

# Amazon-Web-Services

## Exam Questions AIP-C01

AWS Certified Generative AI Developer - Professional



### NEW QUESTION 1

An ecommerce company is developing a generative AI application that uses Amazon Bedrock with Anthropic Claude to recommend products to customers. Customers report that some recommended products are not available for sale on the website or are not relevant to the customer. Customers also report that the solution takes a long time to generate some recommendations.

The company investigates the issues and finds that most interactions between customers and the product recommendation solution are unique. The company confirms that the solution recommends products that are not in the company's product catalog. The company must resolve these issues.

Which solution will meet this requirement?

- A. Increase grounding within Amazon Bedrock Guardrail
- B. Enable Automated Reasoningcheck
- C. Set up provisioned throughput.
- D. Use prompt engineering to restrict the model responses to relevant product
- E. Use streaming techniques such as the InvokeModelWithResponseStream action to reduce perceived latency for the customers.
- F. Create an Amazon Bedrock knowledge bas
- G. Implement Retrieval Augmented Generation RA
- H. Set the PerformanceConfigLatency parameter to optimized.
- I. Store product catalog data in Amazon OpenSearch Servic
- J. Validate the model's product recommendations against the product catalo
- K. Use Amazon DynamoDB to implement response caching.

**Answer: C**

### NEW QUESTION 2

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 base
- B. Use IAM filtering to control access to each knowledge bas
- C. Deploy a supervisor agent to perform natural language intent classification on patient inquirie
- D. Configure the supervisor agent to route queries to specialized collaborator agents to respond to department-specific querie
- E. Configure each specialized collaborator agent to use Retrieval Augmented Generation (RAG) with the agent's department-specific knowledge base.
- F. Create a separate supervisor agent for each departmen
- G. Configure individual collaborator agents to perform natural language intent classification for each specialty domain within each departmen
- H. Integrate each collaborator agent with department-specific knowledge bases onl
- I. Implement manual handoff processes between the supervisor agents.
- J. Isolate data for each department in separate knowledge base
- K. Use IAM filtering to control access to each knowledge bas
- L. Deploy a single general-purpose agen
- M. Configure multiple action groups within the general-purpose agent to perform specific department function
- N. Implement rule-based routing logic within the general-purpose agent instructions.
- O. Implement multiple independent supervisor agents that run in parallel to respond to patient inquiries for each departmen
- P. Configure multiple collaborator agents for each supervisor agen
- Q. Integrate all agents with the same knowledge bas
- R. Use external routing logic to merge responses from multiple supervisor agents.

**Answer: A**

### NEW QUESTION 3

A company has a generative AI (GenAI) application that uses Amazon Bedrock to provide real-time responses to customer queries. The company has noticed intermittent failures with API calls to foundation models (FMs) during peak traffic periods.

The company needs a solution to handle transient errors and provide detailed observability into FM performance. The solution must prevent cascading failures during throttling events and provide distributed tracing across service boundaries to identify latency contributors. The solution must also enable correlation of performance issues with specific FM characteristics.

Which solution will meet these requirements?

- A. Implement a custom retry mechanism with a fixed delay of 1 second between retrie
- B. Configure Amazon CloudWatch alarms to monitor the application's error rates and latency metrics.
- C. Configure the AWS SDK with standard retry mode and exponential backoff with jitte
- D. Use AWS X-Ray tracing with annotations to identify and filter service components.
- E. Implement client-side caching of all FM response
- F. Add custom logging statements in the application code to record API call durations.
- G. Configure the AWS SDK with adaptive retry mod
- H. Use AWS CloudTrail distributed tracing to monitor throttling events.

**Answer: B**

### NEW QUESTION 4

A company deploys multiple Amazon Bedrock-based generative AI (GenAI) applications across multiple business units for customer service, content generation, and document analysis. Some applications show unpredictable token consumption patterns. The company requires a comprehensive observability solution that provides real-time visibility into token usage patterns across multiple models. The observability solution must support custom dashboards for multiple stakeholder groups and provide alerting capabilities for token consumption across all the foundation models that the company's applications use.

Which combination of solutions will meet these requirements with the LEAST operational overhead? (Select TWO.)

- A. Use Amazon CloudWatch metrics as data sources to create custom Amazon QuickSight dashboards that show token usage trends and usage patterns across FMs.

- B. Use CloudWatch Logs Insights to analyze Amazon Bedrock invocation logs for token consumption patterns and usage attribution by application
- C. Create custom queries to identify high-usage scenarios
- D. Add log widgets to dashboards to enable continuous monitoring.
- E. Create custom Amazon CloudWatch dashboards that combine native Amazon Bedrock token and invocation CloudWatch metrics
- F. Set up CloudWatch alarms to monitor token usage thresholds.
- G. Create dashboards that show token usage trends and patterns across the company's FMs by using an Amazon Bedrock zero-ETL integration with Amazon Managed Grafana.
- H. Implement Amazon EventBridge rules to capture Amazon Bedrock model invocation events
- I. Route token usage data to Amazon OpenSearch Serverless by using Amazon Data Firehose
- J. Use OpenSearch dashboards to analyze usage patterns.

**Answer:** CD

#### NEW QUESTION 5

A retail company is using Amazon Bedrock to develop a customer service AI assistant. Analysis shows that 70% of customer inquiries are simple product questions that a smaller model can effectively handle. However, 30% of inquiries are complex return policy questions that require advanced reasoning. The company wants to implement a cost-effective model selection framework to automatically route customer inquiries to appropriate models based on inquiry complexity. The framework must maintain high customer satisfaction and minimize response latency. Which solution will meet these requirements with the LEAST implementation effort?

- A. Create a multi-stage architecture that uses a small foundation model (FM) to classify the complexity of each inquiry
- B. Route simple inquiries to a smaller, more cost-effective model
- C. Route complex inquiries to a larger, more capable model
- D. Use AWS Lambda functions to handle routing logic.
- E. Use Amazon Bedrock intelligent prompt routing to automatically analyze inquiries
- F. Route simple product inquiries to smaller models and route complex return policy inquiries to more capable larger models.
- G. Implement a single-model solution that uses an Amazon Bedrock mid-sized foundation model (FM) with on-demand pricing
- H. Include special instructions in model prompts to handle both simple and complex inquiries by using the same model.
- I. Create separate Amazon Bedrock endpoints for simple and complex inquiries
- J. Implement a rule-based routing system based on keyword detection
- K. Use on-demand pricing for the smaller model and provisioned throughput for the larger model.

