

# Amazon-Web-Services

## Exam Questions AIP-C01

AWS Certified Generative AI Developer - Professional



### NEW QUESTION 1

A company is developing a generative AI (GenAI) application that uses Amazon Bedrock foundation models. The application has several custom tool integrations. The application has experienced unexpected token consumption surges despite consistent user traffic.

The company needs a solution that uses Amazon Bedrock model invocation logging to monitor InputTokenCount and OutputTokenCount metrics. The solution must detect unusual patterns in tool usage and identify which specific tool integrations cause abnormal token consumption. The solution must also automatically adjust thresholds as traffic patterns change.

Which solution will meet these requirements?

- A. Use Amazon CloudWatch Logs to capture model invocation log
- B. Create CloudWatch dashboards for token metric
- C. Configure static CloudWatch alarms with fixed thresholds for each tool integration.
- D. Store model invocation logs in Amazon S3. Use AWS Glue and Amazon Athena to analyze token usage trends.
- E. Use Amazon CloudWatch Logs to capture model invocation log
- F. Create CloudWatch metric filters to extract tool-specific invocation pattern
- G. Apply CloudWatch anomaly detection alarms that automatically adjust baselines for each tool's token metrics.
- H. Store model invocation logs in an Amazon S3 bucket
- I. Use AWS Lambda to process logs in real time
- J. Manually update CloudWatch alarm thresholds based on trends identified by the Lambda function.

**Answer: C**

### NEW QUESTION 2

A company runs a Retrieval Augmented Generation (RAG) application that uses Amazon Bedrock Knowledge Bases to perform regulatory compliance queries. The application uses the RetrieveAndGenerateStream API. The application retrieves relevant documents from a knowledge base that contains more than 50,000 regulatory documents, legal precedents, and policy updates.

The RAG application is producing suboptimal responses because the initial retrieval often returns semantically similar but contextually irrelevant documents. The poor responses are causing model hallucinations and incorrect regulatory guidance. The company needs to improve the performance of the RAG application so it returns more relevant documents.

Which solution will meet this requirement with the LEAST operational overhead?

- A. Deploy an Amazon SageMaker endpoint to run a fine-tuned ranking model
- B. Use an Amazon API Gateway REST API to route request
- C. Configure the application to make requests through the REST API to rerank the results.
- D. Use Amazon Comprehend to classify documents and apply relevance score
- E. Integrate the RAG application's reranking process with Amazon Textract to run document analysis
- F. Use Amazon Neptune to perform graph-based relevance calculations.
- G. Implement a retrieval pipeline that uses the Amazon Bedrock Knowledge Bases Retrieve API to perform initial document retrieval
- H. Call the Amazon Bedrock Rerank API to rerank the result
- I. Invoke the InvokeModelWithResponseStream operation to generate responses.
- J. Use the latest Amazon reranker model through the reranking configuration within Amazon Bedrock Knowledge Base
- K. Use the model to improve document relevance scoring and to reorder results based on contextual assessments.

**Answer: D**

### NEW QUESTION 3

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

**Answer: A**

### NEW QUESTION 4

A publishing company is developing a chat assistant that uses a containerized large language model (LLM) that runs on Amazon SageMaker AI. The architecture consists of an Amazon API Gateway REST API that routes user requests to an AWS Lambda function. The Lambda function invokes a SageMaker AI real-time endpoint that hosts the LLM.

Users report uneven response times. Analytics show that a high number of chats are abandoned after 2 seconds of waiting for the first token. The company wants a solution to ensure that p95 latency is under 800 ms for interactive requests to the chat assistant.

Which combination of solutions will meet this requirement? (Select TWO.)

- A. Enable model preload upon container start
- B. Implement dynamic batching to process multiple user requests together in a single inference pass.
- C. Select a larger GPU instance type for the SageMaker AI endpoint
- D. Set the minimum number of instances to 0. Continue to perform per-request processing
- E. Lazily load model weights on the first request.
- F. Switch to a multi-model endpoint
- G. Use lazy loading without request batching.

