# Amazon-Web-Services

## Exam Questions AIP-C01

AWS Certified Generative AI Developer - Professional

**NEW QUESTION 1**
A company is building a serverless application that uses AWS Lambda functions to help students around the world summarize notes. The application uses Anthropic Claude through Amazon Bedrock. The company observes that most of the traffic occurs during evenings in each time zone. Users report experiencing throttling errors during peak usage times in their time zones.
The company needs to resolve the throttling issues by ensuring continuous operation of the application. The solution must maintain application performance quality and must not require a fixed hourly cost during low traffic periods.
Which solution will meet these requirements?

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

**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 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 inquir
B. Route simple inquiries to a smaller, more cost-effective mode
C. Route complex inquiries to a larger, more capable mode
D. Use AWS Lambda functions to handle routing logic.
E. Use Amazon Bedrock intelligent prompt routing to automatically analyze inquirie
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 pricin
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 inquirie
J. Implement a rule-based routing system based on keyword detectio
K. Use on-demand pricing for the smaller model and provisioned throughput for the larger model.

**Answer:** B

**NEW QUESTION 4**
A company is building a legal research AI assistant that uses Amazon Bedrock with an Anthropic Claude foundation model (FM). The AI assistant must retrieve highly relevant case law documents to augment the FM??s responses. The AI assistant must identify semantic relationships between legal concepts, specific legal terminology, and citations. The AI assistant must perform quickly and return precise results.
Which solution will meet these requirements?

A. Configure an Amazon Bedrock knowledge base to use a default vector search configuratio
B. Use Amazon Bedrock to expand queries to improve retrieval for legal documents based on specific terminology and citations.

C. Use Amazon OpenSearch Service to deploy a hybrid search architecture that combines vector search with keyword searc
D. Apply an Amazon Bedrock reranker model to optimize result relevance.
E. Enable the Amazon Kendra query suggestion feature for end user
F. Use Amazon Bedrock to perform post-processing of search results to identify semantic similarity in the documents and to produce precise results.
G. Use Amazon OpenSearch Service with vector search and Amazon Bedrock Titan Embeddings to index and search legal document
H. Use custom AWS Lambda functions to merge results with keyword-based filters that are stored in an Amazon RDS database.

**Answer:** B


**NEW QUESTION 5**
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 startu
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 endpoin
D. Set the minimum number of instances to 0. Continue to perform per-request processin
E. Lazily load model weights on the first request.
F. Switch to a multi-model endpoin
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 bucke
K. Set the minimum number of instances to 0.

**Answer:** AD


**NEW QUESTION 6**
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 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 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 process 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 simultaneousl
B. Use Amazon CloudWatch metrics to monitor latency, concurrency, and throughpu
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 performanc
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 assistan
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 meanin
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**
A healthcare company is using Amazon Bedrock to build a system to help practitioners make clinical decisions. The system must provide treatment recommendations to physicians based only on approved medical documentation and must cite specific sources. The system must not hallucinate or produce factually incorrect information.
Which solution will meet these requirements with the LEAST operational overhead?

A. Integrate Amazon Bedrock with Amazon Kendra to retrieve approved document
B. Implement custom post-processing to compare generated responses against source documents and to include citations.
C. Deploy an Amazon Bedrock Knowledge Base and connect it to approved clinical source document
D. Use the Amazon Bedrock RetrieveAndGenerate API to return citations from the knowledge base.
E. Use Amazon Bedrock and Amazon Comprehend Medical to extract medical entitie
F. Implement verification logic against a medical terminology database.
G. Use an Amazon Bedrock knowledge base with Retrieve API calls and InvokeModel API calls to retrieve approved clinical source document
H. Implement verification logic to compare against retrieved sources and to cite sources.

**Answer:** B


**NEW QUESTION 9**
A retail company has a generative AI (GenAI) product recommendation application that uses Amazon Bedrock. The application suggests products to customers based on browsing history and demographics. The company needs to implement fairness evaluation across multiple demographic groups to detect and measure bias in recommendations between two prompt approaches. The company wants to collect and monitor fairness metrics in real time. The company must receive an alert if the fairness metrics show a discrepancy of more than 15% between demographic groups. The company must receive weekly reports that compare the performance of the two prompt approaches.
Which solution will meet these requirements with the LEAST custom development effort?