**Answer:** B

#### NEW QUESTION 6

A financial services company is deploying a generative AI (GenAI) application that uses Amazon Bedrock to assist customer service representatives to provide personalized investment advice to customers. The company must implement a comprehensive governance solution that follows responsible AI practices and meets regulatory requirements. The solution must detect and prevent hallucinations in recommendations. The solution must have safety controls for customer interactions. The solution must also monitor model behavior drift in real time and maintain audit trails of all prompt-response pairs for regulatory review. The company must deploy the solution within 60 days. The solution must integrate with the company's existing compliance dashboard and respond to customers within 200 ms. Which solution will meet these requirements with the LEAST operational overhead?

- A. Configure Amazon Bedrock guardrails to apply custom content filters and toxicity detection
- B. Use Amazon Bedrock Model Evaluation to detect hallucination
- C. Store prompt-response pairs in Amazon DynamoDB to capture audit trails and set a TTL
- D. Integrate Amazon CloudWatch custom metrics with the existing compliance dashboard.
- E. Deploy Amazon Bedrock and use AWS PrivateLink to access the application securely
- F. Use AWS Lambda functions to implement custom prompt validation
- G. Store prompt-response pairs in an Amazon S3 bucket and configure S3 Lifecycle policies
- H. Create custom Amazon CloudWatch dashboards to monitor model performance metrics.
- I. Use Amazon Bedrock Agents and Amazon Bedrock Knowledge Bases to ground responses
- J. Use Amazon Bedrock Guardrails to enforce content safety
- K. Use Amazon OpenSearch Service to store and index prompt-response pairs
- L. Integrate OpenSearch Service with Amazon QuickSight to create compliance reports and to detect model behavior drift.
- M. Use Amazon SageMaker Model Monitor to detect model behavior drift
- N. Use AWS WAF to filter content
- O. Store customer interactions in an encrypted Amazon RDS database
- P. Use Amazon API Gateway to create custom HTTP APIs to integrate with the compliance dashboard.

**Answer:** A

#### NEW QUESTION 7

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
- B. Create multiple pipelines in AWS CodePipeline
- C. Configure each pipeline to have its own settings for each FM
- D. Configure the application to invoke the Amazon Bedrock FMs by using the `aws_bedrock.ProvisionedModel.fromProvisionedModelArn()` method.
- E. Create a separate AWS CDK application for each environment
- F. Configure the applications to invoke the Amazon Bedrock FMs by using the `aws_bedrock.FoundationModel.fromFoundationModelId()` method
- G. Create a separate pipeline in AWS CodePipeline for each environment.
- H. Create one AWS CDK application
- I. Configure the application to invoke the Amazon Bedrock FMs by using the `aws_bedrock.FoundationModel.fromFoundationModelId()` method
- J. Create a pipeline in AWS CodePipeline that has a deployment stage for each environment that uses AWS CodeBuild deploy actions.

- K. Create one AWS CDK application for the production environment
- L. Configure the application to invoke the Amazon Bedrock FMs by using the `aws_bedrock.ProvisionedModel.fromProvisionedModelArn()` method
- M. Create a pipeline in AWS CodePipeline
- N. Configure the pipeline to deploy to the production environment by using an AWS CodeBuild deploy action
- O. For the development environment, manually recreate the resources by referring to the production application code.

**Answer: C**

#### NEW QUESTION 8

A GenAI developer is building a Retrieval Augmented Generation (RAG)-based customer support application that uses Amazon Bedrock foundation models (FMs). The application needs to process 50 GB of historical customer conversations that are stored in an Amazon S3 bucket as JSON files. The application must use the processed data as its retrieval corpus. The application's data processing workflow must extract relevant data from customer support documents, remove customer personally identifiable information (PII), and generate embeddings for vector storage. The processing workflow must be cost-effective and must finish within 4 hours.

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

- A. Use AWS Lambda and Amazon Comprehend to process files in parallel, remove PII, and call Amazon Bedrock APIs to generate vector
- B. Configure Lambda concurrency limits and memory settings to optimize throughput.
- C. Create an AWS Glue ETL job to run PII detection scripts on the data
- D. Use Amazon SageMaker Processing to run the HuggingFaceProcessor to generate embeddings by using a pre-trained model
- E. Store the embeddings in Amazon OpenSearch Service.
- F. Deploy an Amazon EMR cluster that runs Apache Spark with user-defined functions (UDFs) that call Amazon Comprehend to detect PII
- G. Use Amazon Bedrock APIs to generate vector
- H. Store outputs in Amazon Aurora PostgreSQL with the pgvector extension.
- I. Implement a data processing pipeline that uses AWS Step Functions to orchestrate a workload that uses Amazon Comprehend to detect PII and Amazon Bedrock to generate embedding
- J. Directly integrate the workflow with Amazon OpenSearch Serverless to store vectors and provide similarity search capabilities.

**Answer: D**

#### NEW QUESTION 9

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 requests
- J. Store requests in an Amazon S3 bucket
- K. Set the minimum number of instances to 0.

**Answer: AD**

#### NEW QUESTION 10

A medical company is creating a generative AI (GenAI) system by using Amazon Bedrock. The system processes data from various sources and must maintain end-to-end data lineage. The system must also use real-time personally identifiable information (PII) filtering and audit trails to automatically report compliance.

Which solution will meet these requirements?

- A. Use AWS Glue Data Catalog to register all data sources and track lineage
- B. Use Amazon Bedrock Guardrails PII filter
- C. Enable AWS CloudTrail logging for all Amazon Bedrock API calls with Amazon S3 integration
- D. Use Amazon Macie to scan stored data for sensitive information and publish findings to Amazon CloudWatch Log
- E. Create CloudWatch dashboards to visualize the findings and generate automated compliance reports.
- F. Use AWS Config to track data source configurations and changes
- G. Use AWS WAF with custom rules to filter PII at the application layer before Amazon Bedrock processes the data
- H. Configure Amazon EventBridge to capture and route audit events to Amazon S3. Use Amazon Comprehend Medical with scheduled AWS Lambda functions to analyze stored outputs for compliance violations.
- I. Use AWS DataSync to replicate data sources to track lineage
- J. Configure Amazon Macie to scan Amazon Bedrock outputs for sensitive information
- K. Use AWS Systems Manager Session Manager to log user interaction
- L. Deploy Amazon Textract with AWS Step Functions workflows to identify and redact PII from generated reports.
- M. Configure Amazon Athena to query data sources to analyze and report on data lineage
- N. Use Amazon CloudWatch custom metrics to monitor PII exposure in Amazon Bedrock responses and establish AWS X-Ray tracing to generate an audit trail
- O. Use an Amazon Rekognition Custom Labels model to detect sensitive information in the data that Amazon Bedrock processes.