- H. Set the minimum number of instances to greater than 0. Enable response streaming.
- I. Switch to Amazon SageMaker Asynchronous Inference for all request
- J. Store requests in an Amazon S3 bucket
- K. Set the minimum number of instances to 0.

**Answer:** AD

#### NEW QUESTION 5

Company configures a landing zone in AWS Control Tower. The company handles sensitive data that must remain within the European Union. The company must use only the eu-central-1 Region. The company uses Service Control Policies (SCPs) to enforce data residency policies. GenAI developers at the company are assigned IAM roles that have full permissions for Amazon Bedrock.

The company must ensure that GenAI developers can use the Amazon Nova Pro model through Amazon Bedrock only by using cross-Region inference (CRI) and only in eu-central-1. The company enables model access for the GenAI developer IAM roles in Amazon Bedrock. However, when a GenAI developer attempts to invoke the model through the Amazon Bedrock Chat/Text playground, the GenAI developer receives the following error:

User arn:aws:sts:123456789012:assumed-role/AssumedDevRole/DevUserName Action: bedrock:InvokeModelWithResponseStream

On resource(s): arn:aws:bedrock:eu-west-3::foundation-model/amazon.nova-pro-v1:0 Context: a service control policy explicitly denies the action

The company needs a solution to resolve the error. The solution must retain the company's existing governance controls and must provide precise access control.

The solution must comply with the company's existing data residency policies.

Which combination of solutions will meet these requirements? (Select TWO.)

- A. Add an AdministratorAccess policy to the GenAI developer IAM role
- B. Extend the existing SCPs to enable CRI for the eu.amazon.nova-pro-v1:0 inference profile
- C. Enable Amazon Bedrock model access for Amazon Nova Pro in the eu-west-3 Region
- D. Validate that the GenAI developer IAM roles have permissions to invoke Amazon Nova Pro through the eu.amazon.nova-pro-v1:0 inference profile on all European Union AWS Regions that can serve the model
- E. Extend the existing SCP to enable CRI for the eu-\* inference profile

**Answer:** BE

#### NEW QUESTION 6

A company has a recommendation system. The system's applications run on Amazon EC2 instances. The applications make API calls to Amazon Bedrock foundation models (FMs) to analyze customer behavior and generate personalized product recommendations.

The system is experiencing intermittent issues. Some recommendations do not match customer preferences. The company needs an observability solution to monitor operational metrics and detect patterns of operational performance degradation compared to established baselines. The solution must also generate alerts with correlation data within 10 minutes when FM behavior deviates from expected patterns.

Which solution will meet these requirements?

- A. Configure Amazon CloudWatch Container Insights for the application infrastructure
- B. Set up CloudWatch alarms for latency threshold
- C. Add custom metrics for token counts by using the CloudWatch embedded metric format
- D. Create CloudWatch dashboards to visualize the data.
- E. Implement AWS X-Ray to trace requests through the application component
- F. Enable CloudWatch Logs Insights for error pattern detection
- G. Set up AWS CloudTrail to monitor all API calls to Amazon Bedrock
- H. Create custom dashboards in Amazon QuickSight.
- I. Enable Amazon CloudWatch Application Insights for the application resource
- J. Create custom metrics for recommendation quality, token usage, and response latency by using the CloudWatch embedded metric format with dimensions for request types and user segment
- K. Configure CloudWatch anomaly detection on the model metric
- L. Establish log pattern analysis by using CloudWatch Logs Insights.
- M. Use Amazon OpenSearch Service with the Observability plugin
- N. Ingest model metrics and logs by using Amazon Kinesis
- O. Create custom Piped Processing Language (PPL) queries to analyze model behavior pattern
- P. Establish operational dashboards to visualize anomalies in real time.

