



Amazon-Web-Services

Exam Questions AIP-C01

AWS Certified Generative AI Developer - Professional

About ExamBible

Your Partner of IT Exam

Found in 1998

ExamBible is a company specialized on providing high quality IT exam practice study materials, especially Cisco CCNA, CCDA, CCNP, CCIE, Checkpoint CCSE, CompTIA A+, Network+ certification practice exams and so on. We guarantee that the candidates will not only pass any IT exam at the first attempt but also get profound understanding about the certificates they have got. There are so many alike companies in this industry, however, ExamBible has its unique advantages that other companies could not achieve.

Our Advances

* 99.9% Uptime

All examinations will be up to date.

* 24/7 Quality Support

We will provide service round the clock.

* 100% Pass Rate

Our guarantee that you will pass the exam.

* Unique Gurantee

If you do not pass the exam at the first time, we will not only arrange FULL REFUND for you, but also provide you another exam of your claim, ABSOLUTELY FREE!

NEW QUESTION 1

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

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

Which solution will meet these requirements?

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

Answer: C

NEW QUESTION 2

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 base
- G. Implement Retrieval Augmented Generation RA
- H. Set the PerformanceConfigLatency parameter to optimized.
- I. Store product catalog data in Amazon OpenSearch Service
- J. Validate the model's product recommendations against the product catalog
- K. Use Amazon DynamoDB to implement response caching.

Answer: C

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 retries
- 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 jitter
- D. Use AWS X-Ray tracing with annotations to identify and filter service components.
- E. Implement client-side caching of all FM responses
- F. Add custom logging statements in the application code to record API call durations.
- G. Configure the AWS SDK with adaptive retry mode
- H. Use AWS CloudTrail distributed tracing to monitor throttling events.

Answer: B

NEW QUESTION 4

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 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 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 6

A company is using Amazon Bedrock to design an application to help researchers apply for grants. The application is based on an Amazon Nova Pro foundation model (FM). The application contains four required inputs and must provide responses in a consistent text format. The company wants to receive a notification in Amazon Bedrock if a response contains bullying language. However, the company does not want to block all flagged responses.

The company creates an Amazon Bedrock flow that takes an input prompt and sends it to the Amazon Nova Pro FM. The Amazon Nova Pro FM provides a response.

Which additional steps must the company take to meet these requirements? (Select TWO.)

- A. Use Amazon Bedrock Prompt Management to specify the required inputs as variables
- B. Select an Amazon Nova Pro F
- C. Specify the output format for the responses
- D. Add the prompt to the prompts node of the flow.
- E. Create an Amazon Bedrock guardrail that applies the hate content filter
- F. Set the filter response to block
- G. Add the guardrail to the prompts node of the flow.
- H. Create an Amazon Bedrock prompt route
- I. Specify an Amazon Nova Pro F
- J. Add the required inputs as variables to the input node of the flow
- K. Add the prompt router to the prompts node
- L. Add the output format to the output node.
- M. Create an Amazon Bedrock guardrail that applies the insults content filter
- N. Set the filter response to detect
- O. Add the guardrail to the prompts node of the flow.
- P. Create an Amazon Bedrock application inference profile that specifies an Amazon Nova Pro F
- Q. Specify the output format for the response in the description
- R. Include a tag for each of the input variables
- S. Add the profile to the prompts node of the flow.

Answer: AD

NEW QUESTION 7

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

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

Which solution will meet these requirements?

- A. Isolate data for each agent by using separate knowledge bases
- B. Use IAM filtering to control access to each knowledge base
- C. Deploy a supervisor agent to perform natural language intent classification on patient inquiries
- D. Configure the supervisor agent to route queries to specialized collaborator agents to respond to department-specific queries
- 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 department
- G. Configure individual collaborator agents to perform natural language intent classification for each specialty domain within each department
- H. Integrate each collaborator agent with department-specific knowledge bases only
- I. Implement manual handoff processes between the supervisor agents.
- J. Isolate data for each department in separate knowledge bases
- K. Use IAM filtering to control access to each knowledge base

- L. Deploy a single general-purpose agent
- 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 department
- P. Configure multiple collaborator agents for each supervisor agent
- Q. Integrate all agents with the same knowledge base
- R. Use external routing logic to merge responses from multiple supervisor agents.

