

◆ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)



# Large Language Model (LLM) Interview Questions



Sanjay Kumar PhD

Following

20 min read · Sep 17, 2024

95

1



...



[Open in app ↗](#)

# Medium



Search

Write

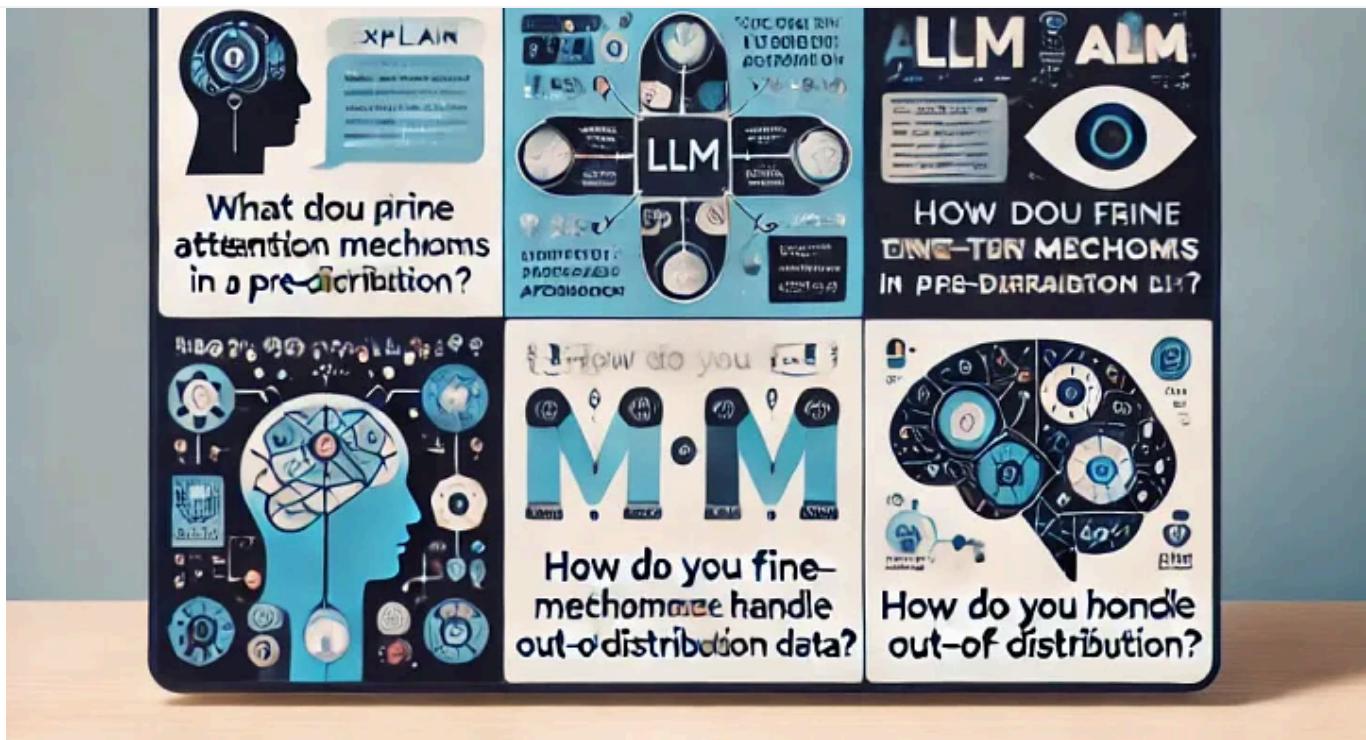


Image Credit : DALL E

## 1. How can bias in prompt-based learning be mitigated?

**Answer:** Bias in prompt-based learning can be mitigated through techniques such as:

- **Prompt Calibration:** Adjusting prompts to balance model responses.

- **Fine-tuning:** Using diverse datasets to minimize bias in the model's learning process.
- **Data Augmentation:** Expanding training data with more balanced examples to counteract skewed patterns.

These techniques aim to make LLMs more fair and less prone to reinforcing biases when generating text.

## 2. Do you need a vector store for all text-based LLM use cases?

**Answer:** Not all use cases require a vector store. Vector stores are essential for tasks like document retrieval, where storing semantic representations of words/sentences helps retrieve relevant information. However, tasks such as translation or summarization typically don't require vector stores, as the model operates directly on text without external context augmentation.