A. Configure an Amazon CloudWatch dashboard to display default metrics from Amazon Bedrock API call
B. Create custom metrics based on model output
C. Set up Amazon EventBridge rules to invoke AWS Lambda functions that perform post-processing analysis on model responses and publish custom fairness metrics.
D. Create the two prompt variants in Amazon Bedrock Prompt Managemen
E. Use Amazon Bedrock Flows to deploy the prompt variants with defined traffic allocatio
F. Configure Amazon Bedrock guardrails to monitor demographic fairnes
G. Set up Amazon CloudWatch alarms on the GuardrailContentSource dimension by using InvocationsIntervened metrics to detect recommendation discrepancy threshold violations.
H. Set up Amazon SageMaker Clarify to analyze model output
I. Publish fairness metrics to Amazon CloudWatc
J. Create CloudWatch composite alarms that combine SageMaker Clarify bias metrics with Amazon Bedrock latency metrics.
K. Create an Amazon Bedrock model evaluation job to compare fairness between the two prompt variant
L. Enable model invocation logging in Amazon CloudWatc
M. Set up CloudWatch alarms for InvocationsIntervened metrics with a dimension for each demographic group.

**Answer:** B


**NEW QUESTION 10**
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 functio
B. Configure the Lambda function to estimate token usage before sending requests to Amazon Bedroc
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 policie
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 failure
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 call
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 10**
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 chang
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 reques
H. Run business logic in the Lambda function to select the appropriate FM for each reques
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 AppConfi
K. Return authorization contexts based on business logi
L. Route requests to model-specific Lambda functions for each Amazon Bedrock FM.

**Answer:** C


**NEW QUESTION 12**
A wildlife conservation agency operates zoos globally. The agency uses various sensors, trackers, and audiovisual recorders to monitor animal behavior. The agency wants to launch a generative AI (GenAI) assistant that can ingest multimodal data to study animal behavior.
The GenAI assistant must support natural language queries, avoid speculative behavioral interpretations, and maintain audit logs for ethical research audits.
Which solution will meet these requirements?

A. Ingest raw videos into Amazon Rekognition to detect animal postures and expression
B. Use Amazon Data Firehose to stream sensor and GPS data into Amazon S3. Prompt an Amazon Bedrock FM using basic templates stored in AWS Systems Manager Parameter Stor
C. Use IAM for access contro
D. Use AWS CloudTrail for audit logging.
E. Use Amazon SageMaker Processing and Amazon Transcribe to pre-process multimodal dat
F. Ingest curated summaries into an Amazon Bedrock Knowledge Base
G. Apply Amazon Bedrock guardrails to restrict speculative output
H. Use AWS AppConfig to manage prompt template
I. Use AWS CloudTrail to log research activity for audits.
J. Use Amazon OpenSearch Serverless to index behavioral logs and telemetr
K. Use Amazon Comprehend to extract entitie
L. Use Amazon Bedrock to answer questions over indexed dat
M. Use IAM for access control and CloudTrail for audit logging.
N. Configure Amazon O Business to federate data across Amazon S3, Amazon Kinesis, and Amazon SageMaker Feature Stor
O. Use EventBridge for ingestion orchestratio
P. Use custom AWS Lambda functions to filter LLM outputs for ethical compliance.

**Answer:** B


**NEW QUESTION 17**
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-Ra
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 plugi
L. Ingest metrics and logs through Amazon Kinesis and analyze behavior with custom queries.

**Answer:** C


**NEW QUESTION 20**
A company uses Amazon Bedrock to build a Retrieval Augmented Generation (RAG) system. The RAG system uses an Amazon Bedrock Knowledge Bases 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 Servic
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 22**
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 latenc
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 chainin
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 ElastiCach
H. Use AWS Lambda functions to pre- process metadata and to trace end-to-end latenc
I. Use AWS X-Ray to identify and remediate performance bottlenecks.
J. Use Amazon Kendra to improve roast log retrieval accurac
K. Store normalized prompt metadata within Amazon DynamoD
L. Use AWS Step Functions to orchestrate multi-step prompts.

