

# **Transformers Interview**

## **Questions and Answers**

### **1.What is a Transformer, and why was it introduced?**

**Answer:**

A **Transformer** is a deep learning model introduced in the paper “*Attention is All You Need*” (*Vaswani et al., 2017*). It was designed to overcome the limitations of RNNs and LSTMs, particularly in handling long-range dependencies, by using the **self-attention mechanism**.

**Key Reasons for its Introduction:**

- **Parallelization:** Unlike RNNs, Transformers process input sequences in parallel, making training more efficient.
- **Handling Long-Term Dependencies:** Self-attention allows the model to capture relationships between distant words in a sequence.
- **Better Scalability:** Transformers scale effectively to larger datasets and bigger models (e.g., BERT, GPT).

## 2. Explain the self-attention mechanism in Transformers.

### Answer:

The **self-attention mechanism** allows each token in a sequence to **weigh the importance** of other tokens before generating an output.

- Attention scores are calculated as:

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

- This determines how much each word contributes to the final representation.

## Steps in Self-Attention:

### 1. Compute Queries (Q), Keys (K), and Values (V):

- Each token is transformed into three vectors: **Q** (query), **K** (key), and **V** (value).
- **Compute Attention Scores:**
- Attention scores are calculated as:

## Apply Softmax and Weight Values:

- Softmax normalizes scores to highlight important tokens.

## Weighted Sum of Values:

- The output representation is generated based on the weighted sum.

### **Why It's Important:**

- Allows for **context-aware embeddings** where words in different contexts get different representations.
- Enables **parallel processing** of sequences.

## **3. What is Multi-Head Attention, and why is it used?**

### **Answer:**

Multi-Head Attention (MHA) extends the **self-attention mechanism** by using multiple attention heads.

### **How it Works:**

- Instead of a **single** set of Q, K, V matrices, MHA **splits the input into multiple heads.**

- Each attention head learns a different aspect of **relationships in the sequence.**
- The outputs of all attention heads are **concatenated and projected** into a final output.

### **Benefits:**

- **Enhances model capacity:** Different attention heads focus on different parts of the sequence.
- **Improves performance** in tasks like **translation and text generation.**

## **4. How does a Transformer handle positional information in sequences?**

### **Answer:**

Unlike RNNs, **Transformers lack inherent sequential order.** To address this, they use **Positional Encoding.**

### **Key Concept:**

- Positional encodings are **added to input embeddings** before processing by the attention mechanism.

- Encodings are computed using **sine and cosine functions**:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

- This helps encode information about **word positions** in a way that is **continuous and differentiable**.

### Why It's Needed:

- Ensures the model understands word **ordering**.
- Helps **capture sequence relationships** in tasks like machine translation.

## 5. What is Layer Normalization, and why is it used in Transformers?

### Answer:

**Layer Normalization (LayerNorm)** is used to **stabilize training** and **speed up convergence**.

## How It Works:

Normalizes activations **across the feature dimension**,

rather than across batch size (like BatchNorm).

## Benefits:

- **Improves training stability.**
- **Reduces internal covariate shift.**
- **Allows deeper architectures** without vanishing/exploding gradients.

## 6. What are the advantages of using Transformers over RNNs and LSTMs?

### Answer:

Feature	Transformers	RNNs/LSTMs
Parallelization	Yes, processes entire sequences at once.	No, sequential processing (slower).
Handling Long Dependencies	Excellent due to self-attention.	Struggles with long sequences (vanishing gradients).
Scalability	Highly scalable (used in GPT, BERT).	Limited scalability due to sequential updates.
Computation Efficiency	More efficient with GPUs/TPUs.	Computationally expensive.

## 7. What is the difference between Encoder and Decoder in Transformers?

**Answer:**

Component	Encoder	Decoder
Role	Processes input data and generates a hidden representation.	Generates output sequence based on the encoder's representation.
Attention Mechanisms	Uses <b>self-attention</b> .	Uses <b>masked self-attention and cross-attention</b> (to encoder outputs).
Application	Used in <b>BERT, T5, etc.</b>	Used in <b>GPT, text generation tasks.</b>

Transformers are composed of two main parts: **Encoders** and **Decoders**.

