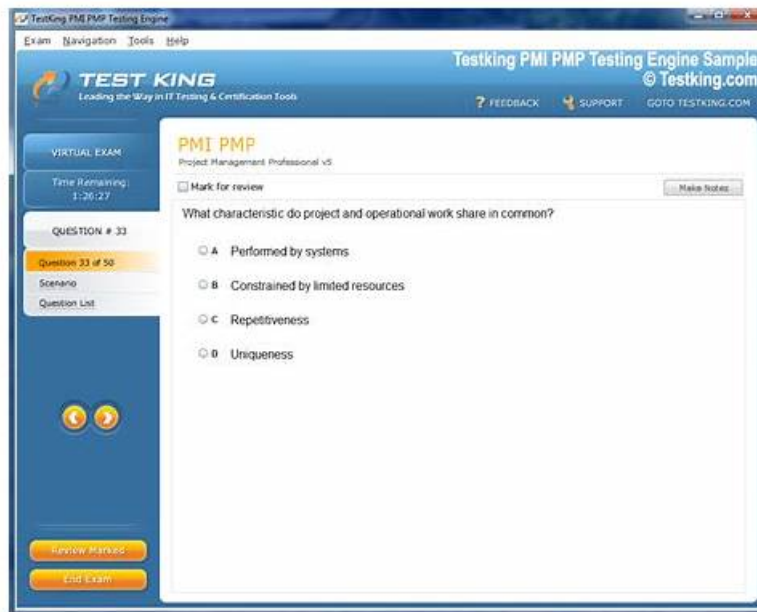


Test Amazon AIP-C01 King & Dumps AIP-C01 Vce



AIP-C01 Exam is an Amazon certification exam and IT professionals who have passed some Amazon certification exams are popular in the IT industry. So more and more people participate in the AIP-C01 certification exam, but the AIP-C01 certification exam is not very simple. If you do not have participated in a professional specialized training course, you need to spend a lot of time and effort to prepare for the exam. But now Prep4SureReview can help you save a lot of your precious time and energy.

The Prep4SureReview Amazon AIP-C01 exam dumps are being offered in three different formats. The names of these formats are AIP-C01 PDF questions file, desktop practice test software, and web-based practice test software. All these three AWS Certified Generative AI Developer - Professional exam dumps formats contain the real Amazon AIP-C01 Exam Questions that will help you to streamline the AIP-C01 exam preparation process.

>> Test Amazon AIP-C01 King <<

Dumps AIP-C01 Vce - AIP-C01 Reliable Test Experience

You can trust top-notch AWS Certified Generative AI Developer - Professional (AIP-C01) exam questions and start preparation with complete peace of mind and satisfaction. The AIP-C01 exam questions are real, valid, and verified by Amazon AIP-C01 certification exam trainers. They work together and put all their efforts to ensure the top standard and relevancy of AIP-C01 Exam Dumps all the time. So we can say that with Amazon AIP-C01 exam questions you will get everything that you need to make the AIP-C01 exam preparation simple, smart, and successful.

Amazon AIP-C01 Exam Syllabus Topics:

Topic	Details
Topic 1	<ul style="list-style-type: none">• Implementation and Integration: This domain focuses on building agentic AI systems, deploying foundation models, integrating GenAI with enterprise systems, implementing FM APIs, and developing applications using AWS tools.
Topic 2	<ul style="list-style-type: none">• AI Safety, Security, and Governance: This domain addresses input• output safety controls, data security and privacy protections, compliance mechanisms, and responsible AI principles including transparency and fairness.
Topic 3	<ul style="list-style-type: none">• Testing, Validation, and Troubleshooting: This domain covers evaluating foundation model outputs, implementing quality assurance processes, and troubleshooting GenAI-specific issues including prompts, integrations, and retrieval systems.

