j$ are given by

$$\alpha_{ij} = \frac{\exp(q_i^\top k_j)}{\sum_{j'=1}^{n} \exp(q_i^\top k_{j'})}.$$

The context vector $c_i$, which the decoder uses to generate the next word, is then a weighted sum of the encoder's hidden states:

$$c_i = \sum_{j=1}^{n} \alpha_{ij} H_j.$$

This mechanism ensures that the translation is contextually accurate and semantically consistent.

Example: In translating the sentence "The cat sat on the mat" to French, the attention mechanism ensures that "cat" is correctly translated to "chat" and "mat" to "tapis," maintaining the semantic integrity of the sentence.

**Evaluation Metrics (BLEU, METEOR)**

Evaluating the quality of machine translation models is essential for understanding their effectiveness. BLEU (Bilingual Evaluation Understudy) and METEOR (Metric for Evaluation of Translation with Explicit ORdering) are two commonly used metrics.

BLEU measures the overlap between the machine-generated translation and one or more reference translations. It is calculated using precision scores for n-grams up to a specified order $N$, typically 4:

$$\text{BLEU} = \text{BP} \cdot \exp\left(\sum_{n=1}^{N} w_n \log p_n\right),$$

where $p_n$ is the precision for n-grams of length $n$, $w_n$ are the weights (typically uniform), and BP is the brevity penalty to penalize translations that are too short.

BLEU is consistent with human judgment to some extent, particularly for high-scoring translations. However, it can be sensitive to exact word matches and may not fully capture the quality of translations that use different wording but convey the same meaning.

METEOR is designed to address some of BLEU's limitations by considering synonymy, stemming, and word order. It combines precision and recall with a penalty for fragmented matches:

$$\text{METEOR} = F_{\text{mean}} \cdot (1 - \text{Penalty}),$$

where $F_{\text{mean}}$ is a weighted harmonic mean of precision and recall, and Penalty accounts for the alignment fragmentation.

METEOR is more sensitive to semantic equivalence and word order, making it a better measure of translation quality in cases where the translation is accurate but uses different words or phrases than the reference.

Example: In evaluating translations from English to French, BLEU might give a high score if the translation is a close word-for-word match with the reference, while METEOR might score it higher if the translation accurately captures the meaning, even with different phrasing.

**Text Generation**

Text generation involves generating coherent and contextually relevant text sequences, often as a continuation of a given prompt. Transformers, particularly those trained with causal language modeling, are highly effective at generating natural language text.

Given a prompt $\mathbf{x} = (x_1, x_2, \ldots, x_n)$, the goal of text generation is to produce a continuation $\mathbf{y} = (y_1, y_2, \ldots, y_m)$ such that the entire sequence $\mathbf{z} = (\mathbf{x}, \mathbf{y})$ is coherent and contextually relevant. The likelihood of the generated sequence is maximized using

$$\mathcal{L}_{\text{gen}}(\theta) = \sum_{i=1}^{m} \log P(y_i \mid x_1, \ldots, x_n, y_1, \ldots, y_{i-1}; \theta).$$

Autoregressive models like transformers ensure coherence in generated text by conditioning each word on the entire preceding context. The chain rule of probability guarantees that the generated sequence $\mathbf{z}$ maintains internal consistency:

$$P(\mathbf{z}; \theta) = \prod_{i=1}^{n+m} P(z_i \mid z_1, \ldots, z_{i-1}; \theta).$$

Example: In a text generation task, a transformer model might be given the prompt "Once upon a time," and generate a coherent continuation like "there was a young princess who lived in a grand castle." The model's ability to generate such text depends on its understanding of narrative structure and language.

## Generative Models

Generative models in NLP aim to learn the underlying distribution of a language to generate new, coherent text that mimics human writing. Transformers, particularly those trained with autoregressive or masked language modeling objectives, are well suited for generative tasks such as text completion, story generation, and dialog systems.

Let $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ be an input sequence and $\mathbf{y} = (y_1, y_2, \ldots, y_m)$ be the generated output sequence. The generative model seeks to maximize the joint probability of the entire sequence $\mathbf{z} = (\mathbf{x}, \mathbf{y})$:

$$P(\mathbf{z}; \theta) = P(x_1, x_2, \ldots, x_n, y_1, y_2, \ldots, y_m; \theta).$$

This joint probability can be factorized using the chain rule of probability:

$$P(\mathbf{z}; \theta) = \prod_{i=1}^{n+m} P(z_i \mid z_1, z_2, \ldots, z_{i-1}; \theta),$$

where $\theta$ represents the parameters of the transformer model.

For a generative model to produce coherent text, it must effectively model the dependencies between tokens. The coherence of the generated text is ensured by maximizing the conditional probabilities in the factorization:

$$\max_{\theta} \prod_{i=1}^{n+m} P(z_i \mid z_1, z_2, \ldots, z_{i-1}; \theta).$$

The ability to capture these dependencies is critical to producing human-like text.

Example: In a story generation task, a transformer-based generative model might be given a prompt like "In a distant land," and generate a continuation such as "a brave

knight set out on a quest to find a hidden treasure." The model's success depends on its understanding of narrative structure and language conventions.

## Evaluation Metrics (Perplexity, ROUGE)

Evaluation metrics are crucial for assessing the quality of generative models and text summarization systems. Perplexity and ROUGE are commonly used metrics, each with specific applications and mathematical foundations.

Perplexity is a measure of how well a probabilistic model predicts a sample. It is commonly used to evaluate language models, where a lower perplexity indicates a better model. Formally, perplexity is defined as

$$\text{Perplexity}(\mathcal{D}; \theta) = \exp\left(-\frac{1}{N} \sum_{i=1}^{N} \log P(x_i \mid x_1, x_2, \ldots, x_{i-1}; \theta)\right),$$

where $\mathcal{D}$ is the dataset and $N$ is the number of tokens in the dataset.

Perplexity can be interpreted as the inverse of the geometric mean per-word likelihood. It measures how surprised the model is by the actual sequence. For an ideal model that perfectly predicts the sequence, perplexity would be 1.

Example: A language model with a perplexity of 10 is, on average, as uncertain about the next word as if it had to choose uniformly among 10 possible words. Lower perplexity values indicate better predictive performance.

ROUGE (Recall-Oriented Understudy for Gisting Evaluation): ROUGE is a set of metrics used to evaluate the quality of text summaries by comparing them to reference summaries. ROUGE-N measures the overlap of n-grams between the generated summary and the reference summary:

$$\text{ROUGE-N} = \frac{\sum_{\text{gram} \in \text{Ref}} \min(\text{Count}_{\text{gen}}(\text{gram}), \text{Count}_{\text{ref}}(\text{gram}))}{\sum_{\text{gram} \in \text{Ref}} \text{Count}_{\text{ref}}(\text{gram})},$$

where $\text{Count}_{\text{gen}}$ and $\text{Count}_{\text{ref}}$ are the counts of the n-gram in the generated and reference summaries, respectively.

ROUGE-N is primarily a recall-based metric, focusing on how much of the reference summary is captured by the generated summary. High ROUGE scores indicate that the generated summary contains a significant portion of the important content from the reference summary.

Example: In text summarization, ROUGE-1 (unigram overlap) and ROUGE-2 (bigram overlap) are commonly used to evaluate how well a model-generated summary captures the essence of a reference summary. High ROUGE scores generally correlate with high-quality summaries.

## Text Summarization

Text summarization is the process of condensing a long piece of text into a shorter version while preserving its key information. Transformers can perform both extractive summarization, which selects important sentences from the original text, and abstractive summarization, which generates new sentences that capture the essence of the original text.

Let $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ be the input document and $\mathbf{y} = (y_1, y_2, \ldots, y_m)$ be the generated summary. The summarization model seeks to maximize the conditional probability:

$$P(\mathbf{y} \mid \mathbf{x}; \theta) = \prod_{i=1}^{m} P(y_i \mid y_1, y_2, \ldots, y_{i-1}, \mathbf{x}; \theta).$$

In extractive summarization, $\mathbf{y}$ is a subset of sentences from $\mathbf{x}$. In abstractive summarization, $\mathbf{y}$ is a newly generated sequence that paraphrases and condenses the information in $\mathbf{x}$.

A summarization model must preserve the essential information from the original document while generating a concise output. Formally, the mutual information between the original document $\mathbf{x}$ and the summary $\mathbf{y}$ should be maximized

$$\text{Mutual Information}(\mathbf{y}; \mathbf{x}) \geq \text{Mutual Information}(\mathbf{y}; \mathbf{x} \setminus \mathbf{y}),$$

ensuring that the summary retains the key content from the document.

Example: For a news article, an extractive summarizer might select sentences that mention the key facts, while an abstractive summarizer might generate a new sentence like "The president announced a new policy today," condensing the article into its core message.

## Extractive and Abstractive Summarization

In extractive summarization, the model identifies and selects key sentences from the original document to form the summary. The challenge lies in selecting sentences that together provide a coherent and informative summary.

Let $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ be the input document consisting of sentences $x_i$. The extractive summarizer selects a subset $\mathbf{y} \subseteq \mathbf{x}$ that maximizes an objective function $f(\mathbf{y}, \mathbf{x}; \theta)$:

$$\mathbf{y}^* = \arg \max_{\mathbf{y} \subseteq \mathbf{x}} f(\mathbf{y}, \mathbf{x}; \theta).$$

The objective function typically measures the informativeness and redundancy of the selected sentences.

Many extractive summarization objectives are submodular, meaning the marginal gain of adding a sentence to the summary decreases as more sentences are added. This property ensures that greedy algorithms can provide near-optimal solutions:

$$f(\mathbf{y} \cup \{x_i\}) - f(\mathbf{y}) \geq f(\mathbf{y}' \cup \{x_i\}) - f(\mathbf{y}')$$

for all $\mathbf{y} \subseteq \mathbf{y}' \subseteq \mathbf{x}$ and $x_i \in \mathbf{x} \setminus \mathbf{y}'$.

In abstractive summarization, the model generates new sentences that capture the main ideas of the document, often rephrasing or condensing the original text. This requires a deeper understanding of the document's content and the ability to generate coherent and fluent text.

Given the input document $\mathbf{x} = (x_1, x_2, \ldots, x_n)$, the abstractive summarizer generates a summary $\mathbf{y} = (y_1, y_2, \ldots, y_m)$ by maximizing the conditional probability:

$$P(\mathbf{y} \mid \mathbf{x}; \theta) = \prod_{i=1}^{m} P(y_i \mid y_1, y_2, \ldots, y_{i-1}, \mathbf{x}; \theta),$$

where $\theta$ represents the model parameters.

Abstractive summarization involves both paraphrasing and compressing the original text. The model learns to generate new text that maintains the semantic content while reducing the length:

$$\text{Mutual Information}(\mathbf{y}; \mathbf{x}) \approx \text{Mutual Information}(\mathbf{y}; \mathbf{x} \setminus \mathbf{y}),$$

ensuring that the summary conveys the essential information from the original document.

Example: In summarizing a research paper, an abstractive model might generate a concise summary like "This study presents a novel algorithm for optimizing network flows," capturing the main contributions of the paper without copying verbatim from the text.

## Evaluation Metrics (ROUGE, BLEU)

The quality of summarization models is typically evaluated using metrics like ROUGE and BLEU, which compare the generated summary to reference summaries.

ROUGE, as previously discussed, measures the overlap of n-grams, sequences of words, and word pairs between the generated summary and the reference summary. ROUGE-1, ROUGE-2, and ROUGE-L (longest common subsequence) are commonly used.

ROUGE-N for n-grams is defined as

$$\text{ROUGE-N} = \frac{\sum_{\text{gram} \in \text{Ref}} \min(\text{Count}_{\text{gen}}(\text{gram}), \text{Count}_{\text{ref}}(\text{gram}))}{\sum_{\text{gram} \in \text{Ref}} \text{Count}_{\text{ref}}(\text{gram})}.$$

ROUGE-N captures the recall of the generated summary, but it can also be adapted to measure precision by swapping the roles of the generated and reference summaries:

$$\text{ROUGE-N}_{\text{Precision}} = \frac{\sum_{\text{gram} \in \text{Gen}} \min(\text{Count}_{\text{gen}}(\text{gram}), \text{Count}_{\text{ref}}(\text{gram}))}{\sum_{\text{gram} \in \text{Gen}} \text{Count}_{\text{gen}}(\text{gram})}.$$

BLEU measures the precision of n-grams in the generated summary against the reference summary, with a brevity penalty to penalize overly short summaries. It is more commonly used for machine translation but can be applied to summarization.

BLEU score is computed as

$$\text{BLEU} = \text{BP} \cdot \exp\left(\sum_{n=1}^{N} w_n \log p_n\right),$$

where $p_n$ is the precision for n-grams of length $n$, $w_n$ are the weights, and BP is the brevity penalty.

Example: In evaluating a summary generated by a transformer model, ROUGE scores might indicate how much of the reference summary's content is captured, while BLEU might provide insight into the fluency and precision of the generated text.

## 5.2 Transformer-Based Models

BERT is one of the most influential models in NLP, setting a new standard for various tasks such as question answering, text classification, and sentiment analysis. This section explores the mathematical foundations of BERT, its architecture, the role of the self-attention mechanism, and the processes of pre-training and fine-tuning that are integral to its success.

### 5.2.1 BERT

BERT is a transformer-based model that captures bidirectional context in text, meaning it considers both the left and right contexts of a word to generate more meaningful representations. This bidirectional approach distinguishes BERT from earlier models, which typically processed text in a unidirectional manner.

**Mathematical Foundation**

BERT is based on the concept of contextualized word embeddings, where the representation of a word depends on its surrounding words. This is achieved through the use of transformers, specifically the encoder part of the transformer architecture.

Let $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ be an input sequence of tokens, where each $x_i$ is a token from a vocabulary $\mathcal{V}$. The goal of BERT is to learn a function $f_\theta$ that maps this sequence to a sequence of contextualized embeddings $\mathbf{h} = (h_1, h_2, \ldots, h_n)$, where $h_i$ represents the embedding of token $x_i$ conditioned on the entire sequence:

$$h_i = f_\theta(x_i \mid x_1, x_2, \ldots, x_{i-1}, x_{i+1}, \ldots, x_n).$$

The function $f_\theta$ is parameterized by the layers of the transformer encoder, which involves multiple layers of self-attention and feedforward networks.

BERT's bidirectional nature is mathematically grounded in the fact that each token's representation $h_i$ is computed by attending to all tokens in the sequence:

$$h_i = \text{SelfAttention}(q_i, K, V),$$

where $q_i$ is the query derived from $x_i$, and $K$ and $V$ are the key and value matrices derived from the entire sequence $\mathbf{x}$. This attention mechanism ensures that $h_i$ captures information from both the left and right contexts of $x_i$.

Example: Consider the sentence "The bank will close at noon." In BERT, the word "bank" will have different embeddings depending on the surrounding words "will" and "close," helping to distinguish between the meanings of "bank" as a financial institution and a riverbank.

**Architecture and Layers**

BERT's architecture consists of multiple transformer encoder layers stacked on top of each other. Each layer comprises a self-attention mechanism followed by a feedforward neural network, with layer normalization and residual connections.

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the input matrix, where $n$ is the sequence length and $d$ is the dimensionality of the embeddings. The output of the $l$-th encoder layer is given by

$$\mathbf{H}^{(l)} = \text{LayerNorm}\left(\mathbf{H}^{(l-1)} + \text{FFN}\left(\text{SelfAttention}\left(\mathbf{H}^{(l-1)}\right)\right)\right),$$

where $\mathbf{H}^{(l-1)}$ is the output of the previous layer and FFN denotes the feedforward network, which applies two linear transformations with a ReLU activation in between:

$$\text{FFN}(\mathbf{h}) = \text{ReLU}(\mathbf{h}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2.$$

The stacking of transformer layers allows for deep contextualization, where information from different parts of the sequence is progressively integrated. The residual

connections help preserve the original input information, while the normalization stabilizes training:

$$\mathbf{H}^{(l)} = \mathbf{H}^{(l-1)} + \mathcal{O}(d^{\frac{3}{2}}),$$

where $\mathcal{O}(d^{\frac{3}{2}})$ represents the complexity of the attention mechanism.

Example: In a 12-layer BERT model (BERT-Base), each word in a sentence is transformed through 12 layers of self-attention and feedforward operations, progressively refining the word representations by incorporating more context from the entire sentence.

## Self-Attention Mechanism in BERT

The self-attention mechanism is a critical component of BERT, allowing the model to weigh the importance of different tokens in a sequence when constructing the representation of a particular token.

The self-attention mechanism computes a weighted sum of values, where the weights (attention scores) are determined by the similarity between the query and key vectors. Formally, for a query vector $q$ and a set of key-value pairs $(K, V)$, the self-attention output is

$$\text{Attention}(q, K, V) = \sum_{i=1}^{n} \alpha_i v_i, \quad \alpha_i = \frac{\exp(q^\top k_i)}{\sum_{j=1}^{n} \exp(q^\top k_j)},$$

where $k_i$ and $v_i$ are the key and value vectors for token $x_i$.

Self-attention enables BERT to capture complex dependencies between words, regardless of their distance in the sequence. The attention scores $\alpha_i$ reflect the model's focus on different parts of the input sequence, leading to context-aware word embeddings:

$$h_i = \sum_{j=1}^{n} \alpha_{ij} v_j.$$

This mechanism allows BERT to handle phenomena like polysemy and syntactic ambiguity effectively.

Example: In the sentence "The animal didn't cross the road because it was too tired," the self-attention mechanism helps BERT understand that "it" refers to "the animal," rather than "the road," by assigning higher attention scores to relevant tokens.

## Pre-training and Fine-tuning

BERT's success is largely due to its two-stage training process: pre-training on large text corpora and fine-tuning on specific downstream tasks.

1. Pre-training: BERT is pre-trained on two tasks—MLM and NSP. The pre-training objective combines the losses from both tasks:

$$\mathcal{L}_{\text{pre-train}}(\theta) = \mathcal{L}_{\text{MLM}}(\theta) + \mathcal{L}_{\text{NSP}}(\theta),$$

where $\mathcal{L}_{\text{MLM}}(\theta)$ is the loss for predicting masked tokens and $\mathcal{L}_{\text{NSP}}(\theta)$ is the loss for predicting sentence order.

2. Fine-tuning: After pre-training, BERT is fine-tuned on specific tasks by adding task-specific layers and minimizing the corresponding loss function $\mathcal{L}_{\text{task}}(\theta)$:

$$\mathcal{L}_{\text{fine-tune}}(\theta) = \mathcal{L}_{\text{task}}(\theta).$$

Fine-tuning typically involves fewer epochs and task-specific data, adapting the pre-trained BERT model to various applications.

The pre-training phase allows BERT to learn transferable representations that generalize well across different NLP tasks. This is mathematically expressed by the low variation in task-specific performance across a wide range of tasks:

$$\text{Var}_{\text{tasks}}(\mathcal{L}_{\text{fine-tune}}(\theta)) \ll \text{Var}_{\text{tasks}}(\mathcal{L}_{\text{from-scratch}}(\theta)),$$

indicating that the fine-tuning loss is much lower for pre-trained models than for models trained from scratch.

Example: BERT is pre-trained on a large corpus like Wikipedia, then fine-tuned on a sentiment analysis task. The model quickly adapts to the new task, leveraging its understanding of language from the pre-training phase to excel in sentiment classification.

**Pre-training Tasks (MLM, NSP)**

BERT's pre-training tasks—MLM and NSP—are designed to enable the model to learn deep contextual relationships in text.

MLM involves randomly masking some tokens in the input sequence and training the model to predict these masked tokens based on the surrounding context. The MLM loss is defined as

$$\mathcal{L}_{\text{MLM}}(\theta) = \mathbb{E}_{(\mathbf{x},\mathbf{m})} \left[ \sum_{i=1}^{n} m_i \log P(x_i \mid \mathbf{x}_{\setminus i}; \theta) \right],$$

where $\mathbf{m}$ is a binary mask vector and $\mathbf{x}_{\setminus i}$ denotes the sequence with the $i$-th token masked.

NSP is designed to help the model understand sentence relationships. The model is given pairs of sentences and trained to predict whether the second sentence logically follows the first:

$$\mathcal{L}_{\text{NSP}}(\theta) = \mathbb{E}_{(\mathbf{s}_1, \mathbf{s}_2), y} \left[ y \log P(\text{True} \mid \mathbf{s}_1, \mathbf{s}_2; \theta) + (1 - y) \log P(\text{False} \mid \mathbf{s}_1, \mathbf{s}_2; \theta) \right],$$

where $y = 1$ if $\mathbf{s}_2$ follows $\mathbf{s}_1$, and $y = 0$ otherwise.

The combination of MLM and NSP tasks in pre-training allows BERT to learn robust representations that are sensitive to both local word context (through MLM) and global sentence relationships (through NSP):

$$\theta_{\text{pre-train}} = \arg \min_{\theta} \left( \mathcal{L}_{\text{MLM}}(\theta) + \mathcal{L}_{\text{NSP}}(\theta) \right).$$

These representations are then fine-tuned for specific tasks, providing a strong foundation for task-specific learning.

Example: MLM enables BERT to predict masked words based on the context, improving its understanding of syntax and semantics. NSP, on the other hand, helps BERT capture discourse-level information, such as the logical sequence of sentences, making it effective for tasks like document classification and passage retrieval.

## Fine-tuning for Downstream Tasks

Fine-tuning is the process by which a pre-trained model like BERT is adapted to a specific task by adjusting its parameters to minimize a task-specific loss function. The key idea is that the pre-trained model already captures general linguistic knowledge, and fine-tuning allows it to specialize in the nuances of the target task.

Let $\mathcal{D}_{\text{task}} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^{N}$ be the dataset for a downstream task, where $\mathbf{x}^{(i)}$ is an input example and $\mathbf{y}^{(i)}$ is the corresponding label or output. The objective of fine-tuning is to find the parameter set $\theta^*$ that minimizes the task-specific loss function $\mathcal{L}_{\text{task}}(\theta)$:

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^{N} \mathcal{L}_{\text{task}}(f_\theta(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}),$$

where $f_\theta(\mathbf{x})$ is the function representing the BERT model with parameters $\theta$.

Under typical conditions, the fine-tuning process converges to a local minimum of the loss function. The rate of convergence and the quality of the solution depend on the choice of optimization algorithm, learning rate, and the structure of the loss landscape:

$$\lim_{t \to \infty} \theta_t = \theta^*, \quad \text{where} \quad \theta_{t+1} = \theta_t - \eta \nabla_\theta \mathcal{L}_{\text{task}}(\theta_t).$$

Example: In fine-tuning BERT for sentiment analysis, the model might be trained to minimize the cross-entropy loss between the predicted sentiment scores and the true sentiment labels. The fine-tuning process adjusts the weights of BERT to optimize its performance on the sentiment analysis dataset.

## Text Classification

Text classification involves assigning a label or category to a given piece of text. BERT excels at text classification tasks due to its ability to capture deep contextual relationships in the text.

Given an input text $\mathbf{x} = (x_1, x_2, \ldots, x_n)$, the goal is to predict a label $y$ from a set of possible labels $\mathcal{Y}$. The BERT model encodes the input into a fixed-dimensional vector $h_{\text{CLS}}$ (corresponding to the [CLS] token) and passes it through a classification layer:

$$P(y \mid \mathbf{x}; \theta) = \text{softmax}(\mathbf{W}h_{\text{CLS}} + \mathbf{b}),$$

where $\mathbf{W}$ and $\mathbf{b}$ are the weights and bias of the classification layer, respectively.

The pre-trained BERT model transforms the input sequence into a high-dimensional space where the different classes are more likely to be linearly separable. This transformation allows the classification layer to achieve high accuracy even with simple linear classifiers:

$$\mathcal{L}_{\text{task}}(\theta) = -\sum_{i=1}^{|\mathcal{Y}|} y_i \log P(y_i \mid \mathbf{x}; \theta).$$

Example: For a spam detection task, BERT can be fine-tuned to classify emails as "spam" or "not spam." The [CLS] token representation serves as a summary of the entire email, enabling the model to make accurate predictions based on the encoded contextual information.

## Named Entity Recognition (NER)

Named Entity Recognition (NER) is the task of identifying and classifying named entities (such as people, organizations, locations) in text. BERT's token-level representations are particularly effective for NER tasks.

In NER, the goal is to assign a label $y_i$ to each token $x_i$ in the input sequence $\mathbf{x} = (x_1, x_2, \ldots, x_n)$. BERT generates a contextualized embedding $h_i$ for each token, and these embeddings are passed through a classification layer:

$$P(y_i \mid x_i, \mathbf{x}; \theta) = \text{softmax}(\mathbf{W}h_i + \mathbf{b}),$$

where $\mathbf{W}$ and $\mathbf{b}$ are the weights and bias specific to the NER task.

BERT's self-attention mechanism enables the model to capture dependencies between tokens, making it particularly effective for NER. The model can distinguish between different entities based on their context:

$$h_i = \text{SelfAttention}(h_i, K, V),$$

where $h_i$ is influenced by other tokens in the sequence, allowing BERT to resolve ambiguities (e.g., distinguishing between "Apple" as a company versus a fruit).

Example: In a sentence like "Apple announced a new product," BERT can accurately classify "Apple" as an organization and "product" as a common noun, using the surrounding context to disambiguate entity types.

### Question Answering

Question Answering (QA) tasks require a model to read a passage and answer questions based on the content. BERT's ability to understand context makes it highly effective for QA.

Given a passage $\mathbf{p} = (p_1, p_2, \ldots, p_m)$ and a question $\mathbf{q} = (q_1, q_2, \ldots, q_n)$, the goal is to find a span in the passage that answers the question. BERT encodes both the passage and the question, and the model predicts the start and end positions $s$ and $e$ of the answer span:

$$P(s, e \mid \mathbf{p}, \mathbf{q}; \theta) = P(s \mid \mathbf{p}, \mathbf{q}; \theta) \cdot P(e \mid \mathbf{p}, \mathbf{q}; \theta),$$

where $P(s \mid \mathbf{p}, \mathbf{q}; \theta)$ and $P(e \mid \mathbf{p}, \mathbf{q}; \theta)$ are obtained from the softmax outputs over the token positions in the passage.

The ability of BERT to predict the correct answer span depends on its capacity to model interactions between the passage and the question. The attention mechanism helps in aligning the passage tokens with the relevant parts of the question:

$$\mathcal{L}_{\text{QA}}(\theta) = -\log P(s \mid \mathbf{p}, \mathbf{q}; \theta) - \log P(e \mid \mathbf{p}, \mathbf{q}; \theta).$$

Example: For the question "Who is the president of the United States?" given the passage "The president of the United States is Joe Biden," BERT can accurately predict "Joe Biden" as the answer span by identifying the relevant context in the passage.

### Applications and Performance

The versatility of BERT extends to a wide range of NLP applications, including sentiment analysis, machine translation, and text summarization. The model's performance is often benchmarked against datasets like GLUE (General Language Understanding Evaluation) and SQuAD (Stanford Question Answering Dataset).

Performance metrics such as accuracy, F1-score, and Exact Match (EM) are used to evaluate BERT's effectiveness across tasks. These metrics are derived from the model's predictions $\hat{\mathbf{y}}$ compared to the true labels $\mathbf{y}$:

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(\hat{y}_i = y_i),$$

$$\text{F1-score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}},$$

$$\text{Exact Match (EM)} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(\hat{\mathbf{y}}_i = \mathbf{y}_i),$$

where $\mathbb{I}$ is the indicator function.

BERT's architecture and pre-training strategy enable it to generalize well across different NLP tasks, reducing the gap between human performance and machine performance on benchmarks:

$$\text{Generalization Error} = \mathbb{E}[\mathcal{L}_{\text{task}}(\theta)] - \mathbb{E}[\mathcal{L}_{\text{task}}(\theta_{\text{human}})],$$

where $\theta_{\text{human}}$ represents the ideal human-level performance parameters.

Example: On the GLUE benchmark, BERT achieves near-human performance across multiple tasks such as natural language inference and sentiment analysis, demonstrating its ability to understand and process complex language patterns.

## 5.2.2   GPT

GPT is an autoregressive model that generates text by predicting the next token in a sequence, conditioned on the previous tokens. Unlike BERT, which is bidirectional, GPT processes text in a unidirectional manner, making it particularly suited for generative tasks such as text completion and story generation.

### Mathematical Foundation

The core idea behind GPT is to model the joint probability distribution of a sequence of tokens $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ in a way that allows the model to generate coherent sequences by sampling from this distribution.

The joint probability of the sequence $\mathbf{x}$ is factorized using the chain rule of probability, which allows GPT to predict each token sequentially:

$$P(\mathbf{x}; \theta) = \prod_{i=1}^{n} P(x_i \mid x_1, x_2, \ldots, x_{i-1}; \theta),$$

where $\theta$ represents the parameters of the model. This factorization enables GPT to generate text in a left-to-right fashion, with each token being conditioned on the previous ones.

GPT's autoregressive nature ensures that the prediction of each token $x_i$ depends solely on the preceding tokens $x_1, x_2, \ldots, x_{i-1}$, which prevents future tokens from influencing the generation of $x_i$:

$$P(x_i \mid x_1, x_2, \ldots, x_{i-1}; \theta) = P(x_i \mid \mathbf{h}_{i-1}),$$

where $\mathbf{h}_{i-1}$ is the hidden state encoding information from the previous tokens.

Example: In a text generation task, GPT might be given the prompt "Once upon a time," and it would generate a continuation like "there was a kingdom ruled by a wise king." The model's ability to produce such coherent text stems from its learned probability distribution over sequences.

## Architecture and Layers

GPT's architecture is built on the transformer's encoder–decoder framework, but it uses only the decoder stack of the transformer. This unidirectional nature of GPT is crucial for its autoregressive capabilities.

Let $\mathbf{H}^{(l)}$ be the hidden state matrix at the $l$-th layer, where $l$ ranges from 1 to $L$, the total number of layers in the model. The computation within each layer can be described by the following equations:

1. Self-Attention:

$$\mathbf{H}^{(l)} = \text{LayerNorm}\left(\mathbf{H}^{(l-1)} + \text{SelfAttention}\left(\mathbf{H}^{(l-1)}\right)\right),$$

where self-attention captures dependencies between tokens in the sequence.

2. Feedforward Network:

$$\mathbf{H}^{(l)} = \text{LayerNorm}\left(\mathbf{H}^{(l)} + \text{FFN}\left(\mathbf{H}^{(l)}\right)\right),$$

where the feedforward network (FFN) applies non-linear transformations to the output of the self-attention mechanism.

The depth of GPT, characterized by the number of layers $L$, contributes to the expressivity of the model. As the number of layers increases, the model can capture

more complex patterns and dependencies in the data, leading to more coherent and contextually relevant text generation:

$$\mathbf{H}^{(L)} = f^{(L)}(f^{(L-1)}(\ldots f^{(1)}(\mathbf{X}))),$$

where $f^{(l)}$ represents the operations at layer $l$.

Example: In GPT-2, a larger version of GPT, the model's depth and size enable it to generate long passages of text that are often indistinguishable from human writing, demonstrating the power of the architecture.

**Causal Self-Attention**

Causal self-attention is a key feature of GPT that ensures each token is only influenced by the preceding tokens, making the model suitable for autoregressive text generation.

In causal self-attention, the attention mechanism is masked to prevent any token from attending to future tokens. Let $\mathbf{Q}$, $\mathbf{K}$, $and\mathbf{V}$ be the query, key, and value matrices, respectively. The attention scores are computed as

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d_k}} + \mathbf{M}\right)\mathbf{V},$$

where $\mathbf{M}$ is a masking matrix that sets the attention scores to negative infinity for future tokens, ensuring that $x_i$ can only attend to tokens $x_1, x_2, \ldots, x_{i-1}$.

The causal nature of self-attention in GPT guarantees that the model's output at each position depends only on the preceding context. This causality is crucial for tasks that require sequential generation, such as text completion:

$$\mathbf{H}_i = \text{SelfAttention}(q_i, K_{\leq i}, V_{\leq i}),$$

where $K_{\leq i}$ and $V_{\leq i}$ are the key and value matrices up to the $i$-th position.

Example: In generating the sentence "The cat sat on the," GPT predicts the next word "mat" by attending only to the preceding words "The cat sat on the," ensuring that the generation is sequential and contextually accurate.

**Pre-training and Fine-tuning**

GPT's training process involves two main stages: unsupervised pre-training on large text corpora and supervised fine-tuning on specific tasks.

1. Pre-training: In the pre-training phase, GPT is trained to predict the next token in a sequence using a language modeling objective. The loss function for pre-training is the negative log-likelihood of the predicted tokens:

$$\mathcal{L}_{\text{pre-train}}(\theta) = -\sum_{i=1}^{n} \log P(x_i \mid x_1, x_2, \ldots, x_{i-1}; \theta).$$

2. Fine-tuning: After pre-training, GPT is fine-tuned on task-specific data. The model is adapted to minimize a task-specific loss function, which could involve classification, generation, or other objectives:

$$\mathcal{L}_{\text{fine-tune}}(\theta) = \mathbb{E}_{(\mathbf{x},\mathbf{y}) \in \mathcal{D}_{\text{task}}} \left[ \mathcal{L}_{\text{task}}(f_\theta(\mathbf{x}), \mathbf{y}) \right],$$

where $\mathcal{D}_{\text{task}}$ is the task-specific dataset.

The representations learned during the pre-training phase are highly transferable to a wide range of tasks. Fine-tuning allows the model to adapt these general representations to the specific requirements of a given task:

$$\theta_{\text{fine-tune}} = \theta_{\text{pre-train}} + \Delta\theta,$$

where $\Delta\theta$ represents the task-specific adjustments made during fine-tuning.