**Answer:** C

#### NEW QUESTION 7

A company upgraded its Amazon Bedrock-powered foundation model (FM) that supports a multilingual customer service assistant. After the upgrade, the assistant exhibited inconsistent behavior across languages. The assistant began generating different responses in some languages when presented with identical questions. The company needs a solution to detect and address similar problems for future updates. The evaluation must be completed within 45 minutes for all supported languages. The evaluation must process at least 15,000 test conversations in parallel. The evaluation process must be fully automated and integrated into the CI/CD pipeline. The solution must block deployment if quality thresholds are not met.

Which solution will meet these requirements?

- A. Create a distributed traffic simulation framework that sends translation-heavy workloads to the assistant in multiple languages simultaneously
- B. Use Amazon CloudWatch metrics to monitor latency, concurrency, and throughput
- C. Run simulations before production releases to identify infrastructure bottlenecks.
- D. Deploy the assistant in multiple AWS Regions with Amazon Route 53 latency-based routing and AWS Global Accelerator to improve global performance
- E. Store multilingual conversation logs in Amazon S3. Perform weekly post-deployment audits to review consistency.
- F. Create a pre-processing pipeline that normalizes all incoming messages into a consistent format before sending the messages to the assistant
- G. Apply rule-based checks to flag potential hallucinations in the output
- H. Focus evaluation on normalized text to simplify testing across languages.
- I. Set up standardized multilingual test conversations with identical meaning
- J. Run the test conversations in parallel by using Amazon Bedrock model evaluation job
- K. Apply similarity and hallucination threshold
- L. Integrate the process into the CI/CD pipeline to block releases that fail.

**Answer:** D

### NEW QUESTION 8

Example Corp provides a personalized video generation service that millions of enterprise customers use. Customers generate marketing videos by submitting prompts to the company's proprietary generative AI (GenAI) model. To improve output relevance and personalization, Example Corp wants to enhance the prompts by using customer-specific context such as product preferences, customer attributes, and business history. The customers have strict data governance requirements. The customers must retain full ownership and control over their own data. The customers do not require real-time access. However, semantic accuracy must be high and retrieval latency must remain low to support customer experience use cases. Example Corp wants to minimize architectural complexity in its integration pattern. Example Corp does not want to deploy and manage services in each customer's environment unless necessary. Which solution will meet these requirements?

- A. Ensure that each customer sets up an Amazon Q Business index that includes the customer's internal data
- B. Ensure that each customer designates Example Corp as a data accessor to allow Example Corp to retrieve relevant content by using a secure API to enrich prompts at runtime.
- C. Use federated search with Model Context Protocol (MCP) by deploying real-time MCP servers for each customer
- D. Retrieve data in real time during prompt generation.
- E. Ensure that each customer configures an Amazon Bedrock knowledge base
- F. Allow cross-account querying so Example Corp can retrieve structured data for prompt augmentation.
- G. Configure Amazon Kendra to crawl customer data source
- H. Share the resulting indexes across accounts so Example Corp can query each customer's Amazon Kendra index to retrieve augmentation data.

**Answer: A**

### NEW QUESTION 9

An ecommerce company operates a global product recommendation system that needs to switch between multiple foundation models (FM) in Amazon Bedrock based on regulations, cost optimization, and performance requirements. The company must apply custom controls based on proprietary business logic, including dynamic cost thresholds, AWS Region-specific compliance rules, and real-time A/B testing across multiple FMs. The system must be able to switch between FMs without deploying new code. The system must route user requests based on complex rules including user tier, transaction value, regulatory zone, and real-time cost metrics that change hourly and require immediate propagation across thousands of concurrent requests. Which solution will meet these requirements?