**Answer: A**

#### NEW QUESTION 10

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 base
- B. Use IAM filtering to control access to each knowledge bas
- C. Deploy a supervisor agent to perform natural language intent classification on patient inquire
- D. Configure the supervisor agent to route queries to specialized collaborator agents to respond to department-specific querie
- E. Configure each specialized collaborator agent to use Retrieval Augmented Generation(RAG) with the agent's department-specific knowledge base.
- F. Create a separate supervisor agent for each departmen
- G. Configure individual collaborator agents to perform natural language intent classification for each specialty domain within each departmen
- H. Integrate each collaborator agent with department-specific knowledge bases onl
- I. Implement manual handoff processes between the supervisor agents.
- J. Isolate data for each department in separate knowledge base
- K. Use IAM filtering to control access to each knowledge bas
- L. Deploy a single general-purpose agen
- M. Configure multiple action groups within the general-purpose agent to perform specific department function
- N. Implement rule-based routing logic in the general-purpose agent instructions.
- O. Implement multiple independent supervisor agents that run in parallel to respond to patient inquiries for each departmen
- P. Configure multiple collaborator agents for each supervisor agen
- Q. Integrate all agents with the same knowledge bas
- R. Use external routing logic to merge responses from multiple supervisor agents.

**Answer: A**

#### NEW QUESTION 14

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 17

A book publishing company wants to build a book recommendation system that uses an AI assistant. The AI assistant will use ML to generate a list of recommended books from the company's book catalog. The system must suggest books based on conversations with customers.

The company stores the text of the books, customers' and editors' reviews of the books, and extracted book metadata in Amazon S3. The system must support low-latency responses and scale efficiently to handle more than 10,000 concurrent users.

Which solution will meet these requirements?

- A. Use Amazon Bedrock Knowledge Bases to generate embedding
- B. Store the embeddings as a vector store in Amazon OpenSearch Servic
- C. Create an AWS Lambda function that queries the knowledge bas
- D. Configure Amazon API Gateway to invoke the Lambda function when handling user requests.
- E. Use Amazon Bedrock Knowledge Bases to generate embedding
- F. Store the embeddings as a vector store in Amazon DynamoD
- G. Create an AWS Lambda function that queries the knowledge bas
- H. Configure Amazon API Gateway to invoke the Lambda function when handling user requests.
- I. Use Amazon SageMaker AI to deploy a pre-trained model to build a personalized recommendation engine for book
- J. Deploy the model as a SageMaker AI endpoint
- K. Invoke the model endpoint by using Amazon API Gateway.
- L. Create an Amazon Kendra GenAI Enterprise Edition index that uses the S3 connector to index the book catalog data stored in Amazon S3. Configure built-in FAQ in the Kendra inde
- M. Develop an AWS Lambda function that queries the Kendra index based on user conversation
- N. Deploy Amazon API Gateway to expose this functionality and invoke the Lambda function.

**Answer: A**

#### NEW QUESTION 18

A university recently digitized a collection of archival documents, academic journals, and manuscripts. The university stores the digital files in an AWS Lake Formation data lake.

The university hires a GenAI developer to build a solution to allow users to search the digital files by using text queries. The solution must return journal abstracts that are semantically similar to a user's query. Users must be able to search the digitized collection based on text and metadata that is associated with the journal abstracts. The metadata of the digitized files does not contain keywords. The solution must match similar abstracts to one another based on the similarity of their text. The data lake contains fewer than 1 million files.

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

- A. Use Amazon Titan Embeddings in Amazon Bedrock to create vector representations of the digitized file
- B. Store embeddings in the OpenSearch Neural plugin for Amazon OpenSearch Service.
- C. Use Amazon Comprehend to extract topics from the digitized file
- D. Store the topics and file metadata in an Amazon Aurora PostgreSQL database
- E. Query the abstract metadata against the data in the Aurora database.
- F. Use Amazon SageMaker AI to deploy a sentence-transformer model
- G. Use the model to create vector representations of the digitized file
- H. Store embeddings in an Amazon Aurora PostgreSQL database that has the pgvector extension.
- I. Use Amazon Titan Embeddings in Amazon Bedrock to create vector representations of the digitized file
- J. Store embeddings in an Amazon Aurora PostgreSQL Serverless database that has the pgvector extension.

**Answer: D**

#### NEW QUESTION 21

A company is building a generative AI (GenAI) application that uses Amazon Bedrock APIs to process complex customer inquiries. During peak usage periods, the application experiences intermittent API timeouts that cause issues such as broken response chunks and delayed data delivery. The application struggles to ensure that prompts remain within token limits when handling complex customer inquiries of varying lengths. Users have reported truncated inputs and incomplete responses. The company has also observed foundation model (FM) invocation failures.

The company needs a retry strategy that automatically handles transient service errors and prevents overwhelming Amazon Bedrock during peak usage periods. The strategy must also adapt to changing service availability and support response streaming and token-aware request handling. Which solution will meet these requirements?

- A. Implement a standard retry strategy that uses a 1-second fixed delay between attempts and a 3-retry maximum for all error
- B. Handle streaming response timeouts by restarting stream
- C. Cap token usage for each session.
- D. Implement an adaptive retry strategy that uses exponential backoff with jitter and a circuit breaker pattern that temporarily disables retries when error rates exceed a predefined threshold
- E. Implement a streaming response handler that monitors for chunk delivery timeout
- F. Configure the handler to buffer successfully received chunks and intelligently resume streaming from the last received chunk when connections are re-established.
- G. Use the AWS SDK to configure a retry strategy in standard mode
- H. Wrap Amazon Bedrock API calls in try-catch blocks that handle timeout exception
- I. Return cached completions for failed streaming request
- J. Enforce a global token limit for all users
- K. Add jitter-based retry logic and lightweight token trimming for each request
- L. Resume broken streams by requesting only missing chunks from the point of failure
- M. Maintain a small in-memory buffer of the most recent chunks.
- N. Set Amazon Bedrock client request timeouts to 30 seconds
- O. Implement client-side load shedding
- P. Buffer partial results and stop new requests when application performance degrades
- Q. Set static token usage caps for all requests
- R. Configure exponential backoff retries, dynamic chunk sizing, and context-aware token limits.