## 8. Explain the role of Masked Self-Attention in Decoders.

**Answer:**

**Masked Self-Attention** is used in Transformer decoders to ensure that **future tokens are not visible** during training.

**How it Works:**

- Applies a **masking operation** to prevent attention to future tokens.
- The attention scores for future tokens are set to **negative infinity** before softmax.
- Ensures the decoder **predicts one token at a time** during training.

### **Why It's Important:**

- Prevents **data leakage** during autoregressive generation.
- Allows the model to **learn sequential dependencies**.

## **9. What are some real-world applications of Transformers?**

### **Answer:**

Transformers are widely used in various domains:

## **Natural Language Processing (NLP):**

- **Machine Translation** (e.g., Google Translate).
- **Text Summarization** (e.g., Pegasus, T5).
- **Chatbots and Conversational AI** (e.g., ChatGPT, Bard).

## **Computer Vision (CV):**

- **Vision Transformers (ViTs)** replace CNNs in image processing.
- Used in **object detection** and **image segmentation**.

## **Multimodal AI:**

- **CLIP** and **DALL·E** combine text and image understanding.

## **Biological & Healthcare Applications:**

- **AlphaFold** (by DeepMind) for **protein structure prediction.**

## **10. What are the challenges of using Transformers?**

### **Answer:**

Despite their advantages, Transformers have some challenges:

- **High Computational Cost:** Requires significant memory and processing power.
- **Training Complexity:** Needs large datasets and powerful hardware (GPUs/TPUs).
- **Inference Latency:** Large models can be slow for real-time applications.
- **Data Hunger:** Requires massive datasets to generalize well.

To address these, models like **Mixture of Experts (MoE)** and **Efficient Transformers (e.g., Linformer, Performer)** have been developed.

## 11. How does the Transformer's computational complexity compare to RNNs and CNNs?

**Answer:**

### Complexity Analysis:

Model	Complexity per Layer	Why?
RNNs	$O(n)$	Sequential processing, cannot parallelize.
CNNs	$O(k \cdot n)$ (for kernel size $k$ )	Captures local dependencies, limited global context.
Transformers	$O(n^2 \cdot d)$ (due to self-attention)	Computes attention for all word pairs simultaneously.

### Optimizations in Transformers:

- **Linformer:** Reduces complexity from  $O(n^2)$  to  $O(n)$ .
- **Performer:** Uses random projections for linear-time attention.

The computational complexity of Transformers is different from **RNNs** and **CNNs** due to the **self-attention mechanism**.

### Complexity Analysis:

## 12. What are some efficient Transformer variants that reduce computational cost?

## Answer:

Variant	Optimization	Complexity Reduction
Linformer	Low-rank approximation of attention matrix.	$O(n)$
Performer	Kernel-based attention instead of softmax.	$O(n)$
Reformer	Uses LSH (Locality Sensitive Hashing) to reduce attention cost.	$O(n \log n)$
Longformer	Uses sparse attention to focus on local context.	$O(n)$
BigBird	Hybrid of global and local attention.	$O(n)$

### 💡 Why this matters?

Efficient Transformers are **critical** for handling **long documents, DNA sequences, or video processing**.

Several Transformer variants have been developed to **reduce computational complexity**:

## 13. How does the Transformer decoder differ from the encoder?

## Answer:

Component	Encoder	Decoder
Self-Attention	Unmasked (can attend to all tokens).	Masked (cannot attend to future tokens).
Cross-Attention	🚫 No cross-attention (operates independently).	✅ Cross-attention (attends to encoder outputs).
Use Case	Processes input text/images.	Generates text sequences.
Examples	BERT (encoder-only).	GPT (decoder-only), T5 (encoder-decoder).

### 💡 Why this matters?

- **Encoder models** (e.g., BERT) are best for **classification** and **embedding tasks**.
- **Decoder models** (e.g., GPT) are best for **text generation**.