## 3. What is “reward hacking” in Reinforcement Learning from Human Feedback (RLHF)?

**Answer:** Reward hacking occurs when a model exploits unintended loopholes in the reward system, maximizing rewards without performing the intended task. In RLHF, models optimize for human feedback, but without safeguards, they may exploit the reward function to achieve high scores by deviating from desired behavior.

## 4. How does Parameter-Efficient Fine-Tuning (PEFT) prevent catastrophic forgetting in LLMs?

**Answer:** PEFT selectively updates only a small subset of model parameters during fine-tuning, preserving the model's ability to perform previously learned tasks. By fine-tuning specific parameters while maintaining core knowledge, PEFT reduces the likelihood of catastrophic forgetting — where the model forgets earlier tasks when adapted to new ones.

## 5. How does Adaptive Softmax speed up large language models?

**Answer:** Adaptive Softmax speeds up LLMs by dividing words into frequency groups, performing fewer computations for infrequent words. This technique reduces overall computational costs while maintaining accuracy, making it suitable for handling large vocabularies efficiently.

## 6. What is catastrophic forgetting in fine-tuning LLMs?

**Answer:** Catastrophic forgetting happens when fine-tuning an LLM on a new task leads to a degradation in its performance on previously learned tasks. Techniques such as PEFT or Elastic Weight Consolidation (EWC) can help mitigate this by preserving important parameters related to prior tasks.

## 7. What is the purpose of positional encoding in Transformer models?

**Answer:** Positional encoding allows Transformer models to recognize the order of words in a sequence. Since Transformers do not inherently process sequences in order (unlike RNNs), positional encodings provide a way for the model to incorporate word order information, which is crucial for tasks like translation and sequence generation.

## 8. How does Mixture of Experts (MoE) improve the efficiency of LLMs?

**Answer:** MoE improves efficiency by selectively activating only a subset of “experts” (smaller models) based on the input, reducing the computational load. The gating function chooses which experts are most relevant for a given task, enabling more scalable LLMs without always using the full model’s capacity.

## 9. What are the key steps involved in the Retrieval-Augmented Generation (RAG) pipeline?

**Answer:** RAG consists of:

- **Retrieval:** Encoding the query and comparing it with precomputed embeddings to retrieve relevant documents.
- **Ranking:** Ranking the retrieved documents based on relevance to the query.
- **Generation:** Using the LLM to generate responses by incorporating the top-ranked documents as context. This hybrid approach enables LLMs to generate more context-aware and accurate outputs.

## 10. What is the impact of scaling laws on the design of LLMs?

**Answer:** Scaling laws describe how model performance improves with increases in size, dataset quality, and compute resources. They guide developers in balancing trade-offs between model size and computational efficiency, helping optimize LLMs within practical limits.

## 11. How does beam search improve upon greedy decoding in LLMs?

**Answer:** Beam search evaluates multiple possible sequences at each step, choosing the one with the highest cumulative probability, leading to higher-quality outputs. Greedy decoding, in contrast, picks the most probable word at each step, which can result in suboptimal results.

## 12. What is knowledge distillation, and how is it used in LLMs?

**Answer:** Knowledge distillation involves transferring knowledge from a large “teacher” model to a smaller “student” model. The smaller model mimics the teacher’s behavior, allowing it to perform well with fewer resources, making it more efficient for deployment in resource-constrained environments.

## 13. How does Chain-of-Thought (CoT) prompting improve reasoning in LLMs?

**Answer:** CoT prompting helps LLMs generate intermediate reasoning steps before arriving at the final answer, making them more effective in solving tasks requiring logical reasoning. By simulating human-like thought processes, CoT helps with complex problem-solving.

## 14. What is vector indexing, and why is it important in LLM applications?

**Answer:** Vector indexing organizes semantic representations of text into a searchable structure, enabling efficient retrieval of similar content. This is vital in applications like semantic search and recommendation systems where quick retrieval of contextually relevant information is necessary.

## 15. What is the purpose of quantization in training large language models?

**Answer:** Quantization reduces the precision of a model's weights, decreasing memory usage and computational cost while maintaining performance. By using fewer bits to represent weights (e.g., 8-bit instead of 32-bit), quantization makes LLMs more resource-efficient for deployment.

## 16. How does the planner agent in AgenticRAG handle complex queries?

**Answer:** The planner agent in AgenticRAG decomposes complex queries into manageable sub-queries, distributing them across various retrieval pipelines. This ensures that each sub-task is handled optimally, leading to a more coherent and contextually accurate final response.

