



Top 50 LLM Interview Questions and Answers for 2025

The ultimate LLM interview prep resource! This guide covers the top 50 LLM interview questions and answers to help you crack the code. | ProjectPro

Get Access To All Artificial Intelligence Projects

[View All Artificial Intelligence Projects](#)



Last Updated: 02 Jan 2025 | BY NISHTHA

The tech world is abuzz with the potential of Large Language Models (LLMs). LLMs have revolutionized natural language processing and understanding. From crafting compelling content to automating complex processes, LLMs are rewriting industry rules. As LLMs redefine the paradigms of the industry, the global LLM market will surge to [USD 32.7 billion](#) by 2030. This exponential growth translates to a booming job

market, which requires a comprehensive grasp of LLM tools and technology. This blog post is your one-stop guide to mastering the LLM interview. We've compiled the top 30 LLM interview questions you're likely to encounter and insightful answers to help you get hired for the AI job roles you've been aspiring to.



Llama2 Project for MetaData Generation using FAISS and RAGs

Downloadable solution code | Explanatory videos | Tech Support

[Start Project](#)

Table of Contents

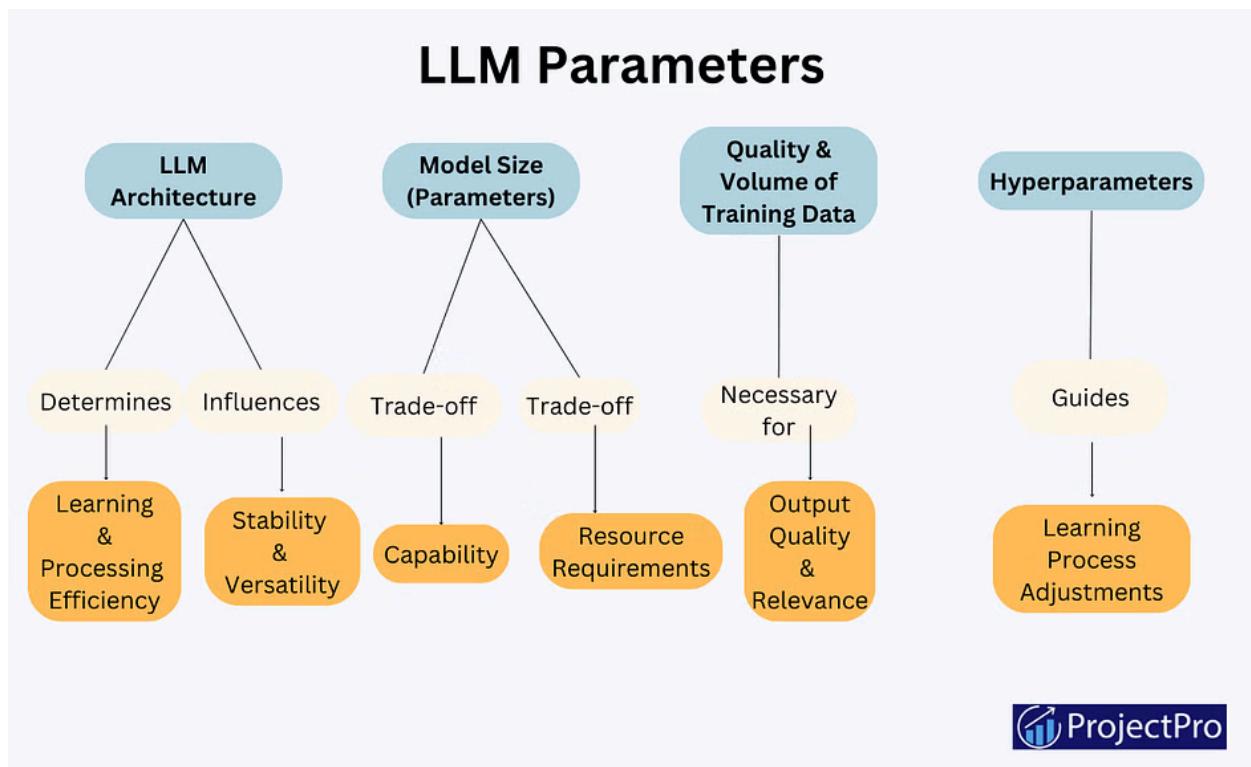
- [Beginner-Level LLM Interview Questions and Answers](#)
- [Intermediate-level LLM / Generative AI Interview Questions and Answers](#)
- [Advanced-level LLM Interview Questions and Answers](#)
- [Open-ended LLM Interview Questions](#)
- [Additional Resources to Prepare for your LLM Interviews](#)
- [A Secret Tip to Ace Large Language Models LLM Interviews](#)
- [Master LLMs with ProjectPro!](#)

Beginner-Level LLM Interview Questions and Answers

Interviews often focus on your grasp of core LLM concepts like [fine-tuning](#), training, architecture, and especially the working mechanism of LLMs. So, this series provides a foundational understanding of [Large Language Models](#) (LLMs) through beginner-friendly interview questions and answers covered below -

1. What are Key LLM Parameters?

LLMs are complex AI models, and their performance is influenced by several key parameters. Understanding these parameters is crucial for effectively developing, training, and utilizing LLMs.



Here's a breakdown of some of the most crucial parameters -

- LLM Architecture:** This refers to the underlying design and structure of the neural network used in the LLM. It determines the model's efficiency in learning and processing information, as well as its stability and versatility in handling different

tasks. Common LLM architectures include transformers, which excel at handling long sequences of text data.

2. **Model Size (Parameters):** This refers to the total number of numerical values (weights and biases) within the LLM's neural network. In simpler terms, it reflects the model's complexity. Generally, larger models with more parameters have a greater capability for complex tasks and can produce more nuanced outputs. However, they also require significantly more computational resources for training and inference.
3. **Quality and Volume of Training Data:** The data an LLM is trained on significantly impacts its output quality and relevance. High-quality, diverse, and relevant data leads to a model that can produce more accurate and reliable outputs. The volume of training data also plays a role, with larger datasets typically leading to better performance, but also requiring more training time and resources.
4. **Hyperparameters:** These are settings that control the LLM's learning process. They don't directly encode knowledge like model parameters, but rather guide how the model learns from the training data. Examples of hyperparameters include learning rate, batch size, and optimizer type. Adjusting these settings allows you to fine-tune the learning process and optimize the model's performance for specific tasks.

Maximize Your Productivity and ROI with ProjectPro



250+ State-of-the-Art End-to-End Projects in Data Engineering, Data Science, Machine Learning, and Cloud.



600 + Hours of Guided and Explanatory Videos by Industry Experts



Personalized Project Paths based on Your Goals



5 New Projects Added Every Month



Deploy Projects to Enterprise Grade Cloud Lab Environment



Unlimited 1:1 Sessions with Top Industry Experts for Project Troubleshooting, Mock Interviews

[Book Free Demo](#)


ProjectPro

2. How do LLMs work?

LLMs process text using word embeddings, a multidimensional representation of words that captures their meaning and relationships to other words. This allows the

transformer model (a deep learning technique) to understand the context and relationships within sentences through an encoder. With this knowledge, the decoder can generate human-like text tailored to the prompt or situation.