Example: GPT can be pre-trained on a vast corpus of web text and then fine-tuned on a smaller dataset for a specific task like summarization or dialog generation. The fine-tuned model benefits from the rich language understanding developed during pre-training.

## Unsupervised Pre-training

GPT's success is largely attributed to its unsupervised pre-training, which allows the model to learn from large amounts of unlabeled text data, capturing the nuances of natural language.

In unsupervised pre-training, the model is trained solely on the task of language modeling, without any task-specific labels. The objective is to maximize the likelihood of the sequence given the previous tokens:

$$\mathcal{L}_{\text{unsupervised}}(\theta) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ -\sum_{i=1}^{n} \log P(x_i \mid x_1, x_2, \ldots, x_{i-1}; \theta) \right],$$

where $\mathcal{D}$ is the pre-training corpus.

Unsupervised pre-training enables GPT to capture a high density of information from the training data, leading to the development of rich, context-aware representations that are beneficial for a wide range of tasks:

$$\text{Information Density} = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \mathbb{I}(x_i; x_1, x_2, \ldots, x_{i-1}),$$

where $\mathbb{I}$ represents the mutual information between tokens.

Example: By pre-training on a massive corpus like the Common Crawl dataset, GPT learns to generate text that is coherent, contextually relevant, and often indistinguishable from human writing, demonstrating the power of unsupervised learning.

### Supervised Fine-tuning

Supervised fine-tuning is the process by which a pre-trained GPT model is adapted to a specific task using labeled data. The goal is to optimize the model's parameters to minimize a task-specific loss function, thereby improving its performance on the task at hand.

Let $\mathcal{D}_{\text{task}} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^{N}$ be the labeled dataset for the task, where $\mathbf{x}^{(i)}$ represents the input text and $\mathbf{y}^{(i)}$ represents the corresponding label or output. The fine-tuning process involves minimizing the task-specific loss function $\mathcal{L}_{\text{task}}(\theta)$ with respect to the model parameters $\theta$:

$$\theta^* = \arg\min_{\theta} \sum_{i=1}^{N} \mathcal{L}_{\text{task}}(f_\theta(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}),$$

where $f_\theta(\mathbf{x})$ represents the output of the GPT model for input $\mathbf{x}$.

Under appropriate conditions, supervised fine-tuning converges to a local minimum of the task-specific loss function. The rate of convergence and the quality of the solution depend on the choice of optimization algorithm, learning rate, and the structure of the loss landscape:

$$\lim_{t \to \infty} \theta_t = \theta^*, \quad \text{where} \quad \theta_{t+1} = \theta_t - \eta \nabla_\theta \mathcal{L}_{\text{task}}(\theta_t).$$

Example: In fine-tuning GPT for a text classification task, the model might be trained to minimize the cross-entropy loss between the predicted class probabilities and the true labels. The fine-tuning process adjusts the model's weights to optimize its performance on the classification dataset.

### Applications and Performance

GPT has been fine-tuned for various NLP applications, including text generation, dialog systems, and story generation. Its performance in these applications is often evaluated using specific metrics that assess the quality and coherence of the generated text.

1. Text Generation: Given a prompt $\mathbf{p}$, the goal of text generation is to produce a coherent and contextually relevant continuation $\mathbf{y} = (y_1, y_2, \ldots, y_m)$. The model generates text by sampling from the conditional distribution:

$$P(\mathbf{y} \mid \mathbf{p}; \theta) = \prod_{i=1}^{m} P(y_i \mid \mathbf{p}, y_1, y_2, \ldots, y_{i-1}; \theta).$$

2. Dialog Systems: In dialog systems, the model generates responses $\mathbf{r}$ to user inputs $\mathbf{u}$ by modeling the conditional probability distribution:

$$P(\mathbf{r} \mid \mathbf{u}; \theta) = \prod_{i=1}^{n} P(r_i \mid \mathbf{u}, r_1, r_2, \ldots, r_{i-1}; \theta),$$

where $r_i$ represents the tokens in the generated response.

3. Story Generation: For story generation, the model generates a narrative based on an initial premise $\mathbf{s}$. The task is to produce a coherent and engaging continuation:

$$P(\mathbf{s}_{\text{cont}} \mid \mathbf{s}; \theta) = \prod_{i=1}^{m} P(s_i \mid \mathbf{s}, s_1, s_2, \ldots, s_{i-1}; \theta),$$

where $\mathbf{s}_{\text{cont}}$ denotes the continuation of the story.

The ability of GPT to maintain contextual coherence in generated text is a result of its autoregressive architecture, where each token is conditioned on the entire preceding context. This ensures that the generated text is both relevant and logically consistent:

$$P(\mathbf{y} \mid \mathbf{p}; \theta) = \prod_{i=1}^{m} P(y_i \mid \mathbf{h}_{i-1}),$$

where $\mathbf{h}_{i-1}$ encodes the history up to the $i$-th token.

Example: In a dialog system, GPT might be fine-tuned to respond to user queries in a conversational manner. Given the input "How is the weather today?" the model could generate a contextually appropriate response like "It's sunny and warm."

## Text Generation

Text generation involves producing coherent and contextually relevant sequences based on a given prompt. GPT excels at this task due to its ability to model long-range dependencies in text.

Given a prompt $\mathbf{p} = (p_1, p_2, \ldots, p_n)$, the goal is to generate a sequence $\mathbf{y} = (y_1, y_2, \ldots, y_m)$ that extends the prompt. The likelihood of the generated sequence is given by

$$P(\mathbf{y} \mid \mathbf{p}; \theta) = \prod_{i=1}^{m} P(y_i \mid \mathbf{p}, y_1, y_2, \ldots, y_{i-1}; \theta).$$

The coherence of the generated text is ensured by maximizing the conditional probabilities of each token in the sequence. The autoregressive nature of GPT allows it to generate text that is consistent with the preceding context:

$$\mathbf{y}^* = \arg\max_{\mathbf{y}} P(\mathbf{y} \mid \mathbf{p}; \theta).$$

Example: Given the prompt "The future of AI is," GPT might generate a continuation like "full of exciting possibilities, with advances in machine learning and robotics shaping our world."

**Dialog Systems**

Dialog systems aim to generate natural and engaging responses to user inputs in a conversation. GPT's ability to model context makes it particularly effective for this task.

For a dialog system, the model generates a response $\mathbf{r} = (r_1, r_2, \ldots, r_n)$ to a user input $\mathbf{u} = (u_1, u_2, \ldots, u_m)$. The likelihood of the response is given by

$$P(\mathbf{r} \mid \mathbf{u}; \theta) = \prod_{i=1}^{n} P(r_i \mid \mathbf{u}, r_1, r_2, \ldots, r_{i-1}; \theta).$$

The effectiveness of GPT in dialog systems stems from its ability to generate responses that are contextually relevant, maintaining the flow of conversation. This is achieved by conditioning each token in the response on both the user input and the preceding tokens in the response:

$$\mathbf{r}^* = \arg\max_{\mathbf{r}} P(\mathbf{r} \mid \mathbf{u}; \theta).$$

Example: In a customer service chatbot, given the user input "I need help with my order," GPT might generate a response like "Sure, I can help with that. Can you provide your order number?"

**Story Generation**

Story generation involves creating coherent and engaging narratives based on an initial premise or a series of prompts. GPT is capable of generating long-form text that maintains narrative consistency.

Given an initial premise $\mathbf{s} = (s_1, s_2, \ldots, s_n)$, the goal is to generate a continuation $\mathbf{s}_{\text{cont}} = (s_{n+1}, s_{n+2}, \ldots, s_m)$ that extends the story. The probability of the continuation is modeled as

$$P(\mathbf{s}_{\text{cont}} \mid \mathbf{s}; \theta) = \prod_{i=n+1}^{m} P(s_i \mid \mathbf{s}, s_{n+1}, s_{n+2}, \ldots, s_{i-1}; \theta).$$

GPT's architecture allows it to generate stories that are both coherent and consistent, ensuring that the generated narrative logically follows from the initial premise:

$$\mathbf{s}_{\text{cont}}^{*} = \arg\max_{\mathbf{s}_{\text{cont}}} P(\mathbf{s}_{\text{cont}} \mid \mathbf{s}; \theta).$$

Example: Given the premise "Once upon a time in a faraway kingdom," GPT might generate a continuation like "there lived a young prince who dreamed of exploring the world beyond his castle."

**Evaluation Metrics and Performance Analysis**

The performance of GPT in tasks such as text generation, dialog systems, and story generation is typically evaluated using metrics that assess the quality, coherence, and relevance of the generated text.

1. Perplexity: Measures the model's ability to predict the next token in a sequence. Lower perplexity indicates better performance:

$$\text{Perplexity} = \exp\left(-\frac{1}{N} \sum_{i=1}^{N} \log P(x_i \mid x_1, x_2, \ldots, x_{i-1}; \theta)\right).$$

2. BLEU (Bilingual Evaluation Understudy): Measures the overlap between the generated text and reference text based on n-gram precision:

$$\text{BLEU} = \text{BP} \cdot \exp\left(\sum_{n=1}^{N} w_n \log p_n\right),$$

where $p_n$ is the precision of n-grams and BP is the brevity penalty.

3. ROUGE (Recall-Oriented Understudy for Gisting Evaluation): Evaluates the quality of summaries by comparing n-grams, word sequences, and word pairs:

$$\text{ROUGE-N} = \frac{\sum_{\text{gram} \in \text{Ref}} \min(\text{Count}_{\text{gen}}(\text{gram}), \text{Count}_{\text{ref}}(\text{gram}))}{\sum_{\text{gram} \in \text{Ref}} \text{Count}_{\text{ref}}(\text{gram})}.$$

Evaluation metrics such as Perplexity, BLEU, and ROUGE are effective for assessing the performance of GPT in generative tasks. These metrics capture different aspects of text quality, such as fluency, precision, and recall:

$$\text{Performance} = \frac{1}{K} \sum_{k=1}^{K} \text{Metric}_k(\hat{\mathbf{y}}, \mathbf{y}),$$

where $\hat{\mathbf{y}}$ and $\mathbf{y}$ are the predicted and reference sequences, respectively.

Example: In evaluating a story generation model, BLEU and ROUGE scores might be used to assess how closely the generated story matches a reference story, while Perplexity would indicate how well the model predicts the continuation of the story.

## 5.3   Advanced Topics in NLP with Transformers

Transformers have revolutionized various aspects of NLP, enabling significant advancements in sequence labeling and semantic tasks. This section explores the application of transformers to sequence labeling tasks such as part-of-speech tagging and chunking, as well as semantic tasks like semantic role labeling. The mathematical foundations of these applications are emphasized, with a focus on Geometry, Symmetry, and Intelligence.

### 5.3.1   Transformers for Sequence Labeling

Sequence labeling involves assigning a label to each token in a sequence. Transformers, with their self-attention mechanisms and ability to capture long-range dependencies, are well suited for sequence labeling tasks.

**Part-of-Speech Tagging**

Part-of-speech (POS) tagging is the task of assigning a grammatical category (such as noun, verb, adjective) to each token in a sentence. Transformers leverage contextualized embeddings to improve the accuracy of POS tagging by considering the entire sentence context.

Given a sequence of tokens $\mathbf{x} = (x_1, x_2, \ldots, x_n)$, the goal is to assign a POS tag $y_i$ to each token $x_i$. The transformer model generates a contextualized embedding $h_i$ for each token, which is then passed through a classification layer to predict the POS tag:

$$P(y_i \mid x_i, \mathbf{x}; \theta) = \text{softmax}(\mathbf{W} h_i + \mathbf{b}),$$

where $\mathbf{W}$ and $\mathbf{b}$ are the weights and bias of the classification layer, respectively.

The self-attention mechanism in transformers enables the model to capture dependencies between tokens, improving the accuracy of POS tagging by generating embeddings that consider the entire sentence context:

$$h_i = \text{SelfAttention}(q_i, K, V),$$

where $q_i$, $K$, and $V$ are the query, key, and value matrices derived from the sequence $\mathbf{x}$.

Example: In the sentence "The quick brown fox jumps over the lazy dog," the word "jumps" is correctly tagged as a verb because the transformer model captures its relationship with surrounding words.

## Chunking

Chunking, also known as shallow parsing, involves grouping sequences of tokens into syntactically correlated parts of a sentence, such as noun phrases (NP) or verb phrases (VP). Transformers excel at chunking by utilizing their ability to model long-range dependencies.

Given a sequence of tokens $\mathbf{x} = (x_1, x_2, \ldots, x_n)$, the goal is to assign a chunk label $c_i$ to each token $x_i$, where $c_i$ indicates the chunk to which the token belongs. The transformer model generates a contextualized embedding $h_i$ for each token and predicts the chunk label:

$$P(c_i \mid x_i, \mathbf{x}; \theta) = \text{softmax}(\mathbf{W}h_i + \mathbf{b}),$$

where $\mathbf{W}$ and $\mathbf{b}$ are the parameters of the classification layer.

Transformers effectively capture long-range dependencies in text, allowing for accurate chunking by modeling the relationships between non-adjacent tokens:

$$h_i = \sum_{j=1}^{n} \alpha_{ij} v_j, \quad \alpha_{ij} = \frac{\exp(q_i^\top k_j)}{\sum_{k=1}^{n} \exp(q_i^\top k_k)},$$

where $\alpha_{ij}$ are the attention weights that determine how much token $x_j$ influences the representation of token $x_i$.

Example: In the sentence "The quick brown fox jumps over the lazy dog," the transformer model can correctly identify "The quick brown fox" as a noun phrase (NP) and "jumps over" as a verb phrase (VP).

## 5.3.2   Transformers for Semantic Tasks

Semantic tasks involve understanding the meaning and relationships of words within a sentence. Transformers are particularly effective at these tasks due to their ability to model complex dependencies and capture semantic nuances.

### Semantic Role Labeling

Semantic role labeling (SRL) is the task of identifying the predicate-argument structure of a sentence, determining the roles that different words play in relation to a predicate (usually a verb). Transformers excel at SRL by capturing the interactions between words and their roles in a sentence.

Given a sentence $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ and a predicate $p$, the goal is to assign a semantic role $r_i$ to each token $x_i$. The transformer model generates a contextualized embedding $h_i$ for each token and predicts the semantic role:

$$P(r_i \mid x_i, \mathbf{x}, p; \theta) = \text{softmax}(\mathbf{W}h_i + \mathbf{b}),$$

where $\mathbf{W}$ and $\mathbf{b}$ are the parameters of the classification layer.

The self-attention mechanism in transformers allows for the modeling of role-specific interactions between tokens, leading to accurate predictions of semantic roles:

$$h_i = \text{SelfAttention}(q_i, K, V),$$

where the query $q_i$ is influenced by the predicate $p$, ensuring that the context of the predicate is considered when determining the role of $x_i$.

Example: In the sentence "She gave him a book," the transformer model can identify "She" as the agent, "him" as the recipient, and "a book" as the theme, correctly assigning the semantic roles based on the predicate "gave."

### Coreference Resolution

Coreference resolution is the task of identifying expressions in a text that refer to the same entity. This task requires understanding the relationships between different parts of a text and is essential for maintaining coherence in text interpretation.

Given a document consisting of tokens $\mathbf{x} = (x_1, x_2, \ldots, x_n)$ and a set of mentions $\mathcal{M} = \{m_1, m_2, \ldots, m_k\}$, where each mention $m_i$ refers to a subsequence of tokens, the goal is to cluster these mentions into groups, where each group represents a distinct entity. The transformer model generates contextualized embeddings $h_i$ for each token, which are used to determine the probability that two mentions $m_i$ and $m_j$ refer to the same entity:

$$P(\text{coref}(m_i, m_j) \mid \mathbf{x}; \theta) = \sigma(\mathbf{W}[h_{m_i}, h_{m_j}] + \mathbf{b}),$$

where $[h_{m_i}, h_{m_j}]$ is the concatenation of the embeddings of the mentions and $\sigma$ is the sigmoid function.

The self-attention mechanism in transformers allows the model to capture relationships between distant tokens, making it effective for resolving coreferences, even when the mentions are far apart in the text:

$$h_{m_i} = \text{SelfAttention}(q_{m_i}, K, V),$$

where $q_{m_i}$ is the query derived from the mention $m_i$, and $K$ and $V$ are the key and value matrices derived from the entire sequence $\mathbf{x}$.

Example: In the text "John said he would come, but he didn't," the transformer model can resolve "he" as referring to "John" by considering the context provided by the entire sentence.

### 5.3.3   Transformers for Structured Prediction

Structured prediction involves making predictions that are interdependent and often require maintaining a certain structure, such as a tree or graph. Transformers have been adapted to handle structured prediction tasks like dependency parsing and constituency parsing.

#### Dependency Parsing

Dependency parsing aims to identify the syntactic structure of a sentence by establishing relationships between words, where each word (except the root) depends on exactly one other word. These relationships form a tree structure, with words as nodes and dependencies as directed edges.

Given a sentence $\mathbf{x} = (x_1, x_2, \ldots, x_n)$, the goal is to construct a dependency tree $\mathcal{T} = (\mathbf{V}, \mathbf{E})$, where $\mathbf{V}$ are the nodes corresponding to the words in the sentence and $\mathbf{E}$ are the directed edges representing dependencies. The transformer model computes a score for each possible edge $e_{ij}$ from word $x_i$ to word $x_j$:

$$s(e_{ij} \mid \mathbf{x}; \theta) = \mathbf{W}_e^\top [h_i; h_j],$$

where $[h_i; h_j]$ is the concatenation of the embeddings of $x_i$ and $x_j$, and $\mathbf{W}_e$ is a weight vector.

The optimal dependency tree $\mathcal{T}^*$ maximizes the sum of the edge scores, subject to the constraint that the resulting structure is a valid tree:

$$\mathcal{T}^* = \arg \max_{\mathcal{T} \in \mathcal{F}(\mathbf{x})} \sum_{(i,j) \in \mathbf{E}} s(e_{ij} \mid \mathbf{x}; \theta),$$

where $\mathcal{F}(\mathbf{x})$ is the set of all valid dependency trees for the sentence $\mathbf{x}$.

Example: In the sentence "The cat chased the mouse," a dependency parser would identify "chased" as the root verb, with "The cat" and "the mouse" as its subjects and objects, respectively.

**Constituency Parsing**

Constituency parsing involves breaking down a sentence into its constituent parts, such as phrases and clauses, which are represented as nodes in a tree. Each node represents a linguistic unit, and the tree structure captures the hierarchical relationships between these units.

Given a sentence $\mathbf{x} = (x_1, x_2, \ldots, x_n)$, the goal is to construct a constituency tree $\mathcal{C} = (\mathbf{V}, \mathbf{E})$, where $\mathbf{V}$ are the nodes representing constituents (e.g., noun phrases, verb phrases) and $\mathbf{E}$ are the edges representing the hierarchical structure. The transformer model computes a score for each possible split point $k$ that divides a span $[i, j]$ into two sub-constituents:

$$s([i, j] \to [i, k], [k+1, j] \mid \mathbf{x}; \theta) = \mathbf{W}_c^\top [h_{[i,k]}; h_{[k+1,j]}],$$

where $[h_{[i,k]}; h_{[k+1,j]}]$ is the concatenation of the embeddings representing the two sub-constituents and $\mathbf{W}_c$ is a weight vector.

The optimal constituency tree $\mathcal{C}^*$ maximizes the sum of the scores of all the splits, subject to the constraint that the resulting structure is a valid binary tree:

$$\mathcal{C}^* = \arg \max_{\mathcal{C} \in \mathcal{F}(\mathbf{x})} \sum_{(i,j) \in \mathbf{E}} s([i, j] \to [i, k], [k+1, j] \mid \mathbf{x}; \theta),$$

where $\mathcal{F}(\mathbf{x})$ is the set of all valid constituency trees for the sentence $\mathbf{x}$.

Example: In the sentence "The quick brown fox jumps over the lazy dog," a constituency parser would identify "The quick brown fox" as a noun phrase and "jumps over the lazy dog" as a verb phrase, capturing the hierarchical structure of the sentence.

## 5.3.4   *Performance Analysis*

Performance analysis of transformer models involves evaluating how model size, computational efficiency, training time, and resource utilization impact the overall effectiveness of the model in NLP tasks.

**Model Size and Computational Efficiency**

Transformer models can vary significantly in size, from smaller models like BERT-Base to massive models like GPT-3. The size of the model influences both its computational requirements and its ability to capture complex patterns in the data.

Let $\mathcal{M}$ be the number of parameters in the model, and let $\mathcal{C}(\mathcal{M})$ represent the computational complexity of the model, typically measured in terms of FLOPs (floating-point operations):

$$\mathcal{C}(\mathcal{M}) = \mathcal{O}(\mathcal{M}^2),$$

where $\mathcal{M}$ is often proportional to the number of layers $L$, the number of heads $H$, and the dimensionality $d$ of the embeddings.

As model size increases, the computational complexity grows quadratically, leading to increased training time and resource requirements. However, larger models often achieve better performance due to their increased capacity to capture complex relationships in the data:

$$\text{Performance} \propto \log(\mathcal{M}),$$

indicating diminishing returns as model size increases.

Example: While GPT-3, with its 175 billion parameters, achieves state-of-the-art performance across many tasks, it also requires enormous computational resources, making it less accessible for smaller scale applications.

**Training Time and Resource Utilization**

The training time and resource utilization of transformer models are critical factors that influence their practical deployment, especially in environments with limited computational resources.

Let $T$ represent the total training time, and let $R$ denote the resources used, such as GPU hours. The training time is a function of model size $\mathcal{M}$, dataset size $\mathcal{D}$, and the batch size $B$:

$$T \propto \frac{\mathcal{M} \cdot \mathcal{D}}{B \cdot \text{Throughput}},$$

where throughput represents the number of examples processed per second.

Optimizing training time involves balancing batch size, learning rate, and resource allocation to minimize $T$ while maximizing model performance. Techniques such as mixed-precision training and distributed training can be used to improve resource utilization:

$$\text{Resource Efficiency} = \frac{\text{Performance}}{T \cdot R}.$$

Training a large transformer model like BERT-Large on a massive dataset may require several days on a high-performance GPU cluster. Techniques like gradient accumulation and model parallelism can be used to reduce the effective training time without compromising model performance.

## 5.4   Future Directions in NLP with Transformers

As transformer models continue to evolve, several key areas of research are emerging that promise to extend their capabilities and address existing challenges. This section explores future directions in NLP with transformers, focusing on scalability, efficiency, multilingual and cross-lingual models, and ethical considerations. Each subsection delves into the mathematical principles that underpin these developments, emphasizing the themes of Geometry, Symmetry, and Intelligence.

### 5.4.1   Scalability and Efficiency

The scalability and efficiency of transformer models are critical as they grow in size and complexity. Researchers are exploring ways to scale transformers to larger datasets and more diverse tasks while maintaining computational efficiency and reducing the resource burden.

Let $\mathcal{M}$ represent the model size, $\mathcal{D}$ the dataset size, and $\mathcal{C}(\mathcal{M}, \mathcal{D})$ the computational complexity required to train the model. The goal is to optimize this complexity while preserving or enhancing model performance. A key challenge is balancing the number of parameters with the model's ability to generalize

$$\mathcal{C}(\mathcal{M}, \mathcal{D}) = \mathcal{O}(\mathcal{M} \cdot \mathcal{D} \cdot \log \mathcal{M}),$$

where the logarithmic factor accounts for optimization techniques such as sparsity and attention pruning.

As transformer models scale, there is a trade-off between the number of parameters $\mathcal{M}$ and the ability to efficiently utilize computational resources. Theoretical results suggest that beyond a certain point, increasing $\mathcal{M}$ yields diminishing returns in performance unless accompanied by innovations in model architecture and training techniques:

$$\lim_{\mathcal{M} \to \infty} \frac{\text{Performance}(\mathcal{M})}{\mathcal{C}(\mathcal{M}, \mathcal{D})} \approx \text{constant}.$$

Research into scalable transformers includes techniques such as sparse attention, where only a subset of attention weights are computed, and model distillation, which compresses larger models into smaller, more efficient versions without significant loss in performance.

## *5.4.2   Multilingual and Cross-lingual Models*

Multilingual and cross-lingual models aim to extend the capabilities of transformers to understand and generate text across multiple languages, facilitating language-agnostic processing and improving performance in low-resource languages.

### Language-Agnostic Models

Language-agnostic models are designed to perform equally well across different languages without relying on language-specific features. This requires the model to learn representations that are invariant to the linguistic characteristics of individual languages.

Given a set of languages $\mathcal{L} = \{L_1, L_2, \ldots, L_k\}$, the goal is to learn a shared representation $h_i$ for each token $x_i$ in a sequence $\mathbf{x}$ such that the representation is invariant to the language:

$$h_i^{L_j} = f_\theta(x_i^{L_j}), \quad \forall L_j \in \mathcal{L},$$

where $f_\theta$ is the shared model architecture parameterized by $\theta$.

The success of language-agnostic models depends on their ability to learn invariant representations that generalize across languages. This invariance can be quantified by measuring the similarity between representations of equivalent sentences across languages:

$$\mathbb{E}_{(x_i^{L_j}, x_i^{L_k}) \in \mathcal{D}} \left[ \| h_i^{L_j} - h_i^{L_k} \|^2 \right] \approx 0, \quad \forall L_j, L_k \in \mathcal{L}.$$

Example: XLM-R (Cross-lingual Language Model—RoBERTa) ([1]) is an example of a multilingual transformer that leverages shared subword representations to achieve strong performance across multiple languages, including low-resource ones.

### Zero-shot and Few-shot Learning

Zero-shot and few-shot learning enable transformer models to perform tasks in new languages or domains with little to no task-specific training data. This is particularly valuable in multilingual contexts where labeled data is scarce for certain languages.

In zero-shot learning, the model is trained on a set of tasks $\mathcal{T}_{\text{train}}$ and is then evaluated on a new task $\mathcal{T}_{\text{test}}$ without any additional training. The goal is to learn a generalizable function $f_\theta$ that can perform well on unseen tasks:

$$\min_\theta \mathbb{E}_{\mathcal{T} \sim \mathcal{T}_{\text{train}}} \left[ \mathcal{L}_{\mathcal{T}}(f_\theta) \right], \quad \text{subject to} \quad \mathbb{E}_{\mathcal{T} \sim \mathcal{T}_{\text{test}}} \left[ \mathcal{L}_{\mathcal{T}}(f_\theta) \right] \leq \epsilon.$$

The ability of a model to generalize in zero-shot scenarios is bounded by the diversity of the tasks in $\mathcal{T}_{\text{train}}$ and the similarity between tasks in $\mathcal{T}_{\text{train}}$ and $\mathcal{T}_{\text{test}}$:

$$\text{Gen. Error} \leq \mathcal{O}\left(\sqrt{\frac{\log |\mathcal{T}_{\text{train}}| + \log(1/\delta)}{N_{\text{train}}}}\right),$$

where $N_{\text{train}}$ is the number of training samples and $\delta$ is the confidence level.

Example: GPT-3 demonstrates remarkable zero-shot learning capabilities, where it can perform tasks like translation or summarization in languages it has not explicitly been trained on, by leveraging the broad knowledge encoded during its pre-training.

### 5.4.3  Ethical Considerations

As transformer models become more integrated into real-world applications, ethical considerations such as bias, fairness, privacy, and security are paramount. Addressing these concerns requires analysis and the development of new methodologies to ensure that models are both effective and responsible.

#### Bias and Fairness

Bias in transformer models can arise from the training data, leading to unfair treatment of certain groups or propagation of harmful stereotypes. Ensuring fairness involves both detecting and mitigating bias through careful model design and evaluation.

Let $\mathcal{D}$ be the dataset used to train the model, and let $f_\theta$ be the model. Bias can be quantified by measuring the disparity in model predictions across different subgroups $\mathcal{G}_i$ within the population:

$$\text{Bias} = \max_{i,j} \left| \mathbb{E}_{x \sim \mathcal{G}_i} \left[ f_\theta(x) \right] - \mathbb{E}_{x \sim \mathcal{G}_j} \left[ f_\theta(x) \right] \right|.$$

A model is considered fair if the disparity in predictions across subgroups is minimized. This can be enforced through regularization techniques that penalize large differences in subgroup performance:

$$\mathcal{L}_{\text{fair}}(\theta) = \mathcal{L}_{\text{task}}(\theta) + \lambda \cdot \text{Bias}(\theta),$$

where $\lambda$ is a regularization parameter controlling the trade-off between task performance and fairness.

Example: Gender bias in language models can manifest in tasks like sentiment analysis or word association, where certain words may be unfairly associated with gendered terms. Techniques like adversarial training or bias regularization can be employed to mitigate such biases.

**Privacy and Security**

Privacy and security concerns arise when transformer models are trained on sensitive data or are deployed in environments where data confidentiality is critical. Ensuring that models do not leak private information or are vulnerable to adversarial attacks is essential.

Privacy in machine learning can be formalized using differential privacy, where the goal is to ensure that the inclusion or exclusion of a single data point $x_i$ in the dataset $\mathcal{D}$ has a limited impact on the model's output:

$$\mathbb{P}[f_\theta(\mathcal{D}) = y] \le e^\epsilon \cdot \mathbb{P}[f_\theta(\mathcal{D} \setminus \{x_i\}) = y] + \delta,$$

where $\epsilon$ is the privacy budget and $\delta$ is a small probability.

A transformer model can be made differentially private by adding noise to the gradients during training or through post-processing of the model outputs. The level of privacy is controlled by the parameter $\epsilon$, with smaller values indicating stronger privacy guarantees:

$$\text{Private Loss}(\theta) = \mathcal{L}_{\text{task}}(\theta) + \frac{1}{2\epsilon^2} \|\theta\|_2^2.$$

Example: In scenarios where transformer models are deployed in healthcare or finance, ensuring that the models are differentially private helps protect sensitive information while still allowing the models to learn useful patterns from the data.

# References

1. Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., Stoyanov, V.: Unsupervised cross-lingual representation learning at scale. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 8440–8451 (2020). https://doi.org/10.18653/v1/2020.acl-main.747. https://aclanthology.org/2020.acl-main.747/
2. Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: North American Chapter of the Association for Computational Linguistics (2019)
3. Radford, A., Narasimhan, K.: Improving language understanding by generative pre-training. OpenAI preprint (2018)
4. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R., Le, Q.V.: Xlnet: Generalized autoregressive pretraining for language understanding. In: Advances in Neural Information Processing Systems, vol. 32 (2019)

# Chapter 6
# Transformers in Computer Vision

## 6.1 Vision Transformers (ViT)

Vision transformers ([2, 5]) have emerged as a powerful alternative to CNNs for image classification and other computer vision tasks. The key idea behind ViT is to treat an image as a sequence of patches, analogous to how text is treated as a sequence of words in natural language processing. This section explores the mathematical formulation of ViT, focusing on the process of embedding image patches into a sequence of vectors that can be processed by the transformer model.

### 6.1.1 Mathematical Formulation

The vision transformer model operates by first converting an input image into a sequence of patches, each of which is then embedded into a vector representation. This sequence of patch embeddings is then fed into a standard transformer model, which processes the sequence and produces a final classification or other task-specific output. The key steps in this process involve image patch embedding, forming patch sequences and constructing the patch embedding layer.

**Image Patch Embedding**

Given an input image $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$, where $H$ is the height, $W$ is the width, and $C$ is the number of channels (e.g., 3 for RGB images), the image is divided into a grid of non-overlapping patches. Each patch $\mathbf{P}$ is of size $P \times P$ pixels, resulting in $N = \frac{H \times W}{P^2}$ patches.

The patch $\mathbf{P}_{i,j}$ at position $(i, j)$ within the image is a sub-matrix of the original image:

$$\mathbf{P}_{i,j} = \mathbf{I}[iP : (i+1)P, jP : (j+1)P, :],$$

where $i$ and $j$ index the vertical and horizontal positions of the patch in the grid.

The dimensionality of each patch $\mathbf{P}_{i,j}$ is $P^2 \times C$. This dimensionality is the number of pixels in the patch multiplied by the number of channels. The patches are typically flattened into vectors:

$$\mathbf{p}_{i,j} = \text{flatten}(\mathbf{P}_{i,j}) \in \mathbb{R}^{P^2 \cdot C}.$$

Example: For an image of size $224 \times 224 \times 3$ and patch size $16 \times 16$, there are $N = \frac{224 \times 224}{16 \times 16} = 196$ patches, each flattened into a vector of length $16 \times 16 \times 3 = 768$.

**Forming Patch Sequences**

Once the image is divided into patches and each patch is flattened into a vector, these vectors are treated as tokens in a sequence, analogous to the words in a sentence processed by a transformer.

The sequence of patch vectors $\mathbf{p}_{i,j}$ is concatenated to form the input sequence $\mathbf{X}$ to the transformer:

$$\mathbf{X} = [\mathbf{p}_1, \mathbf{p}_2, \ldots, \mathbf{p}_N] \in \mathbb{R}^{N \times (P^2 \cdot C)},$$

where each $\mathbf{p}_k$ corresponds to a flattened patch $\mathbf{p}_{i,j}$.

The length of the sequence is $N$, the number of patches, and the dimensionality of each sequence element (token) is $P^2 \cdot C$. This sequence is then input into the transformer model, which expects a sequence of vectors as input.

Example: Continuing with the example of an image with 196 patches, the input sequence $\mathbf{X}$ to the transformer has a length of 196, with each element being a vector of size 768.

**Patch Embedding Layer**

The flattened patch vectors $\mathbf{p}_k$ are then linearly transformed into a lower-dimensional space suitable for processing by the transformer. This is done by the patch embedding layer, which applies a learnable linear transformation to each patch vector.