- A. Deploy an AWS Lambda function that uses environment variables to store routing rules and Amazon Bedrock FM ID
- B. Use the Lambda console to update the environment variables when business requirements change
- C. Configure an Amazon API Gateway REST API to read request parameters to make routing decisions.
- D. Deploy Amazon API Gateway REST API request transformation templates to implement routing logic based on request attribute
- E. Store Amazon Bedrock FM endpoints as REST API stage variable
- F. Update the variables when the system switches between models.
- G. Configure an AWS Lambda function to fetch routing configurations from the AWS AppConfig Agent for each user request
- H. Run business logic in the Lambda function to select the appropriate FM for each request
- I. Expose the FM through a single Amazon API Gateway REST API endpoint.
- J. Use AWS Lambda authorizers for an Amazon API Gateway REST API to evaluate routing rules that are stored in AWS AppConfig
- K. Return authorization contexts based on business logic
- L. Route requests to model-specific Lambda functions for each Amazon Bedrock FM.

**Answer: C**

### NEW QUESTION 10

A specialty coffee company has a mobile app that generates personalized coffee roast profiles by using Amazon Bedrock with a three-stage prompt chain. The prompt chain converts user inputs into structured metadata, retrieves relevant logs for coffee roasts, and generates a personalized roast recommendation for each customer. Users in multiple AWS Regions report inconsistent roast recommendations for identical inputs, slow inference during the retrieval step, and unsafe recommendations such as brewing at excessively high temperatures. The company must improve the stability of outputs for repeated inputs. The company must also improve app performance and the safety of the app's outputs. The updated solution must ensure 99.5% output consistency for identical inputs and achieve inference latency of less than 1 second. The solution must also block unsafe or hallucinated recommendations by using validated safety controls. Which solution will meet these requirements?

- A. Deploy Amazon Bedrock with provisioned throughput to stabilize inference latency
- B. Apply Amazon Bedrock guardrails that have semantic denial rules to block unsafe output
- C. Use Amazon Bedrock Prompt Management to manage prompts by using approval workflows.
- D. Use Amazon Bedrock Agents to manage chains
- E. Log model inputs and outputs to Amazon CloudWatch Log
- F. Use logs from Amazon CloudWatch to perform A/B testing for prompt versions.
- G. Cache prompt results in Amazon ElastiCache
- H. Use AWS Lambda functions to pre-process metadata and to trace end-to-end latency
- I. Use AWS X-Ray to identify and remediate performance bottlenecks.
- J. Use Amazon Kendra to improve roast log retrieval accuracy
- K. Store normalized prompt metadata within Amazon DynamoDB
- L. Use AWS Step Functions to orchestrate multi-step prompts.

**Answer: A**

### NEW QUESTION 10

A company is building a video analysis platform on AWS. The platform will analyze a large video archive by using Amazon Rekognition and Amazon Bedrock. The platform must comply with predefined privacy standards. The platform must also use secure model I/O, control foundation model (FM) access patterns, and provide an audit of who accessed what and when. Which solution will meet these requirements?

- A. Configure VPC endpoints for Amazon Bedrock model API call
- B. Implement Amazon Bedrock guardrails to filter harmful or unauthorized content in prompts and response
- C. Use Amazon Bedrock trace events to track all agent and model invocations for auditing purpose

- D. Export the traces to Amazon CloudWatch Logs as an audit record of model usage
- E. Store all prompts and outputs in Amazon S3 with server-side encryption with AWS KMS keys (SSE-KMS).
- F. Define access control by using IAM with attribute-based access control (ABAC) to map departments to specific permission
- G. Configure VPC endpoints for Amazon Bedrock model API call
- H. Use IAM condition keys to enforce specific GuardrailIdentifier and ModelId value
- I. Configure AWS CloudTrail to capture management and data events for S3 objects and KMS key usage activities
- J. Enable S3 server access logging to record detailed file-level interactions with the video archive
- K. Send all CloudTrail logs to AWS CloudTrail Lake
- L. Set up Amazon CloudWatch alarms to detect and alert on unexpected activity from Amazon Bedrock, Amazon Rekognition, and AWS KMS.
- M. Restrict access to services by using VPC endpoint policies
- N. Use AWS Config to track resource changes and compliance with security rules
- O. Use server-side encryption with AWS KMS keys (SSE-KMS) to encrypt data at rest
- P. Store the model's I/O in separate Amazon S3 buckets
- Q. Enable S3 server access logging to track file-level interactions.
- R. Configure AWS CloudTrail Insights to analyze API call patterns across accounts and detect anomalous activity in Amazon Bedrock, Amazon Rekognition, Amazon S3, and AWS KMS
- S. Deploy Amazon Macie to scan and classify the video archive
- T. Use server-side encryption with AWS KMS keys (SSE-KMS) to encrypt all stored data
- . Configure CloudTrail to capture KMS API usage events for audit purposes
- . Configure Amazon EventBridge rules to process CloudTrail Insights anomalies and Macie findings
- . Use CloudWatch alarms to trigger automated notifications and security responses when potential security issues are detected.

