

Summer of Science

Mid-Term Report 2025

Large Language Models

Koushik Reddy

23B1000

Contents

1	Introduction to LLMs and Machine Learning Basics	3
1.1	Basics of Large Language Models	3
1.2	Applications of Large Language Models	4
1.3	Core Machine Learning Concepts	5
1.3.1	Regression	5
1.3.2	Classification	6
1.3.3	Support Vector Machines (SVMs)	6
1.3.4	Decision Trees	7
1.3.5	Clustering	8
2	Deep Learning Fundamentals	11
2.1	Neural Networks and Architectures	11
2.1.1	Feedforward Neural Networks (FFNN)	11
2.1.2	Convolutional Neural Networks (CNNs)	13
2.2	Optimization Techniques	15
2.2.1	Gradient Descent	15
2.2.2	Stochastic Gradient Descent (SGD)	15
2.2.3	Mini-batch Gradient Descent	15
2.2.4	Momentum	15
2.2.5	RMSProp	16
2.2.6	Adam Optimizer	16
2.3	Regularization Techniques	16
2.3.1	L1 and L2 Regularization	16
2.3.2	Dropout	17
2.3.3	Early Stopping	17
2.3.4	Batch Normalization	17
2.3.5	Data Augmentation	17
3	Natural Language Processing (NLP)	18
3.1	Introduction to Natural Language Processing (NLP)	18
3.1.1	Motivation	18
3.1.2	Historical Evolution	18
3.1.3	Key Components of NLP	18
3.1.4	Challenges in NLP	19
3.1.5	Applications of NLP	19
3.1.6	Role of Deep Learning in NLP	19
3.2	NLP Pipeline	20
3.2.1	Text Collection	20

3.2.2	Text Preprocessing	20
3.2.3	Text Representation / Feature Extraction	20
3.2.4	Model Building	21
3.2.5	Evaluation	21
3.2.6	Deployment and Feedback	21
3.2.7	Summary	21
3.3	Traditional NLP Methods	21
3.3.1	Bag-of-Words (BoW)	21
3.3.2	N-Grams	22
3.3.3	Term Frequency–Inverse Document Frequency (TF-IDF)	22
3.3.4	Limitations of Traditional Methods	23
3.4	Word Embeddings	23
3.4.1	Motivation	23
3.4.2	Word2Vec	23
3.4.3	GloVe (Global Vectors)	24
3.4.4	Properties of Word Embeddings	24
3.4.5	Limitations	25
3.5	Recurrent Neural Networks (RNNs) in NLP	25
3.5.1	Motivation	25
3.5.2	RNN Architecture and Equations	25
3.5.3	Backpropagation Through Time (BPTT)	26
3.5.4	Applications of RNNs in NLP	26
3.5.5	Bidirectional RNNs (BiRNN)	26
3.5.6	Limitations of Vanilla RNNs	26
3.6	Long Short-Term Memory (LSTM)	27
3.6.1	Motivation	27
3.6.2	LSTM Cell Structure	27
3.6.3	LSTM Equations	27
3.6.4	Intuition Behind the Gates	28
3.6.5	Applications in NLP	28
3.6.6	Bidirectional LSTMs	28
3.6.7	Advantages over RNNs	28
3.7	Attention Mechanisms in NLP	29
3.7.1	Motivation	29
3.7.2	General Form of Attention	29
3.7.3	Scoring Functions	29
3.7.4	Types of Attention	29
3.7.5	Self-Attention	30
3.7.6	Multi-Head Attention	30
3.7.7	Applications in NLP	30
3.7.8	Advantages of Attention	30

Chapter 1

Introduction to LLMs and Machine Learning Basics

1.1 Basics of Large Language Models

Large Language Models (LLMs) are a class of deep learning models trained on massive text corpora to understand and generate human-like language. These models belong to the broader category of Natural Language Processing (NLP) systems and have revolutionized fields like conversational AI, search engines, translation, and content generation.

LLMs are based on the **transformer architecture**, which enables them to process input sequences in parallel and capture long-range dependencies using mechanisms like self-attention. Unlike earlier sequence models such as RNNs and LSTMs, transformers allow for more efficient and scalable training.

Some key characteristics of LLMs include:

- **Pretraining on large text corpora:** LLMs are initially trained on vast amounts of text data in an unsupervised manner, learning to predict the next word in a sentence or fill in masked tokens.
- **Fine-tuning for specific tasks:** After pretraining, LLMs can be fine-tuned on domain-specific data for tasks like summarization, question answering, and code generation.
- **Few-shot and zero-shot learning:** With enough parameters, LLMs demonstrate the ability to perform tasks with little or no task-specific training, relying solely on prompts.

Modern LLMs such as GPT, BERT, T5, and LLaMA have shown remarkable capabilities across a wide range of tasks, enabling more natural and intelligent human-computer interaction. However, they also raise challenges related to ethical use, hallucination, and computational cost.

In summary, understanding the foundational ideas behind LLMs is crucial for exploring their applications and limitations in real-world AI systems.

1.2 Applications of Large Language Models

Large Language Models (LLMs) have enabled significant advancements across a wide range of applications by leveraging their ability to understand, generate, and manipulate human language. These models are versatile and adaptable, making them foundational in modern AI systems.

1. Conversational Agents and Chatbots

LLMs are widely used in building chatbots and virtual assistants such as ChatGPT, Siri, and Alexa. They can handle both casual and task-oriented conversations, provide informative responses, and adapt to different contexts based on the prompt history.

2. Text Generation

LLMs are capable of producing coherent and contextually relevant text, making them suitable for applications like:

- Writing assistance tools (e.g., Grammarly, Notion AI)
- Story and content generation
- Email and document drafting

3. Machine Translation

With sufficient training data, LLMs can accurately translate text between multiple languages. Models like Google's mT5 and Meta's NLLB have achieved state-of-the-art results in multilingual translation tasks.

4. Summarization

LLMs can generate concise and coherent summaries of long documents, news articles, or research papers. This is especially useful in academic, legal, and journalistic contexts where digesting large volumes of text is necessary.

