

# Exam Questions Databricks-Generative-AI-Engineer-Associate

Databricks Certified Generative AI Engineer Associate

<https://www.2passeasy.com/dumps/Databricks-Generative-AI-Engineer-Associate/>



#### NEW QUESTION 1

A Generative AI Engineer wants their (inertuned LLMs in their prod Databricks workspace available for testing in their dev workspace as well. All of their workspaces are Unity Catalog enabled and they are currently logging their models into the Model Registry in MLflow. What is the most cost-effective and secure option for the Generative AI Engineer to accomplish their goal?

- A. Use an external model registry which can be accessed from all workspaces
- B. Setup a script to export the model from prod and import it to dev.
- C. Setup a duplicate training pipeline in dev, so that an identical model is available in dev.
- D. Use MLflow to log the model directly into Unity Catalog, and enable READ access in the dev workspace to the model.

**Answer: D**

#### NEW QUESTION 2

A Generative AI Engineer at an automotive company would like to build a question- answering chatbot for customers to inquire about their vehicles. They have a database containing various documents of different vehicle makes, their hardware parts, and common maintenance information. Which of the following components will NOT be useful in building such a chatbot?

- A. Response-generating LLM
- B. Invite users to submit long, rather than concise, questions
- C. Vector database
- D. Embedding model

**Answer: B**

#### NEW QUESTION 3

A Generative AI Engineer is developing a RAG system for their company to perform internal document Q&A for structured HR policies, but the answers returned are frequently incomplete and unstructured. It seems that the retriever is not returning all relevant context. The Generative AI Engineer has experimented with different embedding and response generating LLMs but that did not improve results. Which TWO options could be used to improve the response quality? Choose 2 answers.

- A. Add the section header as a prefix to chunks
- B. Increase the document chunk size
- C. Split the document by sentence
- D. Use a larger embedding model
- E. Fine tune the response generation model

**Answer: AB**

#### NEW QUESTION 4

A Generative AI Engineer is building an LLM to generate article summaries in the form of a type of poem, such as a haiku, given the article content. However, the initial output from the LLM does not match the desired tone or style. Which approach will NOT improve the LLM's response to achieve the desired response?

- A. Provide the LLM with a prompt that explicitly instructs it to generate text in the desired tone and style
- B. Use a neutralizer to normalize the tone and style of the underlying documents
- C. Include few-shot examples in the prompt to the LLM
- D. Fine-tune the LLM on a dataset of desired tone and style

**Answer: B**

#### NEW QUESTION 5

A Generative AI Engineer is building a RAG application that answers questions about internal documents for the company SnoPen AI. The source documents may contain a significant amount of irrelevant content, such as advertisements, sports news, or entertainment news, or content about other companies. Which approach is advisable when building a RAG application to achieve this goal of filtering irrelevant information?

- A. Keep all articles because the RAG application needs to understand non-company content to avoid answering questions about them.
- B. Include in the system prompt that any information it sees will be about SnoPen AI, even if no data filtering is performed.
- C. Include in the system prompt that the application is not supposed to answer any questions unrelated to SnoPen AI.
- D. Consolidate all SnoPen AI related documents into a single chunk in the vector database.

**Answer: C**

#### NEW QUESTION 6

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

**Answer: A**

**NEW QUESTION 7**

Generative AI Engineer at an electronics company just deployed a RAG application for customers to ask questions about products that the company carries. However, they received feedback that the RAG response often returns information about an irrelevant product. What can the engineer do to improve the relevance of the RAG??s response?

- A. Assess the quality of the retrieved context
- B. Implement caching for frequently asked questions
- C. Use a different LLM to improve the generated response
- D. Use a different semantic similarity search algorithm

**Answer:** A

**NEW QUESTION 8**

A Generative AI Engineer needs to design an LLM pipeline to conduct multi-stage reasoning that leverages external tools. To be effective at this, the LLM will need to plan and adapt actions while performing complex reasoning tasks. Which approach will do this?

- A. Train the LLM to generate a single, comprehensive response without interacting with any external tools, relying solely on its pre-trained knowledge.
- B. Implement a framework like ReAct which allows the LLM to generate reasoning traces and perform task-specific actions that leverage external tools if necessary.
- C. Encourage the LLM to make multiple API calls in sequence without planning or structuring the calls, allowing the LLM to decide when and how to use external tools spontaneously.
- D. Use a Chain-of-Thought (CoT) prompting technique to guide the LLM through a series of reasoning steps, then manually input the results from external tools for the final answer.

**Answer:** B

**NEW QUESTION 9**

A Generative AI Engineer has created a RAG application which can help employees retrieve answers from an internal knowledge base, such as Confluence pages or Google Drive. The prototype application is now working with some positive feedback from internal company testers. Now the Generative AI Engineer wants to formally evaluate the system??s performance and understand where to focus their efforts to further improve the system. How should the Generative AI Engineer evaluate the system?