**Answer: B**

#### NEW QUESTION 14

A medical company is building a generative AI (GenAI) application that uses Retrieval Augmented Generation (RAG) to provide evidence-based medical information. The application uses Amazon OpenSearch Service to retrieve vector embeddings. Users report that searches frequently miss results that contain exact medical terms and acronyms and return too many semantically similar but irrelevant documents. The company needs to improve retrieval quality and maintain low end-user latency, even as the document collection grows to millions of documents. Which solution will meet these requirements with the LEAST operational overhead?

- A. Configure hybrid search by combining vector similarity with keyword matching to improve semantic understanding and exact term and acronym matching.
- B. Increase the dimensions of the vector embeddings from 384 to 1536. Use a post-processing AWS Lambda function to filter out irrelevant results after retrieval.
- C. Replace OpenSearch Service with Amazon Kendra
- D. Use query expansion to handle medical acronyms and terminology variants during pre-processing.
- E. Implement a two-stage retrieval architecture in which initial vector search results are re-ranked by an ML model hosted on Amazon SageMaker.

**Answer: A**

#### NEW QUESTION 15

A company is developing a generative AI (GenAI) application that analyzes customer service calls in real time and generates suggested responses for human customer service agents. The application must process 500,000 concurrent calls during peak hours with less than 200 ms end-to-end latency for each suggestion. The company uses existing architecture to transcribe customer call audio streams. The application must not exceed a predefined monthly compute budget and must maintain auto scaling capabilities. Which solution will meet these requirements?

- A. Deploy a large, complex reasoning model on Amazon Bedrock
- B. Purchase provisioned throughput and optimize for batch processing.
- C. Deploy a low-latency, real-time optimized model on Amazon Bedrock
- D. Purchase provisioned throughput and set up automatic scaling policies.
- E. Deploy a large language model (LLM) on an Amazon SageMaker real-time endpoint that uses dedicated GPU instances.
- F. Deploy a mid-sized language model on an Amazon SageMaker serverless endpoint that is optimized for batch processing.

**Answer: B**

#### NEW QUESTION 16

A company is designing a canary deployment strategy for a payment processing API. The system must support automated gradual traffic shifting between multiple Amazon Bedrock models based on real-time inference metrics, historical traffic patterns, and service health. The solution must be able to gradually increase traffic to new model versions. The system must increase traffic if metrics remain healthy and decrease traffic if the performance degrades below acceptable thresholds. The company needs to comprehensively monitor inference latency and error rates during the deployment phase. The company must also be able to halt deployments and revert to a previous model version without any manual intervention. Which solution will meet these requirements?

