

2026 The Best Databricks-Generative-AI-Engineer-Associate–100% Free Reliable Braindumps Questions | Reliable Databricks-Generative-AI-Engineer-Associate Exam Materials



BTW, DOWNLOAD part of ValidTorrent Databricks-Generative-AI-Engineer-Associate dumps from Cloud Storage:
<https://drive.google.com/open?id=1OgJmSjBE9macvrt6raWurtCCcV4fFe3L>

If you can have the certification, you can enter the company you like as well as improve your salary. Databricks-Generative-AI-Engineer-Associate training materials of us can offer you such opportunity, since we have a professional team to compile and verify, therefore Databricks-Generative-AI-Engineer-Associate exam materials are high quality. You can pass the exam just one time. In addition, Databricks-Generative-AI-Engineer-Associate Exam Dumps contain both questions and answers, so that you can have a quick check after practicing. We offer you free update for one year, and the update version for Databricks-Generative-AI-Engineer-Associate exam materials will be sent to your email address automatically.

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">• 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.

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"> • 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.

>> **Databricks-Generative-AI-Engineer-Associate Reliable Braindumps Questions** <<

100% Pass Quiz 2026 Databricks Databricks-Generative-AI-Engineer-Associate – Valid Reliable Braindumps Questions

Are you looking for valid IT exam materials or study guide? You can try our free Databricks Databricks-Generative-AI-Engineer-Associate new exam collection materials. We offer free demo download for our PDF version. You can know several questions of the real test. It can make you master fundamental knowledge quickly. Our Databricks-Generative-AI-Engineer-Associate new exam collection materials are authorized legal products. Our accuracy is nearly 100% pass which will help you clear exam.

Databricks Certified Generative AI Engineer Associate Sample Questions (Q48-Q53):

NEW QUESTION # 48

A team wants to serve a code generation model as an assistant for their software developers. It should support multiple programming languages. Quality is the primary objective.

Which of the Databricks Foundation Model APIs, or models available in the Marketplace, would be the best fit?

- A. Llama2-70b
- B. MPT-7b
- C. BGE-large
- **D. CodeLlama-34B**

Answer: D

Explanation:

For a code generation model that supports multiple programming languages and where quality is the primary objective, CodeLlama-34B is the most suitable choice. Here's the reasoning:

* Specialization in Code Generation: CodeLlama-34B is specifically designed for code generation tasks.

This model has been trained with a focus on understanding and generating code, which makes it particularly adept at handling various programming languages and coding contexts.

* Capacity and Performance: The "34B" indicates a model size of 34 billion parameters, suggesting a high capacity for handling complex tasks and generating high-quality outputs. The large model size typically correlates with better understanding and generation capabilities in diverse scenarios.

* Suitability for Development Teams: Given that the model is optimized for code, it will be able to assist software developers more effectively than general-purpose models. It understands coding syntax, semantics, and the nuances of different programming languages.

* Why Other Options Are Less Suitable:

* A (Llama2-70b): While also a large model, it's more general-purpose and may not be as fine-tuned for code generation as CodeLlama.

* B (BGE-large): This model may not specifically focus on code generation.

* C (MPT-7b): Smaller than CodeLlama-34B and likely less capable in handling complex code generation tasks at high quality.

Therefore, for a high-quality, multi-language code generation application, CodeLlama-34B (option D) is the best fit.

NEW QUESTION # 49

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. 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.
- **C. 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.**
- D. Add a classifier for user queries that predicts which book will best contain the answer. Use this to filter retrieval.
- E. Change embedding models and compare performance.

Answer: A,C

Explanation:

To optimize a chunking strategy for a Retrieval-Augmented Generation (RAG) application, the Generative AI Engineer needs a structured approach to evaluating the chunking strategy, ensuring that the chosen configuration retrieves the most relevant information and leads to accurate and coherent LLM responses.

Here's why C and E are the correct strategies:

Strategy C: Evaluation Metrics (Recall, NDCG)

* Define an evaluation metric: Common evaluation metrics such as recall, precision, or NDCG (Normalized Discounted Cumulative Gain) measure how well the retrieved chunks match the user's query and the expected response.

* Recall measures the proportion of relevant information retrieved.

* NDCG is often used when you want to account for both the relevance of retrieved chunks and the ranking or order in which they are retrieved.