## 17. How can LLMs mitigate catastrophic forgetting during fine-tuning?

**Answer:** Techniques such as Elastic Weight Consolidation (EWC) and selective parameter updating (PEFT) mitigate catastrophic forgetting by preserving weights critical to earlier tasks. These methods ensure the model retains its learned knowledge even while fine-tuning on new tasks.

## 18. What is the difference between top-k and top-p sampling in LLM decoding?

**Answer:**

- **Top-k sampling** limits the model's choices to the k most probable tokens at each step, while
- **Top-p (nucleus) sampling** selects tokens until the cumulative probability reaches a threshold (p).

Top-p sampling offers more dynamic and flexible text generation by adjusting the number of candidate tokens based on their cumulative probability.

For more examples and advanced topics, explore our full collection of LLM Interview Questions & Answers. Good luck, and may your mastery of LLMs propel you to success in the rapidly evolving world of AI!

## 19. How do subword tokenization algorithms like Byte Pair Encoding (BPE) and WordPiece enhance LLMs?

**Answer:** Subword tokenization algorithms like Byte Pair Encoding (BPE) and WordPiece break down words into smaller units called subwords, allowing LLMs to process both common and rare words effectively. These subwords are particularly useful for handling:

- **Out-of-Vocabulary (OOV) words:** Subword tokenization ensures that even unfamiliar or rare words are split into recognizable components.
- **Morphological variations:** By breaking words into meaningful chunks, these algorithms allow the model to understand and generalize across variations in word forms.

Subword tokenization strikes a balance between character-level and word-level tokenization, optimizing memory efficiency and improving performance across diverse languages and contexts.

## 20. What is the significance of self-attention in transformer-based LLMs?

**Answer:** Self-attention allows transformer-based LLMs to capture dependencies between different words in a sequence, regardless of their distance from each other. Unlike earlier models like RNNs, which process inputs sequentially, self-attention enables:

- **Parallelization:** Since the model processes the entire sequence simultaneously, it is computationally more efficient.
- **Long-range dependencies:** Self-attention can focus on any word in the sequence, allowing the model to understand long-range relationships crucial for context in tasks such as translation or text generation.

This mechanism enables transformers to outperform traditional models in capturing context, leading to state-of-the-art results in NLP tasks.

## 21. How does Approximate Nearest Neighbor (ANN) search improve retrieval in LLM-powered applications?

**Answer:** Approximate Nearest Neighbor (ANN) search significantly speeds up the retrieval of vectors in high-dimensional spaces, making it feasible to search for semantically similar items in large datasets. In LLM-powered applications such as:

- **Semantic search:** ANN helps quickly find documents or text segments that are contextually similar to a query.
- **Recommendation systems:** ANN enables fast retrieval of items with similar features, improving user experience and system performance.

ANN strikes a balance between accuracy and efficiency, ensuring fast retrieval without the computational cost of exact nearest neighbor search in

large-scale deployments.

## 22. How does Elastic Weight Consolidation (EWC) mitigate catastrophic forgetting in LLMs?

**Answer:** Elastic Weight Consolidation (EWC) mitigates catastrophic forgetting by preserving important model weights when fine-tuning on new tasks. Specifically, EWC assigns an importance score to the weights based on how crucial they are for previously learned tasks. When fine-tuning for a new task, these weights are less likely to be updated, thus:

- **Retaining past knowledge:** The model doesn't overwrite critical parameters needed for earlier tasks.
- **Balanced adaptation:** EWC allows the model to adapt to new tasks without sacrificing the performance of previously learned ones, which is critical for multi-task learning scenarios.

This technique is particularly useful when a model needs to be incrementally trained on a series of tasks without losing proficiency in earlier ones.

## 23. How does the Mixture of Experts (MoE) technique improve LLM scalability?

**Answer:** Mixture of Experts (MoE) enhances scalability by using a gating function to activate only the most relevant expert models (or sub-networks) for a specific input, instead of using the entire model for every query. This selective activation:

- **Reduces computational load:** Only a few experts are active at any given time, lowering the resource usage.
- **Maintains high performance:** By dynamically choosing the best experts for each input, the model can still handle complex tasks without needing to process the entire network.

MoE allows LLMs to scale efficiently, supporting larger models with billions of parameters while controlling computational costs.