3. How are LLMs typically trained?

The core of LLM training is a transformer-based neural network with billions of parameters. These parameters connect nodes across layers, allowing the model to learn complex relationships. LLMs are trained on massive datasets of high-quality text and code. This data provides the raw material for the model to learn language patterns. During training, the model predicts the next word in a sequence based on the previous ones. It then adjusts its internal parameters to improve its predictions, essentially teaching itself through vast amounts of examples. Once trained, LLMs can be fine-tuned for specific tasks. This involves using smaller datasets to adjust the model's parameters towards a particular application.

There are three main approaches for fine-tuning:

1. **Zero-shot learning:** The base LLM can respond to prompts and requests without specific training, though accuracy may vary.
2. **Few-shot learning:** Providing relevant examples significantly improves the model's performance.
3. **Fine-tuning** is a more intensive version of few-shot learning, in which the model is trained on a larger dataset to optimize performance for a particular task..

[Showcase Your PySpark Expertise with ProjectPro's PySpark Certification Course](#)

Here's what valued users are saying about ProjectPro

I come from a background in Marketing and Analytics and when I developed an interest in Machine Learning algorithms, I did multiple in-class courses from reputed institutions though I got good...

ProjectPro is a unique platform and helps many people in the industry to solve real-life problems with a step-by-step walkthrough of projects. A platform with some fantastic resources to gain hands...

Anand Kumpatla



Ameeruddin

Mohammed

ETL (Abintio) developer at IBM

Sr Data Scientist @ Doubleslash
Software Solutions Pvt Ltd

Not sure what you are looking for?

[View All Projects](#)

4. How is training different from fine-tuning?

Training large language models is like this entire process. You start with the foundation - the basic architecture of the model. Then, you slowly build upon it, brick by brick. This involves feeding the model a massive dataset relevant to the desired task. The model learns by analyzing patterns and relationships within the data, adjusting its internal parameters (like weights and biases) to improve performance. Training takes significant time and computational resources, especially for complex models.

In fine-tuning, you take a model already trained on a large, general dataset. Then, you focus on a specific task by retraining the model on a smaller relevant dataset. Since the model has already learned a lot from the initial training, it adapts to the new task much faster and with less data than training from scratch.

Feature	Training	Fine-Tuning
Starting Point	Building the model architecture from scratch	Pre-trained model with learned parameters
Data Requirements	Large, general dataset	Smaller, task-specific dataset
Time and Resources	High	Lower
Customization	More flexible	Less flexible, focuses on refining existing knowledge.

5. Explain the architecture of large-scale LLMs?

The [architecture of a large language model](#) (LLM) is influenced by several factors, including the specific goals of the model, the computational resources available, and the kind of language processing tasks it's designed to perform. Here's a breakdown of the key components that make up a typical LLM architecture:

1. **Transformer Networks:** At the core of most contemporary LLMs lies the Transformer architecture. This neural network departs from traditional recurrent neural networks ([RNNs](#)) and excels at understanding long-range dependencies within sequences, making it particularly well-suited for language processing tasks. Transformers consist of two sub-components:

Encoder: This section processes the input text, breaking it down into a series of encoded representations, capturing the relationships between words.

Decoder: Here, the model leverages the encoded information from the encoder to generate the output text, one word at a time.

2. **Self-Attention Mechanism:** This ingenious mechanism within the Transformer allows the model to focus on the most relevant parts of the input sequence for a given word or phrase. It attends to different parts of the input text differentially, depending on their importance to the prediction at hand. This capability is crucial for LLMs to grasp the nuances of language and context.

Input Embeddings and Output Decoding

Input Embedding: Before feeding text data into the LLM, word embedding transforms it into numerical representations. This process converts words into vectors, capturing their semantic similarities and relationships.

Output Decoding: Once the LLM has processed the encoded input, it translates the internal representation back into human-readable text through decoding

3. **Model Size and Parameter Count:** The number of parameters (weights and biases) within an LLM significantly impacts its capabilities. Large-scale LLMs often have billions, or even trillions, of parameters, allowing them to learn complex patterns and relationships within language data. However, this also necessitates substantial computational resources for training and running the model.

6. What is Hallucination, and How can it be controlled using Prompt Engineering?

Hallucination refers to the model generating factually incorrect or nonsensical outputs.

Imagine a student confidently presenting a made-up historical event. LLMs can do something similar, filling in gaps in their knowledge with creative fiction.

Prompt engineering is a way to rein in these hallucinations. It involves crafting instructions that guide the LLM towards more reliable responses. Here's how

- Just like you wouldn't ask an essay question without any background, a good prompt sets the scene for the LLM, helping it understand what kind of response is expected.
- Do you want a factual summary or a creative story? Letting the LLM know its objective reduces the chances of it going off on tangents.
- Sometimes, giving the LLM multiple-choice options or a specific format keeps its answer on track.

7. Explain the RAG pipeline and each component.

The [Retrieval-Augmented Generation](#) (RAG) pipeline tackles a fundamental limitation of Large Language Models (LLMs) – their reliance on pre-trained data. RAG injects relevant information from external sources, enhancing the LLM's response accuracy and context. Here's a breakdown of the components:

