

Test Databricks-Generative-AI-Engineer-Associate Prep | Pass4sure Databricks-Generative-AI-Engineer-Associate Exam Prep



What's more, part of that ExamPrepAway Databricks-Generative-AI-Engineer-Associate dumps now are free: <https://drive.google.com/open?id=18vxhkNbkQ458mOF8zQYcU61zGz4vMsbe>

If you are still unsure whether to pursue Databricks Databricks-Generative-AI-Engineer-Associate exam questions for Databricks Databricks Certified Generative AI Engineer Associate exam preparation, you are losing the game at the first stage in a fiercely competitive marketplace. Databricks Databricks-Generative-AI-Engineer-Associate Questions are the best option for becoming Databricks Databricks Certified Generative AI Engineer Associate.

The Databricks-Generative-AI-Engineer-Associate exam is highly competitive and acing it is not a piece of cake for majority of the people. It requires a great skill set and deep knowledge Databricks-Generative-AI-Engineer-Associate Exam Questions. An aspirant achieving Databricks Certified Generative AI Engineer Associate (Databricks-Generative-AI-Engineer-Associate) certificate truly reflects his hard work and consistent struggle. These Databricks-Generative-AI-Engineer-Associate exam practice test a person's true capacities and passing it requires extensive knowledge of each Databricks-Generative-AI-Engineer-Associate topic.

>> Test Databricks-Generative-AI-Engineer-Associate Prep <<

Perfect Databricks - Test Databricks-Generative-AI-Engineer-Associate Prep

Our qualified team of Databricks Databricks Certified Generative AI Engineer Associate study material to improve the quality and to match the changes in the syllabus and pattern shared by Databricks-Generative-AI-Engineer-Associate. Our desktop Databricks Databricks-Generative-AI-Engineer-Associate Practice Exam software is designed for all those candidates who want to learn and practice in the actual Databricks Databricks-Generative-AI-Engineer-Associate exam environment.

Databricks Databricks-Generative-AI-Engineer-Associate Exam Syllabus Topics:

Topic	Details
Topic 1	<ul style="list-style-type: none">Design Applications: The topic focuses on designing a prompt that elicits a specifically formatted response. It also focuses on selecting model tasks to accomplish a given business requirement. Lastly, the topic covers chain components for a desired model input and output.
Topic 2	<ul style="list-style-type: none">Evaluation and Monitoring: This topic is all about selecting an LLM choice and key metrics. Moreover, Generative AI Engineers learn about evaluating model performance. Lastly, the topic includes sub-topics about inference logging and usage of Databricks features.

Topic 3	<ul style="list-style-type: none"> • Data Preparation: Generative AI Engineers covers a chunking strategy for a given document structure and model constraints. The topic also focuses on filter extraneous content in source documents. Lastly, Generative AI Engineers also learn about extracting document content from provided source data and format.
---------	--

Databricks Certified Generative AI Engineer Associate Sample Questions (Q19-Q24):

NEW QUESTION # 19

A Generative AI Engineer is building a production-ready LLM system which replies directly to customers. The solution makes use of the Foundation Model API via provisioned throughput. They are concerned that the LLM could potentially respond in a toxic or otherwise unsafe way. They also wish to perform this with the least amount of effort. Which approach will do this?

- A. Add some LLM calls to their chain to detect unsafe content before returning text
- **B. Host Llama Guard on Foundation Model API and use it to detect unsafe responses**
- C. Add a regex expression on inputs and outputs to detect unsafe responses.
- D. Ask users to report unsafe responses

Answer: B

Explanation:

The task is to prevent toxic or unsafe responses in an LLM system using the Foundation Model API with minimal effort. Let's assess the options.

* Option A: Host Llama Guard on Foundation Model API and use it to detect unsafe responses

* Llama Guard is a safety-focused model designed to detect toxic or unsafe content. Hosting it via the Foundation Model API (a Databricks service) integrates seamlessly with the existing system, requiring minimal setup (just deployment and a check step), and leverages provisioned throughput for performance.

* Databricks Reference: "Foundation Model API supports hosting safety models like Llama Guard to filter outputs efficiently" ("Foundation Model API Documentation," 2023).

* Option B: Add some LLM calls to their chain to detect unsafe content before returning text

* Using additional LLM calls (e.g., prompting an LLM to classify toxicity) increases latency, complexity, and effort (crafting prompts, chaining logic), and lacks the specificity of a dedicated safety model.