Let $\mathbf{W}_{\text{emb}} \in \mathbb{R}^{(P^2 \cdot C) \times D}$ be the weight matrix of the patch embedding layer, where $D$ is the dimensionality of the transformer's input vectors. The embedding for each patch $\mathbf{p}_k$ is computed as

$$\mathbf{e}_k = \mathbf{p}_k \mathbf{W}_{\text{emb}} + \mathbf{b}_{\text{emb}},$$

where $\mathbf{b}_{\text{emb}} \in \mathbb{R}^D$ is the bias term. The sequence of embeddings $\mathbf{E}$ is then

$$\mathbf{E} = [\mathbf{e}_1, \mathbf{e}_2, \ldots, \mathbf{e}_N] \in \mathbb{R}^{N \times D}.$$

The patch embedding layer reduces the dimensionality of each patch vector from $P^2 \cdot C$ to $D$. The choice of $D$ is typically smaller than $P^2 \cdot C$, allowing the model to focus on the most relevant features of each patch.

Example: If the original patch vectors are of size 768 and the transformer input dimensionality $D$ is 512, then the patch embedding layer reduces the dimensionality of each patch to 512, resulting in an embedding sequence $\mathbf{E}$ of size $196 \times 512$.

**Positional Encoding for Images**

Let $\mathbf{E} = [\mathbf{e}_1, \mathbf{e}_2, \ldots, \mathbf{e}_N]$ be the sequence of patch embeddings, where $N$ is the number of patches. Positional encoding adds a vector $\mathbf{p}_i$ to each embedding $\mathbf{e}_i$ to incorporate the positional information. The modified embedding sequence $\mathbf{E}'$ is given by

$$\mathbf{E}' = [\mathbf{e}_1 + \mathbf{p}_1, \mathbf{e}_2 + \mathbf{p}_2, \ldots, \mathbf{e}_N + \mathbf{p}_N],$$

where $\mathbf{p}_i$ is the positional encoding for the $i$-th patch.

The addition of positional encodings ensures that the transformer model can distinguish between different spatial positions of patches, effectively introducing spatial order into the otherwise permutation-invariant architecture. Formally, the positional encoding function $\mathbf{P}$ maps each position $i$ to a unique vector $\mathbf{p}_i$:

$$\mathbf{P} : i \mapsto \mathbf{p}_i \in \mathbb{R}^D,$$

where $D$ is the dimensionality of the embedding space. This mapping preserves the spatial structure of the image within the transformer model.

Example: In an image divided into $14 \times 14$ patches, the positional encoding ensures that patches corresponding to adjacent image regions (e.g., the top-left corner or the bottom-right corner) are recognized as being near each other spatially.

**Sinusoidal Positional Encoding**

Sinusoidal positional encoding is a fixed, deterministic approach that encodes position information using sine and cosine functions of different frequencies. This encoding method is designed to be simple yet effective, capturing relative positions through periodic functions.

For a patch at position $i$ in the sequence, the positional encoding vector $\mathbf{p}_i$ is defined as

$$\mathbf{p}_i[2k] = \sin\left(\frac{i}{10000^{2k/D}}\right), \quad \mathbf{p}_i[2k+1] = \cos\left(\frac{i}{10000^{2k/D}}\right),$$

where $k = 0, 1, \ldots, D/2 - 1$ and $D$ is the dimensionality of the embedding space. This formulation ensures that each dimension of the positional encoding vector captures a different aspect of the position through a unique frequency component.

The sinusoidal positional encoding allows the model to distinguish positions through unique combinations of sine and cosine functions. The periodic nature of the functions introduces a smooth, continuous notion of position that generalizes well to sequences of varying lengths.

Example: For an image patch sequence, the sinusoidal positional encoding will assign unique vectors to each patch, even for images of different sizes, allowing the transformer to consistently interpret spatial relationships.

**Learned Positional Encoding**

Unlike sinusoidal positional encodings, learned positional encodings are parameters that are learned during the training process. Each position $i$ is associated with a learnable vector $\mathbf{p}_i$, which is optimized along with the other model parameters.

Let $\mathbf{P} \in \mathbb{R}^{N \times D}$ be the matrix of positional encodings, where each row $\mathbf{p}_i$ corresponds to the positional encoding of patch $i$. These encodings are initialized randomly and updated through backpropagation:

$$\mathbf{P} \leftarrow \mathbf{P} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{P}},$$

where $\mathcal{L}$ is the loss function and $\eta$ is the learning rate. The learned positional encodings thus adapt to the specific task and dataset.

Learned positional encodings offer greater flexibility compared to sinusoidal encodings, as they can adapt to specific patterns in the data. However, this flexibility comes at the cost of requiring additional parameters and the risk of overfitting, particularly in small datasets.

Example: In large-scale vision tasks, learned positional encodings can capture complex spatial relationships that might be specific to certain datasets, such as the relative positioning of objects in images, thereby improving model performance.

## 6.1.2  ViT Architecture

The vision transformer architecture adapts the standard transformer model, originally designed for sequence-based tasks like natural language processing, to handle visual data. The key components of the ViT architecture include the input embedding layer (which handles image patches), the transformer encoder layers, and the final classification head. Each component plays a crucial role in the model's ability to capture and process visual information effectively.

Let $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$ be an input image, where $H$ is the height, $W$ is the width, and $C$ is the number of channels (e.g., 3 for RGB images). The ViT architecture processes this image as follows:

1. Patch Embedding: The image $\mathbf{I}$ is divided into $N$ non-overlapping patches, each of size $P \times P$. Each patch is then flattened and linearly transformed into an embedding vector. Mathematically, this step is represented as

$$\mathbf{X} = [\mathbf{e}_1, \mathbf{e}_2, \ldots, \mathbf{e}_N],$$

where $\mathbf{e}_i \in \mathbb{R}^D$ is the embedding vector for the $i$-th patch and $D$ is the dimensionality of the embedding space. The transformation from the flattened patch to the embedding vector is a linear mapping:

$$\mathbf{e}_i = \mathbf{P}_i \mathbf{W}_{\text{emb}} + \mathbf{b}_{\text{emb}},$$

where $\mathbf{W}_{\text{emb}} \in \mathbb{R}^{(P^2 \cdot C) \times D}$ and $\mathbf{b}_{\text{emb}} \in \mathbb{R}^D$ ensure that each patch is mapped consistently to the same embedding space.

2. Addition of Positional Encodings: To incorporate spatial information, positional encodings $\mathbf{P}_i$ are added to each patch embedding. This modifies the input sequence to

$$\mathbf{X}' = [\mathbf{e}_1 + \mathbf{p}_1, \mathbf{e}_2 + \mathbf{p}_2, \ldots, \mathbf{e}_N + \mathbf{p}_N],$$

where $\mathbf{p}_i \in \mathbb{R}^D$ represents the positional encoding for the $i$-th patch. The addition of positional encodings ensures that the transformer model can distinguish between patches based on their original positions within the image, preserving spatial relationships during the encoding process.

3. Transformer Encoder Layers: The sequence of patch embeddings with positional encodings is then processed by a series of transformer encoder layers. Each encoder layer consists of multi-head self-attention (MHSA) and feedforward networks (FFNs), followed by layer normalization:

$$\mathbf{Z}^{(l)} = \text{LayerNorm}\left(\mathbf{X}^{(l)} + \text{MHSA}\left(\mathbf{X}^{(l)}\right)\right),$$

$$\mathbf{X}^{(l+1)} = \text{LayerNorm}\left(\mathbf{Z}^{(l)} + \text{FFN}\left(\mathbf{Z}^{(l)}\right)\right),$$

where $\mathbf{X}^{(l)}$ is the input to the $l$-th encoder layer and $\mathbf{X}^{(l+1)}$ is the output. Each encoder layer refines the representations of the patch embeddings by aggregating information across the entire sequence through self-attention and by transforming these representations via feedforward networks. This iterative refinement process allows the model to capture complex, hierarchical features from the image.

4. Classification Head: After passing through multiple transformer encoder layers, the final sequence representation is summarized into a single vector, typically by selecting the representation corresponding to a special "classification token" $\mathbf{z}_{\text{cls}}$. This

vector is then passed through a fully connected layer to produce the final classification logits:

$$\mathbf{y} = \mathrm{softmax}\left(\mathbf{W}_{\mathrm{cls}}\mathbf{z}_{\mathrm{cls}} + \mathbf{b}_{\mathrm{cls}}\right),$$

where $\mathbf{W}_{\mathrm{cls}}$ and $\mathbf{b}_{\mathrm{cls}}$ are the weights and biases of the classification head, and $\mathbf{y}$ represents the predicted probabilities for each class. The classification head ensures that the entire image representation is mapped into a discrete probability distribution over the target classes. This mapping is consistent with the global structure captured by the transformer encoders.

## 6.2  Applications in Computer Vision

ViTs have shown significant potential in various computer vision tasks, including image classification and object detection. This section explores the application of ViTs in these areas, focusing on the mathematical foundations of dataset preparation, model training, and evaluation metrics. By understanding these aspects, we gain deeper insights into the practical deployment of ViTs in real-world vision tasks.

### 6.2.1  Image Classification

Image classification is one of the primary applications of vision transformers. The process involves preparing the dataset, training the model, and evaluating its performance using various metrics. This subsection delves into the mathematical details of each step, ensuring a thorough understanding of how ViTs are applied to classification tasks.

#### Dataset Preparation

Dataset preparation is the foundation of any successful image classification task. It involves curating and preprocessing the data to ensure that it is suitable for training a vision transformer model.

Let $\mathcal{D} = \{(\mathbf{I}_i, y_i)\}_{i=1}^{M}$ represent the dataset, where $\mathbf{I}_i \in \mathbb{R}^{H \times W \times C}$ is an image, $y_i \in \{1, 2, \ldots, K\}$ is the corresponding label, and $M$ is the total number of samples. The dataset preparation involves the following steps:

1. Each image $\mathbf{I}_i$ is normalized to have zero mean and unit variance, ensuring that the pixel values are centered and scaled appropriately:

$$\mathbf{I}_i' = \frac{\mathbf{I}_i - \mu}{\sigma},$$

where $\mu$ and $\sigma$ are the mean and standard deviation of the pixel values, computed across the entire dataset.

2. Data augmentation is applied to increase the diversity of the training data by creating transformed versions of the original images. Common transformations include rotations, flips, and color jittering:

$$\mathcal{T}(\mathbf{I}_i) = \text{Augment}(\mathbf{I}_i),$$

where $\mathcal{T}$ represents the augmentation transformation applied to the image $\mathbf{I}_i$.

Data augmentation increases the effective size of the training dataset, reducing overfitting and improving the model's ability to generalize to unseen data. Formally, let $\mathcal{D}_{\text{aug}}$ represent the augmented dataset. Then, the generalization error $\mathcal{E}_{\text{gen}}$ of the model trained on $\mathcal{D}_{\text{aug}}$ is expected to be lower than that of a model trained on the original dataset $\mathcal{D}$:

$$\mathcal{E}_{\text{gen}}(\mathcal{D}_{\text{aug}}) < \mathcal{E}_{\text{gen}}(\mathcal{D}).$$

Example: For a dataset like CIFAR-10, normalization ensures that all images have consistent pixel value distributions, while data augmentation introduces variations that help the model learn more robust features.

## Training ViT Models

Training a vision transformer involves optimizing the model's parameters using a labeled dataset, typically through a gradient-based optimization method. The process is guided by a loss function that measures the discrepancy between the model's predictions and the true labels.

Let $\theta$ represent the parameters of the vision transformer model. The goal of training is to minimize the loss function $\mathcal{L}(\theta)$, which is typically the cross-entropy loss for classification tasks:

$$\mathcal{L}(\theta) = -\frac{1}{M} \sum_{i=1}^{M} \sum_{k=1}^{K} \mathbf{1}\{y_i = k\} \log P(y_i = k \mid \mathbf{I}_i; \theta),$$

where $P(y_i = k \mid \mathbf{I}_i; \theta)$ is the probability that the model assigns to class $k$ given image $\mathbf{I}_i$, and $\mathbf{1}\{\cdot\}$ is the indicator function.

Optimization: The parameters $\theta$ are updated iteratively using gradient descent:

$$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta),$$

where $\eta$ is the learning rate and $\nabla_\theta \mathcal{L}(\theta)$ is the gradient of the loss function with respect to the model parameters.

Under appropriate conditions on the learning rate $\eta$ and the smoothness of the loss function $\mathcal{L}(\theta)$, gradient descent converges to a local minimum of the loss function, leading to a model that performs well on the training data.

Example: In training a vision transformer on the ImageNet dataset, the cross-entropy loss guides the model to adjust its parameters so that the predicted class probabilities match the true labels as closely as possible.

## Evaluation Metrics (Accuracy, Precision, Recall)

Evaluating the performance of a vision transformer model involves measuring how well the model's predictions align with the true labels. Common evaluation metrics include accuracy, precision, and recall.

Given a test set $\mathcal{D}_{\text{test}} = \{(\mathbf{I}_i, y_i)\}_{i=1}^N$, where $N$ is the number of test samples, the following metrics are defined:

1. Accuracy measures the proportion of correct predictions:

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbf{1}\{\hat{y}_i = y_i\},$$

where $\hat{y}_i$ is the predicted label for image $\mathbf{I}_i$.

2. Precision measures the proportion of correctly predicted positive instances among all predicted positive instances for a given class $k$:

$$\text{Precision}(k) = \frac{\sum_{i=1}^N \mathbf{1}\{\hat{y}_i = k \text{ and } y_i = k\}}{\sum_{i=1}^N \mathbf{1}\{\hat{y}_i = k\}}.$$

3. Recall measures the proportion of correctly predicted positive instances among all actual positive instances for a given class $k$:

$$\text{Recall}(k) = \frac{\sum_{i=1}^N \mathbf{1}\{\hat{y}_i = k \text{ and } y_i = k\}}{\sum_{i=1}^N \mathbf{1}\{y_i = k\}}.$$

There is often a trade-off between precision and recall. High precision may come at the cost of lower recall and vice versa. The F1-score is a metric that combines both:

$$\text{F1-score}(k) = 2 \cdot \frac{\text{Precision}(k) \cdot \text{Recall}(k)}{\text{Precision}(k) + \text{Recall}(k)}.$$

Example: When evaluating a vision transformer on the CIFAR-100 dataset, accuracy gives a general measure of performance, while precision and recall provide more detailed insights into how well the model distinguishes between different classes.

## *6.2.2   Object Detection*

Object detection extends image classification by not only identifying the class of objects in an image but also localizing them within the image using bounding boxes. Applying vision transformers to object detection requires adapting the model to output both class labels and bounding box coordinates.

Let $\mathbf{I}$ be an input image containing $m$ objects, each represented by a class label $y_j$ and a bounding box $\mathbf{b}_j = (x_j, y_j, w_j, h_j)$, where $(x_j, y_j)$ is the center of the bounding box, and $(w_j, h_j)$ are the width and height. The goal of object detection is to predict these labels and bounding boxes.

The vision transformer can be adapted to object detection by introducing additional prediction heads that output the bounding box coordinates and class labels for each detected object:

$$\mathbf{z}_{\text{det}} = \text{MHSA}(\mathbf{X}') + \text{FFN}(\mathbf{z}_{\text{cls}}),$$

where $\mathbf{z}_{\text{det}}$ is the detection token output, which is then passed through separate heads for class prediction and bounding box regression.

The loss function for object detection typically combines a classification loss (e.g., cross-entropy) with a bounding box regression loss (e.g., smooth $L_1$):

$$\mathcal{L}_{\text{det}}(\theta) = \mathcal{L}_{\text{cls}}(\theta) + \lambda \mathcal{L}_{\text{bbox}}(\theta),$$

where $\lambda$ balances the two components of the loss.

Bounding box regression aims to minimize the difference between the predicted and true bounding box coordinates. The smooth $L_1$ loss for bounding box regression is defined as

$$\mathcal{L}_{\text{bbox}}(\theta) = \sum_{j=1}^{m} \text{smooth}_{L_1}(\mathbf{b}_j, \hat{\mathbf{b}}_j),$$

where $\hat{\mathbf{b}}_j$ are the predicted bounding box coordinates.

Example: In tasks like COCO object detection, vision transformers can be trained to accurately detect and localize multiple objects within an image, outputting both class predictions and bounding box coordinates for each object.

**Transformer-Based Detection Models (DETR)**

Detection Transformers (DETR) ( [1]) represent a novel approach to object detection that leverages the power of transformers to directly predict object bounding boxes and class labels from an input image. Unlike traditional detection models that rely on region proposal networks and non-maximum suppression, DETR uses a set-based global loss to handle object detection as a direct set prediction problem.

Let $\mathbf{I}$ be an input image, and let $\{(\mathbf{b}_i, y_i)\}_{i=1}^{N}$ be the set of $N$ predicted objects, where $\mathbf{b}_i = (x_i, y_i, w_i, h_i)$ is the bounding box and $y_i$ is the class label for the $i$-th

object. The DETR model consists of an encoder–decoder transformer architecture, where the encoder processes the image features and the decoder generates object predictions.

1. Object Queries: The decoder operates on a fixed set of learned object queries $\mathbf{Q} \in \mathbb{R}^{N \times D}$, where $N$ is the number of objects to predict, and $D$ is the dimensionality of the query embeddings. Each query corresponds to a potential object in the image.

$$\mathbf{Q}' = \text{Decoder}(\mathbf{Q}, \mathbf{E}),$$

where $\mathbf{E}$ represents the encoded image features.

2. Set-based Prediction: DETR directly predicts a set of bounding boxes and class labels from the decoder outputs. The predictions $\{(\hat{\mathbf{b}}_i, \hat{y}_i)\}_{i=1}^N$ are matched to the ground truth objects $\{(\mathbf{b}_j, y_j)\}_{j=1}^M$ using a bipartite matching loss, ensuring a one-to-one correspondence:

$$\mathcal{L}_{\text{DETR}}(\hat{\mathbf{b}}, \hat{y}; \mathbf{b}, y) = \sum_{i=1}^N \left[ \mathcal{L}_{\text{cls}}(\hat{y}_i, y_{\sigma(i)}) + \mathcal{L}_{\text{bbox}}(\hat{\mathbf{b}}_i, \mathbf{b}_{\sigma(i)}) \right],$$

where $\sigma$ is the optimal assignment that minimizes the total loss, $\mathcal{L}_{\text{cls}}$ is the classification loss, and $\mathcal{L}_{\text{bbox}}$ is the bounding box regression loss.

The optimal assignment $\sigma$ is found by minimizing the total matching cost, ensuring that each predicted object is uniquely paired with a ground truth object. The bipartite matching ensures that the model predicts the correct number of objects and accurately locates them within the image.

Example: In the COCO dataset ([4]), DETR effectively predicts bounding boxes and class labels for multiple objects within an image, handling complex scenes with overlapping objects without relying on traditional post-processing techniques.

**Evaluation Metrics (mAP, IOU)**

Evaluating the performance of object detection models like DETR involves metrics that quantify both the localization accuracy and the classification correctness of the detected objects.

1. Intersection over Union (IoU): The IoU metric measures the overlap between the predicted bounding box $\hat{\mathbf{b}}$ and the ground truth bounding box $\mathbf{b}$. It is defined as

$$\text{IoU}(\hat{\mathbf{b}}, \mathbf{b}) = \frac{\text{Area}(\hat{\mathbf{b}} \cap \mathbf{b})}{\text{Area}(\hat{\mathbf{b}} \cup \mathbf{b})}.$$

IoU is used to determine whether a predicted bounding box is a true positive (IoU above a threshold, typically 0.5) or a false positive.

2. Mean Average Precision (mAP): The mAP metric summarizes the precision–recall curve for all classes. It is computed as the mean of the average precision (AP) scores across all classes:

$$\text{mAP} = \frac{1}{K} \sum_{k=1}^{K} \text{AP}(k),$$

where $K$ is the number of classes and $\text{AP}(k)$ is the area under the precision–recall curve for class $k$.

The mAP metric is sensitive to the IoU threshold. Higher IoU thresholds typically lead to lower mAP scores, as they require more precise localization of objects.

Example: When evaluating DETR on a dataset like Pascal VOC, mAP at different IoU thresholds (e.g., 0.5, 0.75) provides insights into the model's ability to accurately detect and localize objects across various classes.

### 6.2.3  Image Generation

Beyond object detection, vision transformers can also be applied to image generation tasks, often in combination with Generative Adversarial Networks (GANs) ([3]). Integrating transformers with GANs leverages the transformer's ability to model long-range dependencies and the GAN's capability to generate high-quality images.

**Generative Adversarial Networks (GANs) with Transformers**

In a GAN, two neural networks—the generator $G$ and the discriminator $D$—are trained simultaneously. The generator aims to produce realistic images from random noise, while the discriminator attempts to distinguish between real and generated images. When combined with transformers, the architecture can effectively model complex dependencies in image data.

1. Generator: The generator network $G$ maps a random noise vector $\mathbf{z} \sim p(\mathbf{z})$ to a generated image $\hat{\mathbf{I}} = G(\mathbf{z})$. When using transformers, the generator can incorporate self-attention mechanisms to model global image structures.

$$\hat{\mathbf{I}} = \text{TransformerDecoder}(\mathbf{z}, \mathbf{E}),$$

where $\mathbf{E}$ represents the learned embeddings from an encoder.

2. Discriminator: The discriminator network $D$ assigns a probability $D(\mathbf{I})$ to an image $\mathbf{I}$, indicating whether it is real or generated. The discriminator's goal is to maximize the probability of correctly identifying real images while minimizing the probability for generated ones.

$$\mathcal{L}_D = -\mathbb{E}_{\mathbf{I} \sim p_{\text{data}}}[\log D(\mathbf{I})] - \mathbb{E}_{\hat{\mathbf{I}} \sim p_G}[\log(1 - D(\hat{\mathbf{I}}))].$$

The training of GANs is formulated as a minimax game between the generator and discriminator:

$$\min_G \max_D \mathcal{L}_D(G, D).$$

The equilibrium of this game corresponds to the generator producing images that are indistinguishable from real images by the discriminator.

Example: Transformers integrated with GANs have been applied to tasks like image super-resolution and text-to-image generation, where the transformer captures complex patterns and dependencies, leading to more realistic and coherent images.

## ViT for Image Synthesis

Image synthesis involves generating new images from scratch or from certain inputs, such as noise vectors or textual descriptions. Vision transformers, with their ability to model global dependencies and capture complex patterns, can be effectively utilized for this purpose, particularly in generating high-resolution images with coherent structures.

Let $\mathbf{z} \in \mathbb{R}^d$ be a random noise vector sampled from a distribution $p(\mathbf{z})$ (typically Gaussian), which serves as the input to the vision transformer-based generator. The generator $G_{\mathrm{ViT}}$ maps this noise vector to a synthesized image $\hat{\mathbf{I}}$:

$$\hat{\mathbf{I}} = G_{\mathrm{ViT}}(\mathbf{z}),$$

where $G_{\mathrm{ViT}}$ is constructed using a transformer architecture that processes the noise vector through a series of attention layers, ultimately outputting an image.

1. Patch Embedding in Synthesis: Similar to the ViT architecture for classification, the generator splits the image into patches and processes them as a sequence. However, instead of classifying, the transformer layers are used to refine the latent representation until it resembles a real image:

$$\hat{\mathbf{p}}_i = \mathrm{TransformerLayer}(\mathbf{z}_i),$$

where $\hat{\mathbf{p}}_i$ is the synthesized patch after processing through the transformer layer.

2. Reconstruction of Image: The sequence of patches is then concatenated and reshaped to form the synthesized image:

$$\hat{\mathbf{I}} = \mathrm{Reshape}([\hat{\mathbf{p}}_1, \hat{\mathbf{p}}_2, \ldots, \hat{\mathbf{p}}_N]).$$

The expressivity of transformer-based generators allows them to capture and synthesize complex image structures by modeling the global relationships between different parts of the image. The self-attention mechanism within the transformer layers enables the generator to produce images that are globally coherent and exhibit high-level structural consistency.

Example: A ViT-based generator could be used for text-to-image synthesis, where the input noise vector $\mathbf{z}$ is conditioned on textual descriptions. The transformer layers process this information, capturing the relationships between different elements described in the text and generating a coherent image that matches the description.

**Evaluation Metrics (FID, IS)**

Evaluating the quality of synthesized images is crucial for understanding the performance of image synthesis models. Two commonly used metrics are the Fréchet inception distance and the inception score.

1. Fréchet Inception Distance (FID): The FID metric compares the distribution of synthesized images to the distribution of real images by measuring the distance between their feature representations in a pre-trained inception network. Let $\mu_r$, $\Sigma_r$ and $\mu_g$, $\Sigma_g$ be the means and covariances of the feature representations for real and generated images, respectively. The FID is defined as

$$\text{FID}(\mathcal{X}_r, \mathcal{X}_g) = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2}),$$

where Tr denotes the trace of a matrix. A lower FID score indicates that the distribution of synthesized images is closer to the distribution of real images, implying higher quality and more realistic synthesis.

2. Inception Score (IS): The inception score evaluates the quality of generated images based on the confidence of the inception network's classification. It is computed as

$$\text{IS}(\mathcal{X}_g) = \exp\left(\mathbb{E}_{\mathbf{x} \sim \mathcal{X}_g} D_{\text{KL}}(p(y \mid \mathbf{x}) \| p(y))\right),$$

where $D_{\text{KL}}$ is the Kullback–Leibler divergence, $p(y \mid \mathbf{x})$ is the conditional label distribution given the image $\mathbf{x}$, and $p(y)$ is the marginal distribution over all generated images. A high inception score indicates that the generated images are both diverse and high quality, as the network assigns them to specific classes with high confidence.

Example: In assessing a ViT-based generator trained on the CIFAR-10 dataset, a low FID score would suggest that the generated images are similar to the real CIFAR-10 images, while a high IS would indicate that the images are diverse and correspond well to distinct classes within the dataset.

## 6.2.4   Other Applications

Vision transformers also find applications in other areas of computer vision, such as image segmentation and video analysis, where their ability to model long-range dependencies and capture global context proves advantageous.

**Image Segmentation**

Image segmentation involves partitioning an image into multiple segments, each corresponding to a different object or region. Vision transformers can be adapted for segmentation tasks by predicting a segmentation map instead of class labels.

Given an input image $\mathbf{I}$, the goal is to predict a segmentation map $\mathbf{S}$, where each pixel $\mathbf{S}_{ij}$ belongs to a specific class. The vision transformer processes the image patches and outputs class logits for each patch, which are then upsampled to create a pixel-wise segmentation map:

$$\mathbf{S} = \text{Upsample}([\mathbf{l}_1, \mathbf{l}_2, \ldots, \mathbf{l}_N]),$$

where $\mathbf{l}_i$ represents the logits for the $i$-th patch.

The accuracy of segmentation is typically measured using the Intersection over Union (IoU) for each class, similar to object detection. Higher IoU scores indicate more accurate segmentation.

Example: A vision transformer adapted for semantic segmentation on the Pascal VOC dataset would generate a segmentation map where each pixel is labeled according to its corresponding object or region, such as "sky," "car," or "building."

**Video Analysis**

Video analysis involves processing sequences of frames to extract meaningful information, such as detecting actions, tracking objects, or summarizing content. Vision transformers can be extended to video analysis by processing each frame as a sequence of patches and modeling the temporal dependencies between frames.

Let $\mathbf{V} = \{\mathbf{I}_t\}_{t=1}^{T}$ be a video sequence consisting of $T$ frames. Each frame $\mathbf{I}_t$ is divided into patches, and the sequence of patch embeddings is processed by the vision transformer to capture both spatial and temporal dependencies:

$$\mathbf{Z} = \text{TransformerEncoder}(\mathbf{X}_1, \mathbf{X}_2, \ldots, \mathbf{X}_T),$$

where $\mathbf{X}_t$ is the sequence of patch embeddings for frame $t$.

The self-attention mechanism within the transformer layers allows the model to maintain temporal coherence across frames, capturing the dynamics of objects and actions over time.

Example: In video action recognition, a vision transformer could be used to classify actions such as "running" or "jumping" by analyzing the sequence of frames and identifying patterns that correspond to specific actions.

## 6.3  Mathematical Analysis

The success of ViTs in computer vision tasks is not only due to their powerful attention mechanisms but also their mathematical complexity. This section explores the time and space complexity of ViTs, providing a detailed comparison with traditional CNNs. By understanding the mathematical intricacies, we can appreciate the trade-offs involved in using ViTs and how they scale with different tasks.

**Model Complexity**

Model complexity plays a crucial role in determining the feasibility and efficiency of deploying vision transformers in real-world applications. This involves analyzing both the time and space complexity of the model, which directly impact the computational resources required for training and inference.

**Time Complexity of ViT**

The time complexity of vision transformers is primarily determined by the self-attention mechanism, which is the most computationally expensive component of the model. Understanding this complexity requires a detailed analysis of the operations involved in computing the attention scores and aggregating the values.

Given an input image $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$, the image is divided into $N = \frac{H \times W}{P^2}$ patches, each of size $P \times P$. The time complexity of the self-attention mechanism in the ViT is as follows:

1. Self-Attention Complexity: For each layer, the self-attention mechanism computes attention scores between all pairs of patches. The complexity of computing the attention matrix $\mathbf{A}$ is

$$\text{TimeComplexity}_{\text{Attention}} = O(N^2 \cdot D),$$

where $N$ is the number of patches and $D$ is the dimensionality of the embeddings.

2. Feedforward Network Complexity: Each self-attention layer is followed by a feedforward network (FFN), which consists of two linear transformations and a non-linear activation. The time complexity of the FFN is

$$\text{TimeComplexity}_{\text{FFN}} = O(N \cdot D^2).$$

The total time complexity of a vision transformer with $L$ layers is the sum of the time complexities of the self-attention and FFN across all layers:

$$\text{TimeComplexity}_{\text{ViT}} = O(L \cdot (N^2 \cdot D + N \cdot D^2)).$$

This expression shows that the self-attention mechanism dominates the time complexity, particularly as the number of patches $N$ increases.

Example: For a ViT processing an image of size $224 \times 224$ with patch size $16 \times 16$, $N = 196$ patches. If $D = 768$ and the model has $L = 12$ layers, the time complexity is dominated by the self-attention mechanism, scaling quadratically with the number of patches.

## Space Complexity

Space complexity refers to the memory required to store the model parameters, intermediate activations, and gradients during training. Vision transformers, with their large embedding dimensions and multi-head attention mechanisms, have significant space requirements.

1. Parameter Space Complexity: The parameters of a ViT are primarily in the projection matrices for the query, key, value, and output in the self-attention layers, as well as the weights in the FFN. The space complexity for storing the parameters in a single layer is

$$\text{SpaceComplexity}_{\text{Params}} = O(N \cdot D^2 + D^2).$$

The total parameter space complexity for $L$ layers is

$$\text{SpaceComplexity}_{\text{ParamsTotal}} = O(L \cdot (N \cdot D^2 + D^2)).$$

2. Activation Space Complexity: During training, the model must store the activations for each layer, which are necessary for backpropagation. The space complexity for storing the activations is

$$\text{SpaceComplexity}_{\text{Activations}} = O(L \cdot N \cdot D).$$

The total space complexity of a vision transformer is the sum of the space required for storing the parameters and the activations:

$$\text{SpaceComplexity}_{\text{ViT}} = O(L \cdot (N \cdot D^2 + D^2 + N \cdot D)).$$

This expression indicates that the space complexity scales linearly with the number of layers and patches, making ViTs memory-intensive, particularly for large input images and deep models.

Example: For the same ViT with $N = 196$ patches, $D = 768$, and $L = 12$ layers, the space complexity is dominated by the storage of parameters and activations, requiring substantial memory resources during training.

## Comparison with CNNs

CNNs have been the dominant architecture in computer vision for many years, and comparing the complexity of ViTs with CNNs provides insights into the trade-offs between these architectures.

1. Time Complexity of CNNs: In CNNs, the time complexity is determined by the convolutional operations. For a convolutional layer with a filter size $F \times F$ applied to an input of size $H \times W \times C$, the time complexity is

$$\text{TimeComplexity}_{\text{CNN}} = O(H \cdot W \cdot C \cdot F^2).$$

CNNs typically have linear time complexity with respect to the number of pixels in the input image.

2. Space Complexity of CNNs: The space complexity in CNNs is determined by the number of filters and the size of the activations. For a layer with $K$ filters, the space complexity is

$$\text{SpaceComplexity}_{\text{CNN}} = O(K \cdot F^2 \cdot C + H \cdot W \cdot K).$$

Vision Transformers have higher time complexity due to the quadratic scaling of the self-attention mechanism with respect to the number of patches. However, CNNs typically have lower space complexity due to the locality of convolutions, which reduces the number of parameters and intermediate activations.

Example: Comparing a ViT with a ResNet-50 CNN on the same image size, the ViT may require more computational resources due to the attention mechanism, while the CNN may be more efficient in terms of both time and space complexity, particularly for lower resolution images.

## 6.4 Optimization and Training Strategies

The training and optimization of ViTs are crucial for achieving high performance on computer vision tasks. This section explores the mathematical foundations of loss functions, gradient descent and its variants, regularization techniques, and data augmentation strategies. These components are essential for guiding the learning process and ensuring that the model generalizes well to unseen data.

The effectiveness of vision transformers depends on carefully designed optimization strategies that ensure the model converges to a good solution during training. This involves selecting appropriate loss functions, optimizing the parameters using gradient-based methods, applying regularization to prevent overfitting, and augmenting the training data to improve generalization.

### Loss Functions

The choice of loss function is fundamental to training a vision transformer, as it defines the objective that the model seeks to minimize. In classification tasks, the cross-entropy loss is commonly used, while in regression tasks, mean squared error or similar loss functions may be more appropriate.