- A. Use Amazon Bedrock with provisioned throughput to host model version
- B. Configure an Amazon EventBridge rule to invoke an AWS Step Functions workflow when a new model version is released
- C. Configure the workflow to shift traffic in stages, wait for a specified time period, and invoke an AWS Lambda function to check Amazon CloudWatch performance metrics
- D. Configure the workflow to increase traffic if metrics meet thresholds and to trigger a traffic rollback if performance metrics fall below thresholds.
- E. Use AWS Lambda functions to invoke various Amazon Bedrock model versions
- F. Use an Amazon API Gateway HTTP API with stage variables and weighted routing to shift traffic gradually
- G. Use Amazon CloudWatch to monitor performance
- H. Use external logic to adjust traffic and roll back if performance falls below thresholds.
- I. Use Amazon SageMaker AI endpoint variants to represent multiple Amazon Bedrock model versions
- J. Use variant weights to shift traffic
- K. Use Amazon CloudWatch and SageMaker Model Monitor to trigger rollbacks
- L. Use EventBridge to roll back deployments if an anomaly is detected.
- M. Use Amazon OpenSearch Service to track inference logs
- N. Configure OpenSearch Service to invoke an AWS Systems Manager Automation runbook to update Amazon Bedrock model endpoints to shift traffic based on inference logs.

**Answer: A**

#### NEW QUESTION 18

A company uses AWS Lake Formation to set up a data lake that contains databases and tables for multiple business units across multiple AWS Regions. The company wants to use a foundation model (FM) through Amazon Bedrock to perform fraud detection. The FM must ingest sensitive financial data from the data lake. The data includes some customer personally identifiable information (PII).

The company must design an access control solution that prevents PII from appearing in a production environment. The FM must access only authorized data subsets that have PII redacted from specific data columns. The company must capture audit trails for all data access.

Which solution will meet these requirements?

- A. Create a separate dataset in a separate Amazon S3 bucket for each business unit and Region combination
- B. Configure S3 bucket policies to control access based on IAM roles that are assigned to FM training instance
- C. Use S3 access logs to track data access.
- D. Configure the FM to authenticate by using AWS Identity and Access Management roles and Lake Formation permissions based on LF-Tag expression
- E. Define business units and Regions as LF-Tags that are assigned to databases and table
- F. Use AWS CloudTrail to collect comprehensive audit trails of data access.
- G. Use direct IAM principal grants on specific databases and tables in Lake Formation
- H. Create a custom application layer that logs access requests and further filters sensitive columns before sending data to the FM.
- I. Configure the FM to request temporary credentials from AWS Security Token Service
- J. Access the data by using presigned S3 URLs that are generated by an API that applies business unit and Regional filter
- K. Use AWS CloudTrail to collect comprehensive audit trails of data access.

**Answer: B**

#### NEW QUESTION 19

A bank is developing a generative AI (GenAI)-powered AI assistant that uses Amazon Bedrock to assist the bank's website users with account inquiries and financial guidance. The bank must ensure that the AI assistant does not reveal any personally identifiable information (PII) in customer interactions.

The AI assistant must not send PII in prompts to the GenAI model. The AI assistant must not respond to customer requests to provide investment advice. The bank must collect audit logs of all customer interactions, including any images or documents that are transmitted during customer interactions.

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

- A. Use Amazon Macie to detect and redact PII in user inputs and in the model response
- B. Apply prompt engineering techniques to force the model to avoid investment advice topic
- C. Use AWS CloudTrail to capture conversation logs.
- D. Use an AWS Lambda function and Amazon Comprehend to detect and redact PII
- E. Use Amazon Comprehend topic modeling to prevent the AI assistant from discussing investment advice topic
- F. Set up custom metrics in Amazon CloudWatch to capture customer conversations.
- G. Configure Amazon Bedrock guardrails to apply a sensitive information policy to detect and filter PII
- H. Set up a topic policy to ensure that the AI assistant avoids investment advice topic
- I. Use the Converse API to log model invocation
- J. Enable delivery and image logging to Amazon S3.
- K. Use regex controls to match patterns for PII
- L. Apply prompt engineering techniques to avoid returning PII or investment advice topics to customer
- M. Enable model invocation logging, delivery logging, and image logging to Amazon S3.

**Answer: C**

#### NEW QUESTION 21

A financial services company is developing a Retrieval Augmented Generation (RAG) application to help investment analysts query complex financial relationships across multiple investment vehicles, market sectors, and regulatory environments. The dataset contains highly interconnected entities that have multi-hop relationships. Analysts must examine relationships holistically to provide accurate investment guidance. The application must deliver comprehensive answers that capture indirect relationships between financial entities and must respond in less than 3 seconds.