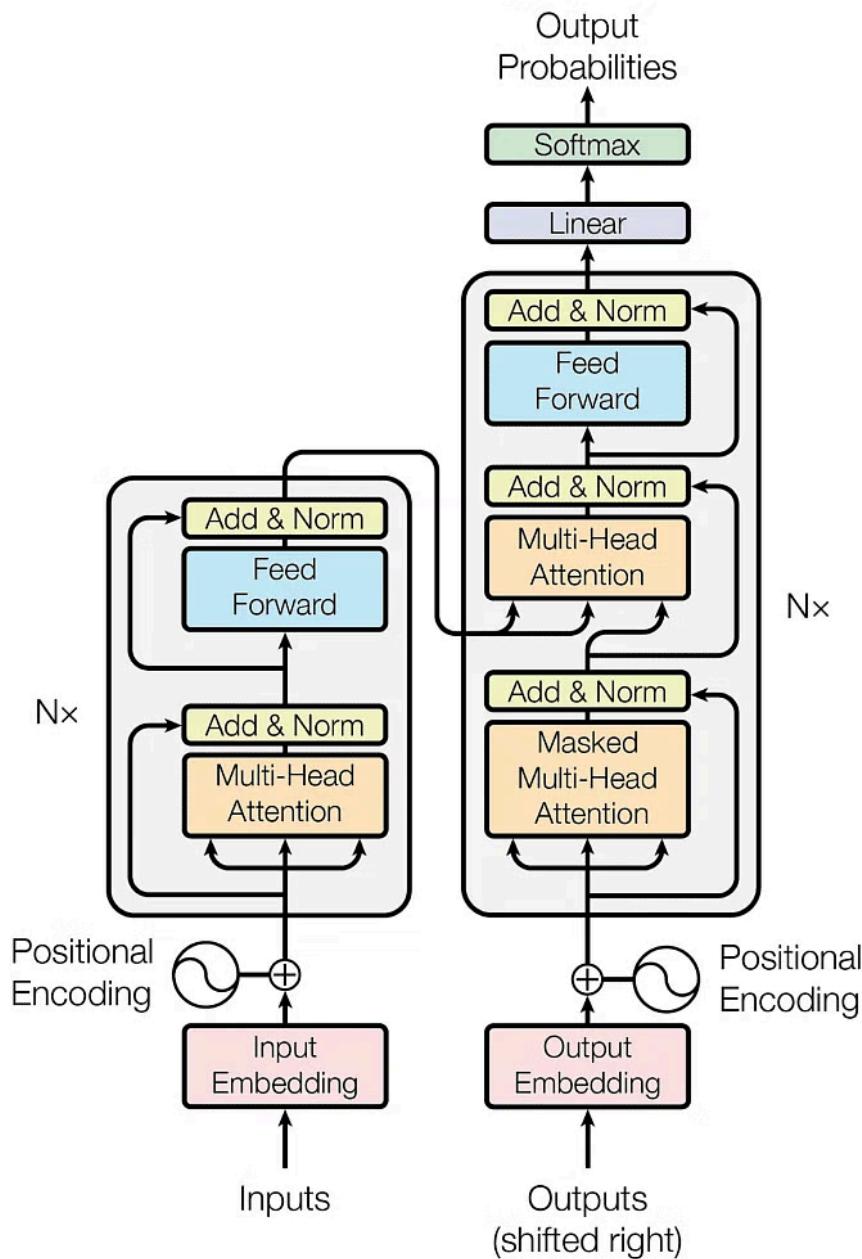
1. **External Data:** RAG utilizes external data sources like databases or documents, transforming this data into numerical representations suitable for the LLM using embedding models. This creates a searchable knowledge base.
2. **Retrieve Relevant Information:** When a user query arrives, RAG converts it into a similar numerical representation and searches the knowledge base for the most relevant data points.
3. **Augment the LLM Prompt:** The retrieved information is strategically added to the user's original query, creating an enriched prompt for the LLM. This provides

additional context for the LLM to generate a more accurate and informative response.

4. **Update External Data:** Maintaining the freshness of external data is crucial. RAG systems employ automated processes to regularly update the information and its corresponding numerical representations.

8. What is the role of transformers in LLM architecture?

Transformers play a critical role in LLM (Large Language Model) architecture by enabling them to understand complex relationships between words in a sentence. This is achieved through self-attention, which lets the model weigh the importance of each word relative to others in the sentence. This is crucial for tasks like machine translation, where a word's meaning depends on the surrounding context.



Source: huggingface.co/

In the encoder-decoder architecture used in translation, the encoder can analyze the entire source sentence simultaneously, while the decoder attends to previously translated words and the complete source sentence to generate the target language. This allows the model to consider word order and context across the entire sentence for accurate translation.

9. What is a token in a Language Model?

Tokens are the basic unit of text that the model processes. Depending on the specific model's design, tokens can be individual words, characters, or phrases. Essentially, they

are the pieces of language the model reads and analyzes to perform various tasks, such as summarizing text or creating new content.

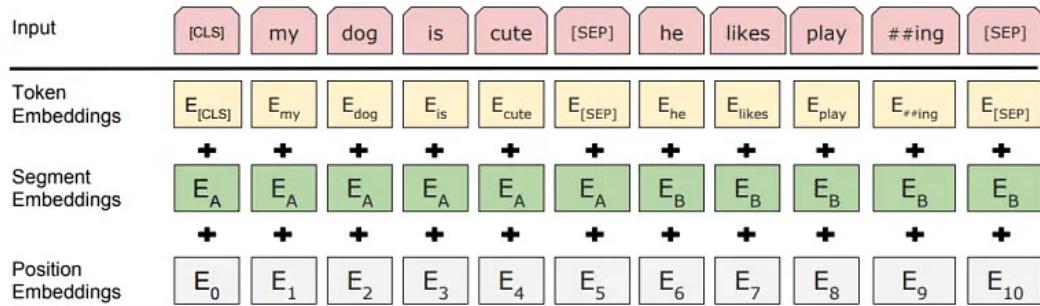
10. Explain the working mechanism of BERT.

BERT, or Bidirectional Encoder Representations from Transformers, leverages the Transformer architecture, specifically its encoder component, to understand contextual relationships between words. The input text is first converted into tokens and enhanced with additional information:

Token embeddings - Special tokens like [CLS] and [SEP] are inserted to mark sentence beginnings and ends.

Segment embeddings - Segment embeddings distinguish between sentences.

Positional embeddings - Positional embeddings indicate each word's location.



Source: huggingface.co/

Unlike traditional models that predict the next word, BERT uses two unique training strategies:

- 1. Masked Language Modeling (MLM):** 15% of words are randomly masked with a [MASK] token. BERT then predicts these masked words based on the context of the remaining unmasked words. This trains the model to understand word meaning based on its surroundings.
- 2. Next Sentence Prediction (NSP):** This technique trains BERT to predict if a given sentence pair is consecutive in the original text. This injects an understanding of sentence relationships, which is crucial for tasks like question answering. During training, BERT is presented with sentence pairs. Half the time, the second sentence truly follows the first. The other half, it's a random sentence entirely. BERT must then classify if the sentences are consecutive, strengthening its ability to grasp connections between sentences.

11. How do you evaluate LLMs?

Evaluating LLMs' performance requires a multifaceted approach considering costs, user experience, and responsible AI practices. Traditional benchmark datasets are limited for LLMs, necessitating real-world user traffic evaluation. This allows for measuring human-like abilities and ensuring a safe and valuable user experience.

Key metrics include:

- GPU utilization (cost estimation).
- Responsible AI (detecting and mitigating risks).
- Performance (latency).
- Utility (user value).

A/B testing is also crucial to measure the impact of LLM features. This includes launch experiments (dark mode, controlled rollout) and post-launch experiments (shadow experiments, regular A/B tests) to optimize the features.

Evaluating LLM's performance requires a multifaceted approach considering costs, user experience, and responsible AI practices.

Intermediate-level LLM / Generative AI Interview Questions and Answers

Mastering LLMs and Generative AI requires understanding critical concepts like components and working mechanisms of the RAG pipeline, prompt engineering, and word2vec for embeddings. This goes beyond encompassing data chunking strategies and mitigating bias in training data. Check out the following interview questions to assess your knowledge of these LLM concepts.

12. Suppose I have a large language model (LLM) and want to tailor it to my specific needs using data. How can I achieve this (prompt

engineering, RAG, fine-tuning, pre-training), and which method is most effective?

There are four main architectural patterns to customize a Large Language Model (LLM) for your needs using data: prompt engineering, Retrieval-Augmented Generation (RAG), fine-tuning, and pre-training. These methods are not competing but complementary, and you can often combine them for optimal results.

