

Exam Questions AIP-C01

AWS Certified Generative AI Developer - Professional

<https://www.2passeasy.com/dumps/AIP-C01/>



NEW QUESTION 1

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

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

Which solution will meet this requirement?

- A. Increase grounding within Amazon Bedrock Guardrail
- B. Enable Automated Reasoningcheck
- C. Set up provisioned throughput.
- D. Use prompt engineering to restrict the model responses to relevant product
- E. Use streaming techniques such as the InvokeModelWithResponseStream action to reduce perceived latency for the customers.
- F. Create an Amazon Bedrock knowledge 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 2

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 Bedrock
- 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 Bedrock
- I. Monitor InvocationLatency, InvocationClientErrors, and InvocationServerErrors metric
- J. Distribute traffic across multiple versions of the same model.

Answer: C

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

Answer: B

NEW QUESTION 4

A financial services company is deploying a generative AI (GenAI) application that uses Amazon Bedrock to assist customer service representatives to provide personalized investment advice to customers. The company must implement a comprehensive governance solution that follows responsible AI practices and meets regulatory requirements.

The solution must detect and prevent hallucinations in recommendations. The solution must have safety controls for customer interactions. The solution must also monitor model behavior drift in real time and maintain audit trails of all prompt-response pairs for regulatory review. The company must deploy the solution within 60 days. The solution must integrate with the company's existing compliance dashboard and respond to customers within 200 ms.

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

- A. Configure Amazon Bedrock guardrails to apply custom content filters and toxicity detection
- B. Use Amazon Bedrock Model Evaluation to detect hallucination
- C. Store prompt-response pairs in Amazon DynamoDB to capture audit trails and set a TTL
- D. Integrate Amazon CloudWatch custom metrics with the existing compliance dashboard.
- E. Deploy Amazon Bedrock and use AWS PrivateLink to access the application securely

- F. Use AWS Lambda functions to implement custom prompt validation
- G. Store prompt-response pairs in an Amazon S3 bucket and configure S3 Lifecycle policies
- H. Create custom Amazon CloudWatch dashboards to monitor model performance metrics.
- I. Use Amazon Bedrock Agents and Amazon Bedrock Knowledge Bases to ground response
- J. Use Amazon Bedrock Guardrails to enforce content safety
- K. Use Amazon OpenSearch Service to store and index prompt-response pairs
- L. Integrate OpenSearch Service with Amazon QuickSight to create compliance reports and to detect model behavior drift.
- M. Use Amazon SageMaker Model Monitor to detect model behavior drift
- N. Use AWS WAF to filter content
- O. Store customer interactions in an encrypted Amazon RDS database
- P. Use Amazon API Gateway to create custom HTTP APIs to integrate with the compliance dashboard.

Answer: A

NEW QUESTION 5

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 configuration
- 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 search
- 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 documents
- 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 6

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

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

Answer: D

NEW QUESTION 7

A company needs a system to automatically generate study materials from multiple content sources. The content sources include document files (PDF files, PowerPoint presentations, and Word documents) and multimedia files (recorded videos). The system must process more than 10,000 content sources daily with peak loads of 500 concurrent uploads. The system must also extract key concepts from document files and multimedia files and create contextually accurate summaries. The generated study materials must support real-time collaboration with version control. Which solution will meet these requirements?

- A. Use Amazon Bedrock Data Automation (BDA) with AWS Lambda functions to orchestrate document file processing
- B. Use Amazon Bedrock Knowledge Bases to process all multimedia
- C. Store the content in Amazon DocumentDB with replication
- D. Collaborate by using Amazon SNS topic subscription
- E. Track changes by using Amazon Bedrock Agents.
- F. Use Amazon Bedrock Data Automation (BDA) with foundation models (FMs) to process document files
- G. Integrate BDA with Amazon Textract for PDF extraction and with Amazon Transcribe for multimedia files
- H. Store the processed content in Amazon S3 with versioning enabled
- I. Store the metadata in Amazon DynamoDB
- J. Collaborate in real time by using AWS AppSync GraphQL subscriptions and DynamoDB.
- K. Use Amazon Bedrock Data Automation (BDA) with Amazon SageMaker AI endpoints to host content extraction and summarization models
- L. Use Amazon Bedrock Guardrails to extract content from all file types
- M. Store document files in Amazon Neptune for time series analysis
- N. Collaborate by using Amazon Bedrock Chat for real-time messaging.
- O. Use Amazon Bedrock Data Automation (BDA) with AWS Lambda functions to process batches of content files
- P. Fine-tune foundation models (FMs) in Amazon Bedrock to classify documents across all content types
- Q. Store the processed data in Amazon ElastiCache (Redis OSS) by using Cluster Mode with sharding
- R. Use Prompt management in Amazon Bedrock for version control.