5. Sentiment Analysis and Text Classification

Businesses use LLMs to extract sentiment or classify text into predefined categories, such as spam detection, review sentiment analysis, or topic categorization in customer feedback.

6. Code Generation and Programming Assistance

LLMs trained on code (e.g., Codex, CodeLLaMA) can write code snippets, debug errors, or explain programming concepts, making them useful in IDEs and platforms like GitHub Copilot.

7. Information Retrieval and Question Answering

LLMs power search systems that go beyond keyword matching by understanding the semantic intent of queries. They can extract answers from documents or generate responses grounded in factual knowledge.

8. Personalization and Recommendation

By analyzing user input and preferences, LLMs can personalize content, recommend products or articles, and generate customized messages or responses in marketing and e-commerce.

9. Education and Tutoring

LLMs can serve as on-demand tutors, explaining complex topics, generating quizzes, or simulating conversations for language learning and exam preparation.

10. Healthcare and Medical Analysis

In medical applications, LLMs assist in summarizing patient records, interpreting clinical notes, and even supporting diagnosis based on textual data from medical literature or electronic health records (EHRs).

Conclusion: The adaptability of LLMs across diverse domains makes them one of the most impactful developments in AI. However, responsible deployment is critical, as misuse or overreliance can lead to misinformation, bias propagation, or ethical concerns.

1.3 Core Machine Learning Concepts

This section explores the foundational concepts in machine learning that underpin many aspects of large language models. We focus on three major paradigms: regression, classification, and clustering.

1.3.1 Regression

Regression is a supervised learning method where the goal is to predict a continuous output variable y given input features $x \in \mathbb{R}^n$. The simplest form is **linear regression**, which models the relationship as:

$$y = w^T x + b$$

where $w \in \mathbb{R}^n$ is the weight vector and $b \in \mathbb{R}$ is the bias.

Given a dataset of m examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$, we aim to minimize the **mean squared error (MSE)** loss:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (w^T x^{(i)} + b - y^{(i)})^2$$

Gradient Descent Update Rules:

$$w \leftarrow w - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (w^T x^{(i)} + b - y^{(i)}) x^{(i)} \quad (1.1)$$

$$b \leftarrow b - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (w^T x^{(i)} + b - y^{(i)}) \quad (1.2)$$

where α is the learning rate.

Linear regression is widely used due to its simplicity, interpretability, and efficiency.

1.3.2 Classification

Classification is another supervised learning task where the goal is to assign a label $y \in \{1, 2, \dots, K\}$ to an input x . For binary classification ($y \in \{0, 1\}$), **logistic regression** is a commonly used algorithm.

It models the probability that $y = 1$ as:

$$P(y = 1|x) = \sigma(w^T x + b)$$

where $\sigma(z) = \frac{1}{1+e^{-z}}$ is the sigmoid function.

Loss Function (Binary Cross-Entropy):

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

where $\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$

Extension to Multi-Class Classification:

For multi-class problems, we use the **softmax function**:

$$P(y = k|x) = \frac{e^{w_k^T x}}{\sum_{j=1}^K e^{w_j^T x}}$$

and the corresponding **categorical cross-entropy** loss:

$$J = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log P(y = k|x^{(i)})$$

1.3.3 Support Vector Machines (SVMs)

Support Vector Machines are powerful supervised learning models used primarily for classification, but also applicable to regression. The key idea is to find the optimal hyperplane that maximally separates the data points of different classes.

Binary Classification Setup: Given a training set of m labeled samples $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$ where $x^{(i)} \in \mathbb{R}^n$, and $y^{(i)} \in \{-1, +1\}$, the goal is to find a decision boundary:

$$w^T x + b = 0$$

such that the margin between the two classes is maximized.

Margin Definition: The margin is defined as the distance between the separating hyperplane and the nearest data point. For a correctly classified point:

$$y^{(i)}(w^T x^{(i)} + b) \geq 1$$

Primal Optimization Problem:

$$\min_{w, b} \frac{1}{2} \|w\|^2 \quad \text{subject to} \quad y^{(i)}(w^T x^{(i)} + b) \geq 1 \quad \forall i$$

This is a convex quadratic optimization problem with linear constraints.

Soft Margin SVM: To handle non-linearly separable data, we introduce slack variables $\xi_i \geq 0$:

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \quad \text{subject to} \quad y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i$$

where C controls the trade-off between margin maximization and error penalty.

Nonlinear SVM (Kernel Trick): To separate non-linearly separable data, we map inputs to a higher-dimensional feature space via $\phi(x)$, and compute dot products using a kernel function:

$$K(x, x') = \phi(x)^T \phi(x')$$

Popular kernels:

- Polynomial: $K(x, x') = (x^T x' + c)^d$
- RBF (Gaussian): $K(x, x') = \exp(-\gamma \|x - x'\|^2)$

Advantages of SVMs:

- Effective in high-dimensional spaces
- Theoretical guarantees (maximum margin)
- Robust to overfitting (with appropriate kernel and regularization)

1.3.4 Decision Trees

Decision Trees (DTs) are supervised learning algorithms used for both classification and regression. They model decisions as a tree-like structure, where each internal node corresponds to a feature test, each branch to an outcome, and each leaf to a class label or numerical value.

Structure:

- **Root node:** The first decision node based on the best splitting feature.
- **Internal nodes:** Contain feature-based decisions.
- **Leaves:** Contain output predictions.

Training Objective: To construct a tree, we recursively split the data to maximize the **purity** of the resulting subsets.

Metrics to Evaluate Splits:

1. Gini Impurity:

$$G(t) = 1 - \sum_{k=1}^K p_k^2$$

where p_k is the fraction of instances of class k in node t .