Which solution will meet these requirements with the LEAST operational overhead?

- A. Use Amazon Bedrock Knowledge Bases with GraphRAG and Amazon Neptune Analytics to store financial data
- B. Analyze multi-hop relationships between entities and automatically identify related information across documents.
- C. Use Amazon Bedrock Knowledge Bases and an Amazon OpenSearch Service vector store to implement custom relationship identification logic that uses AWS Lambda to query multiple vector embeddings in sequence.
- D. Use Amazon OpenSearch Serverless vector search with k-nearest neighbor (k-NN). Implement manual relationship mapping in an application layer that runs on Amazon EC2 Auto Scaling.
- E. Use Amazon DynamoDB to store financial data in a custom indexing system
- F. Use AWS Lambda to query relevant records
- G. Use Amazon SageMaker to generate responses.

**Answer: A**

#### NEW QUESTION 23

An insurance company uses existing Amazon SageMaker AI infrastructure to support a web-based application that allows customers to predict what their insurance premiums will be. The company stores customer data that is used to train the SageMaker AI model in an Amazon S3 bucket. The dataset is growing rapidly. The company wants a solution to continuously re-train the model. The solution must automatically re-train and re-deploy the model to the application when an employee uploads a new customer data file to the S3 bucket.

Which solution will meet these requirements?

- A. Use AWS Glue to run an ETL job on each uploaded file
- B. Configure the ETL job to use the AWS SDK to invoke the SageMaker AI model endpoint
- C. Use real-time inference with the endpoint to re-deploy the model after it is re-trained on the updated customer dataset.
- D. Create an AWS Lambda function and webhook handlers to generate an event when an employee uploads a new file
- E. Configure SageMaker Pipelines to re-deploy the model after it is re-trained on the updated customer dataset
- F. Use Amazon EventBridge to create an event bus

- G. Set the Lambda function event as the source and SageMaker Pipelines as the target.
- H. Create an AWS Step Functions Express workflow with AWS SDK integrations to retrieve the customer data from the S3 bucket when an employee uploads a new file to the S3 bucket.
- I. Use a SageMaker Data Wrangler flow to export the data from the S3 bucket to SageMaker Autopilot.
- J. Use the SageMaker Autopilot to re-deploy the model after it has been re-trained on the updated customer dataset.
- K. Create an AWS Step Functions Standard workflow.
- L. Configure the first state to call an AWS Lambda function to respond when an employee uploads a new file to the S3 bucket.
- M. Use a pipeline in SageMaker Pipelines to re-deploy the model after it has been re-trained on the updated customer dataset.
- N. Use the next state in the workflow to run the pipeline when the first state receives a response.

**Answer: D**

#### NEW QUESTION 25

A company is designing a solution that uses foundation models (FMs) to support multiple AI workloads. Some FMs must be invoked on demand and in real time. Other FMs require consistent high-throughput access for batch processing. The solution must support hybrid deployment patterns and run workloads across cloud infrastructure and on-premises infrastructure to comply with data residency and compliance requirements.

Which combination of steps will meet these requirements? (Select TWO.)

- A. Use AWS Lambda to orchestrate low-latency FM inference by invoking FMs hosted on Amazon SageMaker AI asynchronous endpoints.
- B. Configure provisioned throughput in Amazon Bedrock to ensure consistent performance for high-volume workloads.
- C. Deploy FMs to Amazon SageMaker AI endpoints with support for edge deployment by using Amazon SageMaker Neuron.
- D. Orchestrate the FMs by using AWS Lambda to support hybrid deployment.
- E. Use Amazon Bedrock with auto-scaling to handle unpredictable traffic surges.
- F. Use Amazon SageMaker JumpStart to host and invoke the FMs.