* Databricks Reference: "Ad-hoc LLM checks are less efficient than purpose-built safety solutions" ("Building LLM Applications with Databricks").

* Option C: Add a regex expression on inputs and outputs to detect unsafe responses

* Regex can catch simple patterns (e.g., profanity) but fails for nuanced toxicity (e.g., sarcasm, context-dependent harm), requiring significant manual effort to maintain and update rules.

* Databricks Reference: "Regex-based filtering is limited for complex safety needs" ("Generative AI Cookbook").

* Option D: Ask users to report unsafe responses

* User reporting is reactive, not preventive, and places burden on users rather than the system. It doesn't limit unsafe outputs proactively and requires additional effort for feedback handling.

* Databricks Reference: "Proactive guardrails are preferred over user-driven monitoring" ("Databricks Generative AI Engineer Guide").

Conclusion: Option A (Llama Guard on Foundation Model API) is the least-effort, most effective approach, leveraging Databricks' infrastructure for seamless safety integration.

NEW QUESTION # 20

A Generative AI Engineer is creating an LLM system that will retrieve news articles from the year 1918 and related to a user's query and summarize them. The engineer has noticed that the summaries are generated well but often also include an explanation of how the summary was generated, which is undesirable.

Which change could the Generative AI Engineer perform to mitigate this issue?

- A. Revisit their document ingestion logic, ensuring that the news articles are being ingested properly.
- **B. Provide few shot examples of desired output format to the system and/or user prompt.**
- C. Split the LLM output by newline characters to truncate away the summarization explanation.
- D. Tune the chunk size of news articles or experiment with different embedding models.

Answer: B

Explanation:

To mitigate the issue of the LLM including explanations of how summaries are generated in its output, the best approach is to adjust the training or prompt structure. Here's why Option D is effective:

* Few-shot Learning: By providing specific examples of how the desired output should look (i.e., just the summary without explanation), the model learns the preferred format. This few-shot learning approach helps the model understand not only what content to generate but also how to format its responses.

* Prompt Engineering: Adjusting the user prompt to specify the desired output format clearly can guide the LLM to produce summaries without additional explanatory text. Effective prompt design is crucial in controlling the behavior of generative models.

Why Other Options Are Less Suitable:

* A: While technically feasible, splitting the output by newline and truncating could lead to loss of important content or create awkward breaks in the summary.

* B: Tuning chunk sizes or changing embedding models does not directly address the issue of the model's tendency to generate explanations along with summaries.

* C: Revisiting document ingestion logic ensures accurate source data but does not influence how the model formats its output. By using few-shot examples and refining the prompt, the engineer directly influences the output format, making this approach the most targeted and effective solution.

NEW QUESTION # 21

A Generative AI Engineer is tasked with developing a RAG application that will help a small internal group of experts at their company answer specific questions, augmented by an internal knowledge base. They want the best possible quality in the answers, and neither latency nor throughput is a huge concern given that the user group is small and they're willing to wait for the best answer. The topics are sensitive in nature and the data is highly confidential and so, due to regulatory requirements, none of the information is allowed to be transmitted to third parties.

Which model meets all the Generative AI Engineer's needs in this situation?

- A. Dolly 1.5B
- B. OpenAI GPT-4
- C. Llama2-70B
- **D. BGE-large**

Answer: D

Explanation:

Problem Context: The Generative AI Engineer needs a model for a Retrieval-Augmented Generation (RAG) application that provides high-quality answers, where latency and throughput are not major concerns. The key factors are confidentiality and sensitivity of the data, as well as the requirement for all processing to be confined to internal resources without external data transmission.

Explanation of Options:

* Option A: Dolly 1.5B: This model does not typically support RAG applications as it's more focused on image generation tasks.

* Option B: OpenAI GPT-4: While GPT-4 is powerful for generating responses, its standard deployment involves cloud-based processing, which could violate the confidentiality requirements due to external data transmission.

* Option C: BGE-large: The BGE (Big Green Engine) large model is a suitable choice if it is configured to operate on-premises or within a secure internal environment that meets regulatory requirements.

Assuming this setup, BGE-large can provide high-quality answers while ensuring that data is not transmitted to third parties, thus aligning with the project's sensitivity and confidentiality needs.