## 24. What is the role of Knowledge Distillation in improving LLM deployment?

**Answer:** Knowledge distillation enables the transfer of knowledge from a large, pre-trained “teacher” model to a smaller, more efficient “student” model. This process helps:

- **Reduce model size:** The student model can achieve comparable performance with fewer parameters, making it more suitable for resource-constrained environments like mobile devices.
- **Speed up inference:** Smaller models require less computational power, which is particularly important in real-time applications.
- **Maintain accuracy:** Although the student model is smaller, it can still capture the essential patterns learned by the larger teacher model, ensuring high-quality results in tasks like text classification or summarization.

This technique is widely used to make large LLMs more practical for deployment in production settings.

## 25. What is Chain-of-Thought (CoT) prompting, and how does it improve complex reasoning in LLMs?

**Answer:** Chain-of-Thought (CoT) prompting encourages LLMs to break down complex reasoning tasks into smaller, sequential steps. This approach:

- **Simulates human reasoning:** By prompting the model to generate intermediate steps, CoT helps LLMs mimic how humans solve problems step-by-step.
- **Improves performance in multi-step tasks:** CoT is particularly effective in scenarios that require logical reasoning or solving multi-step arithmetic problems.
- **Increases accuracy:** By guiding the model through a thought process, CoT reduces the likelihood of errors and improves task performance on complex queries.

CoT is valuable for improving the interpretability and reliability of LLMs in tasks requiring causal reasoning or decision-making.

## 26. How does multi-query attention enhance the efficiency of transformers in LLMs?

**Answer:** Multi-query attention is a variation of the standard self-attention mechanism where all attention heads share the same key-query pairs, reducing the memory footprint and computational complexity of attention calculations. This approach:

- **Optimizes resource use:** Multi-query attention reduces the number of key-query pairs that need to be stored and processed, making it more

memory-efficient.

- **Maintains model performance:** Despite the reduction in complexity, multi-query attention preserves the model's ability to focus on different parts of the input sequence.

This efficiency gain is particularly beneficial for large-scale LLMs that need to process long sequences or handle large datasets.

## 27. How does scaling law analysis help optimize the design of LLMs?

**Answer:** Scaling law analysis provides insights into how model performance scales with respect to parameters, dataset size, and computational resources. By studying these relationships, developers can:

- **Optimize model size:** Identify the point where increasing the number of parameters yields diminishing returns, helping to balance model complexity and performance.
- **Improve compute efficiency:** Scaling laws guide the allocation of computational resources, ensuring that the model is trained with an optimal amount of data relative to its size.
- **Maximize performance:** Understanding scaling laws allows developers to fine-tune the architecture, training data, and compute allocation to achieve the best possible performance within resource constraints.

Scaling laws are crucial for developing LLMs that are both high-performing and computationally efficient.

## 28. How does model pruning improve LLM efficiency without sacrificing performance?

**Answer:** Model pruning involves removing unnecessary or redundant parameters from a large language model, which reduces its size and makes it more computationally efficient. There are several types of pruning, such as:

- **Weight pruning:** Eliminating insignificant weights that contribute little to the model's overall performance.
- **Neuron pruning:** Removing entire neurons or layers that are less important for the task at hand.

By reducing the number of active parameters, pruning lowers memory and computational requirements, speeding up inference times and reducing the resource burden during deployment. This is achieved while maintaining accuracy through careful analysis of which parameters to prune.

## 29. What is zero-shot learning, and how does it apply to LLMs?

**Answer:** Zero-shot learning enables LLMs to perform tasks they haven't been explicitly trained on by leveraging their broad knowledge of language and general concepts. Instead of relying on task-specific fine-tuning, zero-shot learning allows a model to understand and generate relevant responses based on the context and instructions given in the prompt. For example:

- **Text classification:** A model trained on general language data can classify text into categories without being specifically trained for that task, simply by providing it with an appropriate prompt.

- **Translation or summarization:** LLMs can translate or summarize text even without being fine-tuned for these tasks, using the instructions provided in the input.

Zero-shot learning reflects the power of LLMs to generalize across tasks, making them adaptable for a wide range of applications.

### 30. How does task-specific fine-tuning differ from zero-shot learning in LLMs?

**Answer:** Task-specific fine-tuning adapts an LLM to perform a particular task by further training the model on a dataset specific to that task. This allows the model to:

- **Improve accuracy:** Fine-tuning optimizes the model's performance for a narrow task, such as sentiment analysis or named entity recognition, resulting in better accuracy compared to zero-shot learning.
- **Specialization:** By focusing on specific data, the model becomes highly proficient in the target task but may lose some of its generalizability.