* Experiment with chunking strategies: Adjusting chunking strategies based on text structure (e.g., splitting by paragraph, chapter, or a fixed number of tokens) allows the engineer to experiment with various ways of slicing the text. Some chunks may better align with the user's query than others.

* Evaluate performance: By using recall or NDCG, the engineer can methodically test various chunking strategies to identify which one yields the highest performance. This ensures that the chunking method provides the most relevant information when embedding and retrieving data from the vector store.

Strategy E: LLM-as-a-Judge Metric

* Use the LLM as an evaluator: After retrieving chunks, the LLM can be used to evaluate the quality of answers based on the chunks provided. This could be framed as a "judge" function, where the LLM compares how well a given chunk answers previous user queries.

* Optimize based on the LLM's judgment: By having the LLM assess previous answers and rate their relevance and accuracy, the engineer can collect feedback on how well different chunking configurations perform in real-world scenarios.

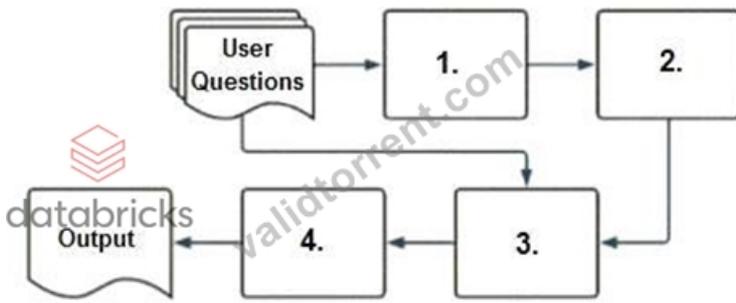
* This metric could be a qualitative judgment on how closely the retrieved information matches the user's intent.

* Tune chunking parameters: Based on the LLM's judgment, the engineer can adjust the chunk size or structure to better align with the LLM's responses, optimizing retrieval for future queries.

By combining these two approaches, the engineer ensures that the chunking strategy is systematically evaluated using both quantitative (recall/NDCG) and qualitative (LLM judgment) methods. This balanced optimization process results in improved retrieval relevance and, consequently, better response generation by the LLM.

NEW QUESTION # 50

A company has a typical RAG-enabled, customer-facing chatbot on its website.



Select the correct sequence of components a user's questions will go through before the final output is returned. Use the diagram above for reference.

- A. 1.response-generating LLM, 2.vector search, 3.context-augmented prompt, 4.embedding model
- **B. 1.embedding model, 2.vector search, 3.context-augmented prompt, 4.response-generating LLM**
- C. 1.response-generating LLM, 2.context-augmented prompt, 3.vector search, 4.embedding model
- D. 1.context-augmented prompt, 2.vector search, 3.embedding model, 4.response-generating LLM

Answer: B

Explanation:

To understand how a typical RAG-enabled customer-facing chatbot processes a user's question, let's go through the correct sequence as depicted in the diagram and explained in option A:

* **Embedding Model (1):**The first step involves the user's question being processed through an embedding model. This model converts the text into a vector format that numerically represents the text. This step is essential for allowing the subsequent vector search to operate effectively.

* **Vector Search (2):**The vectors generated by the embedding model are then used in a vector search mechanism. This search identifies the most relevant documents or previously answered questions that are stored in a vector format in a database.

* **Context-Augmented Prompt (3):**The information retrieved from the vector search is used to create a context-augmented prompt. This step involves enhancing the basic user query with additional relevant information gathered to ensure the generated response is as accurate and informative as possible.

* **Response-Generating LLM (4):**Finally, the context-augmented prompt is fed into a response-generating large language model (LLM). This LLM uses the prompt to generate a coherent and contextually appropriate answer, which is then delivered as the final output to the user.

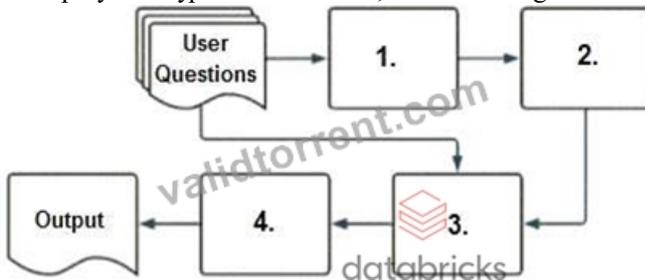
Why Other Options Are Less Suitable:

* B, C, D: These options suggest incorrect sequences that do not align with how a RAG system typically processes queries. They misplace the role of embedding models, vector search, and response generation in an order that would not facilitate effective information retrieval and response generation.