**Answer: BC**

#### NEW QUESTION 27

A company is building a multicloud generative AI (GenAI)-powered secret resolution application that uses Amazon Bedrock and Agent Squad. The application resolves secrets from multiple sources, including key stores and hardware security modules (HSMs). The application uses AWS Lambda functions to retrieve secrets from the sources. The application uses AWS AppConfig to implement dynamic feature gating. The application supports secret chaining and detects secret drift. The application handles short-lived and expiring secrets. The application also supports prompt flows for templated instructions. The application uses AWS Step Functions to orchestrate agents to resolve the secrets and to manage secret validation and drift detection.

The company finds multiple issues during application testing. The application does not refresh expired secrets in time for agents to use. The application sends alerts for secret drift, but agents still use stale data. Prompt flows within the application reuse outdated templates, which cause cascading failures. The company must resolve the performance issues.

Which solution will meet this requirement?

- A. Use Step Functions Map states to run agent workflows in parallel.
- B. Pass updated secret metadata through Lambda function output.
- C. Use AWS AppConfig to version all prompt flows to gate and roll back faulty templates.
- D. Use Amazon Bedrock Agents on Lambda.
- E. Configure Amazon Bedrock guardrails to restrict prompt variations.
- F. Use an inline JSON schema for a single agent's workflow definition to chain tool calls.
- G. Use a centralized Amazon EventBridge pipeline to invoke each agent.
- H. Store intermediate prompts in Amazon DynamoDB.
- I. Resolve agent ordering by using TTL-based backoff and retries.
- J. Use Amazon EventBridge Pipes to invoke resolvers based on Amazon CloudWatch log patterns.
- K. Store response metadata in DynamoDB with TTL and versioned writes.
- L. Use Amazon Q Developer to dynamically generate fallback prompts.

**Answer: A**

#### NEW QUESTION 30

A company is using Amazon Bedrock to develop an AI-powered application that uses a foundation model (FM) that supports cross-Region inference and provisioned throughput. The application must serve users in Europe and North America with consistently low latency. The application must comply with data residency regulations that require European user data to remain within Europe-based AWS Regions.

During testing, the application experiences service degradation when Regional traffic spikes reach service quotas. The company needs a solution that maintains application resilience and minimizes operational complexity.

Which solution will meet these requirements?

- A. Deploy separate Amazon Bedrock instances in North American and European Regions.
- B. Use a custom routing layer that directs traffic based on user location.
- C. Configure Amazon CloudWatch alarms to monitor Regional service usage.
- D. Use Amazon SNS to send email alerts when usage approaches thresholds.
- E. Use Amazon Bedrock cross-Region inference profiles by specifying geographical codes in profile IDs when calling the InvokeModel API.
- F. Configure separate Amazon API Gateway HTTP APIs to direct European and North American users to the appropriate Regional endpoints.
- G. Deploy a multi-Region Amazon API Gateway HTTP API and AWS Lambda functions that implement retry logic to handle throttling.
- H. Configure the Lambda functions to call the FM in the nearest secondary Region when quotas are reached.
- I. Configure provisioned throughput for Amazon Bedrock in multiple Regions.
- J. Implement failover logic in application code to switch Regions when throttling occurs.
- K. Use AWS Global Accelerator to route traffic based on user location.

**Answer: B**

#### NEW QUESTION 31

.....

## **Thank You for Trying Our Product**

### **We offer two products:**

1st - We have Practice Tests Software with Actual Exam Questions

2nd - Questions and Answers in PDF Format

### **AIP-C01 Practice Exam Features:**

- \* AIP-C01 Questions and Answers Updated Frequently
- \* AIP-C01 Practice Questions Verified by Expert Senior Certified Staff
- \* AIP-C01 Most Realistic Questions that Guarantee you a Pass on Your FirstTry
- \* AIP-C01 Practice Test Questions in Multiple Choice Formats and Updatesfor 1 Year

**100% Actual & Verified — Instant Download, Please Click**  
**[Order The AIP-C01 Practice Test Here](#)**