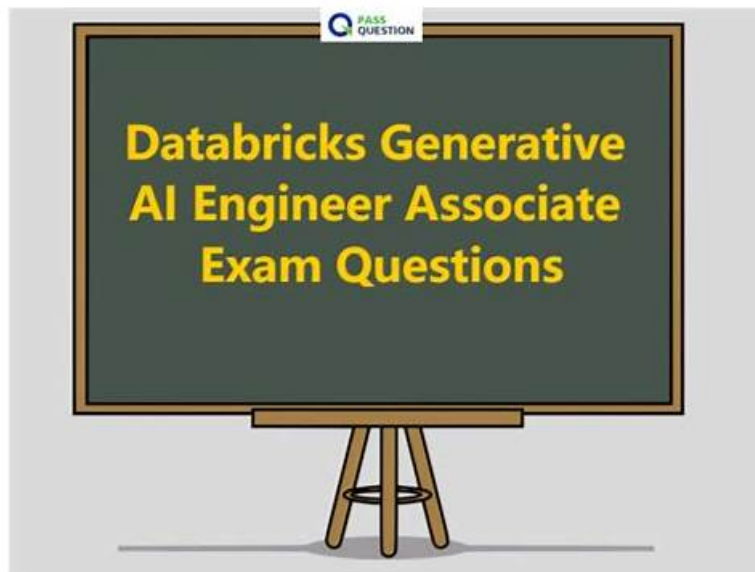


Latest Databricks-Generative-AI-Engineer-Associate Exam Success - Easy and Guaranteed Databricks-Generative-AI-Engineer-Associate Exam Success



DOWNLOAD the newest ExamsReviews Databricks-Generative-AI-Engineer-Associate PDF dumps from Cloud Storage for free: <https://drive.google.com/open?id=1iffD5sYrqJ2VK5oLIThhALyefDe9XsJ3>

If you also need to take the Databricks-Generative-AI-Engineer-Associate exam and want to get the related certification, you can directly select our study materials. We can promise that our Databricks-Generative-AI-Engineer-Associate study question has a higher quality than other study materials in the market. If you want to keep making progress and transcending yourself, we believe that you will harvest happiness and growth. So if you buy and use the Databricks-Generative-AI-Engineer-Associate test dump from our company, we believe that our study materials will make study more interesting and colorful, and it will be very easy for a lot of people to pass their exam and get the related certification if they choose our Databricks-Generative-AI-Engineer-Associate Test Dump and take it into consideration seriously. Now we are willing to introduce the Databricks-Generative-AI-Engineer-Associate exam reference guide from our company to you in order to let you have a deep understanding of our study materials. We believe that you will benefit a lot from our Databricks-Generative-AI-Engineer-Associate study question.

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

Topic	Details
Topic 1	<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.
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">• Application Development: In this topic, Generative AI Engineers learn about tools needed to extract data, Langchain• similar tools, and assessing responses to identify common issues. Moreover, the topic includes questions about adjusting an LLM's response, LLM guardrails, and the best LLM based on the attributes of the application.

Topic 4	<ul style="list-style-type: none"> • Assembling and Deploying Applications: In this topic, Generative AI Engineers get knowledge about coding a chain using a pyfunc mode, coding a simple chain using langchain, and coding a simple chain according to requirements. Additionally, the topic focuses on basic elements needed to create a RAG application. Lastly, the topic addresses sub-topics about registering the model to Unity Catalog using MLflow.
---------	--

>> Databricks-Generative-AI-Engineer-Associate Exam Success <<

Quiz Databricks-Generative-AI-Engineer-Associate - Pass-Sure Databricks Certified Generative AI Engineer Associate Exam Success

Real Databricks-Generative-AI-Engineer-Associate questions in our PDF document can be viewed at any time from any place using your smartphone, tablet, and laptop. If you are busy and don't have time to sit and study for the Databricks Certified Generative AI Engineer Associate Databricks-Generative-AI-Engineer-Associate test, download and use Databricks Databricks-Generative-AI-Engineer-Associate PDF dumps on the go. To pass the Databricks Databricks-Generative-AI-Engineer-Associate exam, it is recommended that you simply use ExamsReviews Databricks-Generative-AI-Engineer-Associate real dumps for a few days.

Databricks Certified Generative AI Engineer Associate Sample Questions (Q53-Q58):

NEW QUESTION # 53

A Generative AI Engineer has been asked to design an LLM-based application that accomplishes the following business objective: answer employee HR questions using HR PDF documentation.