**Answer: B**

#### NEW QUESTION 23

A company runs a generative AI (GenAI)-powered summarization application in an application AWS account that uses Amazon Bedrock. The application architecture includes an Amazon API Gateway REST API that forwards requests to AWS Lambda functions that are attached to private VPC subnets. The application summarizes sensitive customer records that the company stores in a governed data lake in a centralized data storage account. The company has enabled Amazon S3, Amazon Athena, and AWS Glue in the data storage account.

The company must ensure that calls that the application makes to Amazon Bedrock use only private connectivity between the company's application VPC and Amazon Bedrock.

The company's data lake must provide fine-grained column-level access across the company's AWS accounts.

Which solution will meet these requirements?

- A. In the application account, create interface VPC endpoints for Amazon Bedrock runtime
- B. Run Lambda functions in private subnet
- C. Use IAM conditions on inference and data-plane policies to allow calls only to approved endpoints and roles
- D. In the data storage account, use AWS Lake Formation LF-tag-based access control to create table-level and column-level cross-account grants.
- E. Run Lambda functions in private subnet
- F. Configure a NAT gateway to provide access to Amazon Bedrock and the data lake
- G. Use S3 bucket policies and ACLs to manage permissions
- H. Export AWS CloudTrail logs to Amazon S3 to perform weekly reviews.
- I. Create a gateway endpoint only for Amazon S3 in the application account
- J. Invoke Amazon Bedrock through public endpoint
- K. Use database-level grants in AWS Lake Formation to manage data access
- L. Stream AWS CloudTrail logs to Amazon CloudWatch Log
- M. Do not set up metric filters or alarms.
- N. Use VPC endpoints to provide access to Amazon Bedrock and Amazon S3 in the application account
- O. Use only IAM path-based policies to manage data lake access
- P. Send AWS CloudTrail logs to Amazon CloudWatch Log
- Q. Periodically create dashboards and allow public fallback for cross-Region reads to reduce setup time.

**Answer: B**

#### NEW QUESTION 28

A financial services company uses multiple foundation models (FMs) through Amazon Bedrock for its generative AI (GenAI) applications. To comply with a new regulation for GenAI use with sensitive financial data, the company needs a token management solution.

The token management solution must proactively alert when applications approach model-specific token limits. The solution must also process more than 5,000 requests each minute and maintain token usage metrics to allocate costs across business units.

Which solution will meet these requirements?

- A. Develop model-specific tokenizers in an AWS Lambda function
- B. Configure the Lambda function to estimate token usage before sending requests to Amazon Bedrock
- C. Configure the Lambda function to publish metrics to Amazon CloudWatch and trigger alarms when requests approach threshold
- D. Store detailed token usage in Amazon DynamoDB to report costs.
- E. Implement Amazon Bedrock Guardrails with token quota policies
- F. Capture metrics on rejected request
- G. Configure Amazon EventBridge rules to trigger notifications based on Amazon Bedrock Guardrails metric
- H. Use Amazon CloudWatch dashboards to visualize token usage trends across models.
- I. Deploy an Amazon SQS dead-letter queue for failed request
- J. Configure an AWS Lambda function to analyze token-related failures
- K. Use Amazon CloudWatch Logs Insights to generate reports on token usage patterns based on error logs from Amazon Bedrock API responses.
- L. Use Amazon API Gateway to create a proxy for all Amazon Bedrock API calls
- M. Configure request throttling based on custom usage plans with predefined token quota
- N. Configure API Gateway to reject requests that will exceed token limits.

**Answer: A**

#### NEW QUESTION 29

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 31

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 attributes
- E. Store Amazon Bedrock FM endpoints as REST API stage variables
- 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 33

A company is using Amazon Bedrock to develop a customer support AI assistant. The AI assistant must respond to customer questions about their accounts. The AI assistant must not expose personal information in responses. The company must comply with data residency policies by ensuring that all processing occurs within the same AWS Region where each customer is located. The company wants to evaluate how effective the AI assistant is at preventing the exposure of personal information before the company makes the AI assistant available to customers. Which solution will meet these requirements?

- A. Configure a cross-Region Amazon Bedrock guardrail to apply sensitive information filter
- B. Set the guardrail to detect mode during development and testing
- C. Switch to block mode for production deployment.
- D. Configure an Amazon Bedrock guardrail to apply sensitive information filter
- E. Set the guardrail to mask mode during development and testing
- F. Switch to block mode for production deployment
- G. Deploy a copy of the guardrail to each Region where the company operates.

- H. Configure an Amazon Bedrock guardrail to apply content and topic filter
- I. Set the guardrail to detect mode during development, testing, and production
- J. Disable invocation logging for the Amazon Bedrock model.
- K. Configure a cross-Region Amazon Bedrock guardrail to apply a set of content and word filter
- L. Set the guardrail to detect mode during development and testing
- M. Switch to mask mode for production deployment.

**Answer: B**

#### NEW QUESTION 37

A company has a recommendation system running 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 experiences intermittent issues where some recommendations do not match customer preferences. The company needs an observability solution to monitor operational metrics and detect patterns of performance degradation compared to established baselines. The solution must 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 Insight
- B. Set up alarms for latency threshold
- C. Add custom token metrics using the CloudWatch embedded metric format.
- D. Implement AWS X-Ray
- E. Enable CloudWatch Logs Insight
- F. Set up AWS CloudTrail and create dashboards in Amazon QuickSight.
- G. Enable Amazon CloudWatch Application Insight
- H. Create custom metrics for recommendation quality, token usage, and response latency using the CloudWatch embedded metric format with dimensions for request types and user segment
- I. Configure CloudWatch anomaly detection on model metric
- J. Use CloudWatch Logs Insights for pattern analysis.
- K. Use Amazon OpenSearch Service with the Observability plugin
- L. Ingest metrics and logs through Amazon Kinesis and analyze behavior with custom queries.