2. Entropy (Information Gain):

$$H(t) = - \sum_{k=1}^K p_k \log_2 p_k$$

$$IG = H(\text{parent}) - \left[\frac{n_{\text{left}}}{n} H(\text{left}) + \frac{n_{\text{right}}}{n} H(\text{right}) \right]$$

3. Classification Error:

$$E(t) = 1 - \max_k p_k$$

Recursive Tree-Building Process:

- At each node, evaluate all possible feature splits.
- Choose the split that maximizes information gain or minimizes impurity.
- Recurse until:
 - All samples belong to one class
 - Max depth is reached
 - Minimum samples per leaf is met

Pruning: Decision trees tend to overfit. Pruning techniques (e.g., cost-complexity pruning) reduce the size of the tree to improve generalization.

Advantages of Decision Trees:

- Easy to understand and interpret
- Handle both numerical and categorical data
- Require little preprocessing (e.g., no normalization)

Disadvantages:

- Prone to overfitting on noisy data
- Unstable to small data perturbations
- Greedy splits may lead to suboptimal trees

Extensions:

- **Random Forests:** Ensembles of decision trees using bagging.
- **Gradient Boosted Trees:** Sequentially added trees that minimize residuals.

1.3.5 Clustering

Clustering is an unsupervised learning technique where the goal is to group a set of objects such that objects in the same group (called a cluster) are more similar to each other than to those in other groups.

Unlike supervised learning, clustering does not require labeled data. It is useful in exploratory data analysis, anomaly detection, customer segmentation, image compression, etc.

1. K-Means Clustering

K-Means is a centroid-based, iterative algorithm that partitions data into K clusters by minimizing the variance within each cluster.

Objective: Given m data points $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \in \mathbb{R}^n$, and a chosen number of clusters K , the goal is to minimize:

$$J = \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|^2$$

where μ_k is the mean (centroid) of cluster C_k .

Algorithm:

1. Initialize K centroids randomly: μ_1, \dots, μ_K
2. Assign each point x_i to the closest centroid:

$$c_i = \arg \min_k \|x_i - \mu_k\|$$

3. Recompute each centroid as:

$$\mu_k = \frac{1}{|C_k|} \sum_{x_i \in C_k} x_i$$

4. Repeat steps 2–3 until convergence (i.e., assignments no longer change or maximum iterations reached).

Limitations:

- Sensitive to initialization (can lead to suboptimal solutions).
- Assumes spherical and equally sized clusters.
- Requires manual selection of K .

Improvement: Use **K-Means++** for smarter initialization of centroids to improve convergence and quality.

2. Gaussian Mixture Models (GMM)

GMM assumes the data is generated from a mixture of several Gaussian distributions with unknown parameters. It is a **soft clustering** method, where each point belongs to each cluster with a certain probability.

Model:

$$p(x) = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(x|\mu_k, \Sigma_k)$$

where:

- π_k : Mixing coefficient for cluster k ($\sum_k \pi_k = 1$)
- $\mathcal{N}(x|\mu_k, \Sigma_k)$: Multivariate normal distribution

Training via Expectation-Maximization (EM):

1. **E-Step:** Compute posterior probabilities (responsibilities) of cluster k for data point x_i :

$$\gamma_{ik} = \frac{\pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_i|\mu_j, \Sigma_j)}$$

2. **M-Step:** Update parameters:

$$\mu_k = \frac{\sum_{i=1}^m \gamma_{ik} x_i}{\sum_{i=1}^m \gamma_{ik}}, \quad \Sigma_k = \frac{\sum_{i=1}^m \gamma_{ik} (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^m \gamma_{ik}}, \quad \pi_k = \frac{1}{m} \sum_{i=1}^m \gamma_{ik}$$

3. Repeat E and M steps until convergence.

Advantages:

- Probabilistic soft clustering
- Models ellipsoidal clusters

Disadvantages:

- Sensitive to initialization and outliers
- Requires estimating more parameters than K-Means

Conclusion: These foundational techniques form the basis for understanding and building machine learning models. They also provide essential insights into how LLMs learn statistical patterns and structure in large datasets during pretraining and fine-tuning stages.

Chapter 2

Deep Learning Fundamentals

2.1 Neural Networks and Architectures

2.1.1 Feedforward Neural Networks (FFNN)

Feedforward Neural Networks (FFNNs), also known as Multi-Layer Perceptrons (MLPs), are the foundational building blocks of deep learning. They consist of layers of neurons where information flows in one direction—from input to output—without cycles or loops.

Network Architecture:

- **Input Layer:** Receives raw features $x \in \mathbb{R}^n$
- **Hidden Layers:** Each layer transforms input using a linear operation followed by a non-linear activation function
- **Output Layer:** Produces the final prediction, typically with a softmax or sigmoid activation

Forward Pass (Mathematical Formulation):

For a single hidden layer:

$$z^{(1)} = W^{(1)}x + b^{(1)} \quad (\text{Linear transform})$$

$$a^{(1)} = \sigma(z^{(1)}) \quad (\text{Activation})$$

$$z^{(2)} = W^{(2)}a^{(1)} + b^{(2)} \quad (\text{Output logits})$$

$$\hat{y} = f(z^{(2)}) \quad (\text{Output activation: softmax or sigmoid})$$

Where:

- $W^{(l)}$ and $b^{(l)}$: Weights and biases for layer l
- σ : Non-linear activation function (e.g., ReLU, tanh)
- f : Output activation (e.g., softmax for classification)

Loss Function:

For classification tasks with softmax output:

$$\mathcal{L} = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

Where:

- y_i : Ground truth (one-hot vector)
- \hat{y}_i : Predicted probability for class i

Backpropagation: Training the Network

Training is done using Gradient Descent or its variants (SGD, Adam) to minimize the loss. This involves computing the gradient of the loss function with respect to each parameter using the chain rule (backpropagation).

For a layer l , let $\delta^{(l)}$ be the error term (derivative of loss with respect to $z^{(l)}$).

$$\delta^{(2)} = \hat{y} - y \quad (\text{For output layer})$$

$$\delta^{(1)} = (W^{(2)})^T \delta^{(2)} \odot \sigma'(z^{(1)}) \quad (\text{For hidden layer})$$

Gradient updates:

$$\frac{\partial \mathcal{L}}{\partial W^{(l)}} = \delta^{(l)} (a^{(l-1)})^T, \quad \frac{\partial \mathcal{L}}{\partial b^{(l)}} = \delta^{(l)}$$

Update rule:

$$W^{(l)} := W^{(l)} - \eta \frac{\partial \mathcal{L}}{\partial W^{(l)}}, \quad b^{(l)} := b^{(l)} - \eta \frac{\partial \mathcal{L}}{\partial b^{(l)}}$$

Where η is the learning rate.