Answer: A

NEW QUESTION 8

A healthcare company uses Amazon Bedrock to deploy an application that generates summaries of clinical documents. The application experiences inconsistent response quality with occasional factual hallucinations. Monthly costs exceed the company's projections by 40%. A GenAI developer must implement a near real-time monitoring solution to detect hallucinations, identify abnormal token consumption, and provide early warnings of cost anomalies. The solution must require minimal custom development work and maintenance overhead.

Which solution will meet these requirements?

- A. Configure Amazon CloudWatch alarms to monitor InputTokenCount and OutputTokenCount metrics to detect anomalies
- B. Store model invocation logs in an Amazon S3 bucket
- C. Use AWS Glue and Amazon Athena to identify potential hallucinations.
- D. Run Amazon Bedrock evaluation jobs that use LLM-based judgments to detect hallucination
- E. Configure Amazon CloudWatch to track token usage
- F. Create an AWS Lambda function to process CloudWatch metrics
- G. Configure the Lambda function to send usage pattern notifications.
- H. Configure Amazon Bedrock to store model invocation logs in an Amazon S3 bucket
- I. Enable text output logging
- J. Configure Amazon Bedrock guardrails to run contextual grounding checks to detect hallucination
- K. Create Amazon CloudWatch anomaly detection alarms for token usage metrics.
- L. Use AWS CloudTrail to log all Amazon Bedrock API calls
- M. Create a custom dashboard in Amazon QuickSight to visualize token usage patterns
- N. Use Amazon SageMaker Model Monitor to detect quality drift in generated summaries.

Answer: C

NEW QUESTION 9

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 10

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 files
- B. Store embeddings in the OpenSearch Neural plugin for Amazon OpenSearch Service.
- C. Use Amazon Comprehend to extract topics from the digitized files
- 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 files
- 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 files
- J. Store embeddings in an Amazon Aurora PostgreSQL Serverless database that has the pgvector extension.

Answer: D

NEW QUESTION 10

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 role
- 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 permission
- 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 13

An ecommerce company is using Amazon Bedrock to build a generative AI (GenAI) application. The application uses AWS Step Functions to orchestrate a multi-agent workflow to produce detailed product descriptions. The workflow consists of three sequential states: a description generator, a technical specifications validator, and a brand voice consistency checker. Each state produces intermediate reasoning traces and outputs that are passed to the next state. The application uses an Amazon S3 bucket for process storage and to store outputs. During testing, the company discovers that outputs between Step Functions states frequently exceed the 256 KB quota and cause workflow failures. A GenAI Developer needs to revise the application architecture to efficiently handle the Step Functions 256 KB quota and maintain workflow observability. The revised architecture must preserve the existing multi-agent reasoning and acting (ReAct) pattern. Which solution will meet these requirements with the LEAST operational overhead?

- A. Store intermediate outputs in Amazon DynamoDB
- B. Pass only references between state
- C. Create a Map state that retrieves the complete data from DynamoDB when required for each agent's processing step.
- D. Configure an Amazon Bedrock integration to use the S3 bucket URI in the input parameters for large output
- E. Use the ResultPath and ResultSelector fields to route S3 references between the agent steps while maintaining the sequential validation workflow.
- F. Use AWS Lambda functions to compress outputs to less than 256 KB before each agent state
- G. Configure each agent task to decompress outputs before processing and to compress results before passing them to the next state.
- H. Configure a separate Step Functions state machine to handle each agent's processing
- I. Use Amazon EventBridge to coordinate the execution flow between state machines
- J. Use S3 references for the outputs as event data.

Answer: B

NEW QUESTION 15

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 documents
- 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 documents
- 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 entities
- 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 documents
- H. Implement verification logic to compare against retrieved sources and to cite sources.

Answer: B

NEW QUESTION 20

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 Management
 - E. Use Amazon Bedrock Flows to deploy the prompt variants with defined traffic allocation
 - F. Configure Amazon Bedrock guardrails to monitor demographic fairness
 - 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 CloudWatch
 - 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 variants
 - L. Enable model invocation logging in Amazon CloudWatch
 - M. Set up CloudWatch alarms for InvocationsIntervened metrics with a dimension for each demographic group.

Answer: B

NEW QUESTION 21

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 metrics
- 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 25

A financial services company needs to build a document analysis system that uses Amazon Bedrock to process quarterly reports. The system must analyze financial data, perform sentiment analysis, and validate compliance across batches of reports. Each batch contains 5 reports. Each report requires multiple foundation model (FM) calls. The solution must finish the analysis within 10 seconds for each batch. Current sequential processing takes 45 seconds for each batch. Which solution will meet these requirements?