Method	Definition	Primary use case	Data requirements	Advantages	Considerations
Prompt engineering	Crafting specialized prompts to guide LLM behavior	Quick, on-the-fly model guidance	None	Fast, cost-effective, no training required	Less control than fine-tuning
Retrieval augmented generation (RAG)	Combining an LLM with external knowledge retrieval	Dynamic datasets and external knowledge	External knowledge base or database (e.g., vector database)	Dynamically updated context, enhanced accuracy	Increases prompt length and inference computation
Fine-tuning	Adapting a pretrained LLM to specific datasets or domains	Domain or task specialization	Thousands of domain-specific or instruction examples	Granular control, high specialization	Requires labeled data, computational cost
Pretraining	Training an LLM from scratch	Unique tasks or domain-specific corpora	Large datasets (billions to trillions of tokens)	Maximum control, tailored for specific needs	Extremely resource-intensive

Source: Databricks

Prompt engineering is a good starting point if you have limited data and computational power.

If you need more comprehensive responses and have access to some relevant data, consider **RAG**.

If accuracy is paramount and you have a substantial dataset, **fine-tuning** is a powerful option.

Pre-training is ideal for highly specialized applications with extensive domain-specific data and significant computational resources.

13. How to estimate infrastructure requirements for fine-tuning an LLM?

Accurately estimating the infrastructure needed for fine-tuning an LLM requires careful consideration of several factors. The two main drivers are the size and complexity of the model you're working with, and the speed at which you want to complete the training process.

Larger models with billions of parameters will naturally demand more computational power and memory. This translates to needing powerful graphics processing units (GPUs) or tensor processing units (TPUs), potentially in a multi-unit configuration. On the other hand, if achieving faster training times is crucial, you'll need to scale your resources even further. This could involve employing multiple GPUs or TPUs working in parallel or leveraging cloud-based solutions that offer on-demand scalability. It's important to remember that all these choices come with a budget attached. Consider exploring techniques like quantization, which can reduce the model size and consequently the computational demands. This might allow you to train the model on less powerful (and more affordable) hardware, making the fine-tuning process more cost-effective.

14. When should I use Fine-tuning instead of RAG?

Fine-tuning is perfect when you need a laser focus on a specific task. Here are some scenarios where it excels over RAG:

1. Is your task in a niche field like medicine or law? Fine-tuning lets you tailor the model to that specific domain's vocabulary and intricacies.
2. Working with limited task-specific data? Fine-tuning leverages the pre-trained LLM's knowledge as a springboard, improving efficiency and preventing overfitting.

3. Does your task require precise outputs, like sentiment analysis or classification?

Fine-tuning excels at recognizing specific patterns within a domain.

15. What are different decoding strategies for picking output tokens?

Here are the essential decoding strategies for picking output tokens in large language models -

1. **Greedy Search (default):** Picks the most likely token at each step, favoring fluent but potentially repetitive outputs.
2. **Multinomial Sampling:** Randomly selects the next token based on predicted probabilities, encouraging diversity but potentially sacrificing coherence.
3. **Beam Search:** Maintains multiple possible continuations (beams) and expands the most promising ones, balancing fluency and exploration.
4. **Contrastive Search (advanced):** Aims for non-repetitive long outputs by considering the distinctiveness of new tokens compared to previously generated ones.

These strategies can be further tweaked to achieve specific goals with parameters like Temperature, Top-k/Top-p Sampling, and Beam Search parameters.

16. What are some of the aspects to keep in mind while using few-shots prompting?

Few-shot prompting offers a powerful way to steer large language models towards specific tasks. Consider the following tips for effective use -

- While it's called "few-shot," focus on the quality of your examples. Choose clear, concise demonstrations that accurately reflect the desired output format and style.
- Don't just provide isolated examples. Include context in your prompts that ties the examples to the desired task. This helps the model understand the relationship between input and output.
- If your task involves multiple categories, ensure your examples represent each category proportionally. A skewed distribution can lead to the model favoring certain outputs.

- Clearly define the task you want the model to perform within the prompt. This helps the model focus on the specific goal and avoid irrelevant information.

17. What is the difference in embedding short and long content?

Embedding short and long content can have some key differences. For short content, capturing the full meaning can be challenging. With fewer words, there's less context for the embedding model to grasp the nuances of the text. This can lead to issues like similar short pieces being considered more alike simply because of their brevity, rather than actual semantic similarity.

On the other hand, long content offers a richer pool of information for the embedding model. It can analyze the relationships between words throughout the text, taking into account things like sentence structure and topic development. This can create a more accurate and comprehensive embedding that reflects the full meaning of the content. However, extremely long content can also present challenges. Some models might struggle to process the sheer volume of information, potentially losing some detail.

18. How embedding models are used in the context of an LLM application?

Embedding models act like a search engine for the LLM, helping it find the most relevant pieces of information from a vast amount of data. This helps LLMs to deliver more accurate, informative, and helpful responses to user queries. Let's understand with the help of an excellent example mentioned below -

Example - Consider an LLM as a powerful language expert, but it needs a way to understand the meaning behind the user's words. Embedding models come in and convert textual inputs (queries, documents) into numerical representations called embeddings. These embeddings capture the semantic relationships between words and concepts.

LLM applications often leverage retrieval-augmented generation (RAG). In RAG, the embedding model converts both the user query and relevant information from a knowledge base into embeddings. Then, it finds the information in the knowledge base

with the most similar embedding to the user's query. This retrieved information, along with the user's query, becomes context for the LLM. Finally, the LLM uses this context to generate a response that's both relevant and informative.

19. How to use stop sequence in LLMs?

A stop sequence is a specific string of text you instruct the LLM to stop generating text after it encounters. This allows you to define a clear endpoint for the model's output. They work by instructing the LLM to halt text generation upon encountering a specific string of characters. This helps you to define clear stopping points for the model's response, leading to more focused and controlled results.

How do I use stop sequences in the OpenAI API?

Here is what you need to know about the stop sequence parameter in the OpenAI Chat Completions API

Updated over a week ago

Stop sequences are used to make the model stop generating tokens at a desired point, such as the end of a sentence or a list. Using the [Chat Completions API](#), you can specify the `stop` parameter and pass in the sequence. The model response will not contain the stop sequence and you can pass up to four stop sequences.

Simple example:

In this [simple chat example](#), one stop sequence is used, the word "World". The system message and the user message are designed to try to get the model to output "Hello world" but as you will see if you run the example in the playground, the model usually stops after just saying "Hello" since world is a stop sequence.

You can explore additional stop sequence examples using the [OpenAI chat playground](#).

1. What is a key difference between training and fine-tuning a Large Language Model (LLM)?

Training requires less data than fine-tuning

Fine-tuning is done from scratch while training uses pre-trained models

Training uses a large general dataset, while fine-tuning uses a smaller task-specific dataset

Fine-tuning requires more computational resources than training

Previous

1 / 7

Next

Source: Open AI

Example - Setting "Best regards" as the stop sequence ensures the model stops before the closing salutation, keeping the email on point. Similarly, for a numbered list, you can define a unique separator like "===" between list items.