Activation Functions:

- **Sigmoid:** $\sigma(x) = \frac{1}{1+e^{-x}}$
- **Tanh:** $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
- **ReLU:** $\text{ReLU}(x) = \max(0, x)$

ReLU is most widely used in hidden layers due to its non-saturating gradient and efficiency.

Advantages of FFNNs:

- Can approximate any continuous function (Universal Approximation Theorem)
- Scalable to large datasets and tasks

Limitations:

- Require large amounts of labeled data
- Prone to overfitting without regularization (dropout, weight decay)
- Lack temporal memory (handled better by RNNs/LSTMs)

2.1.2 Convolutional Neural Networks (CNNs)

Convolutional Neural Networks (CNNs) are deep learning models specialized for processing data with a grid-like topology, such as images (2D grids of pixels). CNNs leverage spatial hierarchies in data and reduce the number of trainable parameters through shared weights.

Motivation: Traditional FFNNs do not scale well with high-dimensional inputs like images. For a 256×256 RGB image, a fully connected layer would have millions of parameters. CNNs solve this by using local receptive fields and weight sharing.

Key Components of CNN:

- **Convolutional Layer**
- **Activation Function** (typically ReLU)
- **Pooling Layer** (e.g., Max Pooling)
- **Fully Connected Layer**

1. Convolution Operation:

A convolutional layer performs discrete convolution between an input image and a filter (also called a kernel). For a 2D input I and kernel K of size $f \times f$, the convolution operation is defined as:

$$S(i, j) = (I * K)(i, j) = \sum_{m=0}^{f-1} \sum_{n=0}^{f-1} I(i+m, j+n) \cdot K(m, n)$$

This operation is applied over the entire image with a sliding window, producing a feature map.

Hyperparameters:

- **Stride (s):** Step size by which the kernel is moved.
- **Padding (p):** Adding zeros around the border to control the output size.

Output dimension:

$$\text{Output size} = \left\lfloor \frac{N - f + 2p}{s} \right\rfloor + 1$$

2. Activation Function: Commonly used is the ReLU:

$$\text{ReLU}(x) = \max(0, x)$$

ReLU introduces non-linearity and speeds up convergence by avoiding vanishing gradient problems.

3. Pooling Layer:

Pooling downsamples the feature maps, reducing the spatial dimensions and computation. Most common is max pooling.

$$\text{MaxPool}(X) = \max\{x_1, x_2, \dots, x_k\}$$

This retains the most prominent features and makes the network more invariant to small translations.

4. Fully Connected Layer: After several convolution and pooling layers, the output is flattened and passed through one or more fully connected layers for final prediction.

CNN Architecture Example:

- Input: $32 \times 32 \times 3$ image
- Conv layer (5x5 kernel, 6 filters) \rightarrow ReLU $\rightarrow 28 \times 28 \times 6$
- MaxPool (2x2) $\rightarrow 14 \times 14 \times 6$
- Conv layer (5x5, 16 filters) \rightarrow ReLU $\rightarrow 10 \times 10 \times 16$
- MaxPool $\rightarrow 5 \times 5 \times 16$
- Flatten \rightarrow Fully Connected Layer \rightarrow Output

Advantages of CNNs:

- Fewer parameters due to weight sharing
- Captures local spatial features
- Translation invariance due to pooling
- Performs exceptionally well on image and spatial data

Applications of CNNs:

- Image classification (e.g., ImageNet)
- Object detection (e.g., YOLO, Faster R-CNN)
- Image segmentation (e.g., U-Net)
- Facial recognition and medical imaging

Limitations of CNNs:

- Require large labeled datasets and high computation
- Not ideal for sequential or long-range dependent data (handled better by RNNs or Transformers)
- Do not inherently encode positional relationships as well as attention-based models

2.2 Optimization Techniques

Optimization in deep learning is the process of adjusting neural network weights to minimize a cost (loss) function. This is crucial for effective model training.

2.2.1 Gradient Descent

Gradient Descent iteratively updates weights to minimize a cost function $J(\theta)$:

$$\theta := \theta - \eta \nabla_{\theta} J(\theta)$$

- θ : Model parameters (e.g., weights, biases)
- η : Learning rate (controls the step size)
- $\nabla_{\theta} J(\theta)$: Gradient of the loss function with respect to θ

2.2.2 Stochastic Gradient Descent (SGD)

SGD updates weights using a single training example:

$$\theta := \theta - \eta \nabla_{\theta} J(\theta; x^{(i)}, y^{(i)})$$

- $x^{(i)}, y^{(i)}$: Input and label of the i^{th} training sample

2.2.3 Mini-batch Gradient Descent

Mini-batch Gradient Descent updates using a subset of training data:

$$\theta := \theta - \eta \nabla_{\theta} \frac{1}{m} \sum_{i=1}^m J(\theta; x^{(i)}, y^{(i)})$$

- m : Mini-batch size

2.2.4 Momentum

Momentum helps accelerate learning by using a velocity term:

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta), \quad \theta := \theta - v_t$$

- v_t : Velocity vector at time t
- γ : Momentum coefficient (typically 0.9)

2.2.5 RMSProp

RMSProp adapts the learning rate using a moving average of squared gradients:

$$E[g^2]_t = \rho E[g^2]_{t-1} + (1 - \rho)g_t^2, \quad \theta := \theta - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}}g_t$$

- ρ : Decay rate (e.g., 0.9)
- g_t : Gradient at step t
- ϵ : Small value to avoid division by zero

2.2.6 Adam Optimizer

Adam combines ideas from Momentum and RMSProp:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_t \quad (\text{1st moment})$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)g_t^2 \quad (\text{2nd moment})$$

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

$$\theta := \theta - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$$

- β_1, β_2 : Exponential decay rates (default: 0.9 and 0.999)
- m_t, v_t : Estimates of first and second moments of gradient
- \hat{m}_t, \hat{v}_t : Bias-corrected estimates

—

2.3 Regularization Techniques

Regularization prevents overfitting by discouraging complex models.