Thus, the correct sequence is embedding model, vector search, context-augmented prompt, response-generating LLM, which is option A.

NEW QUESTION # 51

A company has a typical RAG-enabled, customer-facing chatbot on its website.



Select the correct sequence of components a user's questions will go through before the final output is returned. Use the diagram above for reference.

- A. 1.response-generating LLM, 2.vector search, 3.context-augmented prompt, 4.embedding model
- **B. 1.embedding model, 2.vector search, 3.context-augmented prompt, 4.response-generating LLM**
- C. 1.response-generating LLM, 2.context-augmented prompt, 3.vector search, 4.embedding model
- D. 1.context-augmented prompt, 2.vector search, 3.embedding model, 4.response-generating LLM

Answer: B

Explanation:

To understand how a typical RAG-enabled customer-facing chatbot processes a user's question, let's go through the correct sequence as depicted in the diagram and explained in option A:

* **Embedding Model (1)**:The first step involves the user's question being processed through an embedding model. This model converts the text into a vector format that numerically represents the text. This step is essential for allowing the subsequent vector search to operate effectively.

* **Vector Search (2)**:The vectors generated by the embedding model are then used in a vector search mechanism. This search identifies the most relevant documents or previously answered questions that are stored in a vector format in a database.

* **Context-Augmented Prompt (3)**:The information retrieved from the vector search is used to create a context-augmented prompt. This step involves enhancing the basic user query with additional relevant information gathered to ensure the generated response is as accurate and informative as possible.

* **Response-Generating LLM (4)**:Finally, the context-augmented prompt is fed into a response- generating large language model (LLM). This LLM uses the prompt to generate a coherent and contextually appropriate answer, which is then delivered as the final output to the user.

Why Other Options Are Less Suitable:

* B, C, D: These options suggest incorrect sequences that do not align with how a RAG system typically processes queries. They misplace the role of embedding models, vector search, and response generation in an order that would not facilitate effective information retrieval and response generation.

Thus, the correct sequence is embedding model, vector search, context-augmented prompt, response- generating LLM, which is option A.

NEW QUESTION # 52

A Generative AI Engineer wants to build an LLM-based solution to help a restaurant improve its online customer experience with bookings by automatically handling common customer inquiries. The goal of the solution is to minimize escalations to human intervention and phone calls while maintaining a personalized interaction. To design the solution, the Generative AI Engineer needs to define the input data to the LLM and the task it should perform.

Which input/output pair will support their goal?

- A. Input: Online chat logs; Output: Cancellation options
- **B. Input: Online chat logs; Output: Buttons that represent choices for booking details**
- C. Input: Customer reviews; Output: Classify review sentiment
- D. Input: Online chat logs; Output: Group the chat logs by users, followed by summarizing each user's interactions

Answer: B

Explanation:

Context: The goal is to improve the online customer experience in a restaurant by handling common inquiries about bookings, minimizing escalations, and maintaining personalized interactions.

Explanation of Options:

* Option A: Grouping and summarizing chat logs by user could provide insights into customer interactions but does not directly address the task of handling booking inquiries or minimizing escalations.

* Option B: Using chat logs to generate interactive buttons for booking details directly supports the goal of facilitating online bookings, minimizing the need for human intervention by providing clear, interactive options for customers to self-serve.

* Option C: Classifying sentiment of customer reviews does not directly help with booking inquiries, although it might provide valuable feedback insights.

* Option D: Providing cancellation options is helpful but narrowly focuses on one aspect of the booking process and doesn't support the broader goal of handling common inquiries about bookings.

Option B best supports the goal of improving online interactions by using chat logs to generate actionable items for customers, helping them complete booking tasks efficiently and reducing the need for human intervention.

NEW QUESTION # 53

.....

It is heartening to announce that all ValidTorrent users will be allowed to capitalize on a free Databricks Databricks-Generative-AI-Engineer-Associate exam questions demo of all three formats of the Databricks Databricks-Generative-AI-Engineer-Associate practice test. It will make them scrutinize how our formats work and what we offer them, for example, the form and pattern of Databricks Databricks-Generative-AI-Engineer-Associate Exam Dumps, and their relevant and updated answers. It is convenient for our consumers to check ValidTorrent Databricks Databricks-Generative-AI-Engineer-Associate exam questions free of charge.