20. What are certain strategies to write good prompts?

Give the LLM specific instructions and provide context, such as setting or characters. Use strong verbs to guide the reaction and keep it concise to avoid overwhelming the model.

 H

HarpOnLife

Nov 2023

Hi everyone,

I'm excited to share some insights on creating effective prompts, a skill that is increasingly essential in various AI applications. This guide aims to help you tailor prompts to your specific needs, whether you're in business, creative fields, or tech.

Effective prompting is key to maximizing the potential of AI tools like GPT-4. The right prompt can significantly influence the output quality and relevance. Here are some strategies:

- 1. Understand Your Objective:** Clearly define what you want to achieve with your prompt. Is it information, creativity, or problem-solving?
- 2. Keep It Clear and Concise:** Avoid overly complex or vague prompts. Clarity leads to better AI responses.
- 3. Context Matters:** Provide enough background for the AI to understand the scenario but avoid unnecessary information.
- 4. Experiment and Iterate:** Don't hesitate to refine your prompts based on the responses you get. Iteration is key to finding the most effective wording.
- 5. Consider Your Audience:** Tailor your prompt based on who will interact with or benefit from the AI's response.
- 6. Evaluate and Adapt:** Continuously assess the effectiveness of your prompts and be ready to adapt as needed.

Source: community.openai.com

21. How do you control LLM hallucinations at different levels?

There are several approaches to controlling LLM hallucinations at different levels. At the foundational level, we can focus on improving the training data. This includes using high-quality, well-sourced information and techniques like few-shot learning to guide the model towards accurate responses with relevant examples.

Secondly, we can address hallucinations during the generation process. Prompt engineering plays a crucial role here. Techniques like chain-of-thought prompting encourage the LLM to show its reasoning steps, revealing inconsistencies or factual errors. Additionally, fine-tuning model parameters like temperature and top-p can influence the randomness and diversity of the output, reducing the likelihood of nonsensical responses.

Finally, post-generation methods can identify and mitigate hallucinations. Consistency checks can analyze the generated text for contradictions or illogical elements.

Additionally, we can leverage self-consistency or voting techniques, where the LLM

generates multiple responses and the most frequent answer (among statistically plausible options) is considered the most reliable.

22. How do we increase accuracy and reliability & make answers verifiable in LLM?

Here are some approaches to improve accuracy and reliability and make answers verifiable in LLMs:

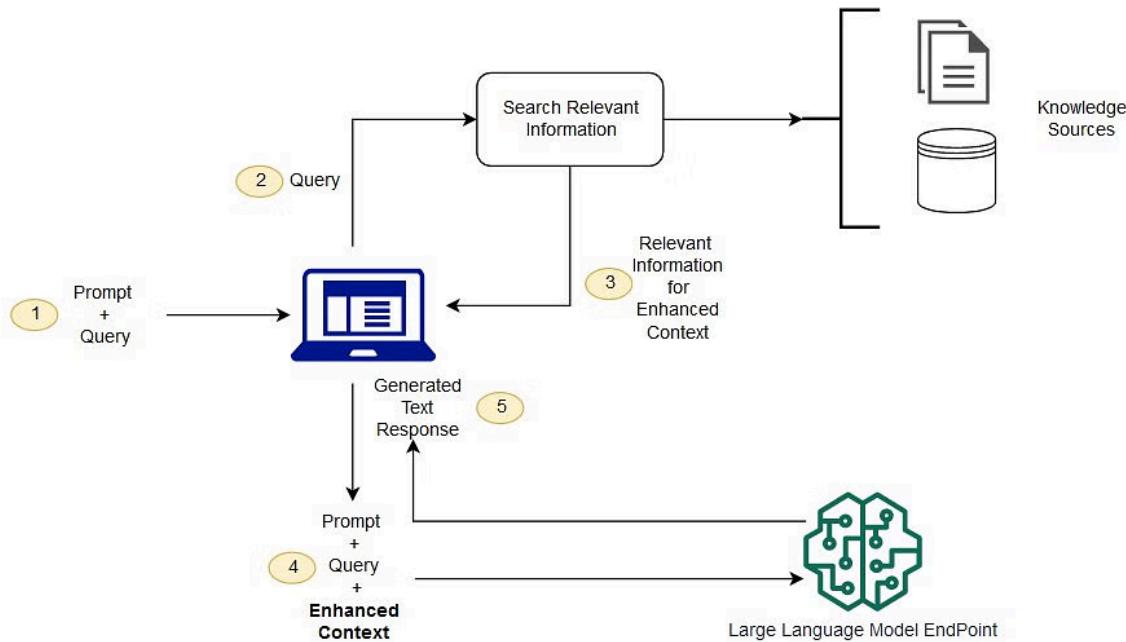
- **LLM-assisted Retrieval** - This technique uses the LLM to generate answers and find supporting documents. The retrieved documents can then be used to verify the accuracy of the generated answer, fostering trust in the LLM's output.
- **Prompt engineering** - How you ask questions can significantly impact the LLM's response. Thus, providing clear and specific prompts can help you guide the LLM toward generating more accurate and relevant answers.
- **Integration with fact-checking APIs** – LLMs can be integrated with third-party APIs designed for fact-checking. This allows the LLM to validate its generated answers against real-world data, improving reliability.
- **Verify-update cycle** – This method involves an iterative process in which the LLM refines its retrieval results based on verification attempts. This loop continues until the retrieved documents convincingly support the answer, enhancing correctness and verifiability.
- **Model parameters and control mechanisms** – Parameters like "number of tokens" and "stop words" can help you influence the LLM's generation process. This allows the model to determine when to stop generating text, potentially reducing the introduction of errors.

23. How does Retrieval augmented generation (RAG) work?

Retrieval-augmented generation (RAG) improves responses from Large Language Models (LLMs) by feeding them relevant external information before they generate text. Here's how it works -

External Data Preparation: First, relevant information from various sources, such as databases or documents, is converted into a format usable by the LLM. This involves converting the data into a numerical representation using embedding models, which creates a searchable knowledge library.

Information Retrieval: When a user asks a question, the system finds related information from the knowledge library. This is done by converting the user query into a numerical representation and matching it with the document embeddings.



Source: aws.amazon.com/

LLM Prompting: The retrieved information is then incorporated into the user's original query to create a more informative prompt for the LLM. This helps the LLM generate a more accurate and relevant response.

Data Updates: To ensure the information stays current, the external data sources and their embeddings are updated regularly. This can be done either continuously or periodically.

24. Can you describe a few different approaches used for chunking text data?

Here are a few different approaches used for chunking text data:

1. **Sentence-based parsing:** This approach breaks down text documents into chunks of complete sentences. It's simple to implement and has low processing costs but

may need to capture the complete context of a thought.