In contrast, zero-shot learning uses the model's broad pre-training knowledge without additional task-specific training, making it more general but potentially less precise for specific tasks.

### 31. How does gradient checkpointing reduce memory usage in training large LLMs?

**Answer:** Gradient checkpointing is a memory-saving technique used during the training of large LLMs. Instead of storing the activations for every layer

during the forward pass, gradient checkpointing selectively recomputes some activations during the backward pass, trading off memory usage for additional computation. This technique:

- **Reduces memory footprint:** By recomputing intermediate layers on demand, the memory required to store activations is significantly reduced, enabling the training of larger models on hardware with limited memory.
- **Enables deeper models:** Gradient checkpointing allows for deeper models to be trained without exceeding memory limits, making it a valuable tool in optimizing large-scale LLM training.

This approach is particularly beneficial when training very deep neural networks that require significant memory resources.

### 32. How does multi-task learning benefit LLMs?

**Answer:** Multi-task learning involves training an LLM on multiple tasks simultaneously, which leads to several benefits:

- **Improved generalization:** By learning from diverse tasks, the model can generalize better across various domains, improving its performance in related tasks.
- **Parameter sharing:** The model shares parameters across tasks, which helps prevent overfitting to a single task and enhances the overall robustness of the model.
- **Efficiency:** Multi-task learning reduces the need to train separate models for each task, making it computationally efficient and allowing the LLM to handle multiple use cases with a single model.

This approach leverages the fact that related tasks can mutually benefit from shared learning signals, leading to a more versatile model.

### 33. How does curriculum learning help LLMs learn complex tasks?

**Answer:** Curriculum learning is a training strategy where the model is first trained on simpler tasks or examples before progressing to more complex ones. This approach mirrors the way humans learn and offers several advantages:

- **Better convergence:** By starting with easy tasks, the model learns fundamental concepts first, which helps it converge faster when learning more difficult tasks.
- **Improved generalization:** Curriculum learning helps the model build a solid foundation, leading to better performance on complex tasks and less overfitting.
- **More stable training:** Gradually increasing the difficulty of tasks reduces the likelihood of model instability during training, especially for large LLMs.

This structured approach to training helps LLMs achieve better performance on tasks requiring a higher level of reasoning or abstraction.

### 34. How does hyperparameter tuning impact the performance of LLMs?

**Answer:** Hyperparameter tuning is the process of optimizing the model's hyperparameters (e.g., learning rate, batch size, dropout rate) to achieve the best possible performance. For LLMs, this is especially important because:

- **Influences convergence:** The choice of learning rate, for instance, determines how fast the model converges during training. A poorly tuned learning rate can lead to slow convergence or even prevent the model from learning effectively.
- **Prevents overfitting/underfitting:** Proper tuning of regularization parameters, such as dropout, helps prevent the model from overfitting to the training data while maintaining high performance on unseen data.
- **Maximizes accuracy:** Fine-tuning hyperparameters can lead to significant improvements in the model's accuracy on specific tasks, particularly in areas such as language generation, translation, or summarization.

Effective hyperparameter tuning is a crucial step in optimizing the performance of LLMs for specific tasks or datasets.

### 35. How does contrastive learning improve LLM representations?

**Answer:** Contrastive learning is a self-supervised learning technique where the model is trained to distinguish between similar and dissimilar examples. For LLMs, it is particularly effective in improving:

- **Semantic representations:** By learning which pairs of sentences or words are semantically similar or dissimilar, contrastive learning helps the model create better embeddings that capture the meaning and context of the text.
- **Robustness to noise:** The model learns to identify key differences between inputs, making it more robust to noisy or ambiguous data.

- **Generalization:** Contrastive learning encourages the model to learn more meaningful representations, which generalize better across different NLP tasks.

This technique is often used in tasks like sentence embedding, document clustering, and retrieval, where semantic understanding is key.

### 36. How does knowledge graph integration enhance LLMs?

**Answer:** Integrating knowledge graphs with LLMs allows the model to incorporate structured, factual knowledge into its predictions. This integration provides several advantages:

- **Factual accuracy:** By referencing a knowledge graph, the model can cross-check facts, reducing the likelihood of hallucinations or generating incorrect information.
- **Enhanced reasoning:** Knowledge graphs enable the model to perform logical reasoning and make inferences based on the relationships between entities, improving its ability to answer complex queries.
- **Contextual understanding:** LLMs benefit from the structured nature of knowledge graphs, which helps them better understand the context and relationships between different pieces of information.