2.3.1 L1 and L2 Regularization

Regularization adds penalty terms to the loss:

- **L1 Regularization (Lasso):**

$$J'(\theta) = J(\theta) + \lambda \sum_i |\theta_i|$$

- **L2 Regularization (Ridge):**

$$J'(\theta) = J(\theta) + \frac{\lambda}{2} \sum_i \theta_i^2$$

- λ : Regularization strength (hyperparameter)

2.3.2 Dropout

Dropout randomly disables neurons during training:

$$\tilde{h}_j = r_j h_j, \quad \text{where } r_j \sim \text{Bernoulli}(p)$$

- h_j : Activation of neuron j
- r_j : Random mask (0 or 1 with probability p)
- p : Keep probability (e.g., 0.5)

2.3.3 Early Stopping

Early stopping halts training when validation performance degrades. It avoids overfitting by selecting the optimal epoch.

2.3.4 Batch Normalization

BatchNorm normalizes activations across the mini-batch:

$$\hat{x}^{(i)} = \frac{x^{(i)} - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}, \quad y^{(i)} = \gamma \hat{x}^{(i)} + \beta$$

- $\mu_{\mathcal{B}}, \sigma_{\mathcal{B}}^2$: Mean and variance of mini-batch
- γ, β : Learnable scale and shift parameters
- ϵ : Small constant for numerical stability

2.3.5 Data Augmentation

This expands the dataset using transformations:

- **Images:** Rotation, flipping, scaling, color jittering
- **Text:** Synonym replacement, random insertion/deletion
- **Audio:** Time shift, noise addition, speed tuning

This improves generalization by increasing sample variety.

Chapter 3

Natural Language Processing (NLP)

3.1 Introduction to Natural Language Processing (NLP)

Natural Language Processing (NLP) is a subfield of Artificial Intelligence (AI) that focuses on enabling machines to understand, interpret, generate, and interact using human languages. It lies at the intersection of linguistics, computer science, and cognitive science, and plays a crucial role in the development of intelligent systems that interact with users in a natural and intuitive manner.

3.1.1 Motivation

Human language is inherently ambiguous, context-dependent, and variable. The goal of NLP is to build models and algorithms that can:

- Parse and understand the structure and meaning of text and speech.
- Extract relevant information from unstructured data.
- Enable tasks like translation, summarization, sentiment analysis, and question answering.

3.1.2 Historical Evolution

- **1950s–1980s:** Rule-based systems and symbolic AI dominated the early NLP era, relying on handcrafted linguistic rules.
- **1990s:** The rise of machine learning techniques led to statistical NLP, including models like Hidden Markov Models (HMMs) and Naive Bayes classifiers.
- **2010s–present:** Deep learning revolutionized NLP, with neural architectures like RNNs, LSTMs, GRUs, and Transformers leading to major breakthroughs.

3.1.3 Key Components of NLP

- **Lexical Analysis:** Breaking down text into tokens (words, punctuation).
- **Syntax Analysis (Parsing):** Analyzing grammatical structure using parse trees.
- **Semantic Analysis:** Interpreting the meaning of words and sentences.

- **Pragmatics and Discourse:** Understanding meaning in context and multi-sentence structures.
- **Language Generation:** Producing grammatically and semantically coherent text.

3.1.4 Challenges in NLP

- Ambiguity in word meaning (e.g., “bank” as a financial institution vs. river bank)
- Contextual dependency and polysemy
- Morphological complexity in different languages
- Sarcasm, irony, and figurative language
- Low-resource languages and multilingual support

3.1.5 Applications of NLP

- Machine Translation (e.g., Google Translate)
- Chatbots and Conversational Agents (e.g., Siri, Alexa)
- Sentiment Analysis (e.g., analyzing reviews or tweets)
- Information Retrieval (e.g., search engines)
- Named Entity Recognition, Part-of-Speech Tagging, Text Summarization

3.1.6 Role of Deep Learning in NLP

Recent advancements like word embeddings (e.g., Word2Vec, GloVe), pre-trained language models (e.g., BERT, GPT), and transformer-based architectures have significantly improved the state-of-the-art in NLP tasks.

In the following sections, we will explore the key models, methodologies, and architectures that have driven progress in NLP, with a special focus on how they are integrated within Large Language Models (LLMs).

3.2 NLP Pipeline

Natural Language Processing (NLP) involves transforming raw text into a form that machines can understand and process. The NLP pipeline refers to the sequence of steps or stages involved in this transformation, from raw input to useful structured output.

3.2.1 Text Collection

This is the initial stage where unstructured text data is gathered from various sources such as websites, books, social media, emails, or logs.

- Web scraping using tools like **BeautifulSoup** or **Scrapy**
- APIs from Twitter, Reddit, etc.
- Predefined corpora (e.g., Brown Corpus, Gutenberg, Wikipedia)

3.2.2 Text Preprocessing

Raw text often contains noise. Preprocessing cleans and standardizes the input.

- **Tokenization:** Splitting text into words, phrases, or sentences.
- **Lowercasing:** Converting all text to lowercase.
- **Stopword Removal:** Removing common non-informative words (e.g., "the", "is").
- **Punctuation Removal:** Eliminating symbols like ".", ",", etc.
- **Stemming:** Reducing words to their base form (e.g., "running" → "run") using algorithms like Porter Stemmer.
- **Lemmatization:** Morphologically reducing words using vocabulary and POS tags (e.g., "better" → "good").

3.2.3 Text Representation / Feature Extraction

After cleaning, text must be represented numerically for machine learning models.

- **Bag of Words (BoW):** Counts word frequency without considering order.
- **TF-IDF (Term Frequency–Inverse Document Frequency):** Weights frequent words lower if they appear in many documents.
- **Word Embeddings:** Vector representation capturing semantic meaning (e.g., Word2Vec, GloVe, FastText).
- **Contextual Embeddings:** Dynamic embeddings that consider context (e.g., BERT, ELMo).