2. **Fixed-size parsing (with overlap):** This approach breaks up a document into chunks based on a fixed number of characters or tokens, with some overlap between chunks. It can capture semantic meaning that spans multiple sentences but requires choosing the chunk size and overlap amount.
3. **Custom code:** This approach uses custom code to create chunks based on patterns within the document's structure. It's useful for documents with known or inferable structures but requires more engineering effort.
4. **Large language model augmentation:** This approach uses large language models to create chunks, such as generating textual representations of images or summaries of tables. It's useful for non-textual data but has high processing costs.
5. **Document layout analysis:** This approach uses libraries and services to extract documents' structure and text. It provides better semantic meaning to content but requires writing code to interact with the model.

25. You are tasked with training an LLM to predict stock market trends. How would you address the challenge of differentiating between factors that cause stock price movements and those that merely correlate with them?

Predicting stock trends with LLMs is tricky. The key challenge is separating factors that cause price movements from those that just correlate. To address this, we can train the LLM on a broader range of data, including news sentiment and social media trends, to capture the broader context. Additionally, causal inference techniques can be implemented during training to help the LLM distinguish true causes from coincidental correlations. And for responsible use, the LLM should be able to explain its predictions by highlighting the factors influencing its reasoning.

26. Why does quantization not decrease accuracy of LLM?

Quantization introduces a slight decrease in accuracy for LLMs. This is because it reduces the precision of the model's parameters, introducing an approximation error. However, quantization offers a compelling trade-off. By using fewer bits to represent parameters, we achieve significant benefits:

- Reduced model size: This allows for deploying LLMs on devices with limited memory.
- Faster inference: Lower precision computations are generally faster to perform.
- Lower power consumption: Smaller models require less energy to run.

The key is that the accuracy loss due to quantization can be minimized. Researchers are constantly developing new techniques, such as targeted quantization for specific parts of the model or mitigating the effects of outliers in activations.

27. What is the difference between vector index, vector DB and vector plugins?

Vector index, vector database, and vector plugins - All of them help with finding similar data points quickly. Vector indexes act like filing systems for high-dimensional data, while vector databases are entire libraries built for this purpose, offering additional features and scalability. Vector plugins are like portable search tools that can be added to existing systems for smaller datasets or initial exploration.

28. Explain different types and challenges associated with filtering in vector DB?

There are two main approaches to filtering in vector databases: pre-filtering and post-filtering. Each method tackles the challenge of efficiently retrieving relevant data while maintaining accuracy. Pre-filtering applies metadata constraints before the actual vector similarity search. This reduces the search space and can be faster. However, it risks excluding relevant data that doesn't perfectly match the filter criteria. Additionally, complex pre-filtering can slow down the query due to the extra processing involved.

Post-filtering performs the vector search first and then filters the results based on metadata. This ensures all potentially relevant data is considered, but it requires filtering out irrelevant results after the search, which can be computationally expensive. To address these challenges, vector databases employ techniques like advanced metadata indexing and parallel processing to speed up filtering. The key lies in finding the right balance between filtering accuracy and search performance. This ensures efficient retrieval of relevant data in vector database queries.

Advanced-level LLM Interview Questions and Answers

Potential employers often test you on advanced LLM concepts to assess your understanding. Check out a few challenging questions and insightful answers below to prepare for your next LLM interview and showcase your expertise!

29. How is word2vec trained from Scratch?

There are several steps involved in training word2vec from scratch. Check them out below -

1. **Data Preprocessing:** First, you'll need a text corpus. This could be news articles, books, or any large text data collection. The text needs to be cleaned and preprocessed. This involves breaking it down into words (tokens) and removing stop words (common words like "the" or "a") that don't hold much meaning.
2. **Defining Window Size:** Next, you define a window size. This determines how many words around a target word (center word) will be considered context. Words within this window are considered "related" to the center word.
3. **Choosing Training Method:** Word2vec has two main training methods: CBOW (Continuous Bag-of-Words) and Skip-gram. CBOW predicts the center word based on its surrounding context words. Skip-gram does the opposite, predicting surrounding words based on a given center word.
4. **Training the Model:** Here, you create a [neural network](#) architecture. Words are converted into one-hot encoded vectors. The network is then trained using the chosen method (CBOW or Skip-gram) to predict the target words based on context

or vice versa. This process iterates through the text corpus, adjusting the network weights to improve prediction accuracy.

5. **Model Evaluation and Use:** After training, you can evaluate the model's performance and analyze the learned word embeddings. These embeddings represent words as numerical vectors, where similar words will have vector representations close to each other in the embedding space. This allows you to perform natural language processing tasks like finding synonyms or measuring word relationships.

30. Imagine you're building a system to automatically categorize customer reviews based on sentiment (positive, negative, neutral). The system needs to understand the meaning of words within the reviews. How can we represent words numerically in a way that captures their meaning and relationships to other words?

Word embeddings are a powerful technique that goes beyond simply assigning a unique number to each word. Word embeddings map words to vectors of real numbers in a high-dimensional space. The key idea is that words with similar meanings will have similar vector representations. Techniques like Word2Vec and GloVe analyze large text corpora to learn these relationships. Words like "king" and "queen" would be closer in vector space compared to "king" and "banana" because of their semantic connection. Thus, this technique will help you represent the sentiment of a review by averaging the word vectors within the review. Reviews with a higher average vector value in the "positive" direction of the embedding space would be classified as positive, and vice versa for negative sentiment.

31. Imagine you're designing software that interacts with its environment to achieve goals. Can you describe the core idea of this program and the different approaches you could take to make it function effectively?

The program you're designing can be classified as an intelligent agent. This means it's a software program that can perceive its environment, reason, and take action to achieve specific goals. It accomplishes this through a continuous loop of:

1. **Perception:** Gathering information about the environment through sensors or user input.
2. **Processing:** Interpreting the perceived information and understanding the current state of the environment.
3. **Decision-Making:** Using reasoning and available knowledge to choose an action that moves it closer to its goal.
4. **Action:** Taking a concrete step in the environment to influence its state.

There are several approaches to make this intelligent agent function effectively:

- Clearly define and represent the program's goals in a way the agent can understand. This could involve decision trees, utility functions, or reinforcement learning rewards.
- Design efficient sensors to gather relevant data from the environment. This could involve cameras, microphones, or API access for software agents.
- Implement algorithms for the agent to process data, make decisions, and plan actions. This could involve rule-based systems, machine-learning models, or a combination of both.
- Allow the agent to learn from its experiences and adapt its behavior over time. This could involve reinforcement learning or other supervised/unsupervised learning techniques.
- Develop methods for the agent to take action in the environment. This could involve controlling robots, manipulating software interfaces, or generating outputs like text or recommendations.

32. How do you mitigate catastrophic forgetting in LLMs ?

Catastrophic forgetting occurs when an LLM learning new information completely forgets or significantly weakens its grasp of previously learned tasks. This is problematic for LLMs aiming to be versatile and adaptable.