While both encoder and decoder have similar architectures, **key differences** exist:

## **14. Explain the concept of Cross-Attention in Transformer decoders.**

**Answer:**

**Cross-attention** allows the decoder to focus on **relevant parts of the encoder's output**.

**How it works:**

- The **Query (Q)** comes from the **decoder**, while **Keys (K) and Values (V)** come from the **encoder**.
- This mechanism **links the encoder and decoder**, allowing the model to use information from the input while generating output.

**Why It's Important?**

- Ensures **better alignment between input and generated output.**
- Critical in **translation models** like **T5 and BART.**

## **15. What are some real-world applications of Transformers outside of NLP?**

### **Answer:**

While Transformers are **dominant in NLP**, they have extended into **other domains**:

### **Computer Vision:**

- **Vision Transformers (ViT):** Replaces CNNs for image classification.
- **DEtection TRansformer (DETR):** Used for object detection.

### **Speech Processing:**

- **Wav2Vec 2.0:** Self-supervised learning for speech recognition.
- **Whisper (OpenAI):** Multi-lingual ASR system.

## **Bioinformatics & Healthcare:**

- **AlphaFold:** Protein structure prediction using attention mechanisms.
- **DNABERT:** Uses BERT for DNA sequence analysis.

### **Why this matters?**

Transformers are shaping **next-gen AI models** across multiple industries.

## **16. What is the role of Feedforward Networks (FFN) in Transformers?**

**Answer:**

- Formula:

$$FFN(x) = \text{ReLU}(W_1x + b_1)W_2 + b_2$$

- Allows **non-linearity** and **feature transformation**.

### **Why It's Important?**

- Adds **capacity** to the model beyond attention mechanisms.
- Ensures **better feature representation** for each token.

The **Feedforward Network (FFN)** is applied **individually** to each token **after attention computation**.

### **Structure:**

- Typically **two dense layers** with an activation function in between.

## **17. How does transfer learning work in Transformers?**

### **Answer:**

Transfer learning in Transformers involves **pretraining on a large dataset** followed by **fine-tuning on a specific task**.

## **Steps:**

### **Pretraining Phase:**

- Models like BERT, GPT, T5 are trained on **massive datasets**.
- Uses **self-supervised tasks** (e.g., masked language modeling, next token prediction).

### **Fine-Tuning Phase:**

- The pretrained model is adapted to **downstream tasks**.
- Requires **less data and computational resources** compared to training from scratch.

## **Why It's Useful?**

- **Generalizes well across domains.**
- Reduces the **need for large task-specific datasets.**

## 18. What are common techniques used to improve Transformer training?

Technique	How It Helps
Layer Normalization	Stabilizes training and prevents exploding gradients.
Gradient Clipping	Prevents large updates that can destabilize learning.
Mixed Precision Training	Uses FP16 to reduce memory consumption.
ZeRO (Zero Redundancy Optimizer)	Distributes model parameters across GPUs to scale large models.
Checkpointing	Saves intermediate results to optimize memory.

### 💡 Why this matters?

These techniques allow **training massive models** like **GPT-4** and **Gemini** efficiently.

### Answer:

Training large Transformer models is challenging. Some key improvements include:

## 19. How does the Transformer model compare to MoE (Mixture of Experts)?

### Answer:

Feature	Transformers	MoE
Computation	Fully activates all layers.	Activates only a few experts per input.
Scalability	Becomes costly at scale.	More efficient for large-scale models.
Example Models	BERT, GPT, T5.	Switch Transformers, GLaM.

💡 **Why this matters?**

Hybrid architectures like **Switch Transformers** are now combining the **best of both worlds**.

Transformers **process all tokens** with **dense computation**, while **Mixture of Experts (MoE)** activates only a subset of **experts**.

## 20. What future improvements can we expect in Transformer models?

**Answer:**

Researchers are working on **more efficient and powerful**

Transformer architectures:

1. **Sparse Transformers:** Reduce quadratic complexity.
2. **Hybrid Architectures:** Combining MoE with Transformers.

**3. Neuromorphic AI:** Adapting Transformers for  
**low-power applications.**

**4. Smaller, Efficient Models:** Reducing **memory and  
inference cost.**

### **Why this matters?**

The future of Transformers is **leaner, faster, and more  
scalable** across various AI domains.