3.2.4 Model Building

The extracted features are used to train models for various NLP tasks.

- **Traditional Models:** Naive Bayes, SVMs, Decision Trees, Logistic Regression.
- **Deep Learning Models:** RNNs, LSTMs, GRUs, CNNs for sequence tasks.
- **Transformers:** Attention-based models (e.g., BERT, GPT) that have revolutionized NLP.

3.2.5 Evaluation

Models are evaluated using appropriate metrics based on the task.

- **Classification:** Accuracy, Precision, Recall, F1 Score
- **Sequence Labeling:** Token-level precision/recall
- **Generation:** BLEU, ROUGE, perplexity

3.2.6 Deployment and Feedback

Once trained, the model can be deployed via APIs or as part of an application. Continuous feedback and retraining help improve the system over time.

- **Deployment Tools:** Flask, FastAPI, Docker, Kubernetes
- **Monitoring:** Logging incorrect predictions, drift detection

3.2.7 Summary

The NLP pipeline ensures that raw text data is transformed in a structured and meaningful way, enabling various applications such as sentiment analysis, chatbots, information retrieval, and machine translation.

3.3 Traditional NLP Methods

Before the rise of deep learning, Natural Language Processing primarily relied on statistical and rule-based methods. These approaches operated on simple vector space representations of text, focusing on frequency and co-occurrence patterns rather than learned contextual semantics.

3.3.1 Bag-of-Words (BoW)

The Bag-of-Words model represents a document as a multiset (bag) of its words, disregarding grammar and word order but keeping multiplicity.

Given a vocabulary $V = \{w_1, w_2, \dots, w_N\}$, each document d is represented as a vector $\vec{v}_d \in \mathbb{R}^N$, where each component $v_{d,i}$ is the count of word w_i in document d .

$$\vec{v}_d = [\text{count}(w_1), \text{count}(w_2), \dots, \text{count}(w_N)]$$

Advantages:

- Simple and efficient to implement.
- Works well with linear models like Naive Bayes and Logistic Regression.

Disadvantages:

- Ignores context and word order.
- Large vocabulary leads to sparse vectors.
- Doesn't capture semantics (e.g., synonyms have unrelated vectors).

3.3.2 N-Grams

N-grams are contiguous sequences of n items from a given text. They are used to model word sequences and local context.

Unigram: $\{w_1, w_2, \dots, w_T\}$, Bigram: $\{(w_1, w_2), (w_2, w_3), \dots, (w_{T-1}, w_T)\}$

Probability Estimation:

For bigrams, the probability of a sentence is approximated using the chain rule:

$$P(w_1, w_2, \dots, w_T) \approx \prod_{t=1}^T P(w_t | w_{t-1})$$

Smoothed Estimation:

$$P(w_t | w_{t-1}) = \frac{\text{count}(w_{t-1}, w_t) + \alpha}{\text{count}(w_{t-1}) + \alpha|V|}$$

Where α is the smoothing factor (Laplace smoothing), and $|V|$ is vocabulary size.

Pros and Cons:

- Captures some local context.
- Still suffers from data sparsity and can't model long-range dependencies.

3.3.3 Term Frequency–Inverse Document Frequency (TF-IDF)

TF-IDF improves upon raw counts by down-weighting common words (like "the", "is") and up-weighting informative ones.

Term Frequency (TF):

$$TF(t, d) = \frac{\text{count of term } t \text{ in document } d}{\text{total terms in } d}$$

Inverse Document Frequency (IDF):

$$IDF(t) = \log \left(\frac{N}{1 + |\{d : t \in d\}|} \right)$$

TF-IDF Score:

$$TF\text{-}IDF(t, d) = TF(t, d) \cdot IDF(t)$$

Advantages:

- Reduces impact of common but uninformative terms.
- Works well for document classification and information retrieval.

Limitations:

- Still based on word frequency, not semantics.
- High-dimensional and sparse representations.

3.3.4 Limitations of Traditional Methods

- No handling of word meaning (e.g., "good" vs. "great").
- Struggles with polysemy and synonymy.
- Feature engineering required for downstream tasks.
- Cannot model context or sequential structure effectively.

Despite their simplicity, traditional NLP methods laid the groundwork for more advanced, neural-based techniques. In the following sections, we explore distributed representations and neural models that address these limitations.

3.4 Word Embeddings

Traditional NLP methods like Bag-of-Words and TF-IDF represent words as discrete indices in high-dimensional sparse vectors, failing to capture semantic relationships or context. **Word embeddings** are dense, low-dimensional continuous vector representations of words that preserve semantic similarity and are learned from large corpora.

3.4.1 Motivation

The goal is to represent each word w in a vocabulary V as a vector $\vec{w} \in \mathbb{R}^d$, where similar words (e.g., "king" and "queen") have nearby vector representations:

$$\text{Similarity}(w_i, w_j) = \cos(\vec{w}_i, \vec{w}_j)$$

3.4.2 Word2Vec

Proposed by Mikolov et al. (2013), Word2Vec includes two architectures:

- **Skip-Gram:** Predicts context words given the center word.
- **CBOW (Continuous Bag-of-Words):** Predicts the center word given context words.

Skip-Gram Model: Given a word sequence w_1, w_2, \dots, w_T , the skip-gram objective maximizes the probability of context words $w_{t \pm j}$ given the center word w_t :

$$\mathcal{L}_{\text{skip-gram}} = \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log P(w_{t+j} | w_t)$$

where c is the context window size. Using softmax:

$$P(w_O | w_I) = \frac{\exp(\vec{v}_{w_O}^\top \vec{v}_{w_I})}{\sum_{w \in V} \exp(\vec{v}_w^\top \vec{v}_{w_I})}$$

To avoid computing the full softmax denominator, techniques like **Negative Sampling** and **Hierarchical Softmax** are used.