Let $\mathcal{D} = \{(\mathbf{I}_i, y_i)\}_{i=1}^M$ be a dataset with $M$ samples, where $\mathbf{I}_i$ is the input image, and $y_i$ is the corresponding label. The model parameters $\theta$ are optimized by minimizing a loss function $\mathcal{L}(\theta)$, defined as

$$\mathcal{L}(\theta) = -\frac{1}{M} \sum_{i=1}^M \sum_{k=1}^K \mathbf{1}\{y_i = k\} \log P(y_i = k \mid \mathbf{I}_i; \theta),$$

where $P(y_i = k \mid \mathbf{I}_i; \theta)$ is the probability that the model assigns to class $k$ given the image $\mathbf{I}_i$ and $\mathbf{1}\{\cdot\}$ is the indicator function.

For binary classification problems, the cross-entropy loss is convex with respect to the model parameters, which ensures that gradient descent will converge to a global minimum if the model is linear. However, for non-linear models like vision transformers, the loss landscape may contain multiple local minima, making optimization more challenging.

Example: For a ViT trained on the CIFAR-10 dataset, the cross-entropy loss guides the model to assign high probabilities to the correct classes for each image, minimizing the overall classification error.

**Regularization Techniques**

Regularization is critical for preventing overfitting, especially in large models like vision transformers, which have a vast number of parameters. Common regularization techniques include weight decay, dropout, and label smoothing.

1. Weight decay (L2 regularization) adds a penalty to the loss function based on the magnitude of the model parameters, encouraging smaller weights and thus simpler models:

$$\mathcal{L}_{\text{reg}}(\theta) = \mathcal{L}(\theta) + \lambda \|\theta\|_2^2,$$

where $\lambda$ is the regularization strength.

2. Dropout randomly sets a fraction of the activations to zero during training, preventing co-adaptation of neurons and improving generalization:

$$\mathbf{h}_{\text{dropout}} = \mathbf{h} \odot \mathbf{m},$$

where $\mathbf{h}$ is the vector of activations, $\mathbf{m}$ is a binary mask with probability $p$ of being zero, and $\odot$ denotes element-wise multiplication.

3. Label smoothing replaces the one-hot encoded labels with a smoothed distribution, reducing the model's confidence in its predictions and improving generalization:

$$\tilde{y}_k = (1 - \epsilon)y_k + \frac{\epsilon}{K},$$

where $\epsilon$ is the smoothing parameter and $K$ is the number of classes.

Regularization techniques like weight decay, dropout, and label smoothing reduce the variance of the model, leading to improved generalization performance. These

techniques are particularly effective in preventing overfitting in deep models like vision transformers.

Example: For a ViT trained on the CIFAR-100 dataset, applying dropout with a rate of 0.1 and weight decay with $\lambda = 0.0005$ can help prevent overfitting, leading to better performance on the test set.

**Data Augmentation Techniques**

Data augmentation is a powerful technique for increasing the effective size of the training dataset by generating new, slightly modified versions of the existing data. This helps the model generalize better by exposing it to a wider variety of data.

Let $\mathcal{T}$ represent a set of augmentation transformations, such as rotations, flips, and color jittering. The augmented dataset $\mathcal{D}_{\text{aug}}$ is generated by applying these transformations to the original dataset $\mathcal{D}$:

$$\mathcal{D}_{\text{aug}} = \{(\mathcal{T}(\mathbf{I}_i), y_i) \mid (\mathbf{I}_i, y_i) \in \mathcal{D}, \mathcal{T} \in \mathcal{T}\}.$$

Data augmentation increases the diversity of the training data, reducing the risk of overfitting and improving the model's generalization performance. The theoretical effect of data augmentation is to smooth the decision boundaries, making the model less sensitive to variations in the input data.

Example: For a vision transformer trained on the SVHN dataset, applying data augmentation techniques like random cropping, horizontal flipping, and color jittering can significantly improve the model's robustness to variations in the input data.

## 6.5  Advanced Topics and Future Directions

ViTs have demonstrated exceptional performance in various computer vision tasks. However, the exploration of hybrid models that combine the strengths of CNNs and transformers offers a promising direction for further advancements. This section explores the mathematical formulation and principles behind hybrid models that integrate CNNs and transformers, providing insights into how these architectures leverage the advantages of both approaches.

**Hybrid Models**

Hybrid models aim to harness the localized feature extraction capabilities of CNNs along with the global attention mechanisms of transformers. This combination seeks to create architectures that are both computationally efficient and capable of capturing complex dependencies across the input data. By combining these two powerful paradigms, hybrid models can potentially outperform either CNNs or transformers alone on a wide range of tasks.

**Combining CNNs and Transformers**

The core idea behind hybrid models is to use CNNs to extract local features from the input data, which are then processed by transformers to capture long-range dependencies and contextual relationships. This approach allows for the benefits of convolutional operations in handling high-resolution data and the strengths of self-attention in modeling global interactions.

1. CNN Feature Extraction: Given an input image $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$, a CNN is first applied to extract a set of feature maps $\mathbf{F} \in \mathbb{R}^{H' \times W' \times C'}$, where $H'$ and $W'$ are the spatial dimensions after downsampling, and $C'$ is the number of output channels:

$$\mathbf{F} = \text{CNN}(\mathbf{I}).$$

The CNN typically consists of multiple convolutional layers, each followed by nonlinear activations and pooling operations, which reduce the spatial dimensions and increase the feature depth.

2. Transformer Processing: The extracted feature maps are then reshaped into a sequence of tokens suitable for processing by a transformer. Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ represent the sequence of tokens, where $N = H' \times W'$ is the number of tokens and $D = C'$ is the token dimensionality:

$$\mathbf{X} = \text{Reshape}(\mathbf{F}).$$

These tokens are then passed through a series of transformer layers, where each layer applies self-attention to model the interactions between different parts of the image:

$$\mathbf{Z} = \text{TransformerEncoder}(\mathbf{X}),$$

where $\mathbf{Z} \in \mathbb{R}^{N \times D}$ is the output of the transformer encoder.

The hybrid model benefits from the computational efficiency of CNNs in early layers, where local features are most important, and the expressive power of transformers in later layers, where capturing global context is crucial. Formally, let $T_{\text{CNN}}$ and $T_{\text{Transformer}}$ denote the time complexities of the CNN and transformer components, respectively. The overall time complexity of the hybrid model is

$$T_{\text{Hybrid}} = T_{\text{CNN}} + T_{\text{Transformer}}.$$

This formulation shows that the hybrid model can be more efficient than a pure transformer model, especially when the CNN effectively reduces the spatial dimensions before passing the features to the transformer.

Example: In a hybrid architecture designed for image classification on the ImageNet dataset, the CNN component might consist of several convolutional layers followed by pooling, reducing the image resolution to $14 \times 14$ patches, which are then processed by a vision transformer. This combination allows the model to leverage the CNN's ability to capture fine details and the transformer's capability to model global relationships.

The formulation of hybrid models involves defining the interaction between the CNN and transformer components, particularly how features are extracted, transformed, and integrated. The key challenge is to ensure that the information flows seamlessly from the convolutional layers to the self-attention layers without losing critical spatial information.

1. Feature Extraction and Tokenization: The feature maps $\mathbf{F}$ extracted by the CNN are tokenized by reshaping them into a sequence $\mathbf{X}$. The tokens $\mathbf{x}_i \in \mathbb{R}^D$ correspond to different spatial regions of the image and are enriched with local features from the CNN:

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_N], \quad \mathbf{x}_i = \text{Flatten}(\mathbf{F}_i),$$

where $\mathbf{F}_i$ represents the $i$-th feature map region.

2. Self-Attention Mechanism: The tokens $\mathbf{X}$ are processed by the self-attention mechanism in the transformer, which computes attention scores between all pairs of tokens to capture global dependencies:

$$\mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{D}}\right),$$

where $\mathbf{Q}$, $\mathbf{K}$, $and\mathbf{V}$ are the query, key, and value matrices derived from $\mathbf{X}$.

3. Output Integration: The output of the transformer encoder $\mathbf{Z}$ is reshaped and combined with any residual CNN features to produce the final prediction $\hat{y}$:

$$\hat{y} = \text{Classifier}(\mathbf{Z}).$$

The combination of CNN and transformer layers preserves both local and global information, enabling the model to make more informed predictions. Let $I_{\text{local}}$ and $I_{\text{global}}$ denote the local and global information captured by the CNN and transformer, respectively. The total information $I_{\text{total}}$ preserved by the hybrid model is

$$I_{\text{total}} = I_{\text{local}} + I_{\text{global}},$$

where $I_{\text{local}}$ is maximized by the CNN's convolutional operations and $I_{\text{global}}$ is maximized by the transformer's self-attention mechanism.

Example: In a hybrid model designed for object detection, the CNN might extract detailed features like edges and textures, while the transformer captures the spatial relationships between different objects in the scene. This synergy allows the model to detect objects with high accuracy and robustness.

# References

1. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part I, pp. 213–229. Springer-Verlag, Berlin, Heidelberg (2020). ISBN 978-3-030-58451-1. https://doi.org/10.1007/978-3-030-58452-8_13

2. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale (2020). arXiv:2010.11929
3. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Advances in Neural Information Processing Systems, vol. 27 (2014)
4. Lin, T.-Y., Maire, M., Belongie, S.J., Hays, J., Perona, P., Ramanan, D., Dollár, P., Lawrence Zitnick, C.: Microsoft coco: Common objects in context. In: European Conference on Computer Vision (ECCV), pp. 740–755 (2014)
5. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 9992–10002 (2021)

# Chapter 7
# Time Series Forecasting with Transformers

## 7.1 Mathematical Modeling of Sequential Data

Time series forecasting involves predicting future values of a sequential data series based on its historical behavior. This section provides a mathematical exploration of traditional time series models, such as autoregressive (AR) models, moving average (MA) models, and their extensions, which lay the foundation for understanding how transformers can be applied to time series forecasting. The analysis will focus on the mathematical formulation of these models, highlighting their strengths and limitations, and motivating the use of more advanced techniques like transformers.

### 7.1.1 Time Series Representation

Time series data is a sequence of observations collected over time, where each observation $y_t$ at time $t$ depends on previous observations. Formally, a time series can be represented as $\{y_t\}_{t=1}^T$, where $T$ is the length of the series. The goal of time series modeling is to capture the underlying structure and dependencies within the series to make accurate predictions about future observations.

**Autoregressive Models**

Autoregressive (AR) models are a fundamental class of time series models where the current value of the series is expressed as a linear combination of its previous values. An AR model of order $p$, denoted as AR($p$), assumes that the value at time $t$ depends on the previous $p$ values:

$$y_t = \phi_1 y_{t-1} + \phi_2 y_{t-2} + \cdots + \phi_p y_{t-p} + \epsilon_t,$$

where $\phi_1, \phi_2, \ldots, \phi_p$ are the autoregressive coefficients and $\epsilon_t$ is a white noise error term with zero mean and constant variance $\sigma^2$.

An AR($p$) model is stationary if the roots of the characteristic equation

$$1 - \phi_1 z - \phi_2 z^2 - \cdots - \phi_p z^p = 0$$

lie outside the unit circle in the complex plane. Stationarity implies that the statistical properties of the time series, such as mean and variance, do not change over time, making the model suitable for forecasting.

Example: Consider a simple AR(1) model $y_t = \phi_1 y_{t-1} + \epsilon_t$. The stationarity condition for this model is $|\phi_1| < 1$. If this condition is satisfied, the series will revert to its mean over time, and predictions will be more stable.

## AR, MA, and ARMA Models

While AR models rely on past values of the series, MA models express the current value as a linear combination of past error terms. An MA model of order $q$, denoted as MA($q$), is given by

$$y_t = \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \cdots + \theta_q \epsilon_{t-q},$$

where $\theta_1, \theta_2, \ldots, \theta_q$ are the moving average coefficients.

An MA($q$) model is invertible if the roots of the characteristic equation

$$1 + \theta_1 z + \theta_2 z^2 + \cdots + \theta_q z^q = 0$$

lie outside the unit circle in the complex plane. Invertibility ensures that the MA model can be uniquely represented as an infinite AR model, providing stability in forecasting.

The autoregressive moving average (ARMA) model combines the AR and MA models to capture both the autoregressive and moving average components in the time series. An ARMA($p, q$) model is defined as

$$y_t = \phi_1 y_{t-1} + \cdots + \phi_p y_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \cdots + \theta_q \epsilon_{t-q}.$$

An ARMA($p, q$) model is both stationary and invertible if the roots of the respective AR and MA characteristic equations lie outside the unit circle. This dual condition ensures that the model is well behaved and suitable for long-term forecasting.

Example: In financial time series, ARMA models are often used to capture both short-term dependencies (via the AR component) and shocks or noise (via the MA component), leading to more accurate predictions.

**ARIMA and Seasonal ARIMA**

The autoregressive integrated moving average (ARIMA) model ([2]) extends ARMA to handle non-stationary time series by introducing a differencing step. An ARIMA($p, d, q$) model is defined as

$$\Delta^d y_t = \phi_1 \Delta^d y_{t-1} + \cdots + \phi_p \Delta^d y_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \cdots + \theta_q \epsilon_{t-q},$$

where $\Delta^d y_t$ represents the $d$-th differenced series used to achieve stationarity.

For an ARIMA($p, d, q$) model, if the differencing order $d$ is chosen such that the differenced series $\{\Delta^d y_t\}$ is stationary, then the ARIMA model can be used to forecast the original non-stationary series.

Seasonal ARIMA (SARIMA) ([8]) models extend ARIMA to capture seasonal patterns in the data. A SARIMA($p, d, q \times P, D, Q, s$) model incorporates both non-seasonal ($p, d, q$) and seasonal ($P, D, Q, s$) components, where $s$ is the seasonal period. The model is given by

$$\Phi(B^s)\Delta_s^D \Delta^d y_t = \Theta(B^s)\epsilon_t,$$

where $\Phi(B^s)$ and $\Theta(B^s)$ are the seasonal AR and MA polynomials, respectively, and $\Delta_s^D$ represents seasonal differencing.

A SARIMA model is stationary and invertible if both the non-seasonal and seasonal components satisfy their respective stationarity and invertibility conditions. This ensures that the model can capture complex seasonal patterns while maintaining stability.

Example: In weather forecasting, SARIMA models are used to predict temperature or precipitation by accounting for both daily fluctuations and seasonal cycles, such as summer and winter patterns.

## 7.1.2 Time Series Representation

Recurrent models form a critical class of methods for modeling sequential data, particularly time series. Unlike traditional autoregressive models, which rely on a fixed number of past observations, recurrent models introduce a more flexible approach by maintaining a hidden state that evolves over time, capturing dependencies of arbitrary length. This section explores the mathematical foundations of RNNs ([5, 10]), Long Short-Term Memory (LSTM) networks ([7]), and Gated Recurrent Units (GRUs) ([3]), emphasizing their structure, dynamics, and the specific problems they address.

**Recurrent Neural Networks**

RNNs are a class of neural networks specifically designed for processing sequential data. The key idea behind RNNs is the incorporation of a hidden state that is updated at each time step, allowing the network to maintain a memory of previous inputs.

Given an input sequence $\{x_t\}_{t=1}^T$, where $x_t$ is the input at time $t$, an RNN computes the hidden state $h_t$ and output $y_t$ at each time step using the following recursive equations:

$$h_t = \sigma_h(W_h x_t + U_h h_{t-1} + b_h),$$

$$y_t = \sigma_y(W_y h_t + b_y),$$

where $h_t \in \mathbb{R}^n$ is the hidden state at time $t$, $W_h \in \mathbb{R}^{n \times m}$ and $U_h \in \mathbb{R}^{n \times n}$ are weight matrices for the input and hidden state, respectively; $b_h \in \mathbb{R}^n$ and $b_y \in \mathbb{R}^o$ are bias vectors; and $\sigma_h$ and $\sigma_y$ are activation functions, typically tanh or ReLU for the hidden state and softmax for the output.

RNNs are prone to the vanishing and exploding gradient problems during back-propagation through time (BPTT). This issue arises because the gradient of the loss with respect to the hidden state at an earlier time step $t$ involves the product of many Jacobian matrices:

$$\frac{\partial \mathcal{L}}{\partial h_t} = \sum_{k=t+1}^{T} \frac{\partial \mathcal{L}}{\partial h_k} \prod_{j=t+1}^{k} \frac{\partial h_j}{\partial h_{j-1}}.$$

If the eigenvalues of $\frac{\partial h_j}{\partial h_{j-1}}$ are not close to 1, the gradients can either vanish (if the eigenvalues are less than 1) or explode (if they are greater than 1).

In tasks requiring long-term dependencies, such as predicting trends over a long horizon, standard RNNs may struggle due to vanishing gradients, leading to poor performance as the model fails to retain important information from earlier in the sequence.

**Long Short-Term Memory (LSTM)**

LSTM networks were introduced to address the limitations of RNNs, particularly the vanishing gradient problem. LSTMs introduce a more complex hidden state structure that includes gating mechanisms, allowing the network to learn when to forget previous information and when to update the hidden state.

An LSTM unit consists of a cell state $c_t$ and three gates: input gate $i_t$, forget gate $f_t$, and output gate $o_t$. The dynamics of an LSTM are governed by the following equations:

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f),$$

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i),$$

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o),$$

$$\tilde{c}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c),$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t,$$

$$h_t = o_t \odot \tanh(c_t),$$

where $f_t, i_t, o_t \in \mathbb{R}^n$ are the forget, input, and output gates, respectively; $c_t \in \mathbb{R}^n$ is the cell state; $\tilde{c}_t \in \mathbb{R}^n$ is the candidate cell state; $\odot$ denotes element-wise multiplication; $\sigma$ is the sigmoid activation function; and tanh is the hyperbolic tangent function.

LSTMs are designed to mitigate the vanishing gradient problem by controlling the flow of information through the cell state using the forget gate $f_t$. The gradient of the loss with respect to the cell state at time $t$ is

$$\frac{\partial \mathcal{L}}{\partial c_t} = \frac{\partial \mathcal{L}}{\partial c_{t+1}} \odot f_{t+1} \cdot \text{(other terms)},$$

where $f_{t+1}$ directly modulates the gradient flow, allowing LSTMs to maintain long-term dependencies.

LSTMs excel in tasks where long-term memory is crucial, such as language modeling or time series forecasting with long periodic patterns. For example, in financial forecasting, LSTMs can capture trends over long periods, such as yearly cycles in stock prices.

**Gated Recurrent Units (GRUs)**

GRUs are a simplified variant of LSTMs, introduced to reduce the complexity while retaining the advantages of gating mechanisms. GRUs combine the forget and input gates into a single update gate and eliminate the separate cell state, directly using the hidden state to carry information.

The GRU dynamics are governed by the following equations:

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z),$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r),$$

$$\tilde{h}_t = \tanh(W_h x_t + U_h(r_t \odot h_{t-1}) + b_h),$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t,$$

where $z_t \in \mathbb{R}^n$ is the update gate, $r_t \in \mathbb{R}^n$ is the reset gate, and $\tilde{h}_t \in \mathbb{R}^n$ is the candidate hidden state.

While GRUs and LSTMs both address the vanishing gradient problem through gating mechanisms, GRUs have fewer parameters and are computationally more efficient. This can lead to faster training and inference times, particularly in applications where model size and computational resources are a concern. The key difference lies in the structure of the gates, with GRUs using fewer gates to control the information flow.

GRUs are often preferred in real-time applications, such as anomaly detection in streaming data, where computational efficiency is critical. Despite their simpler structure, GRUs can perform comparably to LSTMs on many tasks, particularly when the training data is limited.

### 7.1.3   Transformers for Time Series

Transformers, originally developed for natural language processing tasks, have proven to be powerful models for sequential data, including time series. Their ability to model long-range dependencies through attention mechanisms, rather than relying on recurrent structures, makes them particularly suitable for time series forecasting, where capturing both short-term and long-term dependencies is crucial. This section explores the application of transformers to time series data, focusing on the sequence-to-sequence framework, the encoder–decoder architecture, and the modeling of temporal dependencies using attention.

#### Sequence-to-Sequence Models

Sequence-to-sequence (Seq2Seq) models are designed to map an input sequence to an output sequence, making them well suited for tasks like time series forecasting. In the context of time series, the input sequence consists of past observations and the output sequence consists of future predictions.

Let $\mathbf{x} = \{x_1, x_2, \ldots, x_T\}$ be the input sequence of observed values and $\mathbf{y} = \{y_1, y_2, \ldots, y_{T'}\}$ be the output sequence of predicted values. A Seq2Seq model aims to learn a mapping $f : \mathbb{R}^T \to \mathbb{R}^{T'}$ such that

$$\mathbf{y} = f(\mathbf{x}),$$

where $f$ is typically parameterized by a deep neural network, such as a transformer. The goal is to minimize a loss function $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$ that quantifies the difference between the predicted sequence $\hat{\mathbf{y}}$ and the true sequence $\mathbf{y}$.

Given sufficient capacity (number of layers, hidden units, etc.), a transformer-based Seq2Seq model can approximate any continuous mapping from the input sequence $\mathbf{x}$ to the output sequence $\mathbf{y}$. This is a consequence of the universal approximation theorem, which states that neural networks can approximate any continuous function to arbitrary accuracy.

Example: In time series forecasting, a Seq2Seq model might take a sequence of daily temperatures as input and predict the temperature for the next week. The model learns to capture both short-term patterns (e.g., daily fluctuations) and long-term trends (e.g., seasonal variations).

## Encoder–Decoder Framework

The encoder–decoder framework is a common architecture used in Seq2Seq models, particularly in transformers. The encoder processes the input sequence and encodes it into a fixed-dimensional representation, while the decoder takes this representation and generates the output sequence.

Let $\mathbf{H} = \{\mathbf{h}_1, \mathbf{h}_2, \ldots, \mathbf{h}_T\}$ represent the sequence of hidden states produced by the encoder, where each $\mathbf{h}_t$ encodes information from the input sequence up to time $t$. The decoder generates the output sequence $\hat{\mathbf{y}} = \{\hat{y}_1, \hat{y}_2, \ldots, \hat{y}_{T'}\}$ as follows:

$$\hat{y}_t = g(\mathbf{h}_T, \mathbf{y}_{<t}),$$

where $g$ is the decoding function that generates the prediction $\hat{y}_t$ based on the final encoder state $\mathbf{h}_T$ and the previous outputs $\mathbf{y}_{<t}$.

The effectiveness of the encoder–decoder framework relies on the encoder's ability to compress the entire input sequence into a fixed-dimensional representation $\mathbf{h}_T$. This compression can lead to an information bottleneck, where critical information about the input sequence is lost. However, the attention mechanism in transformers alleviates this issue by allowing the decoder to attend to different parts of the input sequence dynamically.

Example: In machine translation, the encoder might process an entire sentence in the source language and produce a context vector $\mathbf{h}_T$ that encapsulates the meaning of the sentence. The decoder then generates the translated sentence in the target language, one word at a time, using this context vector.

## Temporal Dependencies and Attention

One of the key strengths of transformers is their ability to model temporal dependencies in sequential data through the self-attention mechanism. Unlike recurrent models, which process the sequence step by step, transformers use attention to consider all time steps simultaneously, capturing dependencies of varying lengths effectively.

The self-attention mechanism computes a weighted sum of the input sequence, where the weights are determined by the similarity between the elements of the sequence. For an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_T\}$, the attention score between elements $x_i$ and $x_j$ is given by

$$\alpha_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}}\right)}{\sum_{j'=1}^{T} \exp\left(\frac{\mathbf{q}_i \mathbf{k}_{j'}^\top}{\sqrt{d_k}}\right)},$$

where $\mathbf{q}_i = \mathbf{x}_i W_Q$, $\mathbf{k}_j = \mathbf{x}_j W_K$, and $d_k$ is the dimensionality of the keys. The output for each position $i$ is then

$$\mathbf{z}_i = \sum_{j=1}^{T} \alpha_{ij} \mathbf{v}_j,$$

where $\mathbf{v}_j = \mathbf{x}_j W_V$ is the value associated with $x_j$.

The self-attention mechanism in transformers allows the model to capture both local and global dependencies within the input sequence. Formally, the attention mechanism can express any dependency structure within the sequence, making transformers highly flexible for modeling complex temporal relationships.

Example: In time series forecasting, self-attention enables the model to focus on important time steps that are relevant for making accurate predictions, such as previous peaks or trends in the data. This allows the model to learn complex temporal patterns, such as seasonality or trends, without relying on recurrent connections.

**Attention Mechanisms in Time Series**

The ability to model dependencies across different time steps is crucial for accurate time series forecasting. Attention mechanisms, and specifically self-attention, enable transformers to weigh the importance of each time step in the input sequence when making predictions, thereby capturing both local and global patterns.

In a standard self-attention mechanism, the input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_T\}$ is transformed into queries $\mathbf{Q}$, keys $\mathbf{K}$, and values $\mathbf{V}$ using learned weight matrices $W_Q$, $W_K$, and $W_V$:

$$\mathbf{Q} = \mathbf{X} W_Q, \quad \mathbf{K} = \mathbf{X} W_K, \quad \mathbf{V} = \mathbf{X} W_V,$$

where $\mathbf{X} \in \mathbb{R}^{T \times d}$ is the matrix representation of the input sequence, with $T$ as the sequence length and $d$ as the dimensionality of each time step. The self-attention output for each time step $t$ is computed as

$$\mathbf{z}_t = \sum_{j=1}^{T} \alpha_{tj} \mathbf{v}_j,$$

where the attention weights $\alpha_{tj}$ are given by

$$\alpha_{tj} = \frac{\exp\left(\frac{\mathbf{q}_t \mathbf{k}_j^\top}{\sqrt{d_k}}\right)}{\sum_{j'=1}^{T} \exp\left(\frac{\mathbf{q}_t \mathbf{k}_{j'}^\top}{\sqrt{d_k}}\right)}.$$

The self-attention mechanism is capable of capturing dependencies across all time steps in the input sequence. This means that the model can learn to focus on relevant past observations, regardless of their distance in time, making it particularly powerful for time series forecasting where both short-term and long-term dependencies are important.

Example: In financial time series, self-attention can allow the model to weigh recent market trends more heavily while also taking into account significant events from the past, such as economic reports, that may influence future trends.

## Self-Attention in Sequential Data

Self-attention in sequential data offers the advantage of parallelizing the computation across all time steps, in contrast to recurrent models, which process sequences sequentially. This parallelization not only speeds up computation but also allows the model to capture complex interactions between time steps that might be difficult to learn with recurrent structures.

Given the self-attention mechanism described above, the parallel computation across all time steps is expressed as

$$\mathbf{Z} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V},$$

where $\mathbf{Z} \in \mathbb{R}^{T \times d}$ is the matrix of attention outputs for all time steps. The softmax operation ensures that the attention weights for each time step sum to one, thereby normalizing the contribution of each value $\mathbf{v}_j$.

The parallel nature of self-attention allows the model to compute the interactions between all pairs of time steps simultaneously. This parallelization reduces the computational complexity from $O(T^2 \times d)$ for each time step to $O(T^2 \times d)$ overall, making it significantly more efficient than recurrent models for long sequences.

Example: In time series forecasting for large datasets, such as predicting energy consumption across a grid, the parallelization of self-attention enables the model to handle very long sequences efficiently, capturing interactions across different timescales.

## Multi-head Attention for Temporal Data

Multi-head attention extends the self-attention mechanism by applying it multiple times in parallel, with different sets of learned weights for each head. This allows the model to capture different aspects of the temporal data simultaneously, improving its ability to model complex patterns.

In multi-head attention, the input sequence $\mathbf{x}$ is processed by $h$ different attention heads, each with its own set of weight matrices $W_Q^{(i)}$, $W_K^{(i)}$, $W_V^{(i)}$ for $i = 1, 2, \ldots, h$.

The output of each head is

$$\mathbf{z}_t^{(i)} = \sum_{j=1}^{T} \alpha_{tj}^{(i)} \mathbf{v}_j^{(i)},$$

where $\alpha_{tj}^{(i)}$ and $\mathbf{v}_j^{(i)}$ are the attention weights and values for head $i$. The outputs from all heads are concatenated and linearly transformed to produce the final output:

$$\mathbf{z}_t = W_O \left[ \mathbf{z}_t^{(1)} \mid \mathbf{z}_t^{(2)} \mid \ldots \mid \mathbf{z}_t^{(h)} \right],$$

where $W_O$ is a learned weight matrix that combines the outputs of all heads.

Multi-head attention increases the model's capacity by allowing it to attend to different parts of the sequence simultaneously and to capture a richer set of relationships within the data. The different heads can focus on different time steps or learn different patterns, leading to a more comprehensive understanding of the temporal data.

Example: In weather forecasting, multi-head attention can enable the model to simultaneously capture short-term fluctuations (e.g., daily temperature changes) and long-term trends (e.g., seasonal variations), leading to more accurate predictions.

**Positional Encoding for Time Series**

Since transformers lack an inherent sense of the order of the input sequence, positional encoding is crucial for allowing the model to capture temporal information. Positional encoding adds information about the position of each time step to the input data, enabling the model to differentiate between time steps in the sequence.

Positional encodings are added to the input embeddings to incorporate information about the relative or absolute position of each time step. For a time step $t$, the positional encoding PE($t$) is typically defined as

$$\text{PE}(t, 2i) = \sin\left(\frac{t}{10000^{2i/d}}\right), \quad \text{PE}(t, 2i+1) = \cos\left(\frac{t}{10000^{2i/d}}\right),$$

where $d$ is the dimensionality of the input embeddings and $i$ indexes the dimensions. These sinusoidal functions ensure that the encodings are unique for each time step and that the relative distances between time steps are preserved.

Positional encoding ensures that the model can distinguish between different time steps and understand the temporal order of the sequence. The sinusoidal functions used in positional encoding provide a smooth and continuous representation of time, which is crucial for capturing temporal dependencies in time series data.

Example: In financial forecasting, positional encoding allows the transformer to recognize the order of daily stock prices, enabling it to capture trends and patterns that depend on the sequence of observations, such as momentum or reversal patterns.

## 7.2 Applications in Forecasting

Financial time series forecasting is a crucial area where the predictive power of transformers can be leveraged. Financial markets generate vast amounts of sequential data, such as stock prices, trading volumes, and volatility indices, all of which require sophisticated models to predict future trends accurately. This section examines the application of transformers to financial time series, focusing on stock price prediction, portfolio optimization, and volatility modeling. We will explore how the mathematical properties of transformers align with the specific challenges posed by financial data.

### 7.2.1 Financial Time Series

Financial time series are often characterized by non-stationarity, high volatility, and complex dependencies across different assets and time periods. Transformers, with their ability to model long-range dependencies and capture intricate temporal patterns, are well suited to tackle these challenges. The application of transformers to financial time series involves understanding the underlying geometric and probabilistic structures of the data, and how these can be encoded into the model.

**Stock Price Prediction**

Stock price prediction involves forecasting the future prices of stocks based on historical data. This is a fundamental problem in finance, with significant implications for trading strategies, risk management, and portfolio optimization. The challenge lies in capturing the complex, often non-linear dependencies between past prices and future movements.

Let $\{P_t\}_{t=1}^{T}$ be the time series representing the closing prices of a stock over $T$ trading days. The goal is to predict the future price $P_{T+h}$ after a horizon $h$, based on the past prices $\{P_t\}_{t=1}^{T}$. A transformer model can be trained to learn the mapping:

$$\hat{P}_{T+h} = f(\{P_t\}_{t=1}^{T}),$$

where $f$ is a function parameterized by the transformer's weights and $\hat{P}_{T+h}$ is the predicted price.

Under the Efficient Market Hypothesis (EMH) ([6]), prices reflect all available information, making future price movements inherently unpredictable. However, empirical studies show that certain market inefficiencies and patterns, such as momentum or mean reversion, can be exploited for predictive modeling. The transformer, by capturing these patterns through attention mechanisms, can potentially identify and leverage these inefficiencies.

Example: In practice, a transformer might be trained on historical price data from multiple stocks, allowing it to learn common patterns across different assets. By focusing attention on key historical events, such as earnings reports or macroeconomic announcements, the model can make more informed predictions about future price movements.

## Portfolio Optimization

Portfolio optimization involves selecting a combination of financial assets that maximizes expected return for a given level of risk, or, equivalently, minimizes risk for a given level of expected return. The challenge is to forecast the returns and covariances of the assets in the portfolio, often based on historical data.

Let $\mathbf{r} = \{r_1, r_2, \ldots, r_n\}$ be the vector of expected returns for $n$ assets and $\Sigma$ be the covariance matrix of these returns. The objective in mean-variance optimization is to find the portfolio weights $\mathbf{w} = \{w_1, w_2, \ldots, w_n\}$ that solve

$$\min_{\mathbf{w}} \mathbf{w}^\top \Sigma \mathbf{w} - \lambda \mathbf{w}^\top \mathbf{r},$$

where $\lambda$ is a risk-aversion parameter. The role of the transformer in this context is to forecast $\mathbf{r}$ and $\Sigma$ based on historical returns, enabling the optimization of the portfolio.

Markowitz's portfolio theory ([9]) states that for any given expected return $\mathbf{r}^\top \mathbf{w}$, there exists a portfolio that minimizes risk, given by

$$\mathbf{w}^* = \frac{\Sigma^{-1} \mathbf{r}}{\mathbf{r}^\top \Sigma^{-1} \mathbf{r}}.$$

By accurately forecasting returns and covariances using transformers, investors can construct optimal portfolios that align with their risk preferences.