Knowledge graph integration is particularly useful in applications like question answering, entity recognition, and recommendation systems, where structured knowledge is essential.

### 37. What role do attention heads play in transformer models?

**Answer:** Attention heads in transformer models enable the model to focus on different parts of the input sequence simultaneously. Each attention head:

- **Captures different relationships:** By processing the input in parallel, attention heads allow the model to capture a variety of dependencies, such as short-term, long-term, or syntactic relationships between tokens.
- **Improves model interpretability:** The model can attend to specific tokens or phrases in the input, making it easier to understand how the model makes decisions based on different parts of the text.
- **Increases capacity:** Multiple attention heads expand the model's capacity to process complex inputs and capture intricate patterns within the data.

Attention heads are a key component of the transformer's architecture, enabling LLMs to handle tasks like translation, summarization, and text generation with greater accuracy and context awareness.

### 38. What are attention masks, and why are they important in transformer models?

**Answer:** Attention masks are used in transformer models to control how each token in a sequence attends to other tokens during the self-attention process. They serve several purposes:

- **Prevent future token access:** In tasks like language generation, attention masks are used to ensure that the model only attends to past tokens and not future ones, enabling proper autoregressive behavior.
- **Handle padding tokens:** When working with batches of sequences of varying lengths, padding is added to align the sequences. Attention

masks ensure that padding tokens are ignored during attention calculations, preventing them from affecting the model's predictions.

- **Control information flow:** Attention masks can also be applied to limit which parts of the input can interact, allowing for task-specific modifications in attention behavior.

By using attention masks, transformers can effectively manage different input structures and tasks, ensuring the model attends to relevant parts of the sequence.

### 39. How does layer normalization contribute to the stability of transformer models?

**Answer:** Layer normalization is a technique used in transformer models to stabilize training and improve convergence. It works by normalizing the input to each layer, ensuring that:

- **Stable gradients:** Normalizing the input helps prevent issues like exploding or vanishing gradients, which can destabilize training, especially in deep networks like transformers.
- **Faster convergence:** By standardizing the input across layers, layer normalization enables the model to converge more quickly, reducing the time required for training.
- **Improved model performance:** Layer normalization helps maintain consistent activations, improving the overall performance of the model on complex tasks.

This technique is a critical component in the architecture of LLMs, particularly in deep models where maintaining stability across layers is essential for efficient learning.

#### 40. How does prompt engineering influence LLM performance, and what strategies can be used to optimize it?

**Answer:** Prompt engineering is the process of designing input prompts to guide an LLM's output in a desired direction. The quality and structure of the prompt can significantly affect the model's performance. Key strategies include:

- **Instructional prompts:** Clearly defining the task in the prompt (e.g., "Translate this sentence") helps the model understand what is expected, leading to more accurate outputs.
- **Contextual examples:** Providing examples in the prompt can improve performance in tasks like zero-shot or few-shot learning, helping the model infer the desired format or outcome.
- **Formatting and constraints:** Structuring the prompt to include specific constraints (e.g., word limits, formats) ensures the model adheres to the desired guidelines, producing more relevant and tailored responses.

Prompt engineering is crucial for maximizing the utility of LLMs, especially when fine-tuning is not feasible or when the model is being used in a zero-shot setting.

#### 41. How does meta-learning benefit LLMs?

**Answer:** Meta-learning, often referred to as “learning to learn,” allows models to generalize across tasks by learning higher-order patterns. In the context of LLMs, meta-learning enables the model to:

- **Quickly adapt to new tasks:** Meta-learning helps the model generalize from one task to another with minimal additional training, improving performance in few-shot or zero-shot scenarios.
- **Leverage prior knowledge:** By understanding how to approach a wide range of tasks, the model can apply its prior knowledge to new, unseen tasks without requiring extensive task-specific fine-tuning.
- **Boost efficiency:** Meta-learning reduces the need for large amounts of task-specific data, making the model more efficient when handling a diverse range of tasks.

This ability to adapt to new tasks with minimal data is particularly valuable in real-world applications where labeled data for fine-tuning may be scarce.