Topic 4	<ul style="list-style-type: none"> Operational Efficiency and Optimization for GenAI Applications: This domain encompasses cost optimization strategies, performance tuning for latency and throughput, and implementing comprehensive monitoring systems for GenAI applications.
Topic 5	<ul style="list-style-type: none"> Foundation Model Integration, Data Management, and Compliance: This domain covers designing GenAI architectures, selecting and configuring foundation models, building data pipelines and vector stores, implementing retrieval mechanisms, and establishing prompt engineering governance.

Amazon AWS Certified Generative AI Developer - Professional Sample Questions (Q97-Q102):

NEW QUESTION # 97

A healthcare company is using Amazon Bedrock to develop a real-time patient care AI assistant to respond to queries for separate departments that handle clinical inquiries, insurance verification, appointment scheduling, and insurance claims. The company wants to use a multi-agent architecture.

The company must ensure that the AI assistant is scalable and can onboard new features for patients. The AI assistant must be able to handle thousands of parallel patient interactions. The company must ensure that patients receive appropriate domain-specific responses to queries.

Which solution will meet these requirements?

- A. Isolate data for each agent by using separate knowledge bases. Use IAM filtering to control access to each knowledge base. Deploy a supervisor agent to perform natural language intent classification on patient inquiries. Configure the supervisor agent to route queries to specialized collaborator agents to respond to department-specific queries. Configure each specialized collaborator agent to use Retrieval Augmented Generation (RAG) with the agent's department-specific knowledge base.
- B. Implement multiple independent supervisor agents that run in parallel to respond to patient inquiries for each department. Configure multiple collaborator agents for each supervisor agent. Integrate all agents with the same knowledge base. Use external routing logic to merge responses from multiple supervisor agents.
- C. Create a separate supervisor agent for each department. Configure individual collaborator agents to perform natural language intent classification for each specialty domain within each department. Integrate each collaborator agent with department-specific knowledge bases only. Implement manual handoff processes between the supervisor agents.
- D. Isolate data for each department in separate knowledge bases. Use IAM filtering to control access to each knowledge base. Deploy a single general-purpose agent. Configure multiple action groups within the general-purpose agent to perform specific department functions. Implement rule-based routing logic within the general-purpose agent instructions.

Answer: A

Explanation:

Option A is the most appropriate design because it provides scalable multi-agent orchestration, clear domain separation, and strong governance with minimal operational complexity. A supervisor-agent pattern is a standard AWS-recommended approach for multi-agent systems: one agent performs intent classification and routing, while specialized agents handle domain-specific tasks.

Isolating data with separate knowledge bases ensures that each specialized collaborator agent retrieves only the information relevant to its department. This improves response accuracy, reduces hallucinations, and supports privacy controls because clinical content, claims content, and scheduling content can have different access policies. IAM-based filtering ensures that each agent has permission only to the knowledge base it is authorized to use.

Routing patient inquiries through a supervisor agent supports high concurrency and extensibility. New departments or features can be added by introducing new collaborator agents and knowledge bases without redesigning the entire system. Because routing is handled centrally, changes in classification logic do not require updates across many independent supervisors.

Using RAG within each collaborator agent ensures that responses are grounded in department-approved information sources, which is critical in healthcare settings to reduce unsafe or incorrect guidance. This approach also improves performance because each retrieval scope is smaller and more relevant, supporting thousands of parallel interactions.

Option B introduces manual handoffs that do not scale. Option C relies on rule-based routing inside one general agent, which becomes brittle and difficult to govern as complexity grows. Option D mixes all departments into a single knowledge base and merges responses externally, increasing risk of incorrect domain answers and operational overhead.

Therefore, Option A best meets the scalability, correctness, and multi-agent onboarding requirements.

NEW QUESTION # 98

A company is building a serverless application that uses AWS Lambda functions to help students around the world summarize notes. The application uses Anthropic Claude through Amazon Bedrock. The company observes that most of the traffic occurs during evenings in each time zone. Users report experiencing throttling errors during peak usage times in their time zones. The company needs to resolve the throttling issues by ensuring continuous operation of the application. The solution must maintain application performance quality and must not require a fixed hourly cost during low traffic periods. Which solution will meet these requirements?