Example: A transformer model could be used to predict the next month's returns for a set of stocks, taking into account both historical price data and relevant external factors such as interest rates or economic indicators. These predictions could then be fed into a portfolio optimization algorithm to determine the optimal allocation of capital.

## Volatility Modeling

Volatility modeling is critical in finance, as it measures the degree of variation in asset prices and is directly related to risk. Accurate volatility forecasts are essential for pricing derivatives, managing risk, and making informed trading decisions.

Let $\sigma_t^2$ represent the conditional variance of the returns $r_t$. A common model for volatility is the GARCH(1,1) model ([1]), given by

$$\sigma_t^2 = \alpha_0 + \alpha_1 r_{t-1}^2 + \beta_1 \sigma_{t-1}^2,$$

where $\alpha_0 > 0$, $\alpha_1 \geq 0$, and $\beta_1 \geq 0$ are parameters. The transformer can be used to model and forecast volatility by capturing the complex dependencies in the time series of returns, potentially incorporating additional factors such as trading volume or macroeconomic indicators.

For the GARCH(1,1) model to be stationary, the sum of the coefficients $\alpha_1 + \beta_1$ must be less than 1. This ensures that the impact of past shocks on current volatility diminishes over time.

Example: A transformer model could be trained to predict the next day's volatility based on historical price and volume data. By accurately forecasting volatility, the model could assist in pricing options or in adjusting trading strategies to mitigate risk during periods of high market uncertainty.

## 7.2.2  Weather Forecasting

Weather forecasting relies on the analysis of historical meteorological data and the simulation of atmospheric dynamics to predict future weather conditions. This involves a combination of physical models, statistical methods, and machine learning techniques. Transformers, with their capacity to model sequential data and capture complex temporal patterns, offer a powerful approach to improving the accuracy and robustness of weather predictions.

### Short-Term Forecasting

Short-term weather forecasting involves predicting weather conditions over a period ranging from a few hours to a few days. This is essential for planning daily activities, managing energy resources, and issuing early warnings for adverse weather conditions.

Let $\mathbf{x}_t$ represent a vector of meteorological variables (e.g., temperature, pressure, humidity) at time $t$. The objective in short-term forecasting is to predict the state of these variables at a future time $t + h$, where $h$ is the forecast horizon. A transformer model can be trained to learn the mapping:

$$\hat{\mathbf{x}}_{t+h} = f(\{\mathbf{x}_t, \mathbf{x}_{t-1}, \ldots, \mathbf{x}_{t-k}\}),$$

where $k$ is the number of previous time steps considered and $f$ is parameterized by the transformer's weights. The model predicts the future state $\hat{\mathbf{x}}_{t+h}$ based on the sequence of past observations.

The atmosphere is a chaotic system, meaning that small changes in initial conditions can lead to vastly different outcomes over time. According to the Lorenz attractor model, the predictability of weather diminishes rapidly as the forecast horizon

increases. However, within short-term horizons, deterministic models, such as transformers, can provide reasonably accurate forecasts by capturing the local dynamics of the atmosphere.

Example: In practice, a transformer might be trained on historical weather data to predict the next day's temperature and precipitation levels. The model would capture both diurnal patterns and short-term weather phenomena, such as the passage of a cold front, leading to more accurate short-term forecasts.

## Long-Term Climate Predictions

Long-term climate predictions involve forecasting climate variables over extended periods, ranging from months to decades. This is crucial for understanding climate change, planning infrastructure, and developing policies for environmental sustainability.

Let $\mathbf{y}_t$ represent climate variables (e.g., average temperature, sea level, carbon dioxide concentration) over a long-term period. The goal is to predict the state of these variables at a future time $t + H$, where $H$ could be several months or years ahead. A transformer model can be formulated to learn the mapping:

$$\hat{\mathbf{y}}_{t+H} = g(\{\mathbf{y}_t, \mathbf{y}_{t-1}, \ldots, \mathbf{y}_{t-K}\}),$$

where $K$ is the number of previous time steps and $g$ is the function parameterized by the transformer.

Climate data is often non-stationary, with trends and periodicities that evolve over time due to factors such as greenhouse gas emissions and volcanic activity. The Fourier decomposition theorem allows any periodic signal to be expressed as a sum of sinusoidal functions, which can be captured by transformers through positional encoding and attention mechanisms.

Example: In long-term climate prediction, a transformer could be used to forecast global temperature anomalies based on historical data, including natural cycles such as the El Niño-Southern Oscillation (ENSO) and anthropogenic factors. The model would need to account for both seasonal variations and long-term trends, such as global warming.

## Extreme Weather Events Prediction

Predicting extreme weather events, such as hurricanes, tornadoes, and heatwaves, is of paramount importance due to their significant impact on human life and property. These events are rare, but their prediction requires models that can capture sudden, non-linear changes in weather patterns.

Let $\mathbf{z}_t$ represent indicators of extreme weather conditions (e.g., wind speed, atmospheric pressure, temperature gradients). The goal is to predict the likelihood of an

extreme event occurring within a future time window $[t, t + H]$. A transformer can model this as a probabilistic forecasting problem:

$$\hat{P}(E_{t+H}) = p(\mathbf{z}_{t+H} \mid \{\mathbf{z}_t, \mathbf{z}_{t-1}, \ldots, \mathbf{z}_{t-k}\}),$$

where $\hat{P}(E_{t+H})$ is the predicted probability of an extreme event $E$ occurring at time $t + H$.

Extreme Value Theory (EVT) ( [4]) provides a framework for modeling the tails of a distribution, which is essential for predicting rare events. The Generalized Extreme Value (GEV) distribution is often used to model the maximum or minimum of a large sample of independent random variables. Transformers can be integrated with EVT to enhance the prediction of extreme weather events by focusing on the distributional tails.

Example: In the context of hurricane prediction, a transformer could be trained on historical hurricane data, including sea surface temperatures, wind shear, and atmospheric pressure, to predict the likelihood of hurricane formation and its potential intensity. The model could provide early warnings, allowing for better preparedness and risk management.

### *7.2.3  Energy Demand Forecasting*

Energy demand forecasting is essential for efficient grid management, load balancing, and integrating renewable energy sources into the grid. Accurate forecasts help in minimizing costs, reducing waste, and ensuring a stable energy supply.

**Electricity Load Forecasting**

Electricity load forecasting involves predicting the future demand for electricity, which is crucial for energy providers to manage supply efficiently. The demand for electricity varies based on factors such as time of day, weather conditions, and socio-economic activities.

Let $\mathbf{L}_t$ represent the electricity load at time $t$. The goal is to predict the future load $\mathbf{L}_{t+h}$ over a horizon $h$, based on historical load data and possibly exogenous variables such as temperature $\mathbf{T}_t$ and day of the week $\mathbf{D}_t$. A transformer model can be formulated to learn the mapping:

$$\hat{\mathbf{L}}_{t+h} = f(\{\mathbf{L}_t, \mathbf{T}_t, \mathbf{D}_t, \ldots, \mathbf{L}_{t-k}, \mathbf{T}_{t-k}, \mathbf{D}_{t-k}\}),$$

where $k$ is the number of previous time steps considered and $f$ is the function parameterized by the transformer.

Electricity load often exhibits strong seasonal patterns, both daily and weekly. The Fourier decomposition theorem allows the periodic components of the load

to be captured and modeled by the transformer's attention mechanism, particularly when combined with positional encoding.

A transformer model might be trained on historical load data from a utility company, using temperature and calendar information as additional inputs. The model would be able to predict peak load times and help in optimizing energy distribution across the grid.

**Renewable Energy Production Forecasting**

Forecasting renewable energy production, particularly from sources like solar and wind, is critical for integrating these resources into the energy grid. The variability and intermittency of renewable energy sources pose significant challenges for accurate forecasting.

Let $\mathbf{E}_t$ represent the energy produced by a renewable source (e.g., solar or wind) at time $t$. The goal is to predict the future production $\mathbf{E}_{t+h}$ over a horizon $h$, based on historical production data and exogenous variables such as solar irradiance $\mathbf{I}_t$ and wind speed $\mathbf{W}_t$. A transformer model can learn the mapping:

$$\hat{\mathbf{E}}_{t+h} = g(\{\mathbf{E}_t, \mathbf{I}_t, \mathbf{W}_t, \ldots, \mathbf{E}_{t-k}, \mathbf{I}_{t-k}, \mathbf{W}_{t-k}\}),$$

where $g$ is the function parameterized by the transformer.

Renewable energy production is highly intermittent and subject to uncertainties due to changing weather conditions. The transformer's ability to capture long-range dependencies and integrate multiple data sources, such as weather forecasts, makes it well suited for this task.

A transformer could be used to predict the output of a solar power plant by modeling the relationship between historical solar power output, current weather conditions, and forecasts of irradiance and cloud cover. This prediction helps in managing grid stability and planning energy storage.

## 7.2.4  Healthcare Time Series

In healthcare, time series data is generated continuously through patient monitoring systems and public health records. Accurate forecasting of patient conditions and disease outbreaks can significantly improve healthcare outcomes by enabling timely interventions and resource allocation.

**Patient Monitoring**

Patient monitoring involves continuously tracking vital signs and other health indicators to detect early signs of deterioration or improvement. Time series models can predict future health states, allowing for proactive healthcare management.

Let $\mathbf{V}_t$ represent a vector of vital signs (e.g., heart rate, blood pressure, oxygen saturation) at time $t$. The goal is to predict the patient's future state $\mathbf{V}_{t+h}$ over a horizon $h$, based on past observations:

$$\hat{\mathbf{V}}_{t+h} = h(\{\mathbf{V}_t, \mathbf{V}_{t-1}, \ldots, \mathbf{V}_{t-k}\}),$$

where $h$ is a function parameterized by the transformer.

Physiological signals often exhibit complex, non-linear dynamics and can be influenced by various factors such as medication, activity level, and environmental conditions. The transformer's attention mechanism can model these dynamics by focusing on relevant time points in the patient's history.

Example: In an ICU setting, a transformer model could be trained to predict the likelihood of a patient experiencing a critical event, such as a heart attack, within the next few hours. By monitoring vital signs and identifying patterns indicative of deterioration, the model could provide early warnings to medical staff.

**Disease Outbreak Prediction**

Predicting disease outbreaks, such as influenza or COVID-19, involves analyzing time series data from various sources, including hospital records, social media, and environmental factors. Early prediction of outbreaks allows for timely public health interventions.

Let $\mathbf{D}_t$ represent the number of reported cases of a disease at time $t$. The goal is to predict the number of cases $\mathbf{D}_{t+h}$ over a horizon $h$, based on historical case data and possibly other exogenous variables like temperature or mobility data:

$$\hat{\mathbf{D}}_{t+h} = f(\{\mathbf{D}_t, \mathbf{M}_t, \mathbf{E}_t, \ldots, \mathbf{D}_{t-k}, \mathbf{M}_{t-k}, \mathbf{E}_{t-k}\}),$$

where $\mathbf{M}_t$ might represent mobility data and $\mathbf{E}_t$ environmental factors.

Epidemiological models often rely on compartmental models (e.g., SIR) to describe the spread of diseases. Transformers can enhance these models by incorporating high-dimensional data and learning the complex dependencies between different variables that affect disease spread.

Example: During the COVID-19 pandemic, transformer models could be used to predict future case counts based on historical data, mobility patterns, and public health interventions. These predictions could inform decisions on lockdowns, resource allocation, and vaccination strategies.

## 7.3    Advantages and Limitations of Transformers in Time Series

Transformers have become a popular choice for modeling time series data due to their ability to capture long-range dependencies and handle complex temporal patterns. However, like any model, they come with their own set of advantages and limitations. This section examines the mathematical underpinnings of these strengths and weaknesses, providing a clear understanding of where transformers excel and where they may face challenges.

**Advantages of Transformers in Time Series**

1. Modeling Long-Range Dependencies: The self-attention mechanism at the heart of transformers allows them to effectively model long-range dependencies in time series data. Unlike recurrent models, which process sequences sequentially, transformers can directly attend to all time steps in the input sequence simultaneously. This parallel processing enables the capture of relationships between distant time points, which is crucial for many time series applications. For an input sequence $\mathbf{x} = \{x_1, x_2, \ldots, x_T\}$, the self-attention mechanism computes the output as

$$\mathbf{z}_t = \sum_{j=1}^{T} \alpha_{tj} \mathbf{v}_j,$$

where the attention weights $\alpha_{tj}$ are determined by the similarity between the queries and keys, allowing the model to focus on relevant time points regardless of their position in the sequence. The self-attention mechanism enables the transformer to model non-local interactions within the sequence, providing a more comprehensive understanding of the data compared to models that rely solely on local context.

2. Scalability and Parallelism: Transformers are highly scalable due to their parallel processing capabilities. This makes them suitable for large-scale time series datasets, where processing efficiency is critical. The ability to compute attention for all time steps in parallel leads to faster training times compared to sequential models like RNNs or LSTMs. The time complexity of the self-attention mechanism is $O(T^2 \cdot d)$, where $T$ is the sequence length and $d$ is the dimensionality of the model. This contrasts with the $O(T \cdot d^2)$ complexity of LSTMs, where each time step is processed sequentially. The parallel computation of attention allows transformers to handle long sequences more efficiently, reducing the training time and enabling the use of larger datasets.

3. Flexibility in Handling Multivariate Time Series: Transformers can easily handle multivariate time series, where multiple variables are observed over time. The multi-head attention mechanism allows the model to focus on different aspects of the data simultaneously, making it well suited for complex, high-dimensional time series. In multivariate time series, the input can be represented as a matrix $\mathbf{X} \in \mathbb{R}^{T \times d}$,

where each row corresponds to a different time step and each column corresponds to a different variable. The multi-head attention mechanism processes this matrix as

$$\mathbf{Z}^{(i)} = \text{softmax}\left(\frac{\mathbf{Q}^{(i)}(\mathbf{K}^{(i)})^\top}{\sqrt{d_k}}\right)\mathbf{V}^{(i)},$$

for each attention head $i$. The multi-head attention mechanism allows transformers to capture complex interactions between variables, providing a more detailed understanding of the underlying dynamics in multivariate time series.

Example: In the context of energy demand forecasting, transformers can simultaneously model the interactions between temperature, humidity, and electricity load, capturing how these variables jointly influence future demand.

**Limitations of Transformers in Time Series**

1. Computational and Memory Complexity: Despite their scalability, transformers can be computationally and memory-intensive, particularly for long sequences. The quadratic complexity of the self-attention mechanism with respect to the sequence length $T$ can become a bottleneck, especially when $T$ is large. The memory requirement for self-attention scales as $O(T^2 \cdot d)$, which can be prohibitive for long sequences or large batch sizes. This limits the practical application of transformers in scenarios where resources are constrained. For sequences longer than a certain threshold, the memory requirements of transformers may exceed the available resources, necessitating the use of approximations or modifications to the self-attention mechanism.

2. Need for Large Datasets: Transformers typically require large amounts of data to perform well. This is due to the large number of parameters in the model, which need sufficient data for effective training. In scenarios where data is limited, transformers may overfit or fail to generalize. The performance of transformers is directly related to the size of the training dataset. In the low-data regime, the model may struggle to learn meaningful patterns, leading to poor generalization.

3. Lack of Explicit Temporal Structure: Transformers do not have an inherent notion of temporal order, unlike recurrent models. This requires the introduction of positional encoding to provide the model with information about the position of each time step. While effective, this approach may not capture all the nuances of temporal dependencies, particularly in highly non-linear time series. Positional encodings are added to the input embeddings as

$$\mathbf{X}_t = \mathbf{E}_t + \text{PE}(t),$$

where $\text{PE}(t)$ represents the positional encoding. However, this approach may not fully capture the temporal dynamics of the sequence, especially in complex, non-stationary time series. The effectiveness of positional encodings is limited by their ability to represent temporal relationships, particularly in sequences with irregular time intervals or non-stationary patterns.

Example: In healthcare time series, where the data may be sparse or irregular (e.g., patient monitoring with missing data points), transformers may struggle to accurately capture the temporal dynamics without extensive data preprocessing or augmentation.

## 7.4   Optimization and Training Strategies

The effectiveness of transformers in time series forecasting depends not only on the model architecture but also on the optimization and training strategies used. This section explores the mathematical foundations of loss functions, hyperparameter tuning, and regularization techniques that are crucial for training transformers on time series data.

### Loss Functions for Time Series

The choice of loss function plays a critical role in training transformers for time series forecasting. The loss function must be chosen based on the specific objectives of the forecasting task, whether it be minimizing prediction error, capturing uncertainty, or modeling distributional properties.

For a time series forecasting task, let $\hat{y}_t$ be the predicted value and $y_t$ the true value at time $t$. Common loss functions include

1. Mean Squared Error (MSE):

$$\mathcal{L}_{\text{MSE}} = \frac{1}{T} \sum_{t=1}^{T} (\hat{y}_t - y_t)^2,$$

which penalizes larger errors more severely.

2. Mean Absolute Error (MAE):

$$\mathcal{L}_{\text{MAE}} = \frac{1}{T} \sum_{t=1}^{T} |\hat{y}_t - y_t|,$$

which is more robust to outliers.

3. Quantile Loss: For predicting quantiles of the distribution:

$$\mathcal{L}_{\text{Quantile}} = \frac{1}{T} \sum_{t=1}^{T} \begin{cases} \tau(\hat{y}_t - y_t), & \text{if } \hat{y}_t > y_t \\ (1-\tau)(y_t - \hat{y}_t), & \text{otherwise} \end{cases},$$

where $\tau$ is the quantile level.

The choice of loss function should align with the forecasting objective. For example, MSE is optimal for minimizing variance, while quantile loss is suitable for estimating conditional quantiles.

Example: In electricity load forecasting, where extreme errors can be costly, MSE might be preferred to penalize large deviations more heavily.

## Hyperparameter Tuning

Hyperparameter tuning is essential for optimizing transformer performance. Hyperparameters such as learning rate, batch size, and the number of attention heads significantly influence the training dynamics and final model accuracy.

Let $\theta$ represent the hyperparameters of the model. The objective is to minimize the loss function $\mathcal{L}(\theta)$ over a validation set:

$$\theta^* = \arg\min_{\theta} \mathcal{L}(\theta).$$

Common hyperparameter tuning methods include
1. Grid Search: Exhaustively searching over a predefined hyperparameter space.
2. Random Search: Sampling hyperparameters randomly from a distribution.
3. Bayesian Optimization: Using probabilistic models to guide the search for optimal hyperparameters.

Under certain conditions, such as Lipschitz continuity of the loss function with respect to the hyperparameters, Bayesian optimization can converge to the global optimum more efficiently than grid or random search.

In training transformers for disease outbreak prediction, Bayesian optimization could be used to find the optimal learning rate and number of attention heads that minimize the prediction error on a validation set.

## Regularization Techniques

Regularization is crucial for preventing overfitting, especially in high-capacity models like transformers. Regularization techniques add constraints to the optimization problem, promoting simpler models that generalize better to unseen data.

Let $\mathcal{L}(\theta)$ be the original loss function and let $\Omega(\theta)$ be a regularization term. The regularized loss is

$$\mathcal{L}_{\text{reg}}(\theta) = \mathcal{L}(\theta) + \lambda\Omega(\theta),$$

where $\lambda$ is a regularization parameter that controls the trade-off between the original loss and the regularization.

Common regularization techniques include
1. L2 Regularization (Ridge):

$$\Omega(\theta) = \|\theta\|_2^2,$$

which penalizes large weights and promotes smooth solutions.

2. Dropout: Randomly dropping units during training to prevent co-adaptation:

$$\mathbf{h} = \text{Dropout}(\mathbf{h}),$$

where $\mathbf{h}$ represents the hidden states.

3. Label Smoothing: Modifying the target labels to make the model less confident in its predictions:

$$y_t' = (1 - \alpha)y_t + \frac{\alpha}{K},$$

where $\alpha$ is the smoothing parameter and $K$ is the number of classes.

Regularization techniques introduce bias to reduce variance, leading to models that generalize better. The optimal regularization parameter $\lambda^*$ balances this trade-off, minimizing the overall prediction error.

Example: In healthcare time series, where the data may be noisy and sparse, L2 regularization combined with dropout can prevent overfitting, leading to more robust models that perform well on new patients.

## 7.5   Advanced Topics and Future Directions

As the field of time series forecasting with transformers continues to evolve, researchers are exploring advanced topics that extend the capabilities of these models. This section delves into hybrid models that combine transformers with traditional approaches, as well as strategies for improving scalability and efficiency.

### 7.5.1   Hybrid Models

Hybrid models leverage the strengths of both transformers and traditional time series models, creating a synergistic approach that improves forecasting accuracy and robustness. By combining different modeling paradigms, hybrid models can capture a broader range of temporal patterns and dependencies.

**Combining Transformers with Traditional Models**

Traditional time series models, such as ARIMA, Exponential Smoothing, and GARCH, have well-established mathematical foundations and are effective in capturing specific types of patterns, such as seasonality and volatility. However, they often struggle with non-linear dependencies and complex temporal relationships. Transformers, on the other hand, excel in modeling these complex patterns but may lack the interpretability and domain-specific insights provided by traditional models.

Let $\mathbf{y}_t$ represent the time series data, which can be decomposed into a component modeled by a traditional approach $\mathbf{y}_t^{\text{trad}}$ and a component modeled by a transformer $\mathbf{y}_t^{\text{trans}}$. The hybrid model can be expressed as

$$\hat{\mathbf{y}}_t = \alpha \mathbf{y}_t^{\text{trad}} + \beta \mathbf{y}_t^{\text{trans}} + \epsilon_t,$$

where $\alpha$ and $\beta$ are weights that balance the contributions of the traditional and transformer components, and $\epsilon_t$ is the error term.

Under certain conditions, the hybrid model can outperform both the traditional and transformer models individually. If $\alpha$ and $\beta$ are chosen optimally, the hybrid model can minimize the overall prediction error by capturing both linear and non-linear dependencies in the data.

Example: In financial time series forecasting, a hybrid model might combine a GARCH model for capturing volatility with a transformer for modeling complex patterns in price movements. The GARCH component would handle the predictable aspects of volatility, while the transformer would capture more intricate, non-linear relationships.

**Mathematical Formulation of Hybrid Models**

The success of hybrid models depends on the proper integration of the different components. This requires a careful mathematical formulation that ensures the models complement each other rather than compete.

The hybrid model can be formulated as a convex combination of the traditional and transformer components:

$$\hat{\mathbf{y}}_t = \sum_{i=1}^{n} \alpha_i \mathbf{y}_t^{\text{trad},i} + \sum_{j=1}^{m} \beta_j \mathbf{y}_t^{\text{trans},j} + \epsilon_t,$$

where $\mathbf{y}_t^{\text{trad},i}$ represents the output of the $i$-th traditional model and $\mathbf{y}_t^{\text{trans},j}$ represents the output of the $j$-th transformer model. The weights $\alpha_i$ and $\beta_j$ are optimized to minimize a loss function $\mathcal{L}$:

$$\mathcal{L}(\alpha, \beta) = \frac{1}{T} \sum_{t=1}^{T} \|\mathbf{y}_t - \hat{\mathbf{y}}_t\|^2 + \lambda \left( \sum_{i=1}^{n} \|\alpha_i\|^2 + \sum_{j=1}^{m} \|\beta_j\|^2 \right),$$

where $\lambda$ is a regularization parameter to prevent overfitting.

If the loss function is convex, the optimization process will converge to a global minimum, ensuring that the hybrid model effectively balances the contributions of the traditional and transformer components.

Example: In energy demand forecasting, a hybrid model might use an ARIMA model to capture the linear trend and seasonality, combined with a transformer to

model the non-linear relationships between temperature, time of day, and electricity load.

## 7.5.2  Scalability and Efficiency

As transformers are applied to increasingly large and complex time series datasets, scalability and efficiency become critical. This section explores advanced techniques to optimize the computational and memory efficiency of transformers, enabling their application to long sequences and large-scale datasets.

### Efficient Attention Mechanisms for Long Sequences

The self-attention mechanism in transformers, while powerful, has a quadratic complexity with respect to the sequence length, making it computationally expensive for long sequences. Various efficient attention mechanisms have been proposed to address this limitation by approximating or modifying the attention computation.

Let $\mathbf{Q}$, $\mathbf{K}$, $\mathbf{V}$ represent the queries, keys, and values, respectively. The standard attention mechanism computes

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V}.$$

To improve efficiency, approximate methods, such as sparse attention or locality-sensitive hashing (LSH), reduce the complexity by focusing only on a subset of the most relevant key-value pairs.

In sparse attention, only a subset of the attention weights is computed, based on a predefined sparsity pattern:

$$\text{SparseAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \sum_{(i,j)\in\mathcal{S}} \alpha_{ij}\mathbf{v}_j,$$

where $\mathcal{S}$ is the set of non-zero attention pairs.

Efficient attention mechanisms introduce an approximation error $\epsilon$ in the attention computation. The error is bounded by

$$\|\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) - \text{EfficientAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})\| \le \epsilon,$$

where $\epsilon$ depends on the sparsity pattern or approximation technique used.

Example: In weather forecasting, where long sequences of atmospheric data need to be processed, sparse attention mechanisms can significantly reduce the computational burden, allowing transformers to handle longer sequences without sacrificing accuracy.

**Memory Optimization Techniques**

Memory usage is a critical concern when training transformers on large-scale datasets, particularly when working with long sequences. Memory optimization techniques, such as gradient checkpointing and mixed-precision training, can help reduce memory requirements while maintaining model performance.

Memory optimization involves strategically reducing the amount of data stored in memory during training. Gradient checkpointing, for instance, trades off computation for memory by recomputing certain activations during the backward pass instead of storing them:

$$\text{CheckpointedGradient} = \frac{\partial \mathcal{L}}{\partial \theta} = \sum_{t=1}^{T} \frac{\partial \mathcal{L}}{\partial h_t} \cdot \frac{\partial h_t}{\partial \theta},$$

where $h_t$ represents hidden states, and only a subset of $h_t$ is stored, with the others being recomputed as needed.

Mixed-precision training reduces memory usage by storing some model parameters and activations in lower precision (e.g., 16 bit instead of 32 bit) without significantly affecting model accuracy:

$$\theta_{\text{mixed}} = \text{cast}(\theta, \text{float16}) + \text{gradients}(\theta, \text{float32}).$$

Memory optimization techniques introduce a trade-off between memory usage and computational cost. The optimal balance is achieved when the reduction in memory usage leads to only a minimal increase in computation time or a negligible decrease in accuracy.

Example: In healthcare time series analysis, where large patient datasets need to be processed in real time, memory optimization techniques such as gradient checkpointing allow for the training of large transformer models without exceeding memory constraints, making them feasible for deployment in clinical settings.

### 7.5.3 Interpretability

Interpretable models are designed to provide insights into how predictions are made, allowing users to understand the relationships between input variables and model

outputs. In the context of transformers, interpretability can be achieved by analyz-
ing the attention weights, feature importance, and other aspects of the model that
contribute to its decision-making process.

Let $\mathbf{X}$ represent the input time series data, and $\hat{y}$ the predicted output. An inter-
pretable model aims to express $\hat{y}$ as a function of the most relevant features in $\mathbf{X}$,
denoted by $\mathbf{X}_{\text{important}}$. The model can be represented as

$$\hat{y} = f(\mathbf{X}_{\text{important}}) + \epsilon,$$

where $f$ is an interpretable function, such as a linear combination of the selected
features, and $\epsilon$ represents the residual error.

Under the assumption that only a subset of the input features $\mathbf{X}$ is relevant to
the prediction, sparsity-inducing techniques such as Lasso regression can be used to
identify $\mathbf{X}_{\text{important}}$. The solution to the Lasso problem

$$\min_{\beta} \left\{ \frac{1}{2N} \sum_{i=1}^{N} (y_i - \mathbf{X}_i \beta)^2 + \lambda \|\beta\|_1 \right\}$$

promotes sparsity in $\beta$, effectively selecting the most important features.

Example: In financial time series forecasting, an interpretable model might iden-
tify key indicators such as moving averages and momentum as the most relevant
features for predicting stock prices, allowing investors to understand the rationale
behind the model's predictions.

**Explaining Model Predictions**

Explaining model predictions involves providing a clear and detailed understand-
ing of why a model made a particular prediction. For transformers, this often
involves analyzing the attention weights and how they are distributed across the
input sequence.

In a transformer model, the attention mechanism assigns a weight $\alpha_{ij}$ to each
element of the input sequence $\mathbf{X}$. The predicted output can be expressed as a weighted
sum of the input features:

$$\hat{y} = \sum_{i=1}^{T} \alpha_{ij} \mathbf{X}_i,$$

where $\alpha_{ij}$ indicates the importance of feature $\mathbf{X}_i$ at position $i$ for the prediction of
$\hat{y}$. The goal of explainability is to understand how these weights $\alpha_{ij}$ influence the
prediction.

Shapley values, derived from cooperative game theory, provide a method to fairly
attribute the contribution of each feature to the prediction. For a set of features
$\mathbf{X} = \{X_1, X_2, \ldots, X_T\}$, the Shapley value for feature $X_i$ is given by

$$\phi_i = \sum_{S \subseteq \mathbf{X} \setminus \{X_i\}} \frac{|S|!(T - |S| - 1)!}{T!} \left[ f(S \cup \{X_i\}) - f(S) \right],$$

where $f(S)$ is the model prediction based on the subset $S$ of features. Shapley values provide a comprehensive measure of each feature's contribution to the prediction.

Example: In healthcare time series, explaining why a model predicts a certain health outcome for a patient might involve analyzing the attention weights and Shapley values for features like heart rate, blood pressure, and temperature, providing doctors with a clear rationale for the prediction.

## 7.6  Challenges and Future Directions

The application of transformers to time series forecasting presents several challenges that need to be addressed to fully realize their potential. This section explores these challenges, including dealing with missing data and handling non-stationarity, and outlines future research directions in time series forecasting with transformers.

### Dealing with Missing Data

Missing data is a common issue in time series, where observations may be missing due to sensor failures, reporting errors, or other factors. Effective handling of missing data is crucial for maintaining model accuracy and robustness.

Let $\mathbf{X}$ represent the time series data, with missing values denoted by $\mathbf{X}_{\text{missing}}$. The objective is to impute these missing values, $\hat{\mathbf{X}}_{\text{missing}}$, such that the imputed series $\mathbf{X}_{\text{imputed}} = \mathbf{X} \cup \hat{\mathbf{X}}_{\text{missing}}$ preserves the underlying data structure.

Imputation Methods:

1. Mean Imputation:

$$\hat{X}_{t,\text{missing}} = \frac{1}{N} \sum_{i=1}^{N} X_{t,i},$$

where the missing value is replaced with the mean of the observed values.

2. Interpolation:

$$\hat{X}_{t,\text{missing}} = \frac{X_{t-1} + X_{t+1}}{2},$$

which estimates the missing value based on neighboring observations.

3. Model-Based Imputation: Impute using a model trained on the observed data:

$$\hat{X}_{t,\text{missing}} = f(\mathbf{X}_{\text{observed}}),$$

where $f$ is a model such as a transformer or Gaussian process.

The choice of imputation method introduces a trade-off between bias and variance. Mean imputation introduces bias by ignoring the temporal structure, while model-based imputation can reduce bias but increase variance due to model uncertainty.

Example: In energy demand forecasting, where sensors may occasionally fail, model-based imputation using a transformer trained on the observed data can provide accurate estimates of the missing values, ensuring that the forecasting model remains robust.

**Handling Non-Stationarity**

Non-stationarity, where the statistical properties of the time series change over time, poses a significant challenge for time series forecasting models. Transformers, like other models, need to be adapted to handle non-stationary data effectively.

A time series $\mathbf{X}$ is non-stationary if its mean, variance, or auto-correlation structure changes over time. The goal is to transform the series into a stationary form, $\mathbf{X}_{\text{stationary}}$, where these properties remain constant, or to adapt the model to handle the non-stationarity directly.

Methods for Handling Non-Stationarity:

1. Differencing: Apply differencing to remove trends:

$$\mathbf{X}_{\text{diff}} = \mathbf{X}_t - \mathbf{X}_{t-1},$$

which often stabilizes the mean.

2. Transformation: Apply a transformation, such as logarithms or Box-Cox, to stabilize the variance:

$$\mathbf{X}_{\text{transformed}} = \log(\mathbf{X}).$$

3. Model Adaptation: Use models that can handle non-stationarity, such as transformers with time-varying parameters or recurrent models with memory mechanisms.

For a stationary time series, the forecasting error tends to be lower due to the consistent statistical properties. Transforming a non-stationary series into a stationary form before modeling can significantly improve forecasting accuracy.

Example: In weather forecasting, where seasonal patterns and trends can lead to non-stationarity, differencing the data or applying seasonal decomposition before feeding it into a transformer can enhance the model's performance.

**Future Research Directions in Time Series Forecasting**

The field of time series forecasting with transformers is rapidly evolving, with several promising research directions. These include the development of more interpretable models, better methods for handling missing and non-stationary data, and the exploration of new architectures and training techniques that further enhance the performance of transformers in time series applications.

Exploration:

1. Interpretable Transformers: Develop transformers that inherently produce interpretable outputs by incorporating feature importance measures, attention visualization, and model simplification techniques.

2. Robust Handling of Missing Data: Explore new imputation techniques that leverage the full power of deep learning, such as generative models and self-supervised learning, to accurately predict missing values in complex time series.

3. Advanced Non-Stationarity Techniques: Investigate the use of non-linear transformations, wavelet decompositions, and adaptive learning rates to better handle non-stationary time series in transformer models.

As time series forecasting continues to grow in importance, the need for models that are both powerful and interpretable will drive research in hybrid approaches, model explainability, and efficiency. In financial markets, future research might focus on developing interpretable transformers that not only predict market movements with high accuracy but also provide clear explanations of the factors driving these predictions, helping traders and analysts make more informed decisions.