Several techniques address this challenge. One approach is rehearsal-based methods, where the LLM revisits a subset of old data alongside the new data during training. This helps the LLM retain its knowledge of past tasks. Another technique is elastic weight consolidation (EWC), which assigns importance scores to different weights in the LLM's network. More recent advancements involve introducing separate modules for new tasks. This allows the LLM to expand its knowledge base without completely rewriting existing connections. Techniques like progressive neural networks (ProgNet) and optimized fixed expansion layers (OFELs) fall under this category.

33. State the difference between OpenAI functions and LangChain functions.

When it comes to integrating functions with large language models, OpenAI Functions and [LangChain](#) Functions offer different approaches. OpenAI Functions provide a more hands-on experience. You have the flexibility to define custom functions and exert greater control over which functions the AI calls and how it uses them. This comes at the cost of increased development effort, as you'll need to write code to parse the AI's response and manage arguments.

On the other hand, LangChain Functions prioritize ease of use. Using LangChain toolkits can help you avoid writing function definitions and implementations from scratch. However, this comes with limitations. You're restricted to the functions available through chosen agent types or toolkits. Additionally, informing the model about available functions relies on selecting the appropriate agent type within LangChain.

Check out this [Reddit thread](#) for more details on the comparison between OpenAI Functions and LangChain agents.

34. How can you get the best performance and functionality from an LLM system while controlling costs?

Here are five excellent tips to optimize the cost of your LLM system while maintaining good performance and functionality -

1. LLMs are charged by the token (word, punctuation, etc.). Thus, you must focus on clear and specific prompts that get straight to the point and avoid unnecessary information.
2. If you're using an LLM API, consider batching similar requests together. This can improve processing speed and reduce the calls needed, lowering costs.
3. Techniques like LLM Lingua can remove unnecessary words from inputs and outputs, reducing the number of tokens processed and saving money.
4. Limit the conversation history stored in memory for chatbots or interactive systems. This reduces the number of tokens the LLM needs to consider for each response.
5. Different LLMs have varying costs and capabilities. Select an LLM designed explicitly for your tasks to get the most performance for your budget.

35. How do you handle overfitting in LLMs?

Overfitting can lead to the LLM producing nonsensical outputs or mimicking the training data too closely.

We can leverage techniques used in other machine learning models to handle overfitting in LLMs. One approach is data augmentation, where we artificially expand the training data by creating variations of existing examples. This forces the LLM to learn underlying patterns instead of memorizing specific phrases. Additionally, [regularization techniques](#) like dropout can be employed. Dropout randomly deactivate neurons during training, preventing the model from relying too heavily on any single feature and promoting better generalization. These methods can help us train LLMs that perform well on new data and avoid getting stuck on the specifics of the training set.

36. How do you handle long-term dependencies in language models?

Long-term dependencies in language models refer to the challenge of capturing relationships between words that appear far apart in a sequence. Standard neural networks often struggle with this, as information can degrade over time. Here are some standard techniques to tackle long-term dependencies:

- **Recurrent Neural Networks (RNNs):** These networks have internal loops that allow them to store information from previous inputs and use it to influence the processing of current ones. This creates a form of memory that can bridge gaps between distant words.
- **Long Short-Term Memory (LSTM) networks:** A specific type of RNN architecture designed to address the vanishing gradient problem. LSTMs have internal gates that control the flow of information, allowing them to retain relevant information for extended periods and handle dependencies across greater distances. Learn more on this from the words of [Tim Klawa](#) in his post shown below -

Tim Klawa Follow
Head of Product @ F8F | Top AI Voice LinkedIn | Board Member ...

Gated units, like Long Short-Term Memory (LSTM) cells or Gated Recurrent Units (GRUs), are designed to handle long-term dependencies by regulating the flow of information. They use mechanisms called gates to control the retention and forgetting of information in the cell state. These gates can learn which data in a sequence is important to keep or discard, allowing the network to maintain relevant information over long sequences and mitigate the vanishing gradient problem.

Like · 2 Unhelpful

Source: LinkedIn

Image Name: "Handle long-term dependencies in LLM models"

- **Transformer-XL** is a modification of the Transformer architecture, another powerful neural network for language tasks. Transformer-XL addresses long-term dependencies by keeping and utilizing hidden states from previous processing steps when encoding new input segments. This allows it to maintain context over longer sequences.

- **Attention mechanisms:** These techniques allow the model to focus on specific parts of the input sequence most relevant to the current prediction. Attending the informative parts of the sequence, even if they are far away, can help the model learn long-term dependencies.

37. How to improve LLM reasoning if your COT prompt fails?

If your initial Chain-of-Thought (CoT) prompt fails, consider refining it to address vagueness or missing steps. You can also explore other prompting techniques like providing assumptions, background knowledge, or encouraging multi-perspective analysis. Human feedback on the LLM's reasoning process can also be invaluable. Finally, ensure your training data includes diverse examples of the desired reasoning to improve the LLM's capabilities.

38. How do you handle bias in large language models?

Large language models (LLMs) are susceptible to bias because they learn from the data they're trained on. This data can contain inherent biases and misinformation, leading the LLM to reflect those biases in its outputs.

Here are three excellent ways through which we can handle bias in LLMs -

1. **Data Quality:** We can mitigate bias by ensuring high-quality training data. This involves using diverse and representative datasets that reflect the real world. Techniques like data augmentation and filtering can help achieve this.
2. **Mitigating Bias During Training:** New training methods are being developed to address bias during the training process itself. These methods include incorporating logic rules and fine-tuning the model on unbiased datasets for specific tasks.
3. **Identifying and Flagging Bias:** We can develop tools to identify and flag potential bias in LLM outputs. This allows users to be aware of potential biases and make informed decisions.

39. Imagine you are working in a customer service chatbot application. How do you stay up-to-date with new product information and evolving customer inquiries?

I'd continuously learn from a frequently updated knowledge base fed by product launches, manuals, and internal resources. Additionally, I can analyze conversation transcripts and user feedback to identify knowledge gaps and emerging trends. This allows developers to refine my training data. I can seamlessly escalate to human agents who can capture user feedback and update the knowledge base for complex inquiries. [Sentiment analysis](#) can further identify confusion or frustration with new products, prompting targeted knowledge base improvements. This continuous learning loop ensures I provide users with the most up-to-date and relevant information.

40. You are used in a virtual assistant application and encounter a user request you cannot understand or fulfill due to limitations in your training data. How do you handle this situation?

When encountering a user request beyond my expertise in a virtual assistant scenario requires a user-centric approach. I'll acknowledge the request, politely explain my limitations ("I'm sorry, I can't assist you with '{request}' yet"), and offer options to keep the experience positive. This could be learning more from the user ("Would you like to tell me more about your request?") or suggesting alternative assistance ("Can I help you with something else?"). This showcases my ability to understand user intent, communicate clearly, adapt to situations, and continuously learn through user interactions – all vital qualities for a virtual assistant LLM.