Answer: B

NEW QUESTION 8

A healthcare company is developing a document management system that stores medical research papers in an Amazon S3 bucket. The company needs a comprehensive metadata framework to improve search precision for a GenAI application. The metadata must include document timestamps, author information, and research domain classifications.

The solution must maintain a consistent metadata structure across all uploaded documents and allow foundation models (FMs) to understand document context without accessing full content.

Which solution will meet these requirements?

- A. Store document timestamps in Amazon S3 system metadata
- B. Use S3 object tags for domain classification
- C. Implement custom user-defined metadata to store author information.
- D. Set up S3 Object Lock with legal holds to track document timestamp
- E. Use S3 object tags for author information
- F. Implement S3 access points for domain classification.
- G. Use S3 Inventory reports to track timestamp
- H. Create S3 access points for domain classification
- I. Store author information in S3 Storage Lens dashboards.
- J. Use custom user-defined metadata to store author information
- K. Use S3 Object Lock retention periods for timestamp
- L. Use S3 Event Notifications for domain classification.

Answer: A

NEW QUESTION 9

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

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

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

- A. Enable model preload upon container start
- B. Implement dynamic batching to process multiple user requests together in a single inference pass.
- C. Select a larger GPU instance type for the SageMaker AI endpoint
- D. Set the minimum number of instances to 0. Continue to perform per-request processing
- E. Lazily load model weights on the first request.
- F. Switch to a multi-model endpoint
- G. Use lazy loading without request batching.
- H. Set the minimum number of instances to greater than 0. Enable response streaming.
- I. Switch to Amazon SageMaker Asynchronous Inference for all requests
- J. Store requests in an Amazon S3 bucket
- K. Set the minimum number of instances to 0.

Answer: AD

NEW QUESTION 10

A 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 10

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

Which solution will meet these requirements?

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

Answer: A

NEW QUESTION 15

A company is building a generative AI (GenAI) application that produces content based on a variety of internal and external data sources. The company wants to ensure that the generated output is fully traceable. The application must support data source registration and enable metadata tagging to attribute content to its original source. The application must also maintain audit logs of data access and usage throughout the pipeline. Which solution will meet these requirements?

- A. Use AWS Lake Formation to catalog data sources and control access
- B. Apply metadata tags directly in Amazon S3. Use AWS CloudTrail to monitor API activity.
- C. Use AWS Glue Data Catalog to register and tag data source
- D. Use Amazon CloudWatch Logs to monitor access patterns and application behavior.
- E. Store data in Amazon S3 and use object tagging for attribution
- F. Use AWS Glue Data Catalog to manage schema information
- G. Use AWS CloudTrail to log access to S3 buckets.
- H. Use AWS Glue Data Catalog to register all data source
- I. Apply metadata tags to attribute data source
- J. Use AWS CloudTrail to log access and activity across services.

Answer: D

NEW QUESTION 19

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 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 base
- 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 functions
- 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 23

A book publishing company wants to build a book recommendation system that uses an AI assistant. The AI assistant will use ML to generate a list of recommended books from the company's book catalog. The system must suggest books based on conversations with customers. The company stores the text of the books, customers' and editors' reviews of the books, and extracted book metadata in Amazon S3. The system must support low-latency responses and scale efficiently to handle more than 10,000 concurrent users. Which solution will meet these requirements?

- A. Use Amazon Bedrock Knowledge Bases to generate embedding
- B. Store the embeddings as a vector store in Amazon OpenSearch Service
- C. Create an AWS Lambda function that queries the knowledge base
- D. Configure Amazon API Gateway to invoke the Lambda function when handling user requests.
- E. Use Amazon Bedrock Knowledge Bases to generate embedding
- F. Store the embeddings as a vector store in Amazon DynamoDB
- G. Create an AWS Lambda function that queries the knowledge base

- H. Configure Amazon API Gateway to invoke the Lambda function when handling user requests.
- I. Use Amazon SageMaker AI to deploy a pre-trained model to build a personalized recommendation engine for book
- J. Deploy the model as a SageMaker AI endpoint
- K. Invoke the model endpoint by using Amazon API Gateway.
- L. Create an Amazon Kendra GenAI Enterprise Edition index that uses the S3 connector to index the book catalog data stored in Amazon S3. Configure built-in FAQ in the Kendra index
- M. Develop an AWS Lambda function that queries the Kendra index based on user conversation
- N. Deploy Amazon API Gateway to expose this functionality and invoke the Lambda function.