# References

1. Bollerslev, T.: Generalized autoregressive conditional heteroskedasticity. J. Econo. **31**(3), 307–327 (1986)
2. Box, G.E.P., Jenkins, G.M.: Time Series Analysis: Forecasting and Control. Holden-Day (1970)
3. Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder-decoder for statistical machine translation (2014). arXiv preprint arXiv:1406.1078
4. Coles, S.: An Introduction to Statistical Modeling of Extreme Values. Springer (2001)
5. Elman, J.L.: Finding structure in time. Cogn. Sci. **14**(2), 179–211 (1990)
6. Fama, E.F.: Efficient capital markets: A review of theory and empirical work. J. Finance **25**(2), 383–417 (1970)
7. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural Comput. **9**(8), 1735–1780 (1997)
8. Hyndman, R.J., Koehler, A.B., Ord, K., Snyder, R.D.: Forecasting with Exponential Smoothing: The State Space Approach. Springer Science & Business Media (2008)
9. Markowitz, H.: Portfolio selection. J. Finance **7**(1), 77–91 (1952)
10. Rumelhart, D.E., Hinton, G.E., Williams, R.J.: Learning representations by back-propagating errors. Nature **323**(6088), 533–536 (1986)

# Chapter 8
# Signal Analysis and Transformers

## 8.1 Signal Processing

Signal processing is a crucial aspect of many modern technologies, including communications, audio processing, and control systems. It also plays a significant role in the application of transformers to tasks involving time series data, such as speech recognition, financial forecasting, and more. This section provides a mathematical treatment of signal processing in the time domain, focusing on fundamental concepts like time series representation, auto-correlation and cross-correlation, and convolution and filtering.

### 8.1.1 Time-Domain Analysis

Time-domain analysis involves examining signals with respect to time, providing insights into their temporal structure and dynamics. This analysis is foundational for understanding how signals evolve and interact over time, which is essential for tasks involving sequential data.

**Time Series Representation**

A time series is a sequence of data points indexed in time order, often representing how a particular quantity evolves over time. Mathematically, a time series can be defined as a function $x : \mathbb{Z} \to \mathbb{R}$ that maps each time index $t \in \mathbb{Z}$ to a real value $x(t)$.

Let $x(t)$ represent a discrete time signal, where $t$ is an integer representing the time index. The time series can be expressed as

$$x = \{x(t)\}_{t\in\mathbb{Z}} = \{x(t_1), x(t_2), \ldots, x(t_n)\},$$

where $t_1, t_2, \ldots, t_n$ are specific time points and $x(t_i)$ is the value of the signal at time $t_i$.

Example: In financial markets, $x(t)$ could represent the closing price of a stock at time $t$, forming a time series that captures the price fluctuations over a period.

A time series $x(t)$ is said to be stationary if its statistical properties, such as mean, variance, and auto-correlation, do not change over time. Formally, $x(t)$ is stationary if

$$\mathbb{E}[x(t)] = \mu, \quad \text{Var}(x(t)) = \sigma^2, \quad \text{Cov}(x(t), x(t+\tau)) = \gamma(\tau),$$

where $\mu$ is the constant mean, $\sigma^2$ is the constant variance, and $\gamma(\tau)$ is the autocovariance function that depends only on the time difference $\tau$.

Application: In signal processing, stationarity is an important assumption that simplifies the analysis and modeling of time series, particularly in forecasting and filtering.

**Auto-correlation and Cross-correlation**

Auto-correlation and cross-correlation are fundamental tools for analyzing the temporal dependencies within a signal or between multiple signals. These measures capture the degree to which a signal correlates with a lagged version of itself or with another signal.

1. Auto-correlation Function (ACF): The auto-correlation function measures the correlation of a signal with a delayed copy of itself as a function of the delay (or lag) $\tau$. For a signal $x(t)$, the auto-correlation at lag $\tau$ is defined as

$$R_{xx}(\tau) = \mathbb{E}[x(t) \cdot x(t+\tau)] = \lim_{T\to\infty} \frac{1}{2T} \int_{-T}^{T} x(t) \cdot x(t+\tau)\, dt.$$

For a discrete time signal, this becomes

$$R_{xx}(\tau) = \sum_{t=-\infty}^{\infty} x(t) \cdot x(t+\tau).$$

2. Cross-correlation Function (CCF): The cross-correlation function measures the similarity between two signals $x(t)$ and $y(t)$ as a function of the lag $\tau$. It is defined as

$$R_{xy}(\tau) = \mathbb{E}[x(t) \cdot y(t+\tau)] = \lim_{T\to\infty} \frac{1}{2T} \int_{-T}^{T} x(t) \cdot y(t+\tau)\, dt.$$

For discrete time signals:

$$R_{xy}(\tau) = \sum_{t=-\infty}^{\infty} x(t) \cdot y(t+\tau).$$

Properties:

1. Symmetry: The auto-correlation function is symmetric, $R_{xx}(\tau) = R_{xx}(-\tau)$.

2. Maximum at Zero Lag: The maximum value of the auto-correlation function occurs at $\tau = 0$.

3. Cross-correlation Symmetry: Cross-correlation satisfies $R_{xy}(\tau) = R_{yx}(-\tau)$.

Example: In speech signal processing, auto-correlation can be used to identify periodicities in the signal, such as the pitch of a voice, while cross-correlation can measure the similarity between two different audio signals, such as a transmitted and received signal in a communication system.

**Convolution and Filtering**

Convolution is a mathematical operation used to combine two signals, often used in the context of filtering, where a signal is modified by another signal (the filter) to achieve a desired effect, such as smoothing, sharpening, or detecting features.

1. The convolution of two signals $x(t)$ and $h(t)$ is defined as

$$y(t) = (x * h)(t) = \int_{-\infty}^{\infty} x(\tau)h(t - \tau)\, d\tau,$$

where $y(t)$ is the output signal. For discrete time signals, the convolution is given by

$$y[n] = (x * h)[n] = \sum_{m=-\infty}^{\infty} x[m] \cdot h[n - m].$$

2. Filtering is the process of modifying or extracting specific components of a signal. A filter $h(t)$ is applied to a signal $x(t)$ via convolution, producing the filtered signal $y(t)$. Depending on the nature of the filter, this operation can emphasize certain frequencies (e.g., low-pass, high-pass filters) or remove noise.

If a system is linear and time invariant (LTI), the output $y(t)$ of the system to an input $x(t)$ is given by the convolution of the input signal with the system's impulse response $h(t)$:

$$y(t) = (x * h)(t).$$

This property makes convolution a fundamental tool in signal processing for analyzing the response of LTI systems.

Example: In image processing, convolution with a Gaussian filter is used to blur an image, reducing noise and details. In audio processing, convolution with an impulse response can simulate the acoustics of different environments, such as concert halls or small rooms.

## 8.1.2   *Frequency-Domain Analysis*

In frequency-domain analysis, a signal is represented as a sum of sinusoids with different frequencies, amplitudes, and phases. This approach allows for the decomposition of a signal into its constituent frequency components, making it easier to analyze and manipulate the signal in various applications, such as filtering and spectral analysis.

**Fourier Transform**

The Fourier transform is a fundamental mathematical tool that converts a time-domain signal into its frequency-domain representation. It expresses the signal as a sum of sinusoids, each with a specific frequency, amplitude, and phase.

For a continuous-time signal $x(t)$, the Fourier transform $X(f)$ is defined as

$$X(f) = \int_{-\infty}^{\infty} x(t)e^{-j2\pi ft}\, dt,$$

where $f$ represents the frequency and $j$ is the imaginary unit.

The inverse Fourier transform, which reconstructs the time-domain signal from its frequency-domain representation, is given by

$$x(t) = \int_{-\infty}^{\infty} X(f)e^{j2\pi ft}\, df.$$

Parseval's theorem states that the total energy of a signal in the time domain is equal to the total energy in the frequency domain:

$$\int_{-\infty}^{\infty} |x(t)|^2\, dt = \int_{-\infty}^{\infty} |X(f)|^2\, df.$$

This theorem provides a powerful connection between the time and frequency domains, ensuring that energy is conserved during the transformation.

Example: For a sinusoidal signal $x(t) = A\cos(2\pi f_0 t + \phi)$, where $A$ is the amplitude, $f_0$ is the frequency, and $\phi$ is the phase, the Fourier transform is

$$X(f) = \frac{A}{2}\left[\delta(f - f_0) + \delta(f + f_0)\right],$$

indicating that the signal's frequency content is concentrated at $\pm f_0$.

**Discrete Fourier Transform (DFT)**

The Discrete Fourier Transform (DFT) is the discrete counterpart of the Fourier transform, applied to sequences of values, such as digital signals sampled at discrete time intervals. The DFT converts a finite sequence of equally spaced samples of a function into a sequence of coefficients of a finite combination of complex sinusoids.

Given a discrete time signal $x[n]$ of length $N$, the DFT $X[k]$ is defined as

$$X[k] = \sum_{n=0}^{N-1} x[n]e^{-j\frac{2\pi}{N}kn},$$

where $k = 0, 1, 2, \ldots, N - 1$.

The inverse DFT, which reconstructs the time-domain sequence from its frequency-domain representation, is given by

$$x[n] = \frac{1}{N} \sum_{k=0}^{N-1} X[k]e^{j\frac{2\pi}{N}kn}.$$

The exponential functions $e^{-j\frac{2\pi}{N}kn}$ used in the DFT are orthogonal over the interval $n = 0, 1, \ldots, N - 1$:

$$\sum_{n=0}^{N-1} e^{-j\frac{2\pi}{N}kn}e^{j\frac{2\pi}{N}mn} = N\delta[k - m],$$

where $\delta[k - m]$ is the Kronecker delta. This orthogonality property is fundamental to the DFT, as it ensures that the frequency components are independent and can be isolated.

Example: For a digital signal sampled at a rate of $N$ points, the DFT provides the frequency content of the signal in terms of $N$ discrete frequency bins. Each bin represents a specific frequency component of the signal.

**Fast Fourier Transform (FFT)**

The Fast Fourier Transform (FFT) ([2]) is an efficient algorithm for computing the DFT of a sequence. The FFT reduces the computational complexity of the DFT from $O(N^2)$ to $O(N \log N)$, making it feasible to compute the DFT of large sequences.

The FFT algorithm exploits the symmetries in the DFT formula to reduce the number of arithmetic operations required. For a sequence of length $N$, where $N = 2^m$ (a power of two), the FFT recursively divides the DFT into smaller DFTs of even-indexed and odd-indexed elements:

$$X[k] = \sum_{n=0}^{N/2-1} \left[ x[2n]e^{-j\frac{2\pi}{N/2}kn} + x[2n+1]e^{-j\frac{2\pi}{N}kn} \right],$$

where the DFTs of the even and odd parts are computed separately and combined.

The FFT algorithm has a time complexity of $O(N \log N)$, which is significantly more efficient than the $O(N^2)$ complexity of the direct DFT computation. This efficiency makes the FFT the algorithm of choice for frequency-domain analysis in practical applications.

Example: In real-time signal processing, such as audio or video processing, the FFT allows for rapid analysis of the frequency content of signals, enabling tasks like filtering, spectral analysis, and even real-time effects processing.

**Wavelet Transform**

The wavelet transform ([3]) analyzes a signal by decomposing it into shifted and scaled versions of a prototype function called a wavelet. Unlike the Fourier transform, which uses sinusoids as basis functions, the wavelet transform uses wavelets, which are localized in both time and frequency. This localization allows for the analysis of signals with time-varying frequency content.

Let $\psi(t)$ be a mother wavelet, a function that is localized in both time and frequency. The wavelet transform of a signal $x(t)$ is defined as

$$W(a, b) = \frac{1}{\sqrt{|a|}} \int_{-\infty}^{\infty} x(t)\psi^* \left( \frac{t - b}{a} \right) dt,$$

where $a$ is the scale parameter, $b$ is the translation (shift) parameter, and $\psi^*(t)$ denotes the complex conjugate of the mother wavelet.

Scale Parameter $a$: It controls the frequency content of the wavelet. A larger $a$ corresponds to a lower frequency (more stretched wavelet), and a smaller $a$ corresponds to a higher frequency (more compressed wavelet).

Translation Parameter $b$: It controls the time localization of the wavelet, shifting it along the time axis.

The wavelet transform decomposes the signal into wavelets at different scales and translations, providing a multi-resolution analysis of the signal. For a wavelet $\psi(t)$ to be admissible (i.e., suitable for reconstructing the original signal from its wavelet transform), it must satisfy the admissibility condition:

$$C_\psi = \int_{-\infty}^{\infty} \frac{|\hat{\psi}(f)|^2}{|f|} df < \infty,$$

where $\hat{\psi}(f)$ is the Fourier transform of $\psi(t)$. This condition ensures that the wavelet has zero mean and is localized in both time and frequency.

Example: The Morlet wavelet, defined as $\psi(t) = \exp(j2\pi f_0 t)\exp\left(-\frac{t^2}{2\sigma^2}\right)$, is a commonly used wavelet in signal processing. It is a complex sinusoid modulated by a Gaussian window, providing good localization in both time and frequency.

## Continuous Wavelet Transform (CWT)

The Continuous Wavelet Transform (CWT) provides a continuous representation of the signal in both time and scale (frequency). It is particularly useful for analyzing signals with smoothly varying time-frequency content.

The CWT of a signal $x(t)$ is given by

$$W(a, b) = \frac{1}{\sqrt{|a|}} \int_{-\infty}^{\infty} x(t)\psi^*\left(\frac{t-b}{a}\right) dt,$$

where $W(a, b)$ is the wavelet coefficient corresponding to scale $a$ and translation $b$. The CWT provides a continuous mapping of the signal into the time-scale plane.

The original signal $x(t)$ can be reconstructed from its wavelet transform via the inverse CWT:

$$x(t) = \frac{1}{C_\psi} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} W(a, b)\psi\left(\frac{t-b}{a}\right) \frac{da\,db}{a^2},$$

where $C_\psi$ is the admissibility constant. This formula shows that the CWT preserves all the information needed to reconstruct the original signal.

Example: The CWT can be used to analyze the time-frequency characteristics of non-stationary signals, such as chirp signals (signals whose frequency increases or decreases with time). The CWT allows for the detection and characterization of such frequency modulations over time.

## Discrete Wavelet Transform (DWT)

The Discrete Wavelet Transform (DWT) is a sampled version of the CWT, typically implemented using a filter bank approach. The DWT provides a hierarchical decomposition of the signal into different frequency bands, making it computationally efficient and suitable for practical applications.

The DWT is computed by sampling the scale and translation parameters in the CWT. Specifically, let $a = 2^j$ and $b = k2^j$, where $j$ and $k$ are integers. The DWT is then given by

$$W[j, k] = \frac{1}{\sqrt{2^j}} \int_{-\infty}^{\infty} x(t)\psi^*\left(\frac{t-k2^j}{2^j}\right) dt,$$

where $W[j, k]$ represents the wavelet coefficient at scale $2^j$ and translation $k2^j$.

The DWT can be efficiently computed using a multi-level filter bank, where the signal is passed through a series of high-pass and low-pass filters, followed by downsampling. This process decomposes the signal into approximation and detail coefficients at each level.

The DWT provides a multi-resolution analysis of the signal, where the signal is decomposed into approximation and detail components at different levels of resolution. Formally, let $V_j$ and $W_j$ represent the approximation and detail spaces at level $j$. Then

$$V_{j+1} = V_j \oplus W_j,$$

where $\oplus$ denotes the direct sum of the spaces. This decomposition allows for the analysis of the signal at different levels of detail.

Example: In image processing, the DWT is used for image compression, such as in the JPEG2000 standard. The image is decomposed into different frequency bands using the DWT, and the high-frequency components (which often correspond to noise or fine details) can be discarded or quantized more coarsely, leading to efficient compression.

## Spectrogram and Time–Frequency Representation

The spectrogram is a fundamental tool in time–frequency analysis, providing a visual representation of how the frequency content of a signal evolves over time. It is computed by applying the Fourier transform to short, overlapping segments of the signal, resulting in a two-dimensional plot with time on one axis and frequency on the other.

Given a signal $x(t)$, the spectrogram $S(t, f)$ is defined as the magnitude squared of the Short-Time Fourier Transform (STFT) of the signal:

$$S(t, f) = \left| \int_{-\infty}^{\infty} x(\tau) w(\tau - t) e^{-j 2\pi f \tau} \, d\tau \right|^2,$$

where $w(\tau)$ is a window function that is localized in time, such as a Gaussian or Hamming window. The window function is used to isolate segments of the signal, allowing the Fourier transform to capture the frequency content within each segment.

The time-frequency uncertainty principle states that there is a trade-off between the time and frequency resolution of the spectrogram. Formally, for a window function $w(t)$, the product of the time duration $\Delta t$ and the bandwidth $\Delta f$ satisfies

$$\Delta t \cdot \Delta f \geq \frac{1}{4\pi}.$$

This inequality implies that improving time resolution comes at the cost of frequency resolution, and vice versa.

Example: In speech processing, the spectrogram is used to analyze the time-varying frequency content of speech signals, capturing formants and other features that are crucial for tasks such as speech recognition and speaker identification.

### 8.1.3 Advanced Signal Processing Techniques

Here, we discuss the basic tools of Fourier and wavelet analysis to address more complex signals, particularly those with non-linear, non-stationary, or multi-component characteristics.

**Hilbert Transform**

The Hilbert transform is a linear operator that shifts the phase of a signal by $90°$, effectively generating a complex-valued signal from a real-valued signal. It is widely used in the analysis of analytic signals, envelope detection, and instantaneous frequency estimation.

The Hilbert transform $\mathcal{H}[x(t)]$ of a signal $x(t)$ is defined as

$$\mathcal{H}[x(t)] = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{x(\tau)}{t - \tau} \, d\tau.$$

The analytic signal $z(t)$ associated with $x(t)$ is then given by

$$z(t) = x(t) + j\mathcal{H}[x(t)] = A(t)e^{j\phi(t)},$$

where $A(t)$ is the amplitude envelope and $\phi(t)$ is the instantaneous phase of the signal.

The instantaneous frequency $f_i(t)$ of a signal $x(t)$ is the time derivative of the instantaneous phase $\phi(t)$ obtained from the analytic signal:

$$f_i(t) = \frac{1}{2\pi} \frac{d\phi(t)}{dt}.$$

This allows for the analysis of frequency variations over time, which is particularly useful in applications such as radar and communication systems.

Example: The Hilbert transform is used in modulation techniques, where the analytic signal helps separate the carrier and modulation components of a signal, facilitating the extraction of amplitude and phase information.

**Short-Time Fourier Transform**

STFT is a modification of the Fourier transform that allows for the analysis of non-stationary signals by dividing the signal into short segments and applying the Fourier transform to each segment. The STFT provides a time-frequency representation of the signal, capturing how its frequency content evolves over time.

The STFT of a signal $x(t)$ with respect to a window function $w(t)$ is defined as

$$X(t, f) = \int_{-\infty}^{\infty} x(\tau)w(\tau - t)e^{-j2\pi f\tau} \, d\tau.$$

The window function $w(t)$ determines the time resolution of the STFT, with narrower windows providing better time resolution but poorer frequency resolution.

The original signal $x(t)$ can be reconstructed from its STFT using the inverse STFT:

$$x(t) = \frac{1}{w(0)} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} X(\tau, f)w(\tau - t)e^{j2\pi f\tau} \, df \, d\tau.$$

This reconstruction formula ensures that the STFT retains all the information needed to reconstruct the original signal.

Example: The STFT is widely used in music analysis to visualize the harmonic structure of audio signals over time. It allows for the identification of notes, chords, and other musical features by analyzing the time-varying frequency content.

**Empirical Mode Decomposition (EMD)**

EMD is a data-driven technique for decomposing a non-stationary signal into a set of intrinsic mode functions (IMFs), each representing a simple oscillatory mode. EMD is particularly useful for analyzing signals with non-linear and non-stationary characteristics. Given a signal $x(t)$, the EMD process involves iteratively extracting IMFs by identifying local maxima and minima, constructing the upper and lower envelopes, and subtracting the mean envelope from the signal. The first IMF $c_1(t)$ is extracted as

$$c_1(t) = x(t) - \text{mean envelope of } x(t).$$

The residual signal $r_1(t) = x(t) - c_1(t)$ is then used to extract the next IMF, and the process continues until the residual is a monotonic function. The signal is decomposed as

$$x(t) = \sum_{i=1}^{n} c_i(t) + r_n(t),$$

where $c_i(t)$ are the IMFs and $r_n(t)$ is the final residual.

The IMFs $c_i(t)$ obtained through EMD are approximately orthogonal, meaning that their inner product is close to zero:

$$\int_{-\infty}^{\infty} c_i(t)c_j(t)\,dt \approx 0 \quad \text{for } i \neq j.$$

This orthogonality ensures that the IMFs represent distinct oscillatory modes of the signal.

Example: EMD is commonly used in biomedical signal processing, such as in the analysis of electroencephalogram (EEG) signals. It allows for the extraction of different brainwave components (e.g., alpha, beta, delta waves) from raw EEG data, facilitating the study of brain activity.

## 8.2 Transformers in Audio and Speech Processing

The application of transformers in audio and speech processing leverages their powerful sequence-to-sequence modeling capabilities and attention mechanisms to capture complex temporal dependencies in audio signals. This section provides a mathematical formulation of how transformers are applied to audio and speech tasks, focusing on sequence-to-sequence models and the role of attention mechanisms.

### 8.2.1 Mathematical Formulation

Audio and speech processing tasks often involve transforming an input audio signal into a corresponding output sequence, such as in speech recognition, where the goal is to transcribe spoken language into text. Transformers, with their ability to model long-range dependencies and capture contextual information through attention mechanisms, are particularly well suited for these tasks.

**Sequence-to-Sequence Models for Audio**

In sequence-to-sequence (Seq2Seq) models, the goal is to map an input sequence of audio features to an output sequence, such as phonemes, words, or characters. The transformer architecture, originally designed for natural language processing, has been adapted for audio processing by modeling the input audio as a sequence of feature vectors.

Let $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$ represent the input sequence of audio feature vectors, where $\mathbf{x}_t \in \mathbb{R}^d$ is the feature vector at time step $t$ and $T$ is the length of the sequence. The output sequence $\mathbf{y} = (y_1, y_2, \ldots, y_{T'})$ consists of tokens such as

phonemes or words, where $T'$ is the length of the output sequence. The transformer encoder processes the input sequence to produce a sequence of hidden states $\mathbf{H} = (\mathbf{h}_1, \mathbf{h}_2, \ldots, \mathbf{h}_T)$, where each hidden state $\mathbf{h}_t$ captures contextual information from the entire input sequence:

$$\mathbf{H} = \text{Encoder}(\mathbf{x}).$$

The decoder then generates the output sequence by attending to the encoded hidden states and producing one token at a time:

$$y_t = \text{Decoder}(y_{<t}, \mathbf{H}),$$

where $y_{<t}$ denotes the tokens generated so far.

The sequence-to-sequence transformer models the conditional probability of the output sequence given the input sequence as a product of conditional probabilities:

$$P(\mathbf{y} \mid \mathbf{x}) = \prod_{t=1}^{T'} P(y_t \mid y_{<t}, \mathbf{H}),$$

where each conditional probability $P(y_t \mid y_{<t}, \mathbf{H})$ is modeled by the decoder. This formulation captures the dependency of each output token on both the input sequence and the previously generated tokens.

Example: In automatic speech recognition (ASR), the input sequence $\mathbf{x}$ consists of acoustic features extracted from the audio waveform, and the output sequence $\mathbf{y}$ is the transcription of the spoken words. The transformer-based ASR model maps the acoustic features to text, capturing long-range dependencies in both the audio signal and the language model.

**Attention Mechanisms in Audio Processing**

Attention mechanisms are central to the success of transformers in audio processing. They allow the model to dynamically focus on different parts of the input sequence when generating each token of the output sequence, enabling the capture of relevant temporal dependencies and contextual information.

The attention mechanism computes a weighted sum of the encoder hidden states, where the weights are determined by the similarity between the decoder's current state and each encoder hidden state. For an output token $y_t$, the attention mechanism is defined as

$$\mathbf{c}_t = \sum_{i=1}^{T} \alpha_{t,i} \mathbf{h}_i,$$

where $\mathbf{c}_t$ is the context vector and $\alpha_{t,i}$ are the attention weights given by

$$\alpha_{t,i} = \frac{\exp(e_{t,i})}{\sum_{j=1}^{T} \exp(e_{t,j})},$$

with the alignment score $e_{t,i}$ computed as

$$e_{t,i} = \text{Score}(\mathbf{s}_{t-1}, \mathbf{h}_i),$$

where $\mathbf{s}_{t-1}$ is the decoder's state from the previous time step.

The attention mechanism can be interpreted as a soft alignment between the input and output sequences. The attention weights $\alpha_{t,i}$ sum to one and provide a probability distribution over the input sequence, indicating how much each input time step contributes to the generation of the current output token.

Example: In audio-based tasks like speech translation, the attention mechanism allows the model to focus on the relevant portions of the input audio corresponding to the current part of the translation. This enables the model to handle variations in speaking speed, pauses, and other temporal aspects of speech more effectively than traditional methods.

## 8.2.2   Speech Recognition

Speech recognition is the process of converting spoken language into text. In this context, transformers serve as powerful models that can capture the intricate temporal and contextual dependencies in speech, making them ideal for ASR systems.

### Automatic Speech Recognition (ASR) Systems

Traditional ASR systems are typically composed of multiple components, including acoustic models, language models, and pronunciation models. These components work together to transcribe speech into text. The acoustic model maps audio features to phonetic representations, the language model captures the probability of word sequences, and the pronunciation model maps phonetic sequences to words.

Let $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$ be the sequence of acoustic feature vectors extracted from an audio signal, and $\mathbf{y} = (y_1, y_2, \ldots, y_{T'})$ be the corresponding sequence of words or phonemes. The goal of an ASR system is to find the most probable word sequence $\mathbf{y}^*$ given the acoustic features:

$$\mathbf{y}^* = \arg\max_{\mathbf{y}} P(\mathbf{y} \mid \mathbf{x}),$$

where $P(\mathbf{y} \mid \mathbf{x})$ is the posterior probability of the word sequence given the acoustic features.

This probability is typically decomposed using Bayes' theorem into an acoustic model and a language model:

$$P(\mathbf{y} \mid \mathbf{x}) = P(\mathbf{x} \mid \mathbf{y}) P(\mathbf{y}),$$

where $P(\mathbf{x} \mid \mathbf{y})$ is the likelihood from the acoustic model and $P(\mathbf{y})$ is the prior probability from the language model.

In traditional ASR systems, the Viterbi algorithm is used to find the most likely sequence of hidden states (e.g., phonemes) that generate the observed acoustic features. This algorithm optimizes the joint probability by considering the most likely path through the state space, defined by the hidden Markov model (HMM):

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \prod_{t=1}^{T} P(x_t \mid y_t) P(y_t \mid y_{t-1}),$$

where $P(y_t \mid y_{t-1})$ represents the transition probabilities between states.

Example: In a traditional ASR system, the acoustic model might be a Gaussian mixture model (GMM) that estimates the probability distribution of acoustic features given phonemes, while the language model could be an n-gram model that captures the likelihood of word sequences.

### End-to-End ASR Models

End-to-end ASR models, such as those based on transformers, simplify the traditional ASR pipeline by directly mapping acoustic features to text. These models learn to jointly optimize the acoustic, pronunciation, and language models within a single neural network, eliminating the need for separate components.

In an end-to-end ASR model, the transformer encoder processes the input sequence of acoustic features $\mathbf{x}$ to produce a sequence of hidden states $\mathbf{H}$:

$$\mathbf{H} = \text{Encoder}(\mathbf{x}).$$

The decoder then generates the output text sequence $\mathbf{y}$ by attending to the hidden states $\mathbf{H}$:

$$P(\mathbf{y} \mid \mathbf{x}) = \prod_{t=1}^{T'} P(y_t \mid y_{<t}, \mathbf{H}),$$

where each probability $P(y_t \mid y_{<t}, \mathbf{H})$ is modeled by the decoder using attention mechanisms.

A common approach in end-to-end ASR models is Connectionist Temporal Classification (CTC), which allows for alignment-free training by introducing a blank label $\phi$ that represents no output:

$$P(\mathbf{y} \mid \mathbf{x}) = \sum_{\pi \in \text{align}(\mathbf{y})} P(\pi \mid \mathbf{x}),$$

where $\pi$ represents possible alignments of $\mathbf{y}$ with $\mathbf{x}$. The CTC loss function maximizes the probability of the correct label sequence by summing over all valid alignments.

Example: In an end-to-end transformer-based ASR system, the acoustic features might be fed into a transformer encoder, which learns contextual representations of the speech signal. The decoder, equipped with self-attention, generates the transcription by focusing on relevant parts of the encoded signal.

**Evaluation Metrics (WER, CER)**

Evaluating the performance of ASR systems requires metrics that quantify the accuracy of the transcribed text. The two most common metrics are Word Error Rate (WER) and Character Error Rate (CER).

1. Word Error Rate (WER): WER is defined as the ratio of the sum of insertion ($I$), deletion ($D$), and substitution ($S$) errors to the total number of words in the reference transcription $N$:

$$\text{WER} = \frac{S + D + I}{N}.$$

The WER captures how many words in the transcription are incorrect, missing, or extra compared to the reference.

2. Character Error Rate (CER): CER is similar to WER but is calculated at the character level, making it more sensitive to small errors:

$$\text{CER} = \frac{S_c + D_c + I_c}{N_c},$$

where $S_c$, $D_c$, $I_c$, and $N_c$ represent the substitution, deletion, insertion, and total number of characters, respectively.

The WER and CER can be computed using the Levenshtein distance, which measures the minimum number of single-character edits (insertions, deletions, or substitutions) required to change one string into another:

$$\text{Levenshtein}(a, b) = \min\{\text{edit operations to convert } a \text{ to } b\}.$$

This distance directly relates to the error rates, as it quantifies the discrepancy between the predicted and reference transcriptions.

Example: For a speech recognition system transcribing a spoken sentence, WER might capture the overall accuracy of word recognition, while CER would reflect finer

errors such as incorrect letters or misspellings. These metrics guide the optimization of ASR systems by providing a quantitative measure of transcription quality.

### 8.2.3  Audio Classification

Audio classification tasks involve categorizing audio signals into predefined classes. This could range from identifying the genre of a music track to recognizing environmental sounds or identifying speakers. Transformers, with their ability to capture temporal dependencies and contextual information, are well suited for these tasks.

#### Environmental Sound Classification

Environmental sound classification involves identifying and categorizing sounds from the environment, such as footsteps, rain, or traffic noise. The goal is to map an audio signal to one of several predefined classes.

Let $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$ represent the sequence of acoustic features extracted from an audio signal. The task is to predict the class $y$ of the sound, where $y$ belongs to a set of predefined categories $\mathcal{C}$.

The transformer model processes the input sequence through its encoder, generating a sequence of hidden states $\mathbf{H} = (\mathbf{h}_1, \mathbf{h}_2, \ldots, \mathbf{h}_T)$. The final classification decision is made by aggregating these hidden states, often using a pooling mechanism:

$$\hat{y} = \text{softmax}\left(\text{Pooling}(\mathbf{H})\right),$$

where $\hat{y}$ is the predicted probability distribution over the classes, and the pooling operation (e.g., mean or max-pooling) reduces the sequence of hidden states to a single vector.

Global pooling ensures that the classification decision is invariant to the length of the input sequence, making the model robust to variations in the duration of environmental sounds. Formally, if Pooling($\mathbf{H}$) is invariant under permutations of the sequence, then

$$\hat{y}(\mathbf{x}) = \hat{y}(\sigma(\mathbf{x})),$$

for any permutation $\sigma$ of the input sequence, ensuring consistent classification regardless of the sequence order.

Example: In an environmental sound classification task, the transformer might be trained on a dataset containing various environmental sounds. The model learns to recognize patterns in the audio features that are characteristic of different sound classes, such as the frequency patterns associated with rain versus those of traffic noise.

**Music Genre Classification**

Music genre classification aims to categorize music tracks into genres such as jazz, classical, or rock. This task involves capturing both the temporal structure and harmonic content of the music, which transformers can effectively model.

Similar to environmental sound classification, the input is a sequence of acoustic features $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$, extracted from the music track. The transformer processes this sequence to produce hidden states $\mathbf{H} = (\mathbf{h}_1, \mathbf{h}_2, \ldots, \mathbf{h}_T)$, which are then aggregated for classification:

$$\hat{y} = \text{softmax} \left( \text{Pooling}(\mathbf{H}) \right),$$

where $\hat{y}$ represents the predicted genre distribution.

Transformers can capture hierarchical features in music, from low-level acoustic features (e.g., timbre, rhythm) to high-level semantic features (e.g., style, genre). The hierarchical nature of transformers allows the model to build complex representations of the music, making it effective for genre classification:

$$\mathbf{h}_t^l = \text{Layer}^l(\mathbf{h}_t^{l-1}),$$

where $\mathbf{h}_t^l$ represents the hidden state at layer $l$, capturing progressively more abstract features as $l$ increases.

Example: A transformer-based music genre classifier might be trained on a large dataset of music tracks labeled by genre. The model learns to recognize patterns that are indicative of different genres, such as the rhythmic structure of hip-hop or the harmonic richness of classical music.

**Speaker Identification**

Speaker identification involves recognizing and verifying the identity of a speaker based on their voice. This task requires the model to capture unique vocal characteristics that distinguish one speaker from another.