**Answer: C**

#### NEW QUESTION 40

A pharmaceutical company is developing a Retrieval Augmented Generation application that uses an Amazon Bedrock knowledge base. The knowledge base uses Amazon OpenSearch Service as a data source for more than 25 million scientific papers. Users report that the application produces inconsistent answers that cite irrelevant sections of papers when queries span methodology, results, and discussion sections of the papers. The company needs to improve the knowledge base to preserve semantic context across related paragraphs on the scale of the entire corpus of data. Which solution will meet these requirements?

- A. Configure the knowledge base to use fixed-size chunking
- B. Set a 300-token maximum chunk size and a 10% overlap between chunks
- C. Use an appropriate Amazon Bedrock embedding model.
- D. Configure the knowledge base to use hierarchical chunking
- E. Use parent chunks that contain 1,000 tokens and child chunks that contain 200 tokens
- F. Set a 50-token overlap between chunks.
- G. Configure the knowledge base to use semantic chunking
- H. Use a buffer size of 1 and a breakpoint percentile threshold of 85% to determine chunk boundaries based on content meaning.
- I. Configure the knowledge base not to use chunking
- J. Manually split each document into separate files before ingestion
- K. Apply post-processing reranking during retrieval.

**Answer: B**

#### NEW QUESTION 45

A company uses Amazon Bedrock to build a Retrieval Augmented Generation (RAG) system. The RAG system uses an Amazon Bedrock Knowledge Base that is based on an Amazon S3 bucket as the data source for emergency news video content. The system retrieves transcripts, archived reports, and related documents from the S3 bucket. The RAG system uses state-of-the-art embedding models and a high-performing retrieval setup. However, users report slow responses and irrelevant results, which cause decreased user satisfaction. The company notices that vector searches are evaluating too many documents across too many content types and over long periods of time. The company determines that the underlying models will not benefit from additional fine-tuning. The company must improve retrieval accuracy by applying smarter constraints and wants a solution that requires minimal changes to the existing architecture. Which solution will meet these requirements?

- A. Enhance embeddings by using a domain-adapted model that is specifically trained on emergency news content for improved vector similarity.
- B. Migrate to Amazon OpenSearch Service
- C. Use vector fields and metadata filters to define the scope of results retrieval.
- D. Enable metadata-aware filtering within the Amazon Bedrock knowledge base by indexing S3 object metadata.
- E. Migrate to an Amazon Q Business index to perform structured metadata filtering and document categorization during retrieval.

**Answer: C**

#### NEW QUESTION 49

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 search 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 52

A financial technology company is using Amazon Bedrock to build an assessment system for the company's customer service AI assistant. The AI assistant must provide financial recommendations that are factually accurate, compliant with financial regulations, and conversationally appropriate. The company needs to combine automated quality evaluations at scale with targeted human reviews of critical interactions. What solution will meet these requirements?

- A. Configure a pipeline in which financial experts manually score all responses for accuracy, compliance, and conversational quality
- B. Use Amazon SageMaker notebooks to analyze results to identify improvement areas.
- C. Configure Amazon Bedrock evaluations that use Anthropic Claude Sonnet as a judge model to assess response accuracy and appropriateness
- D. Configure custom Amazon Bedrock guardrails to check responses for compliance with financial policies
- E. Add Amazon Augmented AI (Amazon A2I) human reviews for flagged critical interactions.
- F. Create an Amazon Lex bot to manage customer service interaction
- G. Configure AWS Lambda functions to check responses against a static compliance database
- H. Configure intents that call the Lambda function
- I. Add an additional intent to collect end-user reviews.
- J. Configure Amazon CloudWatch to monitor response patterns from the AI assistant
- K. Configure CloudWatch alerts for potential compliance violation
- L. Establish a team of human evaluators to review flagged interactions.

**Answer: B**

#### NEW QUESTION 54

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 permissions
- 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 bucket
- 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 purpose
- . 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 55

A company is creating a workflow to review customer-facing communications before the company sends the communications. The company uses a pre-defined message template to generate the communications and stores the communications in an Amazon S3 bucket. The workflow needs to capture a specific portion from the template and send it to an Amazon Bedrock model. The workflow must store model responses back to the original S3 bucket. Which solution will meet these requirements?

- A. Create a flow in Amazon Bedrock Flow
- B. Configure S3 action nodes at the beginning and end of the flow to retrieve and store the communications and the model response
- C. In the middle of the flow, configure an expression to parse each communication
- D. Configure an agent step to send the parsed input to the model for review.
- E. Create an AWS Step Functions Express workflow state machine
- F. Use an Amazon S3 integration GetObject step to retrieve the original communication
- G. Use an intrinsic function Pass step to parse the communications and to pass the results to an Amazon Bedrock InvokeModel step
- H. Configure an Amazon S3 integration PutObject step to store the model responses back to the S3 bucket.
- I. Create an Amazon Bedrock agent that has an action group
- J. Configure instructions to define how the agent should parse the communication
- K. Configure the action group to retrieve the communications from the S3 bucket, invoke the Amazon Bedrock model, and store the model responses back to the S3 bucket.
- L. Create an Amazon Bedrock agent that has a single action group
- M. Configure three AWS Lambda functions in the action group
- N. Configure the functions to retrieve the communications from the S3 bucket, parse the communications and invoke the Amazon Bedrock model, and store the model responses back to the S3 bucket.

**Answer: A**

#### NEW QUESTION 58

A healthcare company is developing an application to process medical queries. The application must answer complex queries with high accuracy by reducing semantic dilution. The application must refer to domain-specific terminology in medical documents to reduce ambiguity in medical terminology. The application must be able to respond to 1,000 queries each minute with response times less than 2 seconds. Which solution will meet these requirements with the LEAST operational overhead?

- A. Use Amazon API Gateway to route incoming queries to an Amazon Bedrock agent
- B. Configure the agent to use an Anthropic Claude model to decompose queries and an Amazon Titan model to expand queries
- C. Create an Amazon Bedrock knowledge base to store the reference medical documents.
- D. Configure an Amazon Bedrock knowledge base to store the reference medical document
- E. Enable query decomposition in the knowledge base
- F. Configure an Amazon Bedrock flow that uses a foundation model and the knowledge base to support the application.
- G. Use Amazon SageMaker AI to host custom ML models for both query decomposition and query expansion
- H. Configure Amazon Bedrock knowledge bases to store the reference medical document
- I. Encrypt the documents in the knowledge base.
- J. Create an Amazon Bedrock agent to orchestrate multiple AWS Lambda functions to decompose queries
- K. Create an Amazon Bedrock knowledge base to store the reference medical document
- L. Use the agent's built-in knowledge base capabilities
- M. Add deep research and reasoning capabilities to the agent to reduce ambiguity in the medical terminology.