**Answer:** A

**NEW QUESTION 27**
A financial services company is developing a generative AI (GenAI) application that serves both premium customers and standard customers. The application uses AWS Lambda functions behind an Amazon API Gateway REST API to process requests. The company needs to dynamically switch between AI models based on which customer tier each user belongs to. The company also wants to perform A/B testing for new features without redeploying code. The company needs to validate model parameters like temperature and maximum token limits before applying changes.
Which solution will meet these requirements with the LEAST operational overhead?

A. Create AWS Systems Manager Parameter Store parameters for each configuratio
B. Use Lambda functions to poll for parameter update
C. Use Amazon EventBridge events to trigger redeployments when configurations change.
D. Store model configurations in Amazon DynamoDB table
E. Optimize access patterns to retrieve configurations according to customer tie
F. Configure Lambda functions to query DynamoDB at the beginning of each request to determine which model to use.
G. Use AWS AppConfig to manage model configuration
H. Use feature flags to perform A/B testin
I. Define JSON schema validation rules for model parameter
J. Configure Lambda functions to retrieve configurations by using the AWS AppConfig Agent.
K. Create an Amazon ElastiCache (Redis OSS) cluster to store model configuration
L. Set short TTL value
M. Run custom validation logic in Lambda function
N. Use Amazon CloudWatch metrics to monitor configuration usage.

**Answer:** C

**NEW QUESTION 32**
A company developed a multimodal content analysis application by using Amazon Bedrock. The application routes different content types (text, images, and code) to specialized foundation models (FMs).
The application needs to handle multiple types of routing decisions. Simple routing based on file extension must have minimal latency. Complex routing based on content semantics requires analysis before FM selection. The application must provide detailed history and support fallback options when primary FMs fail.
Which solution will meet these requirements?

A. Configure AWS Lambda functions that call Amazon Bedrock FMs for all routing logi
B. Use conditional statements to determine the appropriate FM based on content type and semantics.
C. Create a hybrid solutio
D. Handle simple routing based on file extensions in application cod
E. Handle complex content-based routing by using an AWS Step Functions state machine with JSONata for content analysis and the InvokeModel API for specialized FMs.
F. Deploy separate AWS Step Functions workflows for each content type with routing logic in AWS Lambda function
G. Use Amazon EventBridge to coordinate between workflows when fallback to alternate FMs is required.
H. Use Amazon SQS with different SQS queues for each content typ
I. Configure AWS Lambda consumers that analyze content and invoke appropriate FMs based on message attributes by using Amazon Bedrock with an AWS SDK.

**Answer:** B

**NEW QUESTION 33**
A company wants to select a new FM for its AI assistant. A GenAI developer needs to generate evaluation reports to help a data scientist assess the quality and safety of various foundation models FMs. The data scientist provides the GenAI developer with sample prompts for evaluation. The GenAI developer wants to use Amazon Bedrock to automate report generation and evaluation.
Which solution will meet this requirement?

A. Combine the sample prompts into a single JSON documen
B. Create an Amazon Bedrock knowledge base with the documen
C. Write a prompt that asks the FM to generate a response to each sample promp
D. Use the RetrieveAndGenerate API to generate a report for each model.
E. Combine the sample prompts into a single JSONL documen
F. Store the document in an Amazon S3 bucke
G. Create an Amazon Bedrock evaluation job that uses a judge mode
H. Specify the S3 location as input and a different S3 location as outpu
I. Run an evaluation job for each FM and select the FM as the generator.
J. Combine the sample prompts into a single JSONL documen
K. Store the document in an Amazon S3 bucke
L. Create an Amazon Bedrock evaluation job that uses a judge mode
M. Specify the S3 location as input and Amazon QuickSight as outpu
N. Run an evaluation job for each FM and select the FM as the evaluator.
O. Combine the sample prompts into a single JSON documen
P. Create an Amazon Bedrock knowledge base from the documen
Q. Create an Amazon Bedrock evaluation job that uses the retrieval and response generation evaluation typ
R. Specify an Amazon S3 bucket as the outpu
S. Run an evaluation job for each FM.

**Answer:** B


**NEW QUESTION 37**
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 Kendr
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 38**
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 Bedroc
B. Purchase provisioned throughput and optimize for batch processing.
C. Deploy a low-latency, real-time optimized model on Amazon Bedroc
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 40**
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 rol
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 mode
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 polic
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 polic
K. Use stack sets to deploy the guardrail to each account in the organization.