- A. Use AWS Lambda functions with provisioned concurrency to process each analysis type sequentially
- B. Configure the Lambda function timeouts to 10 seconds
- C. Configure automatic retries with exponential backoff.
- D. Use AWS Step Functions with a Parallel state to invoke separate AWS Lambda functions for each analysis type simultaneously
- E. Configure Amazon Bedrock client timeout
- F. Use Amazon CloudWatch metrics to track execution time and model inference latency.
- G. Create an Amazon SQS queue to buffer analysis requests
- H. Deploy multiple AWS Lambda functions with reserved concurrency
- I. Configure each Lambda function to process different aspects of each report sequentially and then combine the results.
- J. Deploy an Amazon ECS cluster that runs containers that process each report sequentially
- K. Use a load balancer to distribute batch workload
- L. Configure an auto-scaling policy based on CPU utilization.

Answer: B

NEW QUESTION 30

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 expressions
- 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 Store
- C. Use IAM for access control
- D. Use AWS CloudTrail for audit logging.
- E. Use Amazon SageMaker Processing and Amazon Transcribe to pre-process multimodal data
- 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 templates
- I. Use AWS CloudTrail to log research activity for audits.
- J. Use Amazon OpenSearch Serverless to index behavioral logs and telemetry
- K. Use Amazon Comprehend to extract entities
- L. Use Amazon Bedrock to answer questions over indexed data

- 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 35

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 qualit
- 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 appropriatenes
- D. Configure custom Amazon Bedrock guardrails to check responses for compliance with financial policie
- 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 AWSLambda functions to check responses against a static compliance databas
- 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 assistan
- K. Configure CloudWatch alerts for potential compliance violation
- L. Establish a team of human evaluators to review flagged interactions.

Answer: B

NEW QUESTION 37

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 usag
- E. Store all prompts and outputs in Amazon S3 with server-side encryption with AWS KMS keys (SSE-KMS).
- F. Define access control by using IAM with attribute-based access control (ABAC) to map departments to specific permission
- G. Configure VPC endpoints for Amazon Bedrock model API call
- H. Use IAM condition keys to enforce specific GuardrailIdentifier and ModelId value
- I. Configure AWS CloudTrail to capture management and data events for S3 objects and KMS key usage activitie
- J. Enable S3 server access logging to record detailed file-level interactions with the video archive
- K. Send all CloudTrail logs to AWS CloudTrail Lak
- 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 policie
- N. Use AWS Config to track resource changes and compliance with security rule
- O. Use server-side encryption with AWS KMS keys (SSE-KMS) to encrypt data at res
- 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 KM
- S. Deploy Amazon Macie to scan and classify the video archiv
- T. Use server-side encryption with AWS KMS keys (SSE-KMS) to encrypt all stored dat
- . Configure CloudTrail to capture KMS API usage events for audit purpose
- . Configure Amazon EventBridge rules to process CloudTrail Insights anomalies and Macie finding
- . Use CloudWatch alarms to trigger automated notifications and security responses when potential security issues are detected.

Answer: B

NEW QUESTION 38

A financial services company uses an AI application to process financial documents by using Amazon Bedrock. During business hours, the application handles approximately 10,000 requests each hour, which requires consistent throughput.

The company uses the CreateProvisionedModelThroughput API to purchase provisioned throughput. Amazon CloudWatch metrics show that the provisioned capacity is unused while on-demand requests are being throttled. The company finds the following code in the application:

```
response = bedrock_runtime.invoke_model(modelId="anthropic.claude-v2", body=json.dumps(payload))
```

The company needs the application to use the provisioned throughput and to resolve the throttling issues.

Which solution will meet these requirements?

- A. Increase the number of model units (MUs) in the provisioned throughput configuration.
- B. Replace the model ID parameter with the ARN of the provisioned model that the CreateProvisionedModelThroughput API returns.
- C. Add exponential backoff retry logic to handle throttling exceptions during peak hours.
- D. Modify the application to use the invokeModelWithResponseStream API instead of the invokeModel API.

Answer: B

NEW QUESTION 39

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 40

A financial services company wants to develop an Amazon Bedrock application that gives analysts the ability to query quarterly earnings reports and financial statements. The financial documents are typically 5–100 pages long and contain both tabular data and text. The application must provide contextually accurate responses that preserve the relationship between financial metrics and their explanatory text. To support accurate and scalable retrieval, the application must incorporate document segmentation and context management strategies. Which solution will meet these requirements?