**Answer: B**

#### NEW QUESTION 62

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 65

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 67

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 with 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 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 search 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 69

A company uses an organization in AWS Organizations with all features enabled to manage multiple AWS accounts. Employees use Amazon Bedrock across multiple accounts. The company must prevent specific topics and proprietary information from being included in prompts to Amazon Bedrock models. The company must ensure that employees can use only approved Amazon Bedrock models. The company wants to manage these controls centrally. Which combination of solutions will meet these requirements? (Select TWO.)

- A. Create an IAM permissions boundary for each employee's IAM role
- B. Configure the permissions boundary to require an approved Amazon Bedrock guardrail identifier to invoke Amazon Bedrock model
- C. Create an SCP that allows employees to use only approved models.
- D. Create an SCP that allows employees to use only approved model
- E. Configure the SCP to require employees to specify a guardrail identifier in calls to invoke an approved model.
- F. Create an SCP that prevents an employee from invoking a model if a centrally deployed guardrail identifier is not specified in a call to the model
- G. Create a permissions boundary on each employee's IAM role that allows each employee to invoke only approved models.
- H. Use AWS CloudFormation to create a custom Amazon Bedrock guardrail that has a block filtering policy
- I. Use stack sets to deploy the guardrail to each account in the organization.
- J. Use AWS CloudFormation to create a custom Amazon Bedrock guardrail that has a mask filtering policy
- K. Use stack sets to deploy the guardrail to each account in the organization.

**Answer: CD**

#### NEW QUESTION 72

A company uses AWS Lambda functions to build an AI agent solution. A GenAI developer must set up a Model Context Protocol (MCP) server that accesses user information. The GenAI developer must also configure the AI agent to use the new MCP server. The GenAI developer must ensure that only authorized users can access the MCP server.

Which solution will meet these requirements?

- A. Use a Lambda function to host the MCP server
- B. Grant the AI agent Lambda function permission to invoke the Lambda function that hosts the MCP server
- C. Configure the AI agent's MCP client to invoke the MCP server asynchronously.
- D. Use a Lambda function to host the MCP server
- E. Grant the AI agent Lambda function permission to invoke the Lambda function that hosts the MCP server
- F. Configure the AI agent to use the STDIO transport with the MCP server.
- G. Use a Lambda function to host the MCP server
- H. Create an Amazon API Gateway HTTP API that proxies requests to the Lambda function
- I. Configure the AI agent solution to use the Streamable HTTP transport to make requests through the HTTP API
- J. Use Amazon Cognito to enforce OAuth 2.1.
- K. Use a Lambda layer to host the MCP server
- L. Add the Lambda layer to the AI agent Lambda function
- M. Configure the AI agent solution to use the STDIO transport to send requests to the MCP server
- N. In the AI agent's MCP configuration, specify the Lambda layer ARN as the command
- O. Specify the user credentials as environment variables.

**Answer: C**

#### NEW QUESTION 73

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 tables
- 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 75

A financial services company needs to pre-process unstructured data such as customer transcripts, financial reports, and documentation. The company stores the unstructured data in Amazon S3 to support an Amazon Bedrock application.

The company must validate data quality, create auditable metadata, monitor data metrics, and customize text chunking to optimize foundation model (FM) performance.

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

- A. Use Amazon SageMaker Data Wrangler to create a data flow
- B. Configure Amazon CloudWatch metrics and alarms to monitor data quality
- C. Use a custom AWS Lambda function to pre-process the data
- D. Load processed data into Amazon Bedrock.
- E. Set up an AWS Glue crawler to catalog data source
- F. Create AWS Glue ETL jobs to run custom transformation script
- G. Use AWS Glue Data Quality to validate and monitor data quality
- H. Load processed data into Amazon Bedrock.
- I. Use Amazon Comprehend to extract entities
- J. Create an AWS Lambda function to chunk text
- K. Run Amazon Athena to query and validate data quality
- L. Load processed data into Amazon Bedrock.
- M. Create an AWS Step Functions workflow to orchestrate data pre-processing task
- N. Run custom code on Amazon EC2 instance
- O. Use Amazon SageMaker Model Monitor to monitor data quality
- P. Load processed data into Amazon Bedrock.

**Answer: B**

#### NEW QUESTION 80

A financial services company is building a customer support application that retrieves relevant financial regulation documents from a database based on semantic similarity to user queries. The application must integrate with Amazon Bedrock to generate responses. The application must search documents in English, Spanish, and Portuguese. The application must filter documents by metadata such as publication date, regulatory agency, and document type.

The database stores approximately 10 million document embeddings. To minimize operational overhead, the company wants a solution that minimizes management and maintenance effort while providing low-latency responses for real-time customer interactions.

Which solution will meet these requirements?

- A. Use Amazon OpenSearch Serverless to provide vector search capabilities and metadata filtering
- B. Integrate with Amazon Bedrock Knowledge Bases to enable Retrieval Augmented Generation (RAG) using an Anthropic Claude foundation model.
- C. Deploy an Amazon Aurora PostgreSQL database with the pgvector extension
- D. Store embeddings and metadata in table
- E. Use SQL queries for similarity search and send results to Amazon Bedrock for response generation.
- F. Use Amazon S3 Vectors to configure a vector index and non-filterable metadata field
- G. Integrate S3 Vectors with Amazon Bedrock for RAG.
- H. Set up an Amazon Neptune Analytics database with a vector index
- I. Use graph-based retrieval and Amazon Bedrock for response generation.

**Answer: A**

#### NEW QUESTION 81

A company is creating a generative AI (GenAI) application that uses Amazon Bedrock foundation models (FMs). The application must use Microsoft Entra ID to authenticate. All FM API calls must stay on private network paths. Access to the application must be limited by department to specific model families. The company also needs a comprehensive audit trail of model interactions.

Which solution will meet these requirements?