CBOW Model: Given a context window around word w_t , CBOW predicts w_t using the average of context vectors:

$$\mathcal{L}_{\text{CBOW}} = \sum_{t=1}^T \log P(w_t | w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c})$$

3.4.3 GloVe (Global Vectors)

GloVe, proposed by Pennington et al. (2014), learns embeddings by factorizing a word co-occurrence matrix. The core idea is that word meaning can be captured from global co-occurrence statistics.

Let X_{ij} be the number of times word j occurs in the context of word i . The model tries to learn word vectors $\vec{w}_i, \vec{w}_j \in \mathbb{R}^d$ such that:

$$\vec{w}_i^\top \vec{w}_j + b_i + \tilde{b}_j = \log(X_{ij})$$

Loss Function:

$$\mathcal{L}_{\text{GloVe}} = \sum_{i,j=1}^{|V|} f(X_{ij}) \left(\vec{w}_i^\top \vec{w}_j + b_i + \tilde{b}_j - \log X_{ij} \right)^2$$

where $f(X_{ij})$ is a weighting function:

$$f(x) = \begin{cases} \left(\frac{x}{x_{\max}} \right)^\alpha & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases}$$

Typically, $\alpha = 0.75$, $x_{\max} = 100$.

3.4.4 Properties of Word Embeddings

- Capture semantic and syntactic similarity:

$$\text{cosine}(\vec{w}_{\text{king}}, \vec{w}_{\text{queen}}) \approx \text{cosine}(\vec{w}_{\text{man}}, \vec{w}_{\text{woman}})$$

- Enable vector arithmetic:

$$\vec{w}_{\text{king}} - \vec{w}_{\text{man}} + \vec{w}_{\text{woman}} \approx \vec{w}_{\text{queen}}$$

- Dense representations reduce sparsity and generalize better than BoW.

3.4.5 Limitations

- Fixed embeddings: do not adapt based on sentence context.
- Struggle with polysemy (e.g., “bank” as a financial institution vs. river bank).
- Cannot capture compositionality or long-term dependencies.

These limitations led to the development of **contextual embeddings** such as ELMo, BERT, and GPT, which are covered in the next sections.

3.5 Recurrent Neural Networks (RNNs) in NLP

Recurrent Neural Networks (RNNs) are a class of neural networks designed for sequential data, making them highly suitable for Natural Language Processing tasks such as language modeling, text generation, machine translation, and named entity recognition.

3.5.1 Motivation

Most natural language tasks require modeling dependencies between words. Unlike feed-forward networks, RNNs maintain a **hidden state** that captures information from previous time steps, allowing them to model sequences of arbitrary length.

$$\text{Example: } P(w_1, w_2, \dots, w_T) = \prod_{t=1}^T P(w_t | w_1, \dots, w_{t-1})$$

3.5.2 RNN Architecture and Equations

At each time step t , an RNN takes an input vector $\vec{x}_t \in \mathbb{R}^d$ and updates a hidden state $\vec{h}_t \in \mathbb{R}^h$ using the following recurrence relation:

$$\begin{aligned}\vec{h}_t &= \tanh(W_{xh}\vec{x}_t + W_{hh}\vec{h}_{t-1} + \vec{b}_h) \\ \vec{y}_t &= W_{hy}\vec{h}_t + \vec{b}_y\end{aligned}$$

Where:

- W_{xh} is the input-to-hidden weight matrix.
- W_{hh} is the hidden-to-hidden (recurrent) weight matrix.
- W_{hy} is the hidden-to-output weight matrix.
- \vec{h}_0 is typically initialized to zero.

The model is trained to minimize a loss function such as cross-entropy over a sequence of predictions.

3.5.3 Backpropagation Through Time (BPTT)

To train RNNs, the gradients are computed through all time steps via a technique called **Backpropagation Through Time (BPTT)**. The loss is summed over all time steps:

$$\mathcal{L} = \sum_{t=1}^T \ell(\vec{y}_t, \vec{y}_t^{\text{true}})$$

Gradients are backpropagated through the unrolled network. However, this leads to issues with:

- **Vanishing Gradients:** Gradients shrink as they are backpropagated, making it hard to learn long-range dependencies.
- **Exploding Gradients:** Gradients grow exponentially and destabilize training.

3.5.4 Applications of RNNs in NLP

- **Language Modeling:** Predicting the next word given a sequence.
- **Machine Translation:** Encoder-decoder RNNs map source sentences to target sentences.
- **Speech Recognition and Tagging:** Sequence labeling for input data like audio or tokens.
- **Text Generation:** Sampling text one word at a time based on learned distributions.

3.5.5 Bidirectional RNNs (BiRNN)

Standard RNNs process sequences in one direction (past \rightarrow future). **Bidirectional RNNs** use two RNNs: one processes the input forward, and another backward, combining both hidden states:

$$\vec{h}_t = [\vec{h}_t; \overleftarrow{h}_t]$$

This allows access to both past and future context, improving performance on tasks like Named Entity Recognition or POS tagging.

3.5.6 Limitations of Vanilla RNNs

- Cannot retain long-term dependencies due to vanishing gradient problem.
- Sequential computation inhibits parallelization during training.
- Poor at capturing hierarchical sentence structures.

These limitations motivated the development of improved architectures like LSTM (Long Short-Term Memory) and GRU (Gated Recurrent Unit), which we will explore next.

3.6 Long Short-Term Memory (LSTM)

Long Short-Term Memory (LSTM) networks are a special type of Recurrent Neural Networks (RNNs) designed to address the vanishing and exploding gradient problems in standard RNNs. LSTMs are particularly effective at learning long-term dependencies, making them well-suited for many NLP tasks.

3.6.1 Motivation

Traditional RNNs struggle with learning patterns that span long distances in a sequence due to the repeated multiplication of gradients through time. LSTMs solve this by introducing an explicit memory cell and gating mechanisms that control the flow of information.

Key Idea: Maintain a memory cell \vec{c}_t whose updates are regulated by learnable gates.