Given an input sequence of acoustic features $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$ extracted from a speech signal, the goal is to identify the speaker $y$ from a set of known speakers $\mathcal{S}$. The transformer processes the input sequence to produce hidden states $\mathbf{H} = (\mathbf{h}_1, \mathbf{h}_2, \ldots, \mathbf{h}_T)$, which are then used for classification:

$$\hat{y} = \text{softmax} \left( \text{Pooling}(\mathbf{H}) \right).$$

The transformer generates consistent speaker embeddings that are robust to variations in speech content, ensuring that the model accurately identifies the speaker regardless of what is being said:

$$\hat{y}(\mathbf{x}) = \hat{y}(\mathbf{x}'),$$

where $\mathbf{x}$ and $\mathbf{x}'$ are different utterances by the same speaker, and $\hat{y}(\mathbf{x})$ and $\hat{y}(\mathbf{x}')$ are the corresponding predicted speaker identities.

Example: In a speaker identification system, the transformer might be trained on speech samples from multiple speakers. The model learns to recognize speaker-specific characteristics, such as vocal pitch, timbre, and speaking style, which are used to distinguish between different speakers.

### 8.2.4   Speech Synthesis and Enhancement

Speech synthesis and enhancement involve generating or improving speech signals. Transformers are increasingly being used in these tasks due to their ability to model complex temporal dependencies and capture the nuances of natural speech.

#### Text-to-Speech (TTS) Systems

TTS systems convert written text into spoken language. Transformers can be used to model the mapping from text sequences to acoustic features, which are then converted into speech waveforms.

Given an input text sequence $\mathbf{t} = (t_1, t_2, \ldots, t_N)$, where $t_i$ represents a token (e.g., a word or phoneme), the transformer generates a sequence of acoustic features $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$, which are then converted to speech:

$$\mathbf{x} = \text{Decoder}(\mathbf{t}).$$

The generated speech signal $s(t)$ is then reconstructed from the acoustic features using a vocoder or waveform synthesis model.

In TTS systems, the attention mechanism aligns the input text with the generated acoustic features, ensuring that each part of the text is correctly mapped to the corresponding speech segment:

$$\mathbf{c}_t = \sum_{i=1}^{N} \alpha_{t,i} \mathbf{h}_i,$$

where $\alpha_{t,i}$ are the attention weights that align the input text token $t_i$ with the generated acoustic feature $\mathbf{x}_t$.

Example: A transformer-based TTS system might generate natural-sounding speech by learning to map text sequences to acoustic features that capture the prosody, intonation, and rhythm of speech. The attention mechanism ensures that the spoken output corresponds accurately to the input text.

**Speech Enhancement and Noise Reduction**

Speech enhancement aims to improve the quality and intelligibility of speech signals, particularly in noisy environments. Transformers can be used to model the mapping from noisy speech to clean speech, effectively reducing background noise while preserving the speech content.

Given a noisy speech signal $\mathbf{x}_{\text{noisy}} = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$, the goal is to generate a clean speech signal $\mathbf{x}_{\text{clean}} = (\mathbf{x}'_1, \mathbf{x}'_2, \ldots, \mathbf{x}'_T)$. The transformer processes the noisy signal to produce a denoised version:

$$\mathbf{x}_{\text{clean}} = \text{Transformer}(\mathbf{x}_{\text{noisy}}).$$

The attention mechanism in the transformer helps focus on the speech components of the signal while ignoring or suppressing the noise components:

$$\mathbf{c}_t = \sum_{i=1}^{T} \alpha_{t,i} \mathbf{h}_i,$$

where the attention weights $\alpha_{t,i}$ prioritize the parts of the input that contain speech, leading to effective noise suppression.

Example: In a speech enhancement system, a transformer might be trained on pairs of noisy and clean speech signals. The model learns to distinguish between speech and noise, effectively enhancing the speech signal while reducing or eliminating the noise.

## 8.3 Applications and Analysis

The evaluation of audio and speech processing systems, particularly those involving transformers, requires robust performance metrics that quantify the quality and effectiveness of the models. This section focuses on two key metrics: Signal-to-Noise Ratio (SNR) and Perceptual Evaluation of Speech Quality (PESQ). These metrics provide quantitative and perceptual assessments of the audio signals processed by the models, ensuring that the systems meet the desired performance standards.

### 8.3.1 Performance Metrics

Performance metrics in audio and speech processing are essential for evaluating how well a model performs in enhancing, recognizing, or synthesizing speech. These metrics help quantify improvements in signal quality and intelligibility, allowing for comparison and optimization of different models.

**Signal-to-Noise Ratio (SNR)**

SNR is a fundamental metric used to measure the quality of a signal relative to the background noise. It quantifies how much stronger the signal is compared to the noise, providing a direct measure of the effectiveness of noise reduction or speech enhancement techniques.

Given a clean signal $x(t)$ and a noisy signal $y(t) = x(t) + n(t)$, where $n(t)$ represents the noise, the SNR is defined as

$$\text{SNR} = 10 \log_{10} \left( \frac{\sum_{t=1}^{T} x(t)^2}{\sum_{t=1}^{T} n(t)^2} \right) \text{ dB},$$

where $T$ is the duration of the signal. The numerator represents the power of the clean signal and the denominator represents the power of the noise.

If a speech enhancement algorithm successfully reduces the noise $n(t)$ to $n'(t)$, resulting in an enhanced signal $y'(t) = x(t) + n'(t)$, the improvement in SNR, denoted as $\Delta\text{SNR}$, is given by

$$\Delta\text{SNR} = 10 \log_{10} \left( \frac{\sum_{t=1}^{T} n(t)^2}{\sum_{t=1}^{T} n'(t)^2} \right) \text{ dB}.$$

This measure indicates how much the noise has been reduced relative to the original noisy signal.

Example: In a speech enhancement system, SNR is used to evaluate how effectively the system has suppressed background noise while preserving the speech signal. An increase in SNR after processing indicates better noise reduction.

**Perceptual Evaluation of Speech Quality (PESQ)**

While SNR provides a quantitative measure of signal quality, it does not fully capture the perceptual aspects of speech quality as experienced by human listeners. PESQ is a widely used metric that models human auditory perception to evaluate the quality of speech signals.

PESQ compares a degraded speech signal $y(t)$ with a reference clean signal $x(t)$ using a perceptual model that considers factors such as loudness, time alignment, and masking effects. The PESQ score is computed as

$$\text{PESQ} = f(x(t), y(t)),$$

where $f(\cdot)$ represents the perceptual model that transforms the signals into an internal representation, compares them, and produces a quality score. The PESQ score typically ranges from –0.5 to 4.5, with higher scores indicating better speech quality.

The PESQ metric is designed to be consistent with human judgments of speech quality. Mathematically, this means that if a speech enhancement algorithm improves the perceptual quality of speech, the PESQ score should increase accordingly:

$$\text{PESQ}(x(t), y'(t)) > \text{PESQ}(x(t), y(t)),$$

where $y(t)$ is the original degraded signal and $y'(t)$ is the enhanced signal.

Example: In evaluating a TTS system, PESQ can be used to assess how natural and intelligible the generated speech sounds to human listeners. A high PESQ score indicates that the synthetic speech is of high quality and closely resembles natural human speech.

## 8.3.2 Model Complexity

The complexity of transformer models, particularly in terms of time and space, plays a significant role in their scalability and efficiency. Understanding these complexities allows for better design choices and comparisons with traditional models.

### Time Complexity

Time complexity refers to the amount of time required for a transformer model to process an input sequence. For transformers, this is heavily influenced by the attention mechanism, which requires computing pairwise interactions between elements of the input sequence.

Given an input sequence of length $T$, the self-attention mechanism computes attention scores between every pair of elements, leading to a time complexity of

$$\text{Time Complexity} = O(T^2 \cdot d),$$

where $d$ is the dimensionality of the input representations. This quadratic dependence on the sequence length can become a bottleneck, especially for long sequences.

Various strategies have been proposed to reduce the time complexity of transformers, particularly for long sequences. For instance, sparse attention mechanisms and low-rank approximations can reduce the complexity to sub-quadratic levels:

$$\text{Time Complexity (Efficient Transformers)} = O(T \cdot d \cdot \log T).$$

These strategies allow transformers to scale more effectively to long input sequences without sacrificing too much performance.

In speech recognition tasks, where input sequences can be quite long (e.g., a full sentence or paragraph of spoken text), reducing time complexity is crucial for real-time processing. Efficient transformer variants help achieve faster inference times while maintaining accuracy.

## Space Complexity

Space complexity refers to the amount of memory required by a transformer model, which is also influenced by the attention mechanism due to the storage of attention scores and intermediate representations.

For an input sequence of length $T$ and dimensionality $d$, the space complexity of the self-attention mechanism is

$$\text{Space Complexity} = O(T^2 \cdot d + T \cdot d^2).$$

The first term corresponds to the storage of the attention scores and the second term corresponds to the storage of the input and output representations.

To mitigate the space complexity, memory-efficient transformer variants use techniques such as reversible layers, which allow for intermediate activations to be recomputed during backpropagation rather than stored:

$$\text{Space Complexity (Memory-Efficient Transformers)} = O(T \cdot d + d^2).$$

These techniques significantly reduce the memory footprint, enabling the training of larger models on resource-constrained hardware.

In large-scale speech synthesis tasks, where models must generate high-quality audio, reducing space complexity is essential to handle the extensive computations involved, especially when using high-dimensional representations.

## Comparison with Traditional Models

Traditional models for audio and speech processing, such as Hidden Markov Models (HMMs) and Gaussian Mixture Models (GMMs) ([1, 4, 5]), have different complexity profiles compared to transformers. Understanding these differences helps in evaluating the trade-offs between model expressivity and computational efficiency.

For traditional HMM-based ASR systems, the time complexity is primarily linear with respect to the length of the input sequence:

$$\text{Time Complexity (HMM)} = O(T \cdot S^2),$$

where $S$ is the number of states in the HMM. The space complexity is similarly dependent on the number of states and the dimensionality of the feature space.

Transformers, despite their higher time and space complexity, offer greater expressivity due to their ability to model long-range dependencies and context, which traditional models often struggle with. The trade-off is therefore between the computational cost and the representational power:

$$\text{Expressivity (Transformers)} > \text{Expressivity (HMMs)},$$

$$\text{Complexity (Transformers)} > \text{Complexity (HMMs)}.$$

In speaker identification tasks, transformers can leverage contextual information over long sequences, capturing nuances that traditional models might miss. However, this comes at the cost of increased computational resources, necessitating careful consideration of the application requirements.

### 8.3.3 Optimization and Training Strategies

The performance of transformer models in audio and speech processing heavily depends on the optimization and training strategies employed. These include the choice of loss functions, regularization techniques, and data augmentation methods tailored to the specific challenges of audio data.

**Loss Functions for Audio and Speech**

The choice of loss function is critical in training transformer models for audio and speech tasks, as it directly influences the learning dynamics and the model's ability to generalize.

For speech recognition tasks, the most commonly used loss function is the Connectionist Temporal Classification (CTC) loss, which handles sequence alignment issues:

$$\mathcal{L}_{\text{CTC}} = -\log \sum_{\pi \in \text{align}(\mathbf{y})} P(\pi \mid \mathbf{x}),$$

where $\pi$ represents possible alignments of the output sequence $\mathbf{y}$ with the input sequence $\mathbf{x}$.

In speech synthesis, the Mean Squared Error (MSE) between the predicted and target acoustic features is commonly used:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{T} \sum_{t=1}^{T} (x_t - \hat{x}_t)^2,$$

where $\hat{x}_t$ is the predicted feature and $x_t$ is the target.

The CTC loss function is designed to converge even when the alignment between input and output sequences is unknown, making it particularly effective for end-to-end ASR models:

$$\lim_{t \to \infty} \mathcal{L}_{\text{CTC}} = 0 \quad \text{if the model learns the correct alignment.}$$

Example: In training a transformer-based ASR model, CTC loss helps the model learn to align audio frames with corresponding phonemes or words, leading to accurate transcriptions even without explicit alignment information.

## Regularization Techniques

Regularization is crucial for preventing overfitting, especially in large transformer models that are prone to learning spurious patterns in the data.

Dropout is a widely used regularization technique where units in the neural network are randomly dropped during training:

$$\mathbf{h}_t^{(l)} = \text{Dropout}(\mathbf{h}_t^{(l)}),$$

where $\mathbf{h}_t^{(l)}$ is the hidden state at layer $l$. Dropout forces the model to learn more robust features that do not rely on specific neurons.

Weight decay is another regularization technique that penalizes large weights by adding a term to the loss function:

$$\mathcal{L}_{\text{reg}} = \mathcal{L} + \lambda \|\mathbf{W}\|_2^2,$$

where $\mathbf{W}$ represents the model weights and $\lambda$ is a regularization parameter.

Regularization techniques like dropout and weight decay improve the generalization ability of transformer models by preventing overfitting:

$$\text{Generalization Error (Regularized Model)} < \text{Generalization Error (Unregularized Model).}$$

Example: In speech enhancement tasks, regularization helps the model generalize across different noise conditions, ensuring that it performs well not only on the training data but also on unseen noisy environments.

## Data Augmentation for Audio

Data augmentation is a technique used to artificially increase the size of the training dataset by applying transformations to the original data. This is particularly important in audio tasks, where annotated data can be scarce.

Common data augmentation techniques for audio include time-stretching, pitch-shifting, and adding background noise:

$$\mathbf{x}' = \text{Augment}(\mathbf{x}),$$

where $\mathbf{x}$ is the original audio signal and $\mathbf{x}'$ is the augmented version.

For example, time-stretching changes the speed of the audio without altering its pitch:

$$\mathbf{x}'(t) = \mathbf{x}(\alpha t),$$

where $\alpha$ is a time-stretching factor.

Data augmentation enhances the robustness of transformer models by exposing them to a wider variety of input conditions during training:

$$\text{Robustness (Augmented Model)} > \text{Robustness (Non-Augmented Model)}.$$

Example: In music genre classification, applying data augmentation techniques such as pitch-shifting allows the model to recognize genres even when the pitch of the music varies, improving the model's robustness and generalization.

## 8.4   Further Topics

As the field of audio and speech processing continues to evolve, researchers are exploring advanced topics that push the boundaries of current methodologies. This section examines hybrid models that combine transformers with traditional signal processing techniques, as well as the development of efficient transformers specifically tailored for signal processing tasks. These approaches aim to enhance the capabilities of transformers while addressing their computational challenges.

### 8.4.1   Hybrid Models

Hybrid models seek to leverage the strengths of both transformers and traditional signal processing techniques, creating a synergy that can address the limitations of each individual approach. This section provides a mathematical formulation of such hybrid models and explores their potential applications in audio and speech processing.

**Combining Transformers with Traditional Signal Processing Techniques**

Traditional signal processing techniques, such as Fourier transforms, wavelet transforms, and filtering methods, have long been used to analyze and manipulate audio signals. These techniques are well understood, computationally efficient, and effective for specific tasks, but they may lack the representational power needed to capture complex patterns in data. Transformers, on the other hand, excel at modeling long-range dependencies and capturing contextual information but can be computationally intensive.

A hybrid model combines a traditional signal processing technique $\mathcal{T}$ with a transformer model $\mathcal{M}$, where $\mathcal{T}$ preprocesses the input signal $\mathbf{x}$, and $\mathcal{M}$ processes the transformed signal to produce the final output:

$$\mathbf{z} = \mathcal{T}(\mathbf{x}),$$

$$\mathbf{y} = \mathcal{M}(\mathbf{z}),$$

where $\mathbf{z}$ represents the signal processed by $\mathcal{T}$ (e.g., a frequency-domain representation) and $\mathbf{y}$ is the output of the transformer (e.g., a classification or prediction).

For a hybrid model to be effective, the traditional signal processing technique $\mathcal{T}$ must preserve the essential features of the signal that are relevant for the transformer's task. This can be formalized as a preservation theorem, where $\mathbf{z}$ retains the key information from $\mathbf{x}$:

$$\exists f : \mathbf{y} = f(\mathbf{x}) \implies \exists g : \mathbf{y} = g(\mathbf{z}),$$

where $f$ and $g$ are mappings that relate the original signal $\mathbf{x}$ and the transformed signal $\mathbf{z}$ to the final output $\mathbf{y}$.

Example: In speech enhancement, a hybrid model might first apply a wavelet transform to the input speech signal to capture time-frequency characteristics, followed by a transformer that models the temporal dependencies in the wavelet coefficients. This approach leverages the wavelet transform's ability to localize features in time and frequency while utilizing the transformer's capacity for context-aware processing.

**Mathematical Formulation of Hybrid Models**

The mathematical foundation of hybrid models requires a careful integration of signal processing techniques with transformer architectures. The challenge lies in ensuring that the output of the signal processing stage is compatible with the input requirements of the transformer.

Let $\mathbf{x} \in \mathbb{R}^T$ be an input signal (e.g., a time-domain audio signal). A traditional signal processing technique $\mathcal{T}$ transforms $\mathbf{x}$ into a feature space $\mathbb{R}^F$:

$$\mathbf{z} = \mathcal{T}(\mathbf{x}) \in \mathbb{R}^F,$$

where $F$ depends on the nature of the transformation (e.g., $F$ might represent the number of frequency bins in a Fourier transform).

The transformer model $\mathcal{M}$ then processes $\mathbf{z}$ through its layers, each of which is a combination of linear transformations, self-attention mechanisms, and non-linear activations:

$$\mathbf{h}^{(l)} = \text{Attention}\left(\mathbf{W}_q^{(l)}\mathbf{z}, \mathbf{W}_k^{(l)}\mathbf{z}, \mathbf{W}_v^{(l)}\mathbf{z}\right) + \mathbf{b}^{(l)},$$

$$\mathbf{y} = \mathcal{M}(\mathbf{z}) = \mathbf{W}_o\mathbf{h}^{(L)} + \mathbf{b}_o,$$

where $\mathbf{W}_q^{(l)}, \mathbf{W}_k^{(l)}, \mathbf{W}_v^{(l)}, and\,\mathbf{W}_o$ are the weight matrices at the $l$-th layer of the transformer, and $L$ is the total number of layers.

For the hybrid model to function effectively, the output space of the signal processing technique $\mathcal{T}$ must be compatible with the input space of the transformer model $\mathcal{M}$. This compatibility can be expressed as

$$\exists\,\mathcal{T}, \mathcal{M} \text{ such that } \mathcal{M}(\mathcal{T}(\mathbf{x})) = \mathbf{y},$$

where $\mathbf{y}$ is the desired output (e.g., a classification label or enhanced signal). This condition ensures that the features extracted by $\mathcal{T}$ are suitable for the transformer to process.

Example: In music genre classification, a hybrid model might first apply a Fourier transform to the audio signal, converting it into a spectrogram. The transformer then processes the spectrogram, capturing the temporal patterns in the frequency domain that correspond to different musical genres. The Fourier transform reduces the input complexity by focusing on the frequency content, while the transformer captures the genre-specific temporal dependencies.

### 8.4.2 Efficient Transformers for Signal Processing

Transformers, while powerful, are computationally intensive, particularly for long sequences common in audio processing. Efficient transformers aim to reduce the computational burden while maintaining performance, making them more practical for real-time and large-scale applications.

Efficient transformers typically reduce the time and space complexity of the self-attention mechanism by approximating the full attention matrix or by using sparsity in the attention computation. One approach is to replace the full attention mechanism with a sparse attention mechanism:

$$\mathbf{A}_{\text{sparse}} = \text{SparseAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}),$$

where $\mathbf{A}_{\text{sparse}}$ is a sparse matrix that approximates the full attention matrix $\mathbf{A}$ by only considering a subset of the key-value pairs.

The efficiency of sparse attention mechanisms comes at the cost of some approximation error. The error introduced by using sparse attention instead of full attention can be bounded as

$$\|\mathbf{A} - \mathbf{A}_{\text{sparse}}\|_F \leq \epsilon \|\mathbf{A}\|_F,$$

where $\epsilon$ is a small positive constant representing the approximation error and $\|\cdot\|_F$ denotes the Frobenius norm.

Example: In real-time speech recognition systems, efficient transformers might use locality-sensitive hashing (LSH) to approximate the attention mechanism, reducing the time complexity from $O(T^2)$ to $O(T \log T)$. This allows the model to process longer sequences in real time without compromising the accuracy of the recognition.

# References

1. Bishop, C.M.: Pattern Recognition and Machine Learning. Springer (2006)
2. Cooley, J.W., Tukey, J.W.: An algorithm for the machine calculation of complex fourier series. Math. Comput. **19**(90), 297–301 (1965)
3. Daubechies, I.: The wavelet transform, time-frequency localization and signal analysis. IEEE Trans. Inf. Theory **36**(5), 961–1005 (1990)
4. Dempster, A.P., Laird, N.M., Rubin, D.B.: Maximum likelihood from incomplete data via the em algorithm. J. Royal Stat. Soc. Ser. B (Methodological) **39**(1), 1–22 (1977)
5. Rabiner, L.R.: A tutorial on hidden markov models and selected applications in speech recognition. Proceed. IEEE **77**(2), 257–286 (1989). https://doi.org/10.1109/5.18626

# Chapter 9
# Advanced Topics and Future Directions

## 9.1 Memory Optimization

As transformer models grow in complexity and size, memory optimization becomes a critical concern. Efficiently managing memory resources is essential for training large-scale models, ensuring that they remain computationally feasible without sacrificing performance. This section delves into the mathematical principles behind efficient memory storage and retrieval, particularly in the context of large-scale models, and explores the practical implementation of sparse transformers as a means of reducing memory usage.

### 9.1.1 Efficient Storage and Retrieval

Memory efficiency in transformer models involves not only reducing the memory footprint during training but also ensuring that the retrieval of stored information is fast and precise. This requires a careful balance between memory usage, computational complexity, and the fidelity of the stored information.

Consider a transformer model with $L$ layers, each layer producing hidden states $\mathbf{H}^{(l)}$ of size $n \times d$, where $n$ is the sequence length and $d$ is the dimensionality. The total memory required for storing these hidden states is

$$\text{Memory} = L \times n \times d \times \text{sizeof(datatype)},$$

where the datatype could be FP32, FP16, or another precision format.

Efficient Storage Techniques:

1. Gradient checkpointing is a technique that reduces memory usage by selectively storing intermediate activations during the forward pass and recomputing them during the backward pass. Let $\mathcal{C}$ be the set of layers for which activations are stored, and

$\mathcal{R} = L - \mathcal{C}$ be the layers for which activations are recomputed. The total memory required with checkpointing is

$$\text{Memory}_{\text{checkpoint}} = |\mathcal{C}| \times n \times d + \text{RecomputationCost},$$

where RecomputationCost accounts for the additional computational overhead. The memory savings $S_{\text{mem}}$ achieved through checkpointing is

$$S_{\text{mem}} = L \times n \times d - |\mathcal{C}| \times n \times d,$$

showing that significant memory savings can be achieved at the cost of additional recomputation during backpropagation.

2. Quantization reduces the precision of model parameters and activations, thereby reducing the memory required to store them. Quantization maps a high precision value $x$ to a lower precision value $\tilde{x}$:

$$\tilde{x} = \text{Quantize}(x) = \text{round}\left(\frac{x - x_{\text{min}}}{\Delta}\right),$$

where $\Delta$ is the quantization step size and $x_{\text{min}}$ is the minimum value of the range being quantized. If sizeof(FP32) and sizeof(FP16) denote the memory requirements for 32-bit and 16-bit floating-point representations, the memory reduction $R_{\text{quant}}$ through quantization to FP16 is

$$R_{\text{quant}} = \frac{\text{sizeof(FP32)}}{\text{sizeof(FP16)}} = 2.$$

Quantization can halve the memory requirements while introducing only minimal degradation in model performance.

In large-scale transformers like GPT-3, memory optimization techniques such as checkpointing and quantization are crucial for enabling the training of models with billions of parameters. These techniques allow models to be trained on hardware with limited memory resources, ensuring that the models remain scalable and efficient.

### 9.1.2  Implementation in Large-Scale Models

Implementing memory optimization strategies in large-scale transformer models requires careful consideration of the trade-offs between memory usage, computational overhead, and model performance. These strategies must be integrated seamlessly into the training process to ensure that they provide tangible benefits without introducing significant complexities.

Consider a large-scale transformer with parameters $\theta$ distributed across multiple devices. The total memory required for storing the model parameters is

$$\text{Memory}_{\text{params}} = \sum_{i=1}^{L} \text{sizeof}(\theta^{(i)}),$$

where $\theta^{(i)}$ represents the parameters of the $i$-th layer.

Memory Optimization Techniques:

1. Model parallelism can be combined with memory optimization techniques such as checkpointing and quantization to distribute both the computation and memory usage across multiple devices. Let $\mathbf{W}$ be a large weight matrix partitioned across $p$ devices. With quantization applied, the memory required for storing $\mathbf{W}$ on each device is

$$\text{Memory}_{\text{quant}} = \frac{1}{p} \times \frac{\text{sizeof}(\mathbf{W})}{2},$$

where the factor $\frac{1}{2}$ comes from quantization to FP16.

2. Dynamic memory management involves allocating and deallocating memory resources on the fly based on the current requirements of the model, ensuring that memory is used efficiently throughout the training process. Let $M(t)$ represent the memory usage at time $t$ during training. Dynamic memory management seeks to minimize the peak memory usage:

$$M_{\text{peak}} = \max_t M(t),$$

by efficiently reallocating memory as different layers and operations are executed.

Dynamic memory management can significantly reduce the peak memory usage by ensuring that memory is allocated only when needed and deallocated immediately after use. This leads to more efficient utilization of available resources.

In training large-scale models like T5 or BERT-Large, dynamic memory management combined with model parallelism and quantization enables the training of models with over a billion parameters, even on hardware with constrained memory resources. These techniques ensure that the training process remains feasible and efficient, allowing for the exploration of even larger and more powerful models.

### 9.1.3  Sparse Transformers in Practice

Sparse transformers leverage sparsity in the attention mechanism to reduce computational and memory overhead, making them more suitable for practical applications involving long sequences or large datasets. By focusing on the most relevant parts of the input, sparse transformers maintain performance while significantly improving efficiency.

In a sparse transformer, the attention matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, which typically has quadratic complexity, is replaced with a sparse matrix $\mathbf{A}_s$ that has fewer non-zero elements:

$$\mathbf{A}_s = \text{SparseMask}(\mathbf{A}),$$

where SparseMask is a function that retains only a subset of the elements in $\mathbf{A}$, typically those corresponding to the most relevant tokens.

The complexity of sparse attention is reduced from $O(n^2 d)$ to $O(n \cdot k \cdot d)$, where $k$ is the number of tokens each token attends to. If $k \ll n$, the computational savings are significant:

$$\frac{\text{Complexity}_{\text{Sparse}}}{\text{Complexity}_{\text{Dense}}} \approx \frac{k}{n}.$$

In natural language processing tasks, such as document classification or summarization, sparse transformers can efficiently process long documents by focusing attention on the most informative sentences or paragraphs, leading to faster inference times and lower memory usage without sacrificing accuracy.

## 9.2   Transformers in Reinforcement Learning

Reinforcement Learning (RL) involves learning to make decisions by interacting with an environment to maximize cumulative rewards ([5, 8, 10, 12]). The use of transformers in RL is an emerging area that leverages the model's ability to capture long-term dependencies and complex patterns in sequential data. This section explores the mathematical foundations of applying transformers in RL, focusing on state representation, representation learning, and temporal abstractions. These concepts are crucial for building models that can effectively learn from and adapt to dynamic environments.

### 9.2.1   Mathematical Modeling

The integration of transformers into RL frameworks requires a careful mathematical formulation to capture the complexities of sequential decision-making. This involves modeling the state representation, learning meaningful representations of the environment, and incorporating temporal abstractions that allow the model to plan and act over multiple time scales.

### State Representation

In RL, the state $s_t$ at time $t$ encapsulates all relevant information from the environment necessary to make a decision. The transformer model, with its attention mechanisms, provides a powerful tool for encoding this state, particularly when the state is derived from a sequence of observations.

Let $\mathcal{S}$ be the state space and $s_t \in \mathcal{S}$ be the state at time $t$. In a transformer-based RL model, the state is represented as a sequence of past observations $\mathbf{O}_t = (o_1, o_2, \ldots, o_t)$, where each observation $o_i \in \mathcal{O}$ is part of the observation space $\mathcal{O}$. The state representation $\mathbf{s}_t$ is then given by

$$\mathbf{s}_t = f_{\text{transformer}}(\mathbf{O}_t) = \text{LayerNorm}\left(\sum_{i=1}^{t} \alpha_{i,t} \mathbf{o}_i\right),$$

where $\alpha_{i,t}$ are the attention weights computed by the transformer, indicating the importance of each past observation $\mathbf{o}_i$ in determining the current state.

The transformer's state representation is capable of capturing complex dependencies in the observation sequence. Formally, let $\mathcal{H}(\mathbf{s}_t)$ denote the set of possible state representations generated by the transformer. The expressivity theorem states that

$$\mathcal{H}(\mathbf{s}_t) \supseteq \mathcal{H}_{\text{RNN}}(\mathbf{s}_t),$$

where $\mathcal{H}_{\text{RNN}}(\mathbf{s}_t)$ represents the set of state representations possible with RNNs. This indicates that transformers can model all the dependencies captured by RNNs, and more, due to their global attention mechanism.

In an RL task like autonomous driving, where the state is derived from a sequence of visual and sensory inputs, the transformer's ability to attend to relevant observations (e.g., traffic lights, road signs) across the sequence allows for a more nuanced and effective state representation, leading to better decision-making.

**Representation Learning for RL**

Effective representation learning in RL involves discovering compact, meaningful representations of the environment that facilitate decision-making. Transformers excel at learning such representations by capturing the relationships and patterns within the sequence of observations.

Let $\mathcal{E}$ be the environment, and let $\mathcal{T}(\mathcal{E})$ represent the set of trajectories, where each trajectory $\tau \in \mathcal{T}(\mathcal{E})$ consists of a sequence of states and actions: $\tau = \{(s_1, a_1), (s_2, a_2), \ldots, (s_T, a_T)\}$. The transformer's objective in representation learning is to encode the trajectory $\tau$ into a representation $\mathbf{z}$ that captures the essential features for predicting future states and actions:

$$\mathbf{z} = g_{\text{transformer}}(\tau) = \text{MLP}\left(\text{LayerNorm}\left(\sum_{t=1}^{T} \alpha_t \mathbf{s}_t\right)\right),$$

where $\mathbf{s}_t$ is the state at time $t$ and $\alpha_t$ are the attention weights.

Transformers can learn representations that preserve the Markov property, which is critical for RL algorithms. Specifically, for any Markov decision process (MDP)

with state transition dynamics $P(s_{t+1} \mid s_t, a_t)$, the representation $\mathbf{z}$ learned by the transformer satisfies

$$P(s_{t+1} \mid \mathbf{z}, a_t) = P(s_{t+1} \mid s_t, a_t),$$

indicating that the representation $\mathbf{z}$ retains all necessary information for predicting the next state.

In a robotics RL task, where the environment's state is high dimensional (e.g., images or point clouds), transformers can learn compact embeddings of these states that preserve the relevant information for action selection, thus improving the robot's ability to perform tasks like object manipulation or navigation.

### Temporal Abstractions

Temporal abstractions allow the model to operate over different time scales, making decisions that account for both immediate and long-term consequences. Transformers naturally lend themselves to this due to their ability to model long-range dependencies in sequences.

Temporal abstractions in RL involve defining higher-level actions or options that extend over multiple time steps. Let $\mathcal{O}$ represent the set of options, where each option $o \in \mathcal{O}$ is a temporally extended action. The transformer can model the initiation and termination of options through its attention mechanism:

$$\mathbf{q}_o = \sum_{t=1}^{T} \alpha_t^o \mathbf{s}_t,$$

where $\mathbf{q}_o$ is the query vector associated with option $o$ and $\alpha_t^o$ are the attention weights specific to this option.

Transformers can represent temporal abstractions by assigning different attention patterns to different time scales. Formally, let $\mathcal{A}(\mathbf{s}_t)$ represent the set of possible actions, including temporally extended options. The transformer-based policy $\pi$ that incorporates temporal abstractions is defined as

$$\pi(a_t \mid \mathbf{s}_t) = \sum_{o \in \mathcal{O}} \beta_o \pi_o(a_t \mid \mathbf{q}_o),$$

where $\beta_o$ is the probability of selecting option $o$ and $\pi_o$ is the policy associated with that option. This allows the model to plan over multiple time scales effectively.

In hierarchical RL tasks, such as playing a strategy game, transformers can learn to abstract short-term actions into higher-level strategies that span several moves, enabling more strategic planning and better overall performance.

**Policy Learning**

Policy learning in RL involves learning a mapping from states to actions, which guides the agent's behavior in the environment. This mapping, or policy, is often represented as a probability distribution over actions, conditioned on the current state. Transformers, with their ability to model sequential dependencies, can be used to learn such policies, especially in environments where the state representation involves sequences of observations.

Let $\pi(a_t \mid s_t; \theta)$ represent the policy, parameterized by $\theta$, which gives the probability of taking action $a_t$ given state $s_t$ at time $t$. The objective in policy learning is to maximize the expected cumulative reward $\mathbb{E}[\sum_{t=1}^{T} r_t]$, where $r_t$ is the reward received at time $t$.

The transformer's attention mechanism can be used to represent the policy by considering the sequence of past states and actions:

$$\pi(a_t \mid s_t; \theta) = \text{softmax} \left( \text{LayerNorm} \left( \sum_{i=1}^{t} \alpha_{i,t} \mathbf{s}_i \right) \right),$$

where $\mathbf{s}_i$ represents the state at time $i$ and $\alpha_{i,t}$ are the attention weights learned by the transformer.