- A. Configure SAML federation between Microsoft Entra ID and AWS Identity and Access Management
- B. Create department-specific IAM roles that allow only the required ModelId value
- C. Create AWS PrivateLink interface VPC endpoints for Amazon Bedrock runtime service
- D. Enable AWS CloudTrail to capture Amazon Bedrock API calls
- E. Configure Amazon Bedrock model invocation logging to record detailed model interactions.
- F. Create an identity provider (IdP) connection in IAM to authenticate by using Microsoft Entra ID
- G. Assign department permission sets to control access to specific model families
- H. Deploy AWS Lambda functions in private subnets with a NAT gateway for egress to Amazon Bedrock public endpoint
- I. Enable CloudWatch Logs to capture model interactions for auditing purposes.
- J. Create a SAML identity provider (IdP) in IAM to authenticate by using Microsoft Entra ID
- K. Use IAM permissions boundaries to limit department roles' access to specific model families
- L. Configure public Amazon Bedrock API endpoints with VPC routing to maintain private network connectivity
- M. Set up CloudTrail with Amazon S3 Lifecycle rules to manage audit logs of model interactions.
- N. Configure OpenID Connect (OIDC) federation between Microsoft Entra ID and IAM
- O. Use attribute-based access control to map department attributes to specific model access permissions
- P. Apply SCP policies to restrict access to Amazon Bedrock FM families based on department
- Q. Use Microsoft Entra ID's built-in logging capabilities to maintain an audit trail of model interactions.

**Answer: A**

#### NEW QUESTION 86

A media company is launching a platform that allows thousands of users every hour to upload images and text content. The platform uses Amazon Bedrock to process the uploaded content to generate creative compositions.

The company needs a solution to ensure that the platform does not process or produce inappropriate content. The platform must not expose personally identifiable information (PII) in the compositions. The solution must integrate with the company's existing Amazon S3 storage workflow. Which solution will meet these requirements with the LEAST infrastructure management overhead?

- A. Enable the Enhanced Monitoring too
- B. Use an Amazon CloudWatch alarm to filter traffic to the platfor
- C. Use Amazon Comprehend PII detection to pre-process the dat
- D. Create a CloudWatch alarm to monitor for Amazon Comprehend PII detection event
- E. Create an AWS Step Functions workflow that includes an Amazon Rekognition image moderation step.
- F. Use an Amazon API Gateway HTTP API with request validation templates to screen content before storing the uploaded content in Amazon S3. Use Amazon SageMaker AI to build custom content moderation models that process content before sending the processed content to Amazon Bedrock.
- G. Create an Amazon Cognito user pool that uses pre-authentication AWS Lambda functions to run content moderation check
- H. Use Amazon Textract to filter text content and Amazon Rekognition to filter image content before allowing users to upload content to the platform.
- I. Create an AWS Step Functions workflow that uses built-in Amazon Bedrock guardrails to filter conten
- J. Use Amazon Comprehend PII detection to pre-process the conten
- K. Use Amazon Rekognition image moderation.

**Answer: D**

#### NEW QUESTION 90

A bank is building a generative AI (GenAI) application that uses Amazon Bedrock to assess loan applications by using scanned financial documents. The application must extract structured data from the documents. The application must redact personally identifiable information (PII) before inference. The application must use foundation models (FMs) to generate approvals. The application must route low-confidence document extraction results to human reviewers who are within the same AWS Region as the loan applicant.

The company must ensure that the application complies with strict Regional data residency and auditability requirements. The application must be able to scale to handle 25,000 applications each day and provide 99.9% availability.

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

- A. Deploy Amazon Textract and Amazon Augmented AI within the same Region to extract relevant data from the scanned document
- B. Route low-confidence pages to human reviewers.
- C. Use AWS Lambda functions to detect and redact PII from submitted documents before inferenc
- D. Apply Amazon Bedrock guardrails to prevent inappropriate or unauthorized content in model output
- E. Configure Region-specific IAM roles to enforce data residency requirements and to control access to the extracted data.
- F. Use Amazon Kendra and Amazon OpenSearch Service to extract field-level values semantically from the uploaded documents before inference.
- G. Store uploaded documents in Amazon S3 and apply object metadat
- H. Configure IAM policies to store original documents within the same Region as each applican
- I. Enable object tagging for future audits.
- J. Use AWS Glue Data Quality to validate the structured document dat
- K. Use AWS Step Functions to orchestrate a review workflow that includes a prompt engineering step that transforms validated data into optimized prompts before invoking Amazon Bedrock to assess loan applications.
- L. Use Amazon SageMaker Clarify to generate fairness and bias reports based on model scoring decisions that Amazon Bedrock makes.

**Answer: ABD**

#### NEW QUESTION 92

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 fil
- B. Configure the ETL job to use the AWS SDK to invoke the SageMaker AI model endpoin
- 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 fil
- E. Configure SageMaker Pipelines to re-deploy the model after it is re-trained on the updated customer datase
- F. Use Amazon EventBridge to create an event bu
- 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 bucke
- I. Use a SageMaker Data Wrangler flow to export the data from the S3 bucket to SageMaker Autopilo
- 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 workflo
- L. Configure the first state to call an AWS Lambda function to respond when an employee uploads a new file to the S3 bucke
- M. Use a pipeline in SageMaker Pipelines to re-deploy the model after it has been re-trained on the updated customer datase
- N. Use the next state in the workflow to run the pipeline when the first state receives a response.

**Answer: D**

#### NEW QUESTION 95

A company is building an AI advisory application by using Amazon Bedrock. The application will provide recommendations to customers. The company needs the application to explain its reasoning process and cite specific sources for data. The application must retrieve information from company data sources and show step-by-step reasoning for recommendations. The application must also link data claims to source documents and maintain response latency under 3 seconds.

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

- A. Use Amazon Bedrock Knowledge Bases with source attribution enable
- B. Use the Anthropic Claude Messages API with RAG to set high-relevance thresholds for sourcedocument
- C. Store reasoning and citations in Amazon S3 for auditing purposes.
- D. Use Amazon Bedrock with Anthropic Claude models and extended thinkin
- E. Configure a 4,000-token thinking budge

- F. Store reasoning traces and citations in Amazon DynamoDB for auditing purposes.
- G. Configure Amazon SageMaker AI with a custom Anthropic Claude mode
- H. Use the model's reasoning parameter and AWS Lambda to process response
- I. Add source citations from a separate Amazon RDS database.
- J. Use Amazon Bedrock with Anthropic Claude models and chain-of-thought reasoning
- K. Configure custom retrieval tracking with the Amazon Bedrock Knowledge Bases API
- L. Use Amazon CloudWatch to monitor response latency metrics.