- A. Use cosine similarity score to comprehensively evaluate the quality of the final generated answers.
- B. Curate a dataset that can test the retrieval and generation components of the system separately
- C. Use MLflow??s built in evaluation metrics to perform the evaluation on the retrieval and generation components.
- D. Benchmark multiple LLMs with the same data and pick the best LLM for the job.
- E. Use an LLM-as-a-judge to evaluate the quality of the final answers generated.

**Answer:** B

**NEW QUESTION 10**

A Generative AI Engineer is ready to deploy an LLM application written using Foundation Model APIs. They want to follow security best practices for production scenarios. Which authentication method should they choose?

- A. Use an access token belonging to service principals
- B. Use a frequently rotated access token belonging to either a workspace user or a service principal
- C. Use OAuth machine-to-machine authentication
- D. Use an access token belonging to any workspace user

**Answer:** A

**NEW QUESTION 10**

A Generative AI Engineer is using the code below to test setting up a vector store:

```
from databricks.vector_search.client import VectorSearchClient

vsc = VectorSearchClient()

vsc.create_endpoint(
    name="vector_search_test",
    endpoint_type="STANDARD"
)
```

Assuming they intend to use Databricks managed embeddings with the default embedding model, what should be the next logical function call?

- A. vsc.get\_index()
- B. vsc.create\_delta\_sync\_index()
- C. vsc.create\_direct\_access\_index()
- D. vsc.similarity\_search()

**Answer:** B

#### NEW QUESTION 12

A Generative AI Engineer has been asked to build an LLM-based question-answering application. The application should take into account new documents that are frequently published. The engineer wants to build this application with the least cost and least development effort and have it operate at the lowest cost possible.

Which combination of chaining components and configuration meets these requirements?

- A. For the application a prompt, a retriever, and an LLM are required.
- B. The retriever output is inserted into the prompt which is given to the LLM to generate answers.
- C. The LLM needs to be frequently updated with the new documents in order to provide most up-to-date answers.
- D. For the question-answering application, prompt engineering and an LLM are required to generate answers.
- E. For the application a prompt, an agent and a fine-tuned LLM are required.
- F. The agent is used by the LLM to retrieve relevant content that is inserted into the prompt which is given to the LLM to generate answers.

**Answer: A**

#### NEW QUESTION 14

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

**Answer: B**

#### NEW QUESTION 19

A Generative AI Engineer is setting up a Databricks Vector Search that will lookup news articles by topic within 10 days of the date specified. An example query might be "Tell me about monster truck news around January 5th 1992". They want to do this with the least amount of effort.

How can they set up their Vector Search index to support this use case?

- A. Split articles by 10 day blocks and return the block closest to the query.
- B. Include metadata columns for article date and topic to support metadata filtering.
- C. Pass the query directly to the vector search index and return the best articles.
- D. Create separate indexes by topic and add a classifier model to appropriately pick the best index.

**Answer: B**

#### NEW QUESTION 23

A Generative AI Engineer is designing an LLM-powered live sports commentary platform. The platform provides real-time updates and LLM-generated analyses for any users who would like to have live summaries, rather than reading a series of potentially outdated news articles.

Which tool below will give the platform access to real-time data for generating game analyses based on the latest game scores?

- A. DatabricksIQ
- B. Foundation Model APIs
- C. Feature Serving
- D. AutoML

**Answer: C**

#### NEW QUESTION 28

A Generative AI Engineer is developing a chatbot designed to assist users with insurance-related queries. The chatbot is built on a large language model (LLM) and is conversational. However, to maintain the chatbot's focus and to comply with company policy, it must not provide responses to questions about politics.

Instead, when presented with political inquiries, the chatbot should respond with a standard message:

"Sorry, I cannot answer that. I am a chatbot that can only answer questions around insurance."

Which framework type should be implemented to solve this?

- A. Safety Guardrail
- B. Security Guardrail
- C. Contextual Guardrail
- D. Compliance Guardrail

**Answer: A**

#### NEW QUESTION 30

A small and cost-conscious startup in the cancer research field wants to build a RAG application using Foundation Model APIs.

Which strategy would allow the startup to build a good-quality RAG application while being cost-conscious and able to cater to customer needs?

- A. Limit the number of relevant documents available for the RAG application to retrieve from.
- B. Pick a smaller LLM that is domain-specific.
- C. Limit the number of queries a customer can send per day.
- D. Use the largest LLM possible because that gives the best performance for any general queries.