- A. Create custom Amazon CloudWatch metrics to monitor model errors. Set provisioned throughput to a value that is safely higher than the peak traffic observed.
- **B. Enable invocation logging in Amazon Bedrock. Monitor key metrics such as Invocations, InputTokenCount, OutputTokenCount, and InvocationThrottles. Distribute traffic across cross-Region inference endpoints.**
- C. Create custom Amazon CloudWatch metrics to monitor model errors. Set up a failover mechanism to redirect invocations to a backup AWS Region when the errors exceed a specified threshold.
- D. Enable invocation logging in Amazon Bedrock. Monitor InvocationLatency, InvocationClientErrors, and InvocationServerErrors metrics. Distribute traffic across multiple versions of the same model.

Answer: B

Explanation:

Option C is the correct solution because it resolves throttling while preserving performance and avoiding fixed costs during low-traffic periods. Amazon Bedrock supports on-demand inference with usage-based pricing, making it well suited for applications with time-zone-dependent traffic spikes.

Throttling during peak hours typically occurs when inference requests exceed available regional capacity.

Cross-Region inference allows Amazon Bedrock to automatically distribute requests across multiple AWS Regions, reducing contention and preventing throttling without requiring reserved or provisioned capacity.

This approach ensures continuous operation while maintaining low latency for users in different geographic locations.

Invocation logging and native metrics such as InvocationThrottles, InputTokenCount, and OutputTokenCount provide visibility into usage patterns and capacity constraints. Monitoring these metrics enables teams to validate that traffic distribution is working as intended and that performance remains consistent during peak periods.

Option A introduces fixed hourly costs by relying on provisioned throughput, which directly violates the requirement to avoid unnecessary spend during low-traffic periods. Option B introduces regional failover complexity and reactive behavior instead of proactive load distribution. Option D does not address the root cause of throttling, as distributing traffic across model versions within the same Region does not increase available capacity.

Therefore, Option C best aligns with AWS Generative AI best practices for scalable, cost-efficient, global serverless applications.

NEW QUESTION # 99

A company uses an AI assistant application to summarize the company's website content and provide information to customers. The company plans to use Amazon Bedrock to give the application access to a foundation model (FM).

The company needs to deploy the AI assistant application to a development environment and a production environment. The solution must integrate the environments with the FM. The company wants to test the effectiveness of various FMs in each environment. The solution must provide product owners with the ability to easily switch between FMs for testing purposes in each environment.

Which solution will meet these requirements?

- A. Create one AWS CDK application. Create multiple pipelines in AWS CodePipeline. Configure each pipeline to have its own settings for each FM. Configure the application to invoke the Amazon Bedrock FMs by using the `aws_bedrock.ProvisionedModel.fromProvisionedModelArn()` method.
- B. Create one AWS CDK application for the production environment. Configure the application to invoke the Amazon Bedrock FMs by using the `aws_bedrock.ProvisionedModel.fromProvisionedModelArn()` method. Create a pipeline in AWS CodePipeline. Configure the pipeline to deploy to the production environment by using an AWS CodeBuild deploy action. For the development environment, manually recreate the resources by referring to the production application code.
- **C. Create one AWS CDK application. Configure the application to invoke the Amazon Bedrock FMs by using the `aws_bedrock.FoundationModel.fromFoundationModelId()` method. Create a pipeline in AWS CodePipeline that has a deployment stage for each environment that uses AWS CodeBuild deploy actions.**
- D. Create a separate AWS CDK application for each environment. Configure the applications to invoke the Amazon Bedrock FMs by using the `aws_bedrock.FoundationModel.fromFoundationModelId()` method. Create a separate pipeline in AWS CodePipeline for each environment.

Answer: C

Explanation:

Option C best satisfies the requirement for flexible FM testing across environments while minimizing operational complexity and aligning with AWS-recommended deployment practices. Amazon Bedrock supports invoking on-demand foundation models through the FoundationModel abstraction, which allows applications to dynamically reference different models without requiring dedicated provisioned capacity. This is ideal for experimentation and A/B testing in both development and production environments. Using a single AWS CDK application ensures infrastructure consistency and reduces duplication.