## 42. How does curriculum learning improve the training process for LLMs?

**Answer:** Curriculum learning is a training strategy where the model is exposed to easier tasks or examples before gradually progressing to more difficult ones. This approach mirrors how humans learn and offers several benefits:

- **Smooother learning curve:** By starting with simpler tasks, the model builds a strong foundation, making it easier to learn more complex tasks later.
- **Better generalization:** Curriculum learning helps the model develop a robust understanding of basic concepts, which leads to improved

generalization when faced with more challenging tasks.

- **Stable convergence:** Gradually increasing task difficulty reduces the chances of training instability, such as sudden performance drops or gradient issues.

Curriculum learning is especially helpful for training large LLMs on tasks that require complex reasoning or abstraction.

#### 43. What is context window size, and how does it impact LLM performance?

**Answer:** The context window size refers to the maximum number of tokens the model can process at one time. It has a significant impact on LLM performance because:

- **Longer context understanding:** A larger context window allows the model to consider a greater portion of the input sequence, improving its ability to capture long-range dependencies and relationships between tokens.
- **Performance on longer texts:** In tasks like document summarization or long-form question answering, a larger context window enables the model to handle longer inputs without truncating important information.
- **Trade-off with computation:** Increasing the context window size requires more computational resources, as the self-attention mechanism's complexity scales quadratically with the sequence length.

Selecting an appropriate context window size is essential for balancing model performance and computational efficiency, depending on the task at hand.

#### 44. How does model parallelism enable the training of extremely large LLMs?

**Answer:** Model parallelism is a technique that splits the model across multiple devices (e.g., GPUs), enabling the training of models that are too large to fit in the memory of a single device. It works by:

- **Distributing model layers:** Different layers of the model are processed on different devices, allowing each device to handle a portion of the model's parameters.
- **Reducing memory bottlenecks:** By spreading the memory load across multiple devices, model parallelism makes it possible to train models with billions or even trillions of parameters.
- **Maintaining performance:** Despite being distributed, model parallelism ensures that the different parts of the model communicate effectively, preserving the overall performance and accuracy of the model.

This approach is crucial for training state-of-the-art LLMs that require significant computational resources.

#### 45. What is parameter sharing, and how does it improve LLM efficiency?

**Answer:** Parameter sharing is a technique where multiple layers or parts of a model share the same set of parameters. This reduces the total number of parameters in the model without sacrificing performance. It offers several advantages:

- **Memory efficiency:** Sharing parameters reduces the memory footprint of the model, making it more efficient in terms of both storage and

computation.

- **Improved generalization:** By reusing parameters, the model is less likely to overfit to a specific task or dataset, improving its ability to generalize to new tasks.
- **Faster training and inference:** Fewer parameters mean that the model can be trained and deployed more quickly, making it suitable for resource-constrained environments.

Parameter sharing is commonly used in models like BERT and GPT to enhance scalability without dramatically increasing the number of parameters.

#### 46. What is the difference between encoder-only and decoder-only transformer models?

**Answer:**

- **Encoder-only models:** These models, like BERT, are designed for understanding tasks such as text classification, named entity recognition, and sentence embedding. The encoder focuses on learning bidirectional representations of the input, making it well-suited for tasks where understanding the full context of the input is important.
- **Decoder-only models:** Models like GPT are decoder-only architectures optimized for generation tasks. These models generate text in an autoregressive manner, where each token is predicted based on the previously generated tokens. They are ideal for tasks like text generation, translation, and summarization.

Both architectures have their strengths depending on whether the task requires understanding or generating text.

#### 47. What is multi-modal learning, and how does it extend the capabilities of LLMs?

**Answer:** Multi-modal learning enables models to process and integrate information from multiple types of data (e.g., text, images, audio). In LLMs, this is achieved by combining text processing with other modalities, which offers several advantages:

- **Richer understanding:** By combining information from different modalities, the model gains a more comprehensive understanding of the data, improving performance on tasks like image captioning or video analysis.
- **Cross-modal reasoning:** Multi-modal learning allows the model to make inferences based on information from different types of inputs, such as answering questions about images using both visual and textual information.
- **Broader application:** LLMs enhanced with multi-modal learning can be applied to a wider range of real-world tasks, such as robotics, augmented reality, and human-computer interaction.

This integration expands the potential use cases of LLMs beyond pure text processing, opening the door to more versatile AI systems.