**Answer: B**

**NEW QUESTION 33**

A Generative AI Engineer received the following business requirements for an external chatbot. The chatbot needs to know what types of questions the user asks and routes to appropriate models to answer the questions. For example, the user might ask about upcoming event details. Another user might ask about purchasing tickets for a particular event. What is an ideal workflow for such a chatbot?

- A. The chatbot should only look at previous event information
- B. There should be two different chatbots handling different types of user queries.
- C. The chatbot should be implemented as a multi-step LLM workflow
- D. First, identify the type of question asked, then route the question to the appropriate mode
- E. If it??s an upcoming event question, send the query to a text-to-SQL mode
- F. If it??s about ticket purchasing, the customer should be redirected to a payment platform.
- G. The chatbot should only process payments

**Answer: C**

**NEW QUESTION 35**

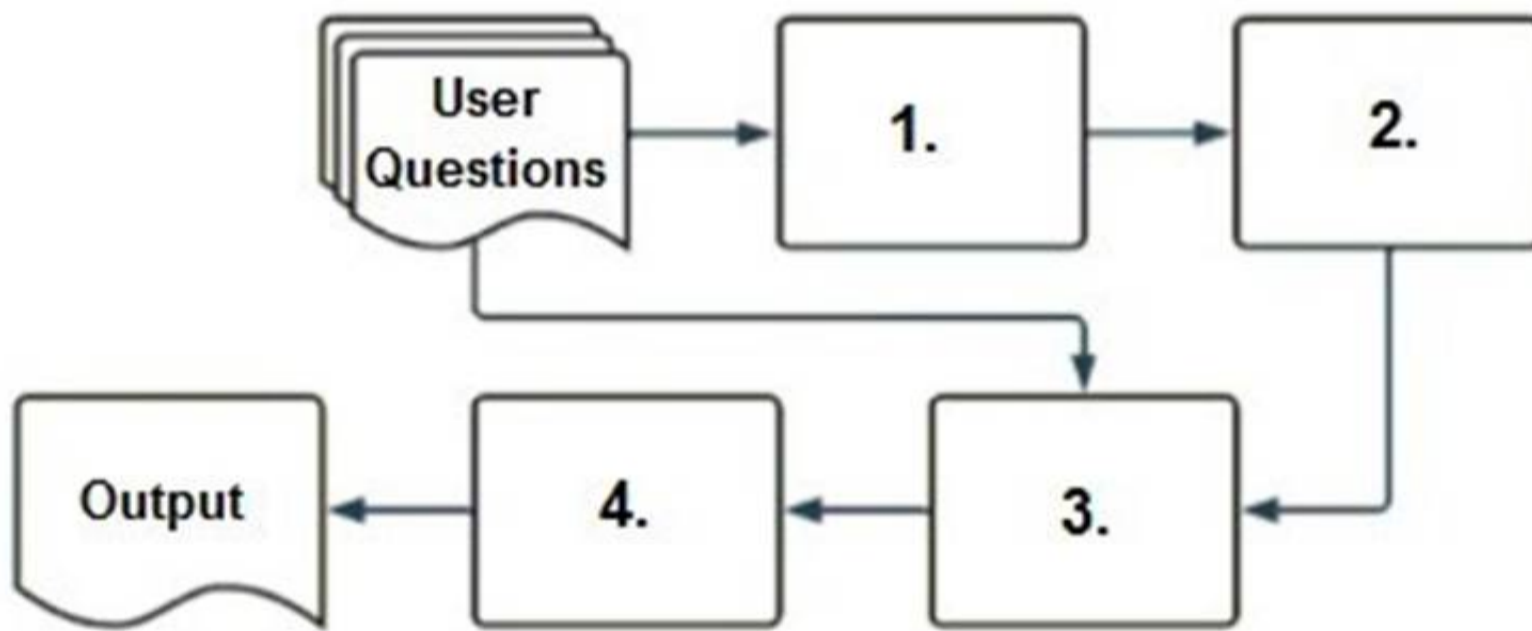
A Generative AI Engineer is building a Generative AI system that suggests the best matched employee team member to newly scoped projects. The team member is selected from a very large team. The match should be based upon project date availability and how well their employee profile matches the project scope. Both the employee profile and project scope are unstructured text. How should the Generative AI Engineer architect their system?

- A. Create a tool for finding available team members given project date
- B. Embed all project scopes into a vector store, perform a retrieval using team member profiles to find the best team member.
- C. Create a tool for finding team member availability given project dates, and another tool that uses an LLM to extract keywords from project scope
- D. Iterate through available team members?? profiles and perform keyword matching to find the best available team member.
- E. Create a tool to find available team members given project date
- F. Create a second tool that can calculate a similarity score for a combination of team member profile and the project scop
- G. Iterate through the team members and rank by best score to select a team member.
- H. Create a tool for finding available team members given project date
- I. Embed team profiles into a vector store and use the project scope and filtering to perform retrieval to find the available best matched team members.