Environment-specific configuration, such as selecting different foundation model IDs, can be externalized through parameters, context variables, or environment-specific configuration files. This allows product owners to easily switch between FMs in each environment without modifying application logic.

A single AWS CodePipeline with distinct deployment stages for development and production is an AWS best practice for multi-environment deployments. It enforces consistent build and deployment steps while still allowing environment-level customization. AWS CodeBuild deploy actions enable automated, repeatable deployments, reducing manual errors and improving governance. Option A increases complexity by introducing multiple pipelines and relies on provisioned models, which are not necessary for FM evaluation and experimentation. Provisioned throughput is better suited for predictable, high-volume production workloads rather than frequent model switching.

Option B creates unnecessary operational overhead by duplicating CDK applications and pipelines, making long-term maintenance more difficult.

Option D directly conflicts with infrastructure-as-code best practices by manually recreating development resources, which increases configuration drift and reduces reliability.

Therefore, Option C provides the most flexible, scalable, and AWS-aligned solution for testing and switching foundation models across development and production environments.

NEW QUESTION # 100

A company has deployed an AI assistant as a React application that uses AWS Amplify, an AWS AppSync GraphQL API, and Amazon Bedrock Knowledge Bases. The application uses the GraphQL API to call the Amazon Bedrock RetrieveAndGenerate API for knowledge base interactions. The company configures an AWS Lambda resolver to use the RequestResponse invocation type.

Application users report frequent timeouts and slow response times. Users report these problems more frequently for complex questions that require longer processing.

The company needs a solution to fix these performance issues and enhance the user experience.

Which solution will meet these requirements?

- A. Increase the timeout value of the Lambda resolver. Implement retry logic with exponential backoff.
- B. Update the application to send an API request to an Amazon SQS queue. Update the AWS AppSync resolver to poll and process the queue.
- C. Use AWS Amplify AI Kit to implement streaming responses from the GraphQL API and to optimize client-side rendering.
- D. Change the RetrieveAndGenerate API to the InvokeModelWithResponseStream API. Update the application to use an Amazon API Gateway WebSocket API to support the streaming response.

Answer: C

Explanation:

Option A is the best solution because it directly addresses both observed problems: user-perceived latency and resolver timeouts that occur more frequently for complex prompts. In the current design, an AWS AppSync Lambda resolver is configured with synchronous RequestResponse behavior. That means the client receives nothing until the entire retrieval and generation workflow completes. For longer-running knowledge base queries, this increases the likelihood of hitting request time limits in the synchronous path and creates a poor user experience because the UI appears stalled.

Using AWS Amplify AI Kit to implement streaming responses allows the application to return partial output incrementally as the model produces tokens. This improves perceived responsiveness because users can see the answer forming immediately, even when the full response takes longer. Streaming also reduces the impact of variable model latency and retrieval time because the client no longer waits for a single final payload before rendering content. From a troubleshooting perspective, streaming makes it easier to distinguish "slow generation" from "no response," and it provides faster feedback during testing of complex questions.

Option B is not sufficient because increasing timeouts and adding retries can worsen load and cost while still producing a stalled UI experience. Retries also risk duplicating requests to the knowledge base and can amplify token usage. Option C introduces an awkward polling model for GraphQL interactions and adds significant operational complexity, while not inherently improving interactivity. Option D adds major architectural changes by replacing the knowledge base RetrieveAndGenerate call path with a different streaming invocation API and introducing a WebSocket layer, which is unnecessary when the goal is primarily to fix timeouts and improve UX within the existing AppSync and Amplify design.

Therefore, streaming through Amplify AI Kit is the most effective and lowest-friction improvement.

Thought for 24s

NEW QUESTION # 101

A company provides a service that helps users from around the world discover new restaurants. The service has 50 million monthly active users. The company wants to implement a semantic search solution across a database that contains 20 million restaurants and 200 million reviews. The company currently stores the data in a PostgreSQL database.