Under certain conditions, transformers can represent the optimal policy $\pi^*(a_t \mid s_t)$ that maximizes the expected cumulative reward. Specifically, given sufficient capacity, the transformer can approximate the policy function to arbitrary precision, capturing the necessary temporal dependencies in the sequence of states and actions.

In a game-playing environment like chess, where the state is represented by the sequence of moves leading up to the current position, a transformer-based policy can effectively model the complex dependencies between moves and select actions that maximize the chance of winning.

**Policy Gradient Methods**

Policy gradient methods are a class of RL algorithms that optimize the policy directly by computing gradients of the expected reward with respect to the policy parameters. These methods are particularly well suited for environments with continuous or high-dimensional action spaces.

The objective in policy gradient methods is to maximize the expected cumulative reward $J(\theta) = \mathbb{E}_{\pi_\theta}[\sum_{t=1}^{T} r_t]$ with respect to the policy parameters $\theta$. The policy gradient is given by

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=1}^{T} \nabla_\theta \log \pi(a_t \mid s_t; \theta) R_t \right],$$

where $R_t = \sum_{k=t}^{T} r_k$ is the return starting from time $t$.

Under the assumption that the policy $\pi(a_t \mid s_t; \theta)$ is differentiable with respect to $\theta$, and that the learning rate $\eta$ satisfies certain conditions (e.g., $\sum \eta_t = \infty$, $\sum \eta_t^2 < \infty$), the policy gradient method converges to a local maximum of $J(\theta)$.

In a continuous control task like robotic arm manipulation, where the actions are continuous valued (e.g., joint angles), policy gradient methods using transformers can effectively learn policies that maximize the reward, such as successfully grasping and moving objects.

## Value-Based Methods

Value-based methods in RL involve learning a value function that estimates the expected return from a given state or state–action pair. The policy is then derived indirectly from the value function, typically by selecting actions that maximize the estimated value.

The value function $V^\pi(s_t)$ under policy $\pi$ is defined as the expected return starting from state $s_t$:

$$V^\pi(s_t) = \mathbb{E}_\pi \left[ \sum_{k=t}^{T} r_k \mid s_t \right].$$

Similarly, the action-value function $Q^\pi(s_t, a_t)$ is defined as the expected return starting from state $s_t$ and taking action $a_t$:

$$Q^\pi(s_t, a_t) = \mathbb{E}_\pi \left[ \sum_{k=t}^{T} r_k \mid s_t, a_t \right].$$

The optimal policy $\pi^*$ is derived by selecting actions that maximize $Q^\pi$:

$$\pi^*(a_t \mid s_t) = \arg \max_{a_t} Q^\pi(s_t, a_t).$$

Transformers can be used to approximate the value function $V^\pi(s_t)$ or action-value function $Q^\pi(s_t, a_t)$. The Bellman optimality equation, which defines the relationship between the value functions at consecutive time steps, holds for the transformer-based approximation:

$$Q^\pi(s_t, a_t) = r_t + \gamma \mathbb{E}_{s_{t+1} \sim P} \left[ \max_{a_{t+1}} Q^\pi(s_{t+1}, a_{t+1}) \right],$$

where $\gamma$ is the discount factor and $P$ is the state transition probability.

In an environment like Atari games, where the state is represented by a sequence of frames, transformers can be used to approximate the value function, guiding the agent to take actions that maximize the score.

## 9.2.2  *Applications and Analysis*

The application of transformers in RL spans various domains, from games and robotics to finance and healthcare. The flexibility of transformers in modeling sequential data makes them particularly powerful in environments where decisions depend on long-term planning and context.

Transformers in RL provide several advantages, including the ability to model long-range dependencies in state sequences and the flexibility to handle various types of action spaces. The mathematical foundations explored in this chapter, including policy gradients and value-based methods, underscore the potential of transformers to revolutionize RL by offering more robust and scalable solutions.

Applications:

1. Games: Transformers can be used to model strategies in complex games like Go or StarCraft, where the state space is vast, and decisions depend on the entire sequence of past moves.

2. Robotics: In robotic tasks requiring precision and adaptability, transformers help in learning policies that can handle diverse scenarios by modeling the sequential nature of sensor inputs and actions.

3. Finance: Transformers can be employed in trading algorithms, where the state is represented by historical price sequences, and actions correspond to buy/sell decisions. The model learns policies that maximize profit while managing risk.

**Transformers for Model-Free RL**

Model-free RL involves learning policies or value functions directly from interaction with the environment, without explicitly modeling the environment's dynamics. Transformers in this context can be used to enhance policy learning and value estimation by processing sequences of states, actions, and rewards.

In model-free RL, the objective is to learn a policy $\pi(a_t \mid s_t; \theta)$ or a value function $V^{\pi}(s_t; \theta)$ that maximizes the expected cumulative reward. The transformer, with its attention mechanism, can be used to process a history of observations $\mathbf{O}_t = (o_1, o_2, \ldots, o_t)$ to produce a state representation $\mathbf{s}_t$:

$$\mathbf{s}_t = f_{\text{transformer}}(\mathbf{O}_t) = \text{LayerNorm}\left(\sum_{i=1}^{t} \alpha_{i,t} \mathbf{o}_i\right),$$

where $\alpha_{i,t}$ are attention weights.

The policy or value function can then be defined as

$$\pi(a_t \mid s_t; \theta) = \text{softmax}(g_{\pi}(\mathbf{s}_t; \theta))$$

or

$$V^{\pi}(s_t; \theta) = g_V(\mathbf{s}_t; \theta),$$

where $g_\pi$ and $g_V$ are functions parameterized by $\theta$ that map the state representation to action probabilities or value estimates.

Given a sufficiently expressive transformer model, the policy gradient methods or value iteration methods converge to a locally optimal policy or value function. The transformer's ability to model long-term dependencies ensures that it can learn from sequences of observations that span multiple time steps, capturing relevant information for decision-making.

In a robotic control task where actions must be selected based on a sequence of sensor readings, transformers can process these readings as a sequence and learn a policy that maps the processed sequence to control actions, leading to improved performance in tasks such as navigation or manipulation.

### Transformers for Model-Based RL

Model-based RL involves learning a model of the environment's dynamics, which is then used to plan and make decisions. Transformers can be used in model-based RL to learn complex environment dynamics by processing sequences of states, actions, and rewards.

In model-based RL, the environment is typically modeled by a transition function $P(s_{t+1} \mid s_t, a_t)$ and a reward function $R(s_t, a_t)$. The transformer can be employed to model these functions by processing sequences of states and actions:

$$\hat{s}_{t+1} = f_{\text{transformer}} \left( (s_1, a_1), (s_2, a_2), \dots, (s_t, a_t) \right),$$

where $\hat{s}_{t+1}$ is the predicted next state.

The planning process involves using the learned model to simulate future trajectories and optimize the policy accordingly:

$$\pi^*(a_t \mid s_t) = \arg \max_\pi \mathbb{E}_{\hat{P}} \left[ \sum_{k=t}^{T} R(\hat{s}_k, \pi(\hat{s}_k)) \right],$$

where $\hat{P}$ is the learned transition model.

The accuracy of the environment model learned by the transformer directly impacts the performance of the model-based RL. If the transformer can accurately predict the environment dynamics, the resulting policy will closely approximate the optimal policy.

In a simulated environment like a video game, where the environment's dynamics are complex and involve long sequences of actions leading to rewards, transformers can be used to model these dynamics and plan effective strategies, improving performance in the game.

**Case Studies in Reinforcement Learning**

Case studies demonstrate the practical application of transformers in various RL settings, providing insights into their effectiveness and versatility. These studies highlight how transformers can be adapted to different RL challenges, from simple control tasks to complex decision-making scenarios.

Each case study involves the application of transformers in a specific RL environment, with the following steps:

1. State Representation: Using transformers to encode the sequence of observations or states.

2. Policy Learning or Planning: Applying policy gradient or value-based methods in model-free RL, or planning with a learned model in model-based RL.

3. Evaluation: Assessing the performance of the transformer-based RL model against traditional methods.

Case Studies:

1. Autonomous Driving: Transformers are used to encode sequences of sensor data (e.g., LIDAR, camera images) and learn policies that allow a vehicle to navigate complex environments safely and efficiently.

2. Strategic Games: In games like chess or Go, transformers model the sequence of moves and plan strategies that consider long-term consequences, leading to high-level play.

3. Industrial Control: Transformers are applied to process sequences of sensor readings and control signals in manufacturing processes, optimizing production efficiency and reducing downtime.

**Performance Metrics and Evaluation**

Evaluating the performance of transformer-based RL models involves an analysis of various metrics, including cumulative reward, stability, generalization, and computational efficiency. These metrics are crucial for understanding the strengths and limitations of transformers in RL.

Performance metrics in RL can be broadly categorized as follows:

1. Cumulative Reward: The total reward accumulated over a sequence of actions, defined as

$$G_T = \sum_{t=1}^{T} r_t,$$

where $r_t$ is the reward at time $t$.

2. Stability: The consistency of the learned policy across different training runs, often measured by the variance in cumulative rewards.

3. Generalization: The ability of the learned policy to perform well on unseen environments or tasks, evaluated by testing the model on variations of the training environment.

4. Computational Efficiency: The time and resources required to train the model, often evaluated by the number of iterations or the time taken to converge.

Given that transformers can model complex dependencies in RL tasks, the performance improvement (in terms of cumulative reward) from using transformers is bounded by the accuracy of the state or environment model they learn. Formally, if $\epsilon$ represents the model error, the improvement in cumulative reward is bounded by

$$\Delta G_T \leq \frac{1}{1 - \gamma} \cdot \epsilon,$$

where $\gamma$ is the discount factor.

In an RL task like robot arm control, performance metrics such as cumulative reward and stability can be used to evaluate how well the transformer-based policy handles different objects and tasks, with generalization measured by testing on new, unseen objects.

## 9.3 Convergence of Transformer Models: A Dynamical Systems Perspective

The convergence properties of transformer models can be analyzed using concepts from dynamical systems theory. By viewing the training dynamics of transformers as trajectories in a high-dimensional phase space, we can gain insights into their stability, convergence behavior, and the existence of fixed points. This section introduces the fundamental concepts of dynamical systems and applies them to the study of transformer models.

### 9.3.1 Introduction to Dynamical Systems

Dynamical systems theory ([2, 4, 6, 11]) provides a mathematical framework for analyzing the behavior of systems that evolve over time according to a set of deterministic rules. In the context of transformer models, the training process can be seen as a dynamical system where the parameters of the model evolve under the influence of gradients derived from the loss function.

A dynamical system can be described by a set of differential equations:

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{f}(\mathbf{x}(t), t),$$

where $\mathbf{x}(t) \in \mathbb{R}^n$ represents the state of the system at time $t$ and $\mathbf{f} : \mathbb{R}^n \times \mathbb{R} \to \mathbb{R}^n$ is a function that defines the evolution of the system.

In the context of transformers, $\mathbf{x}(t)$ can represent the model parameters at training iteration $t$ and $\mathbf{f}(\mathbf{x}(t), t)$ corresponds to the gradient of the loss function with respect to the parameters:

$$\mathbf{f}(\mathbf{x}(t)) = -\nabla_\theta \mathcal{L}(\theta(t)),$$

where $\mathcal{L}(\theta(t))$ is the loss function at iteration $t$ and $\theta(t)$ represents the model parameters.

## Phase Space and Trajectories

The phase space of a dynamical system is a mathematical space in which all possible states of the system are represented. Each point in this space corresponds to a unique state of the system, and the trajectory of the system represents the path taken by the system as it evolves over time.

Let $\mathcal{X} \subset \mathbb{R}^n$ be the phase space of the dynamical system, where $n$ is the dimensionality of the state vector $\mathbf{x}(t)$. The trajectory $\gamma(t)$ of the system is a curve in $\mathcal{X}$ defined by the solution to the differential equation:

$$\gamma(t) = \mathbf{x}(t),$$

with initial condition $\mathbf{x}(0) = \mathbf{x}_0$.

In transformer training, the trajectory of the model parameters $\theta(t)$ in the phase space corresponds to the sequence of parameter updates driven by gradient descent:

$$\theta(t + 1) = \theta(t) - \eta \nabla_\theta \mathcal{L}(\theta(t)),$$

where $\eta$ is the learning rate. Under certain conditions, such as Lipschitz continuity of $\mathbf{f}(\mathbf{x}(t))$, the trajectory $\gamma(t)$ is uniquely determined by the initial condition $\mathbf{x}_0$. This implies that the evolution of the transformer parameters during training is well defined and predictable, given the initial parameter values and the loss landscape. Consider a simple transformer model with a small number of parameters. The phase space in this case might be a low-dimensional space, allowing us to visualize the trajectory of the parameters during training. As the training progresses, the trajectory will ideally move toward a region of lower loss, corresponding to better model performance.

## Fixed Points and Stability

Fixed points in a dynamical system are states where the system remains unchanged over time, meaning that once the system reaches a fixed point, it will stay there. In the context of transformers, fixed points correspond to model parameters where the gradient of the loss function is zero, indicating that the model has reached a state of equilibrium.

A fixed point $\mathbf{x}^*$ of the dynamical system is defined by the condition:

$$\mathbf{f}(\mathbf{x}^*) = 0.$$

In transformer training, a fixed point $\theta^*$ occurs when the gradient of the loss function vanishes

$$\nabla_\theta \mathcal{L}(\theta^*) = 0.$$

The stability of a fixed point can be analyzed by examining the eigenvalues of the Jacobian matrix $\mathbf{J}(\mathbf{x}^*)$ of the dynamical system at the fixed point. The Jacobian matrix is given by

$$\mathbf{J}(\mathbf{x}^*) = \left.\frac{\partial \mathbf{f}}{\partial \mathbf{x}}\right|_{\mathbf{x}=\mathbf{x}^*}.$$

A fixed point $\mathbf{x}^*$ is stable if all the eigenvalues of $\mathbf{J}(\mathbf{x}^*)$ have negative real parts. If any eigenvalue has a positive real part, the fixed point is unstable, meaning that small perturbations will cause the system to move away from the fixed point.

In the training of a transformer, a stable fixed point corresponds to a set of parameters where the model has converged to a local minimum of the loss function. If the fixed point is unstable, the training process may diverge, or the model may oscillate around the fixed point without converging.

## 9.3.2   Lyapunov Exponents and Stability Analysis

Lyapunov exponents are critical tools in analyzing the stability of dynamical systems. They measure the rate at which nearby trajectories converge or diverge in phase space, providing a quantitative assessment of stability. In the context of transformer models, Lyapunov exponents help us understand how perturbations in model parameters evolve during training, influencing the convergence and stability of the learning process.

### Lyapunov Functions and Stability Criteria

Lyapunov functions are scalar functions used to prove the stability of a dynamical system. A Lyapunov function $V(\mathbf{x})$ decreases along trajectories of the system, indicating that the system is moving toward a stable equilibrium.

A Lyapunov function $V : \mathbb{R}^n \to \mathbb{R}$ for a dynamical system $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$ satisfies the following conditions:

1. $V(\mathbf{x}) > 0$ for all $\mathbf{x} \neq \mathbf{x}^*$ and $V(\mathbf{x}^*) = 0$, where $\mathbf{x}^*$ is an equilibrium point.
2. The time derivative of $V(\mathbf{x})$ along the trajectories of the system is non-positive:

$$\dot{V}(\mathbf{x}) = \nabla V(\mathbf{x}) \cdot \mathbf{f}(\mathbf{x}) \leq 0.$$

If $\dot{V}(\mathbf{x}) < 0$ for all $\mathbf{x} \neq \mathbf{x}^*$, the equilibrium point $\mathbf{x}^*$ is globally asymptotically stable.

If there exists a Lyapunov function $V(\mathbf{x})$ for the dynamical system $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$, then the equilibrium point $\mathbf{x}^*$ is stable. If $\dot{V}(\mathbf{x}) < 0$, the equilibrium is globally asymptotically stable.

In the context of transformer models, consider a simplified gradient descent scenario where the parameter update rule is $\theta_{t+1} = \theta_t - \eta \nabla_\theta \mathcal{L}(\theta_t)$. A potential Lyapunov function is the loss function itself $V(\theta) = \mathcal{L}(\theta)$, as it decreases along trajectories of gradient descent:

$$\dot{V}(\theta) = \nabla_\theta \mathcal{L}(\theta) \cdot \dot{\theta} = -\eta \|\nabla_\theta \mathcal{L}(\theta)\|^2 \leq 0.$$

This indicates that the system is moving toward a minimum of the loss function, contributing to the stability of the learning process.

**Applications to Neural Networks**

In neural networks, Lyapunov exponents can be used to assess the stability of learning dynamics. A positive Lyapunov exponent indicates chaotic behavior, where small perturbations in the input or parameter space lead to exponentially diverging trajectories, making learning unpredictable. Negative exponents suggest that trajectories converge, indicating stable learning dynamics.

The Lyapunov exponent $\lambda$ for a trajectory $\mathbf{x}(t)$ is defined as

$$\lambda = \lim_{t \to \infty} \frac{1}{t} \ln \frac{\|\delta \mathbf{x}(t)\|}{\|\delta \mathbf{x}(0)\|},$$

where $\delta \mathbf{x}(t)$ is an infinitesimal perturbation in the initial conditions. In the context of a transformer model, this can be interpreted as the sensitivity of the model's parameters to small changes during training.

For a neural network, including transformers, if all Lyapunov exponents are negative, the training process is stable, leading to convergence toward a local or global minimum of the loss function. If any Lyapunov exponent is positive, the system may exhibit chaotic behavior, potentially leading to non-convergence or unpredictable outcomes.

In a transformer model with many layers, the gradients can become highly sensitive to small perturbations due to the depth of the model. By calculating the Lyapunov exponents, we can determine whether the model is likely to exhibit stable learning or if it is prone to chaotic dynamics, which could explain phenomena like exploding or vanishing gradients.

### *9.3.3   Chaos Theory and Non-linear Dynamics*

Chaos theory ([1, 3, 7, 9, 11]) studies the behavior of dynamical systems that are highly sensitive to initial conditions, a property known as the "butterfly effect." In machine learning, particularly in deep learning models like transformers, understanding chaotic dynamics can provide insights into the complexities of the learning process and the potential for irregular behavior.

**Bifurcations and Attractors**

Bifurcations occur when a small change in the system's parameters causes a sudden qualitative change in its behavior. Attractors are sets toward which a system tends to evolve, regardless of the initial conditions. These concepts are crucial for understanding how transformer models behave under different training regimes and parameter settings.

A bifurcation in a dynamical system occurs when a parameter $\mu$ is varied, leading to a qualitative change in the structure of its phase space. The system's equation is generally of the form:

$$\frac{d\mathbf{x}}{dt} = \mathbf{f}(\mathbf{x}, \mu).$$

An attractor is a set $\mathcal{A} \subset \mathcal{X}$ in phase space such that trajectories starting near $\mathcal{A}$ remain close to $\mathcal{A}$ as $t \to \infty$. Common types of attractors include fixed points, limit cycles, and strange attractors.

When a bifurcation occurs, the stability of fixed points or periodic orbits can change. For example, a supercritical pitchfork bifurcation leads to the emergence of a stable fixed point, while a subcritical bifurcation can lead to the sudden appearance of chaotic dynamics.

In the training of transformers, a bifurcation could occur when adjusting hyperparameters such as the learning rate or model depth. A small change in these parameters might lead to a transition from stable learning (convergence) to chaotic learning (non-convergence), where the model's performance becomes highly sensitive to initial conditions or noise.

**Implications for Learning Dynamics**

Understanding bifurcations and attractors in the context of transformer models allows us to predict and control the learning dynamics. By identifying the parameter regions associated with stable attractors, we can ensure that the training process remains stable and convergent.

Consider a transformer model where the training dynamics are influenced by a set of hyperparameters $\mu = (\mu_1, \mu_2, \ldots, \mu_m)$. The phase space of the model's parameters $\theta$ evolves according to

$$\frac{d\theta}{dt} = -\nabla_\theta \mathcal{L}(\theta, \mu).$$

A bifurcation occurs when a critical value $\mu_c$ is reached, causing a sudden change in the nature of the fixed points or attractors in the phase space. The system may transition from a regime with a stable fixed point to one with a strange attractor, characterized by chaotic dynamics.

In the context of neural networks, including transformers, the existence of a strange attractor can lead to highly irregular and sensitive learning dynamics, where small changes in the initialization or training data can lead to vastly different outcomes. By contrast, stable attractors ensure robust learning dynamics, where the model reliably converges to a good solution.

Consider the phenomenon of mode collapse in GANs (Generative Adversarial Networks), which can be viewed as the model being trapped in a low-dimensional attractor within the phase space. Similarly, in transformers, certain configurations of hyperparameters might lead to attractors that cause the model to consistently underperform, suggesting the need for careful tuning and regularization.

## 9.4 Convergence in Transformer Training

The convergence of gradient descent algorithms during transformer training is crucial for ensuring that the model reaches a minimum of the loss function efficiently and effectively. Understanding the convergence properties, particularly the rates of convergence and the impact of learning rates, allows us to optimize the training process, ensuring stability and accelerating progress toward an optimal solution. This section explores the mathematical foundations underlying these aspects of convergence in transformer training.

### 9.4.1 Convergence of Gradient Descent

Gradient descent is the foundational optimization algorithm used in training transformers, where the model parameters are iteratively updated in the direction of the negative gradient of the loss function. The convergence of gradient descent depends on the nature of the loss function, the choice of learning rate, and the structure of the parameter space.

Let $\mathcal{L}(\theta)$ be the loss function, where $\theta \in \mathbb{R}^n$ represents the model parameters. The gradient descent update rule is given by

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta \mathcal{L}(\theta_t),$$

where $\eta > 0$ is the learning rate and $t$ denotes the iteration step.

The objective is to ensure that the sequence $\{\theta_t\}$ converges to a point $\theta^*$ where the gradient of the loss function vanishes:

$$\nabla_\theta \mathcal{L}(\theta^*) = 0.$$

Assume that the loss function $\mathcal{L}(\theta)$ is convex and differentiable, and that the gradient $\nabla_\theta \mathcal{L}(\theta)$ is Lipschitz continuous with constant $L$:

$$\|\nabla_\theta \mathcal{L}(\theta_1) - \nabla_\theta \mathcal{L}(\theta_2)\| \leq L\|\theta_1 - \theta_2\| \quad \text{for all } \theta_1, \theta_2 \in \mathbb{R}^n.$$

Then, if the learning rate $\eta$ satisfies $0 < \eta < \frac{2}{L}$, the sequence $\{\theta_t\}$ generated by gradient descent converges to a minimum $\theta^*$:

$$\mathcal{L}(\theta_{t+1}) \leq \mathcal{L}(\theta_t) - \frac{\eta}{2}\|\nabla_\theta \mathcal{L}(\theta_t)\|^2.$$

This result guarantees that, under appropriate conditions, the gradient descent algorithm will reduce the loss function at each step, eventually converging to a minimum.

In the training of a transformer model, if the loss function (e.g., cross-entropy) satisfies the convexity and Lipschitz continuity conditions, and if the learning rate is chosen within the prescribed bounds, gradient descent will reliably decrease the loss and lead to convergence. This ensures that the model parameters move closer to the optimal values with each iteration.


**Analyzing Convergence Rates**

The rate at which gradient descent converges to the minimum of the loss function is of paramount importance in practical applications. Convergence rates depend on the properties of the loss function and the choice of learning rate.

The convergence rate of gradient descent can be classified as linear, sublinear, or superlinear, depending on how the loss function decreases over iterations. For strongly convex loss functions, the convergence rate is linear, meaning that the error decreases exponentially fast:

$$\|\theta_t - \theta^*\| \leq C \cdot \rho^t,$$

where $C > 0$ is a constant and $0 < \rho < 1$ is the convergence rate.

For general convex functions, the convergence is typically sublinear, often of the form:

$$\mathcal{L}(\theta_t) - \mathcal{L}(\theta^*) \leq \frac{D}{\sqrt{t}},$$

where $D > 0$ is a constant. This indicates that the convergence slows down as the algorithm progresses.

If the loss function $\mathcal{L}(\theta)$ is strongly convex with parameter $m > 0$ and has a Lipschitz-continuous gradient with constant $L > 0$, then the gradient descent algorithm with learning rate $\eta = \frac{2}{L+m}$ has a linear convergence rate:

$$\|\theta_t - \theta^*\| \leq \left(1 - \frac{2m}{L+m}\right)^t \|\theta_0 - \theta^*\|.$$

This result shows that the distance to the optimal parameters decreases exponentially with each iteration, leading to fast convergence.

In a transformer model, if the loss function is strongly convex (which may occur in simpler settings or with regularization), the parameters will converge linearly to the optimum. This rapid convergence is particularly beneficial in large-scale models where training time is a critical factor.

**Impact of Learning Rates**

The learning rate $\eta$ plays a crucial role in determining both the speed and stability of convergence in gradient descent. Too large a learning rate can cause divergence, while too small a learning rate can result in slow convergence.

The choice of learning rate affects the step size of parameter updates. If the learning rate is too large, the updates can overshoot the minimum, causing oscillations or divergence. If the learning rate is too small, the updates become too conservative, leading to slow progress.

For a given learning rate $\eta$, the update rule is

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta \mathcal{L}(\theta_t).$$

The stability of this update can be analyzed by examining the Taylor expansion of the loss function around the minimum $\theta^*$:

$$\mathcal{L}(\theta_t) \approx \mathcal{L}(\theta^*) + \frac{1}{2}(\theta_t - \theta^*)^\top H(\theta_t - \theta^*),$$

where $H$ is the Hessian matrix of second derivatives. The condition for stability requires that the eigenvalues of the matrix $I - \eta H$ lie within the unit circle, leading to the constraint on the learning rate:

$$0 < \eta < \frac{2}{\lambda_{\max}},$$

where $\lambda_{\max}$ is the largest eigenvalue of the Hessian matrix.

For quadratic loss functions, the optimal learning rate $\eta^*$ that minimizes the convergence time is given by

$$\eta^* = \frac{2}{\lambda_{\max} + \lambda_{\min}},$$

where $\lambda_{\min}$ is the smallest eigenvalue of the Hessian. This learning rate balances fast convergence with stability.

In transformer training, adjusting the learning rate dynamically can lead to improved convergence. Techniques such as learning rate schedules or adaptive learning rates (e.g., Adam optimizer) are often used to optimize the learning process, ensuring that the model converges efficiently without the risk of divergence.

## 9.4.2   Stability of Learned Representations

In the context of transformers, stability refers to the consistency of the learned representations across different layers of the model and under various perturbations, such as changes in input data or model parameters. Stable representations are desirable as they indicate that the model is robust to small changes, leading to better generalization and reliability.

### Empirical versus Theoretical Stability

Empirical stability is observed through experiments and simulations, where the behavior of learned representations is analyzed under different conditions. Theoretical stability, on the other hand, is derived from mathematical models and analyses that predict the conditions under which learned representations remain consistent and reliable.

Let $\mathbf{h}^{(l)}$ denote the representation learned at layer $l$ of the transformer model. The stability of this representation can be quantified by measuring the sensitivity of $\mathbf{h}^{(l)}$ to perturbations in the input or parameters:

$$\Delta \mathbf{h}^{(l)} = \|\mathbf{h}^{(l)}(\theta + \delta\theta) - \mathbf{h}^{(l)}(\theta)\|,$$

where $\delta\theta$ represents a small perturbation in the model parameters $\theta$.

The stability of the representation $\mathbf{h}^{(l)}$ is characterized by the condition:

$$\frac{\Delta \mathbf{h}^{(l)}}{\|\delta\theta\|} \leq \kappa,$$

where $\kappa > 0$ is a constant that bounds the sensitivity. A smaller value of $\kappa$ indicates higher stability.

If the transformation from layer $l - 1$ to layer $l$ in the transformer is Lipschitz continuous with constant $L_l$, then the learned representation $\mathbf{h}^{(l)}$ is stable under small perturbations, with the stability bound given by

$$\frac{\Delta \mathbf{h}^{(l)}}{\|\delta \theta\|} \leq L_l \|\mathbf{h}^{(l-1)}\|,$$

where $\mathbf{h}^{(l-1)}$ is the representation from the previous layer.

In practice, empirical studies might involve adding Gaussian noise to the input or parameters and observing how the learned representations change across layers. If the representations remain consistent despite these perturbations, the model is considered to have high empirical stability. The theoretical stability can be verified by ensuring that the transformations between layers satisfy the Lipschitz continuity condition.

**Stability Across Layers**

Stability across layers is crucial for ensuring that the learned representations do not degrade as the data passes through multiple layers of the transformer. Each layer in the transformer performs a linear or non-linear transformation that ideally preserves or enhances the relevant features of the input.

Consider a transformer model with $L$ layers. The overall stability of the learned representation $\mathbf{h}^{(L)}$ at the final layer depends on the cumulative stability of each preceding layer. The total sensitivity to perturbations can be expressed as

$$\Delta \mathbf{h}^{(L)} = \prod_{l=1}^{L} L_l \|\mathbf{h}^{(0)}\|,$$

where $\mathbf{h}^{(0)}$ is the input representation and $L_l$ is the Lipschitz constant for layer $l$.

If each layer $l$ in the transformer satisfies $L_l \leq 1$, the overall model exhibits stable representations, with the total sensitivity being bounded by

$$\|\mathbf{h}^{(L)}\| \leq \|\mathbf{h}^{(0)}\|.$$

This implies that the depth of the transformer does not amplify perturbations, leading to stable and reliable learned representations.

In deep transformers, where the number of layers $L$ is large, ensuring that each layer has a Lipschitz constant $L_l \leq 1$ becomes crucial to prevent the model from becoming unstable. Techniques such as layer normalization and regularization are often employed to maintain stability across layers, ensuring that the final representation retains the essential features of the input.

**Dynamical Systems in Practice**

Dynamical systems theory provides a framework for analyzing and simulating the training dynamics of transformer models. By viewing the training process as a dynamical system, we can gain insights into the behavior of the model under different training regimes, leading to better optimization strategies and model design.

**Numerical Simulations**

Numerical simulations are used to study the behavior of dynamical systems in transformer training. These simulations involve discretizing the continuous training dynamics and iterating over time to observe the evolution of model parameters and learned representations.

Consider the discrete time dynamical system representing the training process:

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta \mathcal{L}(\theta_t),$$

where $t$ represents the iteration step and $\eta$ is the learning rate. Numerical simulations involve iterating this update rule over a large number of steps to observe the trajectory of $\theta_t$ in the parameter space.

Under the conditions of convexity and Lipschitz continuity of the gradient, the numerical simulation of gradient descent converges to a minimum $\theta^*$ of the loss function, with the convergence rate determined by the learning rate $\eta$.

Simulations of transformer training might involve varying the learning rate, initialization, or other hyperparameters to observe their impact on convergence and stability. These simulations provide empirical insights that complement the theoretical analysis, helping to identify optimal training strategies.

**Applications in Real-World Scenarios**

Applying dynamical systems theory to real-world transformer training scenarios allows us to predict and control the behavior of the model under practical conditions. This includes adjusting hyperparameters, understanding the impact of training data variability, and ensuring robust generalization.

In real-world applications, the dynamical system governing transformer training is often influenced by factors such as noise in the data, non-convexity of the loss function, and stochasticity in the optimization process. The practical stability of the model can be assessed by analyzing how these factors affect the trajectory of $\theta_t$ and the resulting learned representations.

If the dynamical system representing transformer training is robust to small perturbations in data and parameters, the model is likely to generalize well to unseen data. Formally, if $\theta_t$ remains within a bounded region $\mathcal{B} \subset \mathbb{R}^n$ under perturbations, the learned representations will exhibit robustness.

In a real-world scenario such as natural language processing, where the data may be noisy or incomplete, applying dynamical systems analysis helps in designing models that maintain stability and performance despite these challenges. This could involve using regularization techniques, adaptive learning rates, or robust optimization methods.

# References

1. Alligood, K.T., Sauer, T.D., Yorke, J.A.: Chaos: An Introduction to Dynamical Systems. Springer (1996)
2. Arnold, V.I.: Geometrical Methods in the Theory of Ordinary Differential Equations. Springer-Verlag (1988)
3. Gleick, J.: Chaos: Making a New Science. Penguin Books (1987)
4. Guckenheimer, J., Holmes, P.: Nonlinear Oscillations, Dynamical Systems, and Bifurcations of Vector Fields. Springer (2002)
5. Kaelbling, L.P., Littman, M.L., Moore, A.W.: Reinforcement learning: a survey. J. Artif. Int. Res. **4**(1), 237–285 (1996). ISSN 1076-9757
6. Katok, A., Hasselblatt, B.: Introduction to the Modern Theory of Dynamical Systems. Cambridge University Press (1997)
7. Lorenz, E.N.: Deterministic nonperiodic flow. J. Atmospheric Sci. **20**(2), 130–141 (1963)
8. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M.A., Fidjeland, A.K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D.: Human-level control through deep reinforcement learning. Nature **518**(7540), 529–533 (2015)
9. Ott, E.: Chaos in Dynamical Systems. Cambridge University Press (2002)
10. 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.P., Leach, M., Kavukcuoglu, K., Graepel, T., Hassabis, D.: Mastering the game of go with deep neural networks and tree search. Nature **529**(7587), 484–489 (2016)
11. Strogatz, S.H.: Nonlinear Dynamics and Chaos: With Applications to Physics. Chemistry, and Engineering. Westview Press, Biology (1994)
12. Sutton, R.S., Barto, A.G.: Reinforcement Learning: An Introduction. MIT Press (2018)