Which set of high level tasks should the Generative AI Engineer's system perform?

- A. Calculate averaged embeddings for each HR document, compare embeddings to user query to find the best document. Pass the best document with the user query into an LLM with a large context window to generate a response to the employee.
- B. Use an LLM to summarize HR documentation. Provide summaries of documentation and user query into an LLM with a large context window to generate a response to the user.
- C. Split HR documentation into chunks and embed into a vector store. Use the employee question to retrieve best matched chunks of documentation, and use the LLM to generate a response to the employee based upon the documentation retrieved.
- D. Create an interaction matrix of historical employee questions and HR documentation. Use ALS to factorize the matrix and create embeddings. Calculate the embeddings of new queries and use them to find the best HR documentation. Use an LLM to generate a response to the employee question based upon the documentation retrieved.

Answer: C

Explanation:

To design an LLM-based application that can answer employee HR questions using HR PDF documentation, the most effective approach is option D. Here's why:

* **Chunking and Vector Store Embedding:**HR documentation tends to be lengthy, so splitting it into smaller, manageable chunks helps optimize retrieval. These chunks are then embedded into a vector store(a database that stores vector representations of text). Each chunk of text is transformed into an embedding using a transformer-based model, which allows for efficient similarity-based retrieval.

* **Using Vector Search for Retrieval:**When an employee asks a question, the system converts their query into an embedding as well. This embedding is then compared with the embeddings of the document chunks in the vector store. The most semantically similar chunks are retrieved, which ensures that the answer is based on the most relevant parts of the documentation.

* **LLM to Generate a Response:**Once the relevant chunks are retrieved, these chunks are passed into the LLM, which uses them as context to generate a coherent and accurate response to the employee's question.

* **Why Other Options Are Less Suitable:**

* **A (Calculate Averaged Embeddings):** Averaging embeddings might dilute important information. It doesn't provide enough granularity to focus on specific sections of documents.

* **B (Summarize HR Documentation):** Summarization loses the detail necessary for HR-related queries, which are often specific. It would likely miss the mark for more detailed inquiries.

* **C (Interaction Matrix and ALS):** This approach is better suited for recommendation systems and not for HR queries, as it's focused on collaborative filtering rather than text-based retrieval.

Thus, option D is the most effective solution for providing precise and contextual answers based on HR documentation.

NEW QUESTION # 54

A Generative AI Engineer is building a system that will answer questions on currently unfolding news topics. As such, it pulls information from a variety of sources including articles and social media posts. They are concerned about toxic posts on social media causing toxic outputs from their system.

Which guardrail will limit toxic outputs?

- A. Log all LLM system responses and perform a batch toxicity analysis monthly.
- **B. Use only approved social media and news accounts to prevent unexpected toxic data from getting to the LLM.**
- C. Implement rate limiting
- D. Reduce the amount of context items the system will include in consideration for its response.

Answer: B

Explanation:

The system answers questions on unfolding news topics using articles and social media, with a concern about toxic outputs from toxic inputs. A guardrail must limit toxicity in the LLM's responses. Let's evaluate the options.

Option A: Use only approved social media and news accounts to prevent unexpected toxic data from getting to the LLM. Curating input sources (e.g., verified accounts) reduces exposure to toxic content at the data ingestion stage, directly limiting toxic outputs.

This is a proactive guardrail aligned with data quality control.

Databricks Reference: "Control input data quality to mitigate unwanted LLM behavior, such as toxicity" ("Building LLM Applications with Databricks," 2023).

Option B: Implement rate limiting

Rate limiting controls request frequency, not content quality. It prevents overload but doesn't address toxicity in social media inputs or outputs.

Databricks Reference: Rate limiting is for performance, not safety: "Use rate limits to manage compute load" ("Generative AI Cookbook").

Option C: Reduce the amount of context items the system will include in consideration for its response. Reducing context might limit exposure to some toxic items but risks losing relevant information, and it doesn't specifically target toxicity. It's an indirect, imprecise fix.

Databricks Reference: Context reduction is for efficiency, not safety: "Adjust context size based on performance needs" ("Databricks Generative AI Engineer Guide").

Option D: Log all LLM system responses and perform a batch toxicity analysis monthly. Logging and analyzing responses is reactive, identifying toxicity after it occurs rather than preventing it. Monthly analysis doesn't limit real-time toxic outputs.

Databricks Reference: Monitoring is for auditing, not prevention: "Log outputs for post-hoc analysis, but use input filters for safety" ("Building LLM-Powered Applications").