The solution must support complex natural language queries and return results for at least 95% of queries within 500 ms. The solution must maintain data freshness for restaurant details that update hourly. The solution must also scale cost-effectively during peak usage periods.

Which solution will meet these requirements with the LEAST development effort?

- A. Keep the restaurant data in PostgreSQL and implement a pgvector extension. Use a foundation model (FM) in Amazon Bedrock to generate vector embeddings from restaurant data. Store the vector embeddings directly in PostgreSQL. Create an AWS Lambda function to convert natural language queries to vector representations by using the same FM. Configure the Lambda function to perform similarity searches within the database.
- B. Migrate the restaurant data to an Amazon Bedrock knowledge base by using a custom ingestion pipeline. Configure the knowledge base to automatically generate embeddings from restaurant information. Use the Amazon Bedrock Retrieve API with built-in vector search capabilities to query the knowledge base directly by using natural language input.
- C. Migrate the restaurant data to Amazon OpenSearch Service. Use a foundation model (FM) in Amazon Bedrock to generate vector embeddings from restaurant descriptions, reviews, and menu items. When users submit natural language queries, convert the queries to embeddings by using the same FM. Perform k-nearest neighbors (k-NN) searches to find semantically similar results.
- D. Migrate the restaurant data to Amazon OpenSearch Service. Implement keyword-based search rules that use custom analyzers and relevance tuning to find restaurants based on attributes such as cuisine type, feature, and location. Create Amazon API Gateway HTTP API endpoints to transform user queries into structured search parameters.

Answer: B

Explanation:

Option D requires the least development effort because it uses a managed retrieval workflow that bundles the most time-consuming parts of semantic search: embedding generation, vector indexing, and natural language retrieval. With an Amazon Bedrock knowledge base, the application does not need to implement and operate separate services to (1) generate embeddings for hundreds of millions of records, (2) store and manage vectors, (3) build query-time embedding conversion logic, and (4) implement k-NN search orchestration.

Instead, the knowledge base is configured to automatically create embeddings during ingestion, and the application queries it using the Amazon Bedrock Retrieve API, which accepts natural language input and performs the vector search as a managed capability. The performance requirement (95% of queries within 500 ms) is best served by a purpose-built vector search backend rather than running similarity search directly inside a transactional PostgreSQL system at this scale.

A knowledge base is designed for retrieval patterns and can be backed by scalable vector stores, which helps meet latency goals under heavy concurrency. The hourly freshness requirement maps naturally to ingestion updates: the pipeline can re-ingest updated restaurant details on a schedule so the knowledge base remains current without building custom re-embedding workflows in application code.

Cost-effective scaling during peak periods is also easier with a managed retrieval layer because scaling the retrieval workload is separated from the operational database. This avoids overprovisioning PostgreSQL for peak semantic-search traffic and reduces the engineering effort to tune performance, sharding, indexing, and retry logic.

Options B and C can work, but they require the team to build and maintain embedding pipelines, query embedding generation, vector index management, and operational scaling strategies. Option A does not provide semantic search because it relies on keyword-based matching rather than embeddings.

NEW QUESTION # 102

.....

Amazon AIP-C01 frequently changes the content of the AWS Certified Generative AI Developer - Professional (AIP-C01) exam. Therefore, to save your valuable time and money, we keep a close eye on the latest updates. Furthermore, Prep4SureReview also offers free updates of AIP-C01 exam questions for up to 365 days after buying AWS Certified Generative AI Developer - Professional (AIP-C01) dumps. We guarantee that nothing will stop you from earning the esteemed Amazon Certification Exam on your first attempt if you diligently prepare with our Amazon in AIP-C01 real exam questions.

Dumps AIP-C01 Vce: <https://www.prep4surereview.com/AIP-C01-latest-braindumps.html>

- AIP-C01 Exam Flashcards AIP-C01 Quiz AIP-C01 Download Demo Download AIP-C01 for free by

