

1. Transformer Architecture

Topic: Explain the Encoder-Decoder structure and Self-Attention mechanism.

The Transformer architecture, introduced in the paper "Attention Is All You Need," revolutionized Natural Language Processing (NLP) by dispensing with recurrence and convolutions entirely, relying instead on a mechanism called Self-Attention. This architecture is built upon a standard Encoder-Decoder structure.

The Encoder-Decoder Structure: The Encoder maps an input sequence of symbol representations (like words in a sentence) to a sequence of continuous representations. It consists of a stack of identical layers, each containing two sub-layers: a multi-head self-attention mechanism and a position-wise fully connected feed-forward network. The encoder's job is to "understand" the input by creating a rich, contextual representation of it.

The Decoder also consists of a stack of identical layers. In addition to the two sub-layers found in the encoder, the decoder inserts a third sub-layer that performs multi-head attention over the output of the encoder stack. This allows the decoder to focus on relevant parts of the input sentence while generating the output sequence one element at a time. Crucially, the decoder is masked to ensure that predictions for position \$i\$ can depend only on the known outputs at positions less than i.

The Self-Attention Mechanism: Self-attention is the core engine that allows the model to weigh the importance of different words in a sentence regardless of their distance from each other. It relates different positions of a single sequence to compute a representation of the sequence. Mathematically, this is often described using Queries (Q), Keys (K), and Values (V). The input vector is multiplied by weight matrices to create these three vectors. The attention score is calculated by taking the dot product of the Query with the Key, scaling it, and applying a softmax function to obtain weights. These weights are then applied to the Values:

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

This mechanism allows the model to look at other words in the input sentence to better encode the current word. For example, in "The animal didn't cross the street because it was too tired," self-attention associates "it" with "animal," allowing for nuanced understanding of context.

2. Backpropagation

Topic: Detail the chain rule and how gradients update weights in a neural network.

Backpropagation, or "backward propagation of errors," is the fundamental algorithm used to train neural networks. It is the method by which the network adjusts its internal parameters (weights and biases) to minimize the difference between its predicted output and the actual target output.

The Chain Rule: At the heart of backpropagation is the calculus concept known as the Chain Rule. A neural network is essentially a composite function—a chain of functions nested within each other (layers). To find out how a change in a weight deep inside the network affects the final error (Loss function L), we must calculate the partial derivative of the error with respect to that weight.

The chain rule states that the derivative of a composite function is the product of the derivatives

of the constituent functions. If L depends on y , which depends on x , which depends on w , the gradient is calculated as:

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial x} \cdot \frac{\partial x}{\partial w}$$

During the forward pass, the input data flows through the network to generate a prediction. During the backward pass, the algorithm computes the gradient of the loss function with respect to the weights by working backward from the output layer to the input layer, propagating the error gradient through the network layer by layer.

Updating Weights: Once the gradients for all weights have been computed, the network uses an optimization algorithm, typically Gradient Descent or a variant like Adam, to update the weights. The gradient points in the direction of the steepest increase of the loss function. To minimize the loss, we move the weights in the opposite direction—down the slope. The update rule for a weight w is given by:

$$w_{new} = w_{old} - \eta \cdot \frac{\partial L}{\partial w}$$

Here, η represents the learning rate, a hyperparameter that controls the size of the step taken during the update. If the step is too small, training is slow; if it is too large, the model might overshoot the minimum. Through iterative cycles of forward passes, error calculation, backpropagation, and weight updates, the network gradually converges on a state where the loss is minimized.

3. RAG Systems

Topic: Explain Retrieval-Augmented Generation, vector databases, and semantic search.

Retrieval-Augmented Generation (RAG) is a framework designed to improve the output of Large Language Models (LLMs) by connecting them to external knowledge sources. Standard LLMs are limited to the data they were trained on, which can quickly become outdated. RAG addresses this by retrieving relevant, up-to-date information from a designated knowledge base and injecting it into the user's prompt before the model generates a response. This grounding process significantly reduces hallucinations and improves factual accuracy.

To enable this retrieval, RAG systems typically utilize Vector Databases. These databases store information not as text, but as numerical vectors (embeddings). When data is ingested, it is converted into high-dimensional vectors. When a user queries the system, the vector database identifies stored data chunks that are numerically close to the query vector. The system then pulls these matching chunks to provide context for the LLM. The system relies entirely on these vector matches to decide what information is relevant to the user.

4. GANs

Topic: Explain the Generator vs Discriminator dynamic and training challenges.

Generative Adversarial Networks (GANs) are a class of machine learning frameworks designed to generate new data instances that resemble your training data. The architecture is unique because it pits two neural networks against each other in a competitive game.

The first network is the Generator. Its objective is to create "fake" data samples, such as images or audio, starting from random noise. It tries to produce output that is convincing enough to pass as real data.

The second network is the Discriminator. It acts as a binary classifier. It receives both real data (from the actual dataset) and fake data (produced by the Generator). Its goal is to correctly distinguish between the two, labeling them as either "Real" or "Fake." During training, the Generator tries to maximize the probability of the Discriminator making a mistake, while the Discriminator tries to minimize that probability. This adversarial process forces the Generator to create increasingly realistic outputs over time to deceive its opponent.