**Answer: A**

#### NEW QUESTION 97

An elevator service company has developed an AI assistant application by using Amazon Bedrock. The application generates elevator maintenance recommendations to support the company's elevator technicians. The company uses Amazon Kinesis Data Streams to collect the elevator sensor data. New regulatory rules require that a human technician must review all AI-generated recommendations. The company needs to establish human oversight workflows to review and approve AI recommendations. The company must store all human technician review decisions for audit purposes. Which solution will meet these requirements?

- A. Create a custom approval workflow by using AWS Lambda functions and Amazon SQS queues for human review of AI recommendation
- B. Store all review decisions in Amazon DynamoDB for audit purposes.
- C. Create an AWS Step Functions workflow that has a human approval step that uses the waitForResource API to pause execution
- D. After a human technician completes a review, use an AWS Lambda function to call the SendTaskSuccess API with the approval decision
- E. Store all review decisions in Amazon DynamoDB.
- F. Create an AWS Glue workflow that has a human approval step
- G. After the human technician review, integrate the application with an AWS Lambda function that calls the SendTaskSuccess API
- H. Store all human technician review decisions in Amazon DynamoDB.
- I. Configure Amazon EventBridge rules with custom event patterns to route AI recommendations to human technicians for review
- J. Create AWS Glue jobs to process human technician approval queue
- K. Use Amazon ElastiCache to cache all human technician review decisions.

**Answer: B**

#### 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. Migrate the restaurant data to Amazon OpenSearch Service
- B. 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
- C. Create Amazon API Gateway HTTP API endpoints to transform user queries into structured search parameters.
- D. Migrate the restaurant data to Amazon OpenSearch Service
- E. Use a foundation model (FM) in Amazon Bedrock to generate vector embeddings from restaurant descriptions, reviews, and menu items
- F. When users submit natural language queries, convert the queries to embeddings by using the same FM
- G. Perform k-nearest neighbors (k-NN) searches to find semantically similar results.
- H. Keep the restaurant data in PostgreSQL and implement a pgvector extension
- I. Use a foundation model (FM) in Amazon Bedrock to generate vector embeddings from restaurant data
- J. Store the vector embeddings directly in PostgreSQL
- K. Create an AWS Lambda function to convert natural language queries to vector representations by using the same FM
- L. Configure the Lambda function to perform similarity searches within the database.
- M. Migrate the restaurant data to an Amazon Bedrock knowledge base by using a custom ingestion pipeline
- N. Configure the knowledge base to automatically generate embeddings from restaurant information
- O. Use the Amazon Bedrock Retrieve API with built-in vector search capabilities to query the knowledge base directly by using natural language input.

**Answer: D**

#### NEW QUESTION 105

A GenAI developer is evaluating Amazon Bedrock foundation models (FMs) to enhance a Europe-based company's internal business application. The company has a multi-account landing zone in AWS Control Tower. The company uses Service Control Policies (SCPs) to allow its accounts to use only the eu-north-1 and eu-west-1 Regions. All customer data must remain in private networks within the approved AWS Regions.

The GenAI developer selects an FM based on analysis and testing and hosts the model in the eu-central-1 Region and the eu-west-3 Region. The GenAI developer must enable access to the FM for the company's employees. The GenAI developer must ensure that requests to the FM are private and remain within the same Regions as the FM.

Which solution will meet these requirements?

- A. Deploy an AWS Lambda function that is exposed by a private Amazon API Gateway REST API to a VPC in eu-north-1. Create a VPC endpoint for the selected FM in eu-central-1 and eu-west-3. Extend existing SCPs to allow employees to use the FM
- B. Integrate the REST API with the business application.
- C. Deploy the FM on Amazon EC2 instances in eu-north-1. Deploy a private Amazon API Gateway REST API in front of the EC2 instance
- D. Configure an Amazon Bedrock VPC endpoint
- E. Integrate the REST API with the business application.
- F. Configure the FM to use cross-Region inference through a Europe-scoped endpoint
- G. Configure an Amazon Bedrock VPC endpoint
- H. Extend existing SCPs to allow employees to use the FM through inference profiles in Europe-based Regions where the FM is available
- I. Use an inference profile to integrate Amazon Bedrock with the business application.
- J. Deploy the FM in Amazon SageMaker in eu-north-1. Configure a SageMaker VPC endpoint
- K. Extend existing SCPs to allow employees to use the SageMaker endpoint
- L. Integrate the FM in SageMaker with the business application.

**Answer: C**

**NEW QUESTION 109**

A company is developing a customer communication platform that uses an AI assistant powered by an Amazon Bedrock foundation model (FM). The AI assistant summarizes customer messages and generates initial response drafts.

The company wants to use Amazon Comprehend to implement layered content filtering. The layered content filtering must prevent sharing of offensive content, protect customer privacy, and detect potential inappropriate advice solicitation. Inappropriate advice solicitation includes requests for unethical practices, harmful activities, or manipulative behaviors.

The solution must maintain acceptable overall response times, so all pre-processing filters must finish before the content reaches the FM.

Which solution will meet these requirements?

- A. Use parallel processing with asynchronous API call
- B. Use toxicity detection for offensive content
- C. Use prompt safety classification for inappropriate advice solicitation
- D. Use personally identifiable information (PII) detection without redaction.
- E. Use custom classification to build an FM that detects offensive content and inappropriate advice solicitation
- F. Apply personally identifiable information (PII) detection as a secondary filter only when messages pass the custom classifier.
- G. Deploy a multi-stage process
- H. Configure the process to use prompt safety classification first, then toxicity detection on safe prompts only, and finally personally identifiable information (PII) detection in streaming mode
- I. Route flagged messages through Amazon EventBridge for human review.
- J. Use toxicity detection with thresholds configured to 0.5 for all categories
- K. Use parallel processing for both prompt safety classification and personally identifiable information (PII) detection with entity redaction
- L. Apply Amazon CloudWatch alarms to filter metrics.

**Answer: D**

**NEW QUESTION 113**

.....

## **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](#)**