**Answer: D**

**NEW QUESTION 38**

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

**Answer: A**

**NEW QUESTION 40**

A Generative AI Engineer is tasked with deploying an application that takes advantage of a custom MLflow Pyfunc model to return some interim results. How should they configure the endpoint to pass the secrets and credentials?

- A. Use spark.conf.set ()
- B. Pass variables using the Databricks Feature Store API
- C. Add credentials using environment variables
- D. Pass the secrets in plain text

**Answer: C**

**NEW QUESTION 41**

When developing an LLM application, it??s crucial to ensure that the data used for training the model complies with licensing requirements to avoid legal risks.

Which action is NOT appropriate to avoid legal risks?

- A. Reach out to the data curators directly before you have started using the trained model to let them know.
- B. Use any available data you personally created which is completely original and you can decide what license to use.
- C. Only use data explicitly labeled with an open license and ensure the license terms are followed.
- D. Reach out to the data curators directly after you have started using the trained model to let them know.

**Answer:** D

#### NEW QUESTION 42

A Generative AI Engineer is deciding between using LSH (Locality Sensitive Hashing) and HNSW (Hierarchical Navigable Small World) for indexing their vector database. Their top priority is semantic accuracy.

Which approach should the Generative AI Engineer use to evaluate these two techniques?

- A. Compare the cosine similarities of the embeddings of returned results against those of a representative sample of test inputs.
- B. Compare the Bilingual Evaluation Understudy (BLEU) scores of returned results for a representative sample of test inputs.
- C. Compare the Recall-Oriented-Understudy for Gisting Evaluation (ROUGE) scores of returned results for a representative sample of test inputs.
- D. Compare the Levenshtein distances of returned results against a representative sample of test inputs.

**Answer:** A

#### NEW QUESTION 46

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. BGE-large
- D. Llama2-70B

**Answer:** C

#### NEW QUESTION 49

A Generative AI Engineer is building a RAG application that will rely on context retrieved from source documents that are currently in PDF format. These PDFs can contain both text and images. They want to develop a solution using the least amount of lines of code.

Which Python package should be used to extract the text from the source documents?

- A. flask
- B. beautifulsoup
- C. unstructured
- D. numpy

**Answer:** B

#### NEW QUESTION 53

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

**Answer:** CE

#### NEW QUESTION 58

A Generative AI Engineer has a provisioned throughput model serving endpoint as part of a RAG application and would like to monitor the serving endpoint's incoming requests and outgoing responses. The current approach is to include a micro-service in between the endpoint and the user interface to write logs to a remote server.

Which Databricks feature should they use instead which will perform the same task?

- A. Vector Search
- B. Lakeview
- C. DBSQL
- D. Inference Tables

Answer: D

**NEW QUESTION 63**

A Generative AI Engineer is developing a RAG application and would like to experiment with different embedding models to improve the application performance. Which strategy for picking an embedding model should they choose?

- A. Pick an embedding model trained on related domain knowledge
- B. Pick the most recent and most performant open LLM released at the time
- C. pick the embedding model ranked highest on the Massive Text Embedding Benchmark (MTEB) leaderboard hosted by HuggingFace
- D. Pick an embedding model with multilingual support to support potential multilingual user questions

Answer: A

**NEW QUESTION 68**

A Generative AI Engineer just deployed an LLM application at a digital marketing company that assists with answering customer service inquiries. Which metric should they monitor for their customer service LLM application in production?

- A. Number of customer inquiries processed per unit of time
- B. Energy usage per query
- C. Final perplexity scores for the training of the model
- D. HuggingFace Leaderboard values for the base LLM

Answer: A

**NEW QUESTION 70**

.....

## THANKS FOR TRYING THE DEMO OF OUR PRODUCT

Visit Our Site to Purchase the Full Set of Actual Databricks-Generative-AI-Engineer-Associate Exam Questions With Answers.

We Also Provide Practice Exam Software That Simulates Real Exam Environment And Has Many Self-Assessment Features. Order the Databricks-Generative-AI-Engineer-Associate Product From:

<https://www.2passeasy.com/dumps/Databricks-Generative-AI-Engineer-Associate/>

## Money Back Guarantee

### **Databricks-Generative-AI-Engineer-Associate Practice Exam Features:**

- \* Databricks-Generative-AI-Engineer-Associate Questions and Answers Updated Frequently
- \* Databricks-Generative-AI-Engineer-Associate Practice Questions Verified by Expert Senior Certified Staff
- \* Databricks-Generative-AI-Engineer-Associate Most Realistic Questions that Guarantee you a Pass on Your FirstTry
- \* Databricks-Generative-AI-Engineer-Associate Practice Test Questions in Multiple Choice Formats and Updatesfor 1 Year