**Answer:** CD


**NEW QUESTION 43**
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 serve
B. Grant the AI agent Lambda functions permission to invoke the Lambda function that hosts the MCP serve
C. Configure the AI agent??s MCP client to invoke the MCP server asynchronously.
D. Use a Lambda function to host the MCP serve
E. Grant the AI agent Lambda functions permission to invoke the Lambda function that hosts the MCP serve
F. Configure the AI agent to use the STDIO transport with the MCP server.
G. Use a Lambda function to host the MCP serve
H. Create an Amazon API Gateway HTTP API that proxies requests to the Lambda functio
I. Configure the AI agent solution to use the Streamable HTTP transport to make requests through the HTTP AP
J. Use Amazon Cognito to enforce OAuth 2.1.
K. Use a Lambda layer to host the MCP serve
L. Add the Lambda layer to the AI agent Lambda function
M. Configure the agentic AI solution to use the STDIO transport to send requests to the MCP serve
N. In the AI agent??s MCP configuration, specify the Lambda layer ARN as the comman
O. Specify the user credentials as environment variables.

**Answer:** C


**NEW QUESTION 47**
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 PostgreSQL.
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 Servic
B. Implement keyword-based search rules that use custom analyzers and relevance tuning to find restaurants based on attributes such as cuisine type, features, and locatio
C. Create Amazon API Gateway HTTP API endpoints to transform user queries into structured search parameters.
D. Migrate the restaurant data to Amazon OpenSearch Servic
E. Use a foundation model (FM) in Amazon Bedrock to generate vector embeddings from restaurant descriptions, reviews, and menu item
F. When users submit natural language queries, convert the queries to embeddings by using the same F
G. Perform k-nearest neighbors (k-NN) searches to find semantically similar results.
H. Keep the restaurant data in PostgreSQL and implement a pgvector extensio
I. Use a foundation model (FM) in Amazon Bedrock to generate vector embeddings from restaurant dat
J. Store the vector embeddings directly in PostgreSQ
K. Create an AWS Lambda function to convert natural language queries to vector representations by using the same F
L. Configure the Lambda function to perform similarity searches within the database.
M. Migrate restaurant data to an Amazon Bedrock knowledge base by using a custom ingestion pipelin
N. Configure the knowledge base to automatically generate embeddings from restaurant informatio
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:** B


**NEW QUESTION 48**
A company is building a generative AI (GenAI) application that processes financial reports and provides summaries for analysts. The application must run two compute environments. In one environment, AWS Lambda functions must use the Python SDK to analyze reports on demand. In the second environment, Amazon EKS containers must use the JavaScript SDK to batch process multiple reports on a schedule. The application must maintain conversational context throughout multi-turn interactions, use the same foundation model (FM) across environments, and ensure consistent authentication.
Which solution will meet these requirements?

A. Use the Amazon Bedrock InvokeModel API with a separate authentication method for each environmen
B. Store conversation states in Amazon DynamoD
C. Use custom I/O formatting logic for each programming language.
D. Use the Amazon Bedrock Converse API directly in both environments with a common authentication mechanism that uses IAM role
E. Store conversation states in Amazon ElastiCach
F. Create programming language-specific wrappers for model parameters.
G. Create a centralized Amazon API Gateway REST API endpoint that handles all model interactions by using the InvokeModel AP
H. Store interaction history in application process memory in each Lambda function or EKS containe
I. Use environment variables to configure model parameters.
J. Use the Amazon Bedrock Converse API and IAM roles for authenticatio
K. Pass previous messages in the request messages array to maintain conversational contex
L. Use programming language-specific SDKs to establish consistent API interfaces.

**Answer:** D


**NEW QUESTION 49**
A company is using Amazon Bedrock to develop an AI-powered application that uses a foundation model 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 Region
B. Use a custom routing layer that directs traffic based on user locatio
C. Configure Amazon CloudWatch alarms to monitor Regional service usag
D. Use Amazon SNS to send email alerts to the company when usage approaches specified thresholds.
E. Use Amazon Bedrock cross-Region inference profiles by specifying geographical codes in profile IDs when the application calls the InvokeModel AP

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 throttlin
H. Configure the Lambda functions to call the foundation model in the nearest secondary Region when the application reaches service quotas in the primary Regio
I. Use intelligent routing to ensure compliance with data residency requirements.
J. Configure provisioned throughput for Amazon Bedrock in multiple Region
K. Implement failover logic in the application code to switch between Regions when throttling occur
L. Use AWS Global Accelerator to route traffic to the appropriate endpoints based on user location.