- A. Use a direct model invocation approach that uses Anthropic Claude to process each financial document as a single input
- B. Use fine-tuned prompts that instruct the model to parse tables and text separately.
- C. Use Amazon Bedrock Knowledge Bases to create a Retrieval Augmented Generation (RAG) application that retrieves relevant information from contextually chunked sections of financial document
- D. Segment documents based on their structural layout
- E. Include citations that reference the original source materials.
- F. Deploy an Amazon Bedrock agent that has an action group that calls custom AWS Lambda functions to analyze financial document
- G. Configure the Lambda functions to perform fixed-size chunking when a user submits a query about financial metrics.
- H. Create one specialized Amazon Bedrock application that is optimized for structured data
- I. Create a second application that is optimized for unstructured data
- J. Configure each application to use a tailored chunking strategy that is suited to the application's content type
- K. Implement logic to link queries to the appropriate sources.

Answer: B

NEW QUESTION 41

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 46

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 49

An ecommerce company is developing a generative AI (GenAI) solution that uses Amazon Bedrock with Anthropic Claude to recommend products to customers. Customers report that some recommended products are not available for sale or are not relevant. Customers also report long response times for some recommendations. The company confirms that most customer interactions are unique and that the solution recommends products not present in the product catalog. Which solution will meet this requirement?

- A. Increase grounding within Amazon Bedrock Guardrail
- B. Enable automated reasoning check
- C. Set up provisioned throughput.
- D. Use prompt engineering to restrict model responses to relevant product
- E. Use streaming inference to reduce perceived latency.
- F. Create an Amazon Bedrock Knowledge Bases and implement Retrieval Augmented Generation (RAG). Set the PerformanceConfigLatency parameter to optimized.
- G. Store product catalog data in Amazon OpenSearch Service
- H. Validate model recommendations against the catalog
- I. Use Amazon DynamoDB for response caching.

Answer: C

NEW QUESTION 51

A pharmaceutical company is developing a Retrieval Augmented Generation (RAG) 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 55

A healthcare company is using Amazon Bedrock to build a Retrieval Augmented Generation (RAG) application that helps practitioners make clinical decisions. The application must achieve high accuracy for patient information retrievals, identify hallucinations in generated content, and reduce human review costs. Which solution will meet these requirements?

- A. Use Amazon Comprehend to analyze and classify RAG responses and to extract medical entities and relationships
- B. Use AWS Step Functions to orchestrate automated evaluation
- C. Configure Amazon CloudWatch metrics to track entity recognition confidence score
- D. Configure CloudWatch to send an alert when accuracy falls below specified thresholds.
- E. Implement automated large language model (LLM)-based evaluations that use a specialized model that is fine-tuned for medical content to assess all responses
- F. Deploy AWS Lambda functions to parallelize evaluation
- G. Publish results to Amazon CloudWatch metrics that track relevance and factual accuracy.
- H. Configure Amazon CloudWatch Synthetics to generate test queries that have known answers on a regular schedule, and track model success rate
- I. Set up dashboards that compare synthetic test results against expected outcomes.
- J. Deploy a hybrid evaluation system that uses an automated LLM-as-a-judge evaluation to initially screen responses and targeted human reviews for edge cases
- K. Use a built-in Amazon Bedrock evaluation to track retrieval precision and hallucination rates.

Answer: D

NEW QUESTION 57

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

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

Answer: A

NEW QUESTION 61

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 64

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 68

A financial services company is creating a Retrieval Augmented Generation (RAG) application that uses Amazon Bedrock to generate summaries of market activities. The application relies on a vector database that stores a small proprietary dataset with a low index count. The application must perform similarity searches. The Amazon Bedrock model's responses must maximize accuracy and maintain high performance.

The company needs to configure the vector database and integrate it with the application. Which solution will meet these requirements?

- A. Launch an Amazon MemoryDB cluster and configure the index by using the Flat algorithm
- B. Configure a horizontal scaling policy based on performance metrics.
- C. Launch an Amazon MemoryDB cluster and configure the index by using the Hierarchical Navigable Small World (HNSW) algorithm
- D. Configure a vertical scaling policy based on performance metrics.
- E. Launch an Amazon Aurora PostgreSQL cluster and configure the index by using the Inverted File with Flat Compression (IVFFlat) algorithm
- F. Configure the instance class to scale to a larger size when the load increases.
- G. Launch an Amazon DocumentDB cluster that has an IVFFlat index and a high probe value
- H. Configure connections to the cluster as a replica set
- I. Distribute reads to replica instances.