* Option D: Llama2-70B: Similar to GPT-4, unless specifically set up for on-premises use, it generally relies on cloud-based services, which might risk confidential data exposure.

Given the sensitivity and confidentiality concerns, BGE-large is assumed to be configurable for secure internal use, making it the optimal choice for this scenario.

NEW QUESTION # 22

A Generative AI Engineer is designing a RAG application for answering user questions on technical regulations as they learn a new sport.

What are the steps needed to build this RAG application and deploy it?

- A. User submits queries against an LLM -> Ingest documents from a source -> Index the documents and save to Vector Search -> LLM retrieves relevant documents -> LLM generates a response -> Evaluate model -> Deploy it using Model

Serving

- B. Ingest documents from a source -> Index the documents and saves to Vector Search -> User submits queries against an LLM -> LLM retrieves relevant documents -> Evaluate model -> LLM generates a response -> Deploy it using Model Serving
- C. Ingest documents from a source -> Index the documents and save to Vector Search -> Evaluate model -> Deploy it using Model Serving
- D. Ingest documents from a source -> Index the documents and save to Vector Search -> User submits queries against an LLM -> LLM retrieves relevant documents -> LLM generates a response -> Evaluate model -> Deploy it using Model Serving

Answer: D

Explanation:

The Generative AI Engineer needs to follow a methodical pipeline to build and deploy a Retrieval-Augmented Generation (RAG) application. The steps outlined in option B accurately reflect this process:

* Ingest documents from a source: This is the first step, where the engineer collects documents (e.g., technical regulations) that will be used for retrieval when the application answers user questions.

* Index the documents and save to Vector Search: Once the documents are ingested, they need to be embedded using a technique like embeddings (e.g., with a pre-trained model like BERT) and stored in a vector database (such as Pinecone or FAISS). This enables fast retrieval based on user queries.

* User submits queries against an LLM: Users interact with the application by submitting their queries.

These queries will be passed to the LLM.

* LLM retrieves relevant documents: The LLM works with the vector store to retrieve the most relevant documents based on their vector representations.

* LLM generates a response: Using the retrieved documents, the LLM generates a response that is tailored to the user's question.

* Evaluate model: After generating responses, the system must be evaluated to ensure the retrieved documents are relevant and the generated response is accurate. Metrics such as accuracy, relevance, and user satisfaction can be used for evaluation.

* Deploy it using Model Serving: Once the RAG pipeline is ready and evaluated, it is deployed using a model-serving platform such as Databricks Model Serving. This enables real-time inference and response generation for users.

By following these steps, the Generative AI Engineer ensures that the RAG application is both efficient and effective for the task of answering technical regulation questions.

NEW QUESTION # 23

A Generative AI Engineer interfaces with an LLM with prompt/response behavior that has been trained on customer calls inquiring about product availability. The LLM is designed to output "In Stock" if the product is available or only the term "Out of Stock" if not. Which prompt will work to allow the engineer to respond to call classification labels correctly?

- A. You will be given a customer call transcript where the customer inquires about product availability. Respond with "In Stock" if the product is available or "Out of Stock" if not.
- B. Respond with "Out of Stock" if the customer asks for a product.
- C. Respond with "In Stock" if the customer asks for a product.
- D. You will be given a customer call transcript where the customer asks about product availability. The outputs are either "In Stock" or "Out of Stock". Format the output in JSON, for example: `{"call_id": "123", "label": "In Stock"}`.

Answer: D

Explanation:

* Problem Context: The Generative AI Engineer needs a prompt that will enable an LLM trained on customer call transcripts to classify and respond correctly regarding product availability. The desired response should clearly indicate whether a product is "In Stock" or "Out of Stock," and it should be formatted in a way that is structured and easy to parse programmatically, such as JSON.

* Explanation of Options:

* Option A: Respond with "In Stock" if the customer asks for a product. This prompt is too generic and does not specify how to handle the case when a product is not available, nor does it provide a structured output format.

* Option B: This option is correctly formatted and explicit. It instructs the LLM to respond based on the availability mentioned in the customer call transcript and to format the response in JSON.

This structure allows for easy integration into systems that may need to process this information automatically, such as customer service dashboards or databases.

* Option C: Respond with "Out of Stock" if the customer asks for a product. Like option A, this prompt is also insufficient as it only covers the scenario where a product is unavailable and does not provide a structured output.