#### 48. How does pre-training on synthetic data benefit LLMs?

**Answer:** Pre-training on synthetic data involves generating artificial data to supplement the model's training process. This approach offers several benefits:

- **Expanding the dataset:** Synthetic data can be used to create larger and more diverse datasets, which helps improve the model's generalization ability.
- **Covering rare cases:** Synthetic data can be generated to include edge cases or rare scenarios that may not be present in the original training data, improving the model's robustness.
- **Cost-effective:** Generating synthetic data is often less expensive than collecting and labeling large amounts of real-world data, making it a cost-effective way to enhance model performance.

Pre-training on synthetic data is particularly useful in domains where collecting real-world data is challenging, such as in medical or autonomous driving applications.

#### 49. What is the role of contextual embeddings in LLMs, and how do they differ from static embeddings?

**Answer:** Contextual embeddings are dynamic word representations generated by models like transformers that change based on the surrounding words in a sentence. In contrast, static embeddings like Word2Vec or GloVe assign a fixed representation to each word, regardless of its context. The key advantages of contextual embeddings include:

- **Capturing polysemy:** Contextual embeddings can represent multiple meanings of a word based on its context, whereas static embeddings

struggle with words that have different meanings in different sentences.

- **Better performance:** Since they take into account the full context of a word within a sentence, contextual embeddings lead to better performance in tasks like sentiment analysis, machine translation, and question answering.
- **Greater flexibility:** Contextual embeddings adapt to the input, making them more effective in capturing nuanced relationships between words and sentences.

LLMs like BERT and GPT rely on contextual embeddings to achieve state-of-the-art performance in NLP tasks.

Llm

Rags

AI

Data Science



Written by Sanjay Kumar PhD

944 followers · 443 following

Following ▾

Data Science | Machine Learning | AI Product | GenAI | RAG | LLM | AI Agents | NLP | Analytics | Data Engineering | Deep Learning | Statistics

## Responses (1)



Hlgsagar

What are your thoughts?



Mabdullahhalhasib him

Oct 27, 2024

...

Great Article.



2 Reply

## More from Sanjay Kumar PhD



Sanjay Kumar PhD

### Statistics for Data Science Interview Questions and Answers

1. What are the key topics in statistics that are often tested in interviews?

Jan 1 662 32



...



Sanjay Kumar PhD

### Microsoft Azure Interview Questions and Answers

Core Azure Concepts & Cloud Computing Basics

Feb 27 51



...



In Artificial Intelligence in Plain... by Sanjay Kumar...

Sanjay Kumar PhD

## Top 100 AI Agent Interview Questions

What is an AI agent, and how does it work?

Jun 1 67 2

...

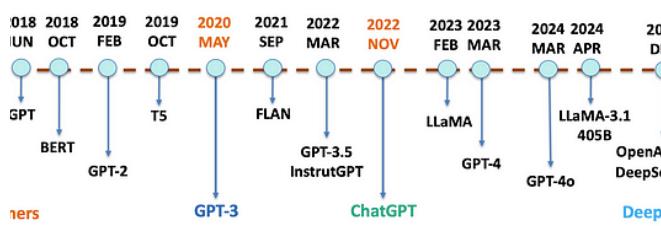
Dec 26, 2024 61 2

...

See all from Sanjay Kumar PhD

## Recommended from Medium

### A Brief History of LLMs



### JP Morgan

Interview Questions



LM LM Po

 In Coding Odyssey by Shivam Srivastava

## A Brief History of LLMs

From Transformers (2017) to DeepSeek-R1 (2025)

 Feb 11  140 



...



 Sanjay Kumar PhD

## RAG Interview Questions and Answers

Q1. What is Retrieval-Augmented Generation (RAG)? A: RAG is a hybrid approach that...

Jan 10  16 



...

## JP Morgan Java Developer Interview

Java Lead Interview Experience

 Jan 24  538 



...



 In Write A Catalyst by Adarsh Gupta

## How I Study Consistently While Working a 9–5 Full-Time Job

No, I don't wake up at 5 AM. And yes, I have a life.

 Apr 21  6.2K 



...



LangChain

 Souvik Majumder

## LangChain Interview Questions & Answers

Detailed list of LangChain interview questions & answers, including general topics & those...



 In AI-ML Interview Playbook by Sajid Khan

## Crack Your Next AI/ML/GenAI Interview: The Ultimate Prep Guid...

Master ML fundamentals, system design, GenAI workflows, and interview strategies in...

Feb 18  4

•••



Jun 3



5



1



•••

---

[See more recommendations](#)