Answer: B

NEW QUESTION 73

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

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

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

Answer: A

NEW QUESTION 76

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 reasoning
- 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 EventBridge
- 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 management
- F. Use AWS Step Functions and Amazon SQS for orchestration
- G. Store agent state in Amazon DynamoDB.
- H. Deploy the reasoning logic as a container on Amazon ECS behind API Gateway
- I. Use Amazon Aurora to store memory and identity data.
- J. Build a custom RAG pipeline by using Amazon Kendra and Amazon Bedrock
- K. Use AWS Lambda to orchestrate tool invocation
- L. Store agent state in Amazon S3.

Answer: AB

NEW QUESTION 80

A legal research company has a Retrieval Augmented Generation (RAG) application that uses Amazon Bedrock and Amazon OpenSearch Service. The application stores 768-dimensional vector embeddings for 15 million legal documents, including statutes, court rulings, and case summaries. The company's current chunking strategy segments text into fixed-length blocks of 500 tokens. The current chunking strategy often splits contextually linked information such as legal arguments, court opinions, or statute references across separate chunks. Researchers report that generated outputs frequently omit key context or cite outdated legal information.

Recent application logs show a 40% increase in response times. The p95 latency metric exceeds 2 seconds. The company expects storage needs for the application to grow from 90 GB to 360 GB within a year.

The company needs a solution to improve retrieval relevance and system performance at scale.

Which solution will meet these requirements?

- A. Increase the embedding vector dimensionality from 768 to 4,096 without changing the existing chunking or pre-processing strategy.
- B. Replace dynamic retrieval with static, pre-written summaries that are stored in Amazon S3. Use Amazon CloudFront to serve the summaries to reduce compute demand and improve predictability.
- C. Update the chunking strategy to use semantic boundaries such as complete legal arguments, clauses, or sections rather than fixed token limit
- D. Regenerate vector embeddings to align with the new chunk structure.
- E. Migrate from OpenSearch Service to Amazon DynamoDB
- F. Implement keyword-based indexes to enable faster lookups for legal concepts.

Answer: C

NEW QUESTION 81

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 recommendations
- 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 85

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 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 the 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: D

NEW QUESTION 87

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 88

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 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 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 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 90

A company is developing a generative AI (GenAI) application by using Amazon Bedrock. The application will analyze patterns and relationships in the company's data. The application will process millions of new data points daily across AWS Regions in Europe, North America, and Asia before storing the data in Amazon S3. The application must comply with local data protection and storage regulations. Data residency and processing must occur within the same continent. The application must also maintain audit trails of the application's decision-making processes and provide data classification capabilities.

Which solution will meet these requirements?

- A. Deploy the application in each Region with local IAM policie
- B. Use Amazon Bedrock cross-Region inference to distribute the workloa
- C. Use Amazon CloudWatch to log AI decision-making processe
- D. Manually track compliance certifications across Regions.
- E. Use SCPs with AWS Organizations to manage location-specific permission

- F. Use AWS CloudTrail immutable logs to audit decision-making processes
- G. Import a custom model into Amazon Bedrock and deploy the model to each Region.
- H. Use Amazon S3 Object Lock with Region-specific S3 bucket policies
- I. Pre-process the data points within the Region based on geographic origin before sending the data points to Amazon Bedrock
- J. Use Amazon Macie to classify the data
- K. Use AWS CloudTrail immutable logs to audit the decision-making processes.
- L. Create separate AWS accounts for each Region with individual compliance frameworks
- M. Use Amazon SageMaker AI with custom monitoring
- N. Create manual compliance reports for each regulatory jurisdiction.

Answer: C

NEW QUESTION 91

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 93

A company is designing a solution that uses foundation models (FMs) to support multiple AI workloads. Some FMs must be invoked on demand and in real time. Other FMs require consistent high-throughput access for batch processing.

The solution must support hybrid deployment patterns and run workloads across cloud infrastructure and on-premises infrastructure to comply with data residency and compliance requirements.

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

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

Answer: BC

NEW QUESTION 95

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

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

Which solution will meet these requirements?

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

Answer: B

NEW QUESTION 96

.....

Relate Links

100% Pass Your AIP-C01 Exam with Exam Bible Prep Materials

<https://www.exambible.com/AIP-C01-exam/>

Contact us

We are proud of our high-quality customer service, which serves you around the clock 24/7.

Viste - <https://www.exambible.com/>