**Answer:** B

NEW QUESTION 50
An ecommerce company is building an internal platform to develop generative AI applications by using Amazon Bedrock foundation models (FMs). Developers need to select models based on evaluations that are aligned to ecommerce use cases. The platform must display accuracy metrics for text generation and summarization in dashboards. The company has custom ecommerce datasets to use as standardized evaluation inputs.
Which combination of steps will meet these requirements with the LEAST operational overhead? (Select TWO.)

A. Import the datasets to an Amazon S3 bucke
B. Provide appropriate IAM permissions and cross-origin resource sharing (CORS) permissions to give the evaluation jobs access to the datasets.
C. Import the datasets to an Amazon S3 bucke
D. Provide appropriate IAM permissions and a VPC endpoint configuration to give the evaluation jobs access to the datasets.
E. Configure an AWS Lambda function to create model evaluation jobs on a schedule in the Amazon Bedrock consol
F. Provide the URI of the S3 bucket that contains the datasets as an inpu
G. Configure the evaluation jobs to measure the real world knowledge (RWK) score for text generation and BERTScore for summarizatio
H. Configure a second Lambda function to check the status of the jobs and publish custom logs to Amazon CloudWatc
I. Create a custom Amazon CloudWatch Logs Insights dashboard.
J. Use Amazon SageMaker Clarify on a schedule to create model evaluation job
K. Useopen source frameworks to create and run standardized evaluation
L. Publish results to Amazon CloudWatch namespace
M. Use an AWS Lambda function to check the status of the jobs and publish custom logs to Amazon CloudWatc
N. Create a custom Amazon CloudWatch Logs Insights dashboard.
O. Run an Amazon SageMaker AI notebook job on a schedule by using the fmvelos or ragas framework to run evaluations that use the datasets in the S3 bucke
P. Write Python code in the notebook that makes direct InvokeModel API calls to the FMs and processes their responses for evaluatio
Q. Publish job status and results to Amazon CloudWatch Logs to measure the real world knowledge (RWK) score for text generation and toxicity for summarization as metrics for accurac
R. Create a custom CloudWatch Logs Insights dashboard.

**Answer:** BC

NEW QUESTION 51
A company is using AWS Lambda and REST APIs to build a reasoning agent to automate support workflows. The system must preserve memory across interactions, share relevant agent state, and support event-driven invocation and synchronous invocation. The system must also enforce access control and session-based permissions.
Which combination of steps provides the MOST scalable solution? (Select TWO.)

A. Use Amazon Bedrock AgentCore to manage memory and session-aware reasonin
B. Deploy the agent with built-in identity support, event handling, and observability.
C. Register the Lambda functions and REST APIs as actions by using Amazon API Gateway and Amazon EventBridg
D. Enable Amazon Bedrock AgentCore to invoke the Lambda functions and REST APIs without custom orchestration code.
E. Use Amazon Bedrock Agents for reasoning and conversation managemen
F. Use AWS Step Functions and Amazon SQS for orchestratio
G. Store agent state in Amazon DynamoDB.
H. Deploy the reasoning logic as a container on Amazon ECS behind API Gatewa
I. Use Amazon Aurora to store memory and identity data.
J. Build a custom RAG pipeline by using Amazon Kendra and Amazon Bedroc
K. Use AWS Lambda to orchestrate tool invocation
L. Store agent state in Amazon S3.

**Answer:** AB

NEW QUESTION 53
A company is developing a customer support application that uses Amazon Bedrock foundation models (FMs) to provide real-time AI assistance to the company??s employees. The application must display AI-generated responses character by character as the responses are generated. The application needs to support thousands of concurrent users with minimal latency. The responses typically take 15 to 45 seconds to finish.
Which solution will meet these requirements?