3.6.2 LSTM Cell Structure

An LSTM cell contains:

- **Forget Gate** \vec{f}_t
- **Input Gate** \vec{i}_t
- **Output Gate** \vec{o}_t
- **Cell State** \vec{c}_t
- **Hidden State** \vec{h}_t

Inputs:

$\vec{x}_t \in \mathbb{R}^d$ (input vector), $\vec{h}_{t-1}, \vec{c}_{t-1} \in \mathbb{R}^h$ (previous hidden and cell states)

3.6.3 LSTM Equations

$$\begin{aligned}\vec{f}_t &= \sigma(W_f \vec{x}_t + U_f \vec{h}_{t-1} + \vec{b}_f) \quad (\text{Forget Gate}) \\ \vec{i}_t &= \sigma(W_i \vec{x}_t + U_i \vec{h}_{t-1} + \vec{b}_i) \quad (\text{Input Gate}) \\ \vec{c}_t &= \tanh(W_c \vec{x}_t + U_c \vec{h}_{t-1} + \vec{b}_c) \quad (\text{Cell Candidate}) \\ \vec{c}_t &= \vec{f}_t \odot \vec{c}_{t-1} + \vec{i}_t \odot \vec{c}_t \quad (\text{New Cell State}) \\ \vec{o}_t &= \sigma(W_o \vec{x}_t + U_o \vec{h}_{t-1} + \vec{b}_o) \quad (\text{Output Gate}) \\ \vec{h}_t &= \vec{o}_t \odot \tanh(\vec{c}_t) \quad (\text{New Hidden State})\end{aligned}$$

Where:

- σ is the sigmoid activation function.
- \tanh is the hyperbolic tangent function.
- \odot denotes element-wise multiplication.
- W_*, U_*, b_* are learnable parameters.

3.6.4 Intuition Behind the Gates

- **Forget Gate \vec{f}_t** : Decides what portion of the previous memory to forget.
- **Input Gate \vec{i}_t** : Determines how much of the new candidate memory to add.
- **Cell Update \vec{c}_t** : Combines old and new memory.
- **Output Gate \vec{o}_t** : Controls what part of the memory to output as hidden state.

3.6.5 Applications in NLP

LSTMs are widely used in:

- **Text Classification**
- **Language Modeling**
- **Named Entity Recognition (NER)**
- **Machine Translation (with Encoder–Decoder frameworks)**
- **Speech Recognition**

3.6.6 Bidirectional LSTMs

Like RNNs, LSTMs can be made bidirectional:

$$\vec{h}_t = [\vec{h}_t; \overleftarrow{h}_t]$$

where \vec{h}_t is from forward LSTM and \overleftarrow{h}_t is from backward LSTM.

3.6.7 Advantages over RNNs

- Better at capturing long-term dependencies.
- Mitigates vanishing gradient problem.
- More expressive due to gating mechanisms.

8. Limitations

- Computationally intensive due to multiple gates.
- Still sequential—cannot fully parallelize across time.
- Less effective than Transformers on long sequences.

These limitations motivated the development of attention mechanisms and Transformer architectures, which are covered in the next section.

3.7 Attention Mechanisms in NLP

Attention mechanisms have revolutionized Natural Language Processing by enabling models to dynamically focus on relevant parts of an input sequence when generating an output. Unlike traditional sequence models that compress entire sequences into a fixed-size vector, attention allows flexible and context-aware alignment between input and output tokens.

3.7.1 Motivation

In tasks such as machine translation or summarization, not all input words are equally important for producing a given output word. Traditional RNNs and LSTMs often suffer from the bottleneck of encoding the entire sequence into a single context vector.

Attention addresses this by allowing the model to look back at the entire input sequence.

3.7.2 General Form of Attention

Given a query \vec{q} , a set of key-value pairs $\{(\vec{k}_i, \vec{v}_i)\}_{i=1}^n$, attention computes a weighted sum of the values:

$$\text{Attention}(\vec{q}, K, V) = \sum_{i=1}^n \alpha_i \vec{v}_i$$

Where the weights α_i are computed via a similarity function between the query and each key:

$$\alpha_i = \frac{\exp(\text{score}(\vec{q}, \vec{k}_i))}{\sum_{j=1}^n \exp(\text{score}(\vec{q}, \vec{k}_j))}$$

3.7.3 Scoring Functions

Common choices for the score function:

- Dot-product: $\text{score}(\vec{q}, \vec{k}_i) = \vec{q}^\top \vec{k}_i$
- Additive (Bahdanau): $\text{score}(\vec{q}, \vec{k}_i) = \vec{v}_a^\top \tanh(W_q \vec{q} + W_k \vec{k}_i)$

3.7.4 Types of Attention

- **Additive Attention (Bahdanau et al., 2015):** Introduced in neural machine translation using encoder-decoder models with alignment scores.
- **Dot-Product Attention (Luong et al., 2015):** Uses the dot product between query and key vectors for simplicity and efficiency.
- **Scaled Dot-Product Attention:** Used in Transformers, scales the dot product by $\sqrt{d_k}$ to stabilize gradients:

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

3.7.5 Self-Attention

In self-attention, the queries, keys, and values all come from the same source sequence:

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

Where:

- $X \in \mathbb{R}^{n \times d}$: sequence of token embeddings
- $W^Q, W^K, W^V \in \mathbb{R}^{d \times d_k}$: learned projections

Self-attention captures dependencies between all tokens in the sequence regardless of their position.

3.7.6 Multi-Head Attention

To allow the model to attend to information from different representation subspaces, multiple attention heads are used:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

Each head is computed as:

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

This enhances the model's ability to capture various aspects of similarity and relevance across positions.

3.7.7 Applications in NLP

- **Machine Translation:** Alignment between source and target sentences.
- **Text Summarization:** Focusing on key parts of a document.
- **Question Answering:** Attending to relevant context for accurate answers.
- **Transformers:** The foundation of models like BERT, GPT, T5, etc.

3.7.8 Advantages of Attention

- Removes fixed-length bottleneck in sequence encoding.
- Enables parallelization unlike RNNs.
- Improves long-range dependency modeling.

Attention paved the way for the development of the **Transformer** architecture, which completely replaces recurrence with stacked self-attention layers.