Open-ended LLM Interview Questions

41. How do you determine the best vector database for your needs?
42. How to choose the right LLM model for your use case ?
43. How would you decide ideal search similarity metrics for the use case?
44. In what scenarios do you fine-tune an LLM model and how do you make this decision ?
45. How do I improve the reasoning ability of my LLM through prompt engineering ?
46. What are different ways you can define stopping criteria in a large language model?
47. List a few metrics that you have used to evaluate an LLM?
48. Compare different Vector index and given a scenario, how would you decide which vector index you would use for a project?
49. LLMs are constantly evolving, with new capabilities emerging. What do you see as the most promising future application of LLMs, and what challenges do you foresee in bringing this application to reality?
50. What architecture patterns do you see when you want to customize your LLM with proprietary data?

Additional Resources to Prepare for your LLM Interviews

[Youssef Hosni](#), (Founder & Author @To Data & Beyond), recently shared some insightful Large Language Model (LLM) interview questions in a LinkedIn post. Check them out below:



Youssef Hosni • Following

Data Scientist | AI Researcher | Founder & Author @ To Data & Beyond

9mo •

...

Top Large Language Models (LLMs) Interview Questions & Answers

1. NLP Basics Interview Questions & Answers

1.1. Describe the process of tokenization in NLP. Why is it important, and what challenges can arise during tokenization?

1.2. How do you handle the problem of bias in NLP models, and what techniques can be used to mitigate it?

1.3. What are the key evaluation metrics for assessing the performance of NLP models, especially in tasks like text classification and machine translation?

1.4. Explain the concept of word embeddings, and how are they used in NLP models like Word2Vec and GloVe.

1.5. Explain the concept of transfer learning in NLP and its importance in building effective NLP models.

1.6. Can you discuss some common techniques for handling out-of-vocabulary (OOV) words in NLP tasks?

1.7. How do you choose the appropriate architecture or model size for a specific NLP task?

1.8. What are the differences between supervised, unsupervised, and semi-supervised learning in NLP, and when would you use each approach?

1.9. In the context of NLP, what are some common techniques for text data preprocessing and cleaning?

1.10. Describe some popular libraries and frameworks used for working with NLP tasks, such as spaCy, NLTK, or Hugging Face Transformers.

1.11. Explain the concept of sequence-to-sequence models in NLP and their applications, especially in tasks like language translation and text summarization.

2. Transformer-Based Interview Questions &Answers

2.1. Explain the Transformer architecture and its significance in NLP.

2.2. What is the purpose of attention mechanisms in Transformer models, and how do they work?

2.3. Can you provide an overview of the attention mechanism's evolution, from the original Transformer to more recent variations like BERT and GPT models?

3. Large Language Models Questions & Answers

3.1. Can you explain how a large language model like GPT-3 works, and what are some of its key applications and limitations?

3.2. What is the difference between pre-training and fine-tuning in the context of large language models like BERT or GPT-3?

3.3. What are the limitations of large language models like GPT-3, and how can those limitations be addressed?

Source: LinkedIn

You can also check out this GitHub Repo - [Devinterview-io/llms-interview-questions](https://github.com/Devinterview-io/llms-interview-questions) that compiles over 60 essential LLM interview questions, helping you prepare for your next data science challenge.

A Secret Tip to Ace Large Language Models LLM Interviews

Acing Large Language Model (LLM) interviews goes beyond just theoretical knowledge. The secret tip to impress potential employers is gaining hands-on experience by working on real-world general AI and [LLM projects](#). This demonstrates your practical understanding of the field and your ability to apply your knowledge to solve complex problems.

Check out the following compilation of Gen AI / LLM projects and Learning paths that you must consider exploring to stand out in a competitive interview setting -

1. [LLM Roadmap](#)
2. [Generative AI Learning Path](#)
3. [LLM Project to Build and Fine Tune a Large Language Model](#)
4. [Build a Langchain Streamlit Chatbot for EDA using LLMs](#)
5. [Langchain Project for Customer Support App in Python](#)
6. [Llama2 Project for MetaData Generation using FAISS and RAGs](#)
7. [AI Video Summarization Project using Mixtral, Whisper, and AWS](#)
8. [Learn to Build Generative Models Using PyTorch Autoencoders](#)

Working on these projects will solidify your understanding of LLMs and allow you to showcase your practical skills and problem-solving abilities to potential employers. Imagine discussing a project where you tackled a specific challenge using LLMs – it demonstrates theoretical knowledge and the ability to apply it effectively. Consider including these projects in your portfolio or resume to set yourself apart and increase your chances of getting shortlisted.

Master LLMs with ProjectPro!

When it comes to building projects using different LLM architectures, applying LLM knowledge to solve practical problems in areas like [text summarization](#), machine translation, and chatbot development, and learning how to fine-tune pre-trained LLMs for your chosen project domain, ProjectPro is all that you need! [ProjectPro's LLM roadmap](#) offers a strategic advantage in helping you traverse through the complexities of LLM job interviews. With ProjectPro's curated LLM projects, candidates not only

showcase practical expertise but also demonstrate readiness to tackle real-world challenges.

[PREVIOUS](#)[NEXT](#)

Your Data Skills Need to Get Stronger
And We Have The Projects Ready

GET ACCESS TO SOLVED PROJECTS

The banner features a teal background with white text. On the right side, there is a 3D-style illustration of three people interacting with a large smartphone. The phone's screen displays a globe and various data visualization elements like bars and clouds. One person is sitting on a cloud, another is standing next to the phone, and a third is running towards it. There are also small icons of a notepad and a pencil.

About the Author



Nishtha

Nishtha is a professional Technical Content Analyst at ProjectPro with over three years of experience in creating high-quality content for various industries. She holds a bachelor's degree in...

[Meet The Author >](#)

Project Categories

Machine Learning Projects

Data Science Projects

Deep Learning Projects

Projects

Walmart Sales Forecasting Data Science Project

BigMart Sales Prediction ML Project

Music Recommender System Project

[Big Data Projects](#)[Apache Hadoop Projects](#)[Apache Spark Projects](#)[Show more](#)[Credit Card Fraud Detection Using Machine Learning](#)[Resume Parser Python Project for Data Science](#)[Time Series Forecasting Projects](#)[Show more](#)

Blogs

[Machine Learning Projects for Beginners with Source Code](#)[Data Science Projects for Beginners with Source Code](#)[Big Data Projects for Beginners with Source Code](#)[IoT Projects for Beginners with Source Code](#)[Data Analyst vs Data Scientist](#)[Data Science Interview Questions and Answers](#)[Show more](#)

Certification Courses

[Practical MLOps Course](#)[Data Engineering Course](#)[AWS Data Engineering Course](#)[Azure Data Engineering Course](#)[GCP Data Engineering Course](#)[PySpark Course](#)[Snowflake Course](#)[Show more](#)

Tutorials

PCA in Machine Learning Tutorial

PySpark Tutorial

Hive Commands Tutorial

MapReduce in Hadoop Tutorial

Apache Hive Tutorial -Tables

Linear Regression Tutorial

Show more

ProjectPro

© 2025 Iconiq Inc.

About us

Contact us

Privacy policy

User policy

Write for ProjectPro