A. Configure an Amazon API Gateway WebSocket API with an AWS Lambda integratio
B. Configure the WebSocket API to invoke the Amazon Bedrock InvokeModelWithResponseStream API and stream partial responses through WebSocket connections.
C. Configure an Amazon API Gateway REST API with an AWS Lambda integratio
D. Configure the REST API to invoke the Amazon Bedrock standard InvokeModel API and implement frontend client-side polling every 100 ms for complete response chunks.
E. Implement direct frontend client connections to Amazon Bedrock by using IAM user credentials and the InvokeModelWithResponseStream API without any intermediate gateway or proxy layer.
F. Configure an Amazon API Gateway HTTP API with an AWS Lambda integratio
G. Configure the HTTP API to cache complete responses in an Amazon DynamoDB table and serve the responses through multiple paginated GET requests to frontend clients.

**Answer:** A

**NEW QUESTION 54**
A financial services company is developing a customer service AI assistant application that uses a foundation model (FM) in Amazon Bedrock. The application must provide transparent responses by documenting reasoning and by citing sources that are used for Retrieval Augmented Generation (RAG). The application must capture comprehensive audit trails for all responses to users. The application must be able to serve up to 10,000 concurrent users and must respond to each customer inquiry within 2 seconds.
Which solution will meet these requirements with the LEAST operational overhead?

A. Enable tracing for Amazon Bedrock Agent
B. Configure structured prompts that direct the FM to provide evidence presentation
C. Integrate Amazon Bedrock Knowledge Bases with data sources to enable RA
D. Configure the application to reference and cite authoritative conten
E. Deploy the application in a Multi-AZ architectur
F. Use Amazon API Gateway and AWS Lambda functions to scale the applicatio
G. Use Amazon CloudFront to provide low- latency delivery.
H. Enable tracing for Amazon Bedrock agent
I. Integrate a custom RAG pipeline with Amazon OpenSearch Service to retrieve and cite source
J. Configure structured prompts to present retrieved evidenc
K. Deploy the application behind an Amazon API Gateway REST AP
L. Use AWS Lambda functions and Amazon CloudFront to scale the application and to provide low latenc
M. Store logs in Amazon S3 and use AWS CloudTrail to capture audit trails.
N. Use Amazon CloudWatch to monitor latency and error rate
O. Embed model prompts directly in the application backend to cite source
P. Store application interactions with users in Amazon RDS for audits.
Q. Store generated responses and supporting evidence in an Amazon S3 bucke
R. Enable versioning on the bucket for audit
S. Use AWS Glue to catalog retrieved document
T. Process the retrieved documents in Amazon Athena to generate periodic compliance reports.

**Answer:** A


**NEW QUESTION 58**
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 F
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 endpoin
E. Integrate the REST API with the business application.
F. Configure the FM to use cross-Region inference through a Europe-scoped endpoin
G. Configure an Amazon Bedrock VPC endpoin
H. Extend existing SCPs to allow employees to use the FM through inference profiles in Europe-based Regions where the FM is availabl
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 endpoin
K. Extend existing SCPs to allow employees to use the SageMaker endpoin
L. Integrate the FM in SageMaker with the business application.

**Answer:** C


**NEW QUESTION 60**
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 conten
C. Use prompt safety classification for inappropriate advice solicitatio
D. Use personally identifiable information (PII) detection without redaction.
E. Use custom classification to build an FM that detects offensive content and inappropriate advice solicitatio
F. Apply personally identifiable information (PII) detection as a secondary filter only when messages pass the custom classifier.
G. Deploy a multi-stage proces
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 mod
I. Route flagged messages through Amazon EventBridge for human review.
J. Use toxicity detection with thresholds configured to 0.5 for all categorie
K. Use parallel processing for both prompt safety classification and personally identifiable information (PII) detection with entity redactio
L. Apply Amazon CloudWatch alarms to filter metrics.

**Answer:** D

**NEW QUESTION 62**

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 paralle
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 onl
E. Configure Amazon Bedrock guardrails to restrict prompt variatio
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 agen
H. Store intermediate prompts in Amazon DynamoD
I. Resolve agent ordering by using TTL-based backoff and retries.
J. Use Amazon EventBridge Pipes to invoke resolvers based on Amazon CloudWatch log pattern
K. Store response metadata in DynamoDB with TTL and versioned write
L. Use Amazon Q Developer to dynamically generate fallback prompts.

**Answer:** A

**NEW QUESTION 64**
......

# Thank You for Trying Our Product

## We offer two products:

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

2nd - Questons 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