Conclusion: Option A is the most effective guardrail, proactively filtering toxic inputs from unverified sources, which aligns with Databricks' emphasis on data quality as a primary safety mechanism for LLM systems.

NEW QUESTION # 55

A Generative AI Engineer at a legal firm is designing a RAG system to analyze historical legal cases. The system needs to process millions of court opinions and legal documents, already organized by time and topic, to track how interpretations of specific laws have evolved over time. All of these documents are in plain-text. The engineer needs to choose a chunking method that would most effectively preserve continuity and the temporal nature of the cases. Which method do they choose?

- **A. Implement windowed summarization with overlapping chunks.**
- B. Implement paragraph level embeddings with each chunk.
- C. Implement a hierarchical tree structure, like RAPTOR, to group similar legal concepts.
- D. Implement sentence level embeddings with each chunk tagged with the time to enable metadata filtering.

Answer: A

Explanation:

In the context of legal document analysis where the "evolution of interpretation" is the primary goal, preserving narrative continuity is paramount. Windowed summarization with overlapping chunks is the most effective method for this use case. Overlapping (e.g., 10-15% of the chunk size) ensures that sentences or concepts split at the boundary of one chunk are preserved in the next, preventing the loss of critical context that often occurs in legal jargon. Furthermore, windowed summarization allows the system to condense long-form court opinions into manageable parts while maintaining the chronological "thread" of the argument. While sentence-level

embeddings with metadata (D) are useful for filtering, they often lack the sufficient context required to understand the nuances of a legal ruling. A windowed approach provides the LLM with enough surrounding text to understand the "why" behind a legal evolution, rather than just the "when."

NEW QUESTION # 56

After changing the response generating LLM in a RAG pipeline from GPT-4 to a model with a shorter context length that the company self-hosts, the Generative AI Engineer is getting the following error:

□ What TWO solutions should the Generative AI Engineer implement without changing the response generating model? (Choose two.)

- A. Use a smaller embedding model to generate
- B. Retrain the response generating model using ALiBi
- C. Decrease the chunk size of embedded documents
- D. Reduce the maximum output tokens of the new model
- E. Reduce the number of records retrieved from the vector database

Answer: C,E

Explanation:

Problem Context: After switching to a model with a shorter context length, the error message indicating that the prompt token count has exceeded the limit suggests that the input to the model is too large.

Explanation of Options:

Option A: Use a smaller embedding model to generate - This wouldn't necessarily address the issue of prompt size exceeding the model's token limit.

Option B: Reduce the maximum output tokens of the new model - This option affects the output length, not the size of the input being too large.

Option C: Decrease the chunk size of embedded documents - This would help reduce the size of each document chunk fed into the model, ensuring that the input remains within the model's context length limitations.

Option D: Reduce the number of records retrieved from the vector database - By retrieving fewer records, the total input size to the model can be managed more effectively, keeping it within the allowable token limits.

Option E: Retrain the response generating model using ALiBi - Retraining the model is contrary to the stipulation not to change the response generating model.

Options C and D are the most effective solutions to manage the model's shorter context length without changing the model itself, by adjusting the input size both in terms of individual document size and total documents retrieved.

NEW QUESTION # 57

A Generative AI Engineer has created a RAG application to look up answers to questions about a series of fantasy novels that are being asked on the author's web forum. The fantasy novel texts are chunked and embedded into a vector store with metadata (page number, chapter number, book title), retrieved with the user's query, and provided to an LLM for response generation. The Generative AI Engineer used their intuition to pick the chunking strategy and associated configurations but now wants to more methodically choose the best values.

Which TWO strategies should the Generative AI Engineer take to optimize their chunking strategy and parameters? (Choose two.)

- A. Choose an appropriate evaluation metric (such as recall or NDCG) and experiment with changes in the chunking strategy, such as splitting chunks by paragraphs or chapters. Choose the strategy that gives the best performance metric.
- B. Create an LLM-as-a-judge metric to evaluate how well previous questions are answered by the most appropriate chunk. Optimize the chunking parameters based upon the values of the metric.
- C. Change embedding models and compare performance.
- D. Add a classifier for user queries that predicts which book will best contain the answer. Use this to filter retrieval.
- E. Pass known questions and best answers to an LLM and instruct the LLM to provide the best token count. Use a summary statistic (mean, median, etc.) of the best token counts to choose chunk size.

Answer: A,B

NEW QUESTION # 58

.....