Answer: A

NEW QUESTION 26

A company has a recommendation system. The system's applications run on Amazon EC2 instances. The applications make API calls to Amazon Bedrock foundation models (FMs) to analyze customer behavior and generate personalized product recommendations. The system is experiencing intermittent issues. Some recommendations do not match customer preferences. The company needs an observability solution to monitor operational metrics and detect patterns of operational performance degradation compared to established baselines. The solution must also generate alerts with correlation data within 10 minutes when FM behavior deviates from expected patterns. Which solution will meet these requirements?

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

Answer: C

NEW QUESTION 31

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 34

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

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

L. Integrate the process into the CI/CD pipeline to block releases that fail.

Answer: D

NEW QUESTION 38

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 41

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 43

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 OpenSearch to federate data across Amazon S3, Amazon Kinesis, and Amazon SageMaker Feature Store
- O. Use EventBridge for ingestion orchestration
- P. Use custom AWS Lambda functions to filter LLM outputs for ethical compliance.

Answer: B

NEW QUESTION 45

A company is developing a generative AI (GenAI)-powered customer support application that uses Amazon Bedrock foundation models (FMs). The application must maintain conversational context across multiple interactions with the same user. The application must run clarification workflows to handle ambiguous user queries. The company must store encrypted records of each user conversation to use for personalization. The application must be able to handle thousands of concurrent users while responding to each user quickly.

Which solution will meet these requirements?

- A. Use an AWS Step Functions Express workflow to orchestrate conversation flow
- B. Invoke AWS Lambda functions to run clarification logic
- C. Store conversation history in Amazon RDS and use session IDs as the primary key.
- D. Use an AWS Step Functions Standard workflow to orchestrate clarification workflow
- E. Include Wait for a Callback patterns to manage the workflow
- F. Store conversation history in Amazon DynamoDB
- G. Purchase on-demand capacity and configure server-side encryption.
- H. Deploy the application by using an Amazon API Gateway REST API to route user requests to an AWS Lambda function to update and retrieve conversation context
- I. Store conversation history in Amazon S3 and configure server-side encryption
- J. Save each interaction as a separate JSON file.
- K. Use AWS Lambda functions to call Amazon Bedrock inference API
- L. Use Amazon SQS queues to orchestrate clarification step
- M. Store conversation history in an Amazon ElastiCache (Redis OSS) cluster
- N. Configure encryption at rest.

Answer: B

NEW QUESTION 48

A company uses Amazon Bedrock to implement a Retrieval Augmented Generation (RAG)-based system to serve medical information to users. The company needs to compare multiple chunking strategies, evaluate the generation quality of two foundation models (FMs), and enforce quality thresholds for deployment. Which Amazon Bedrock evaluation configuration will meet these requirements?

- A. Create a retrieve-only evaluation job that uses a supported version of Anthropic Claude Sonnet as the evaluator model
- B. Configure metrics for context relevance and context coverage
- C. Define deployment thresholds in a separate CI/CD pipeline.
- D. Create a retrieve-and-generate evaluation job that uses custom precision-at-k metrics and an LLM-as-a-judge metric with a scale of 1–5. Include each chunking strategy in the evaluation dataset
- E. Use a supported version of Anthropic Claude Sonnet to evaluate responses from both FMs.
- F. Create a separate evaluation job for each chunking strategy and FM combination
- G. Use Amazon Bedrock built-in metrics for correctness and completeness
- H. Manually review scores before deployment approval.
- I. Set up a pipeline that uses multiple retrieve-only evaluation jobs to assess retrieval quality
- J. Create separate evaluation jobs for both FMs that use Amazon Nova Pro as the LLM-as-a-judge model
- K. Evaluate based on faithfulness and citation precision metrics.

Answer: B

NEW QUESTION 49

A specialty coffee company has a mobile app that generates personalized coffee roast profiles by using Amazon Bedrock with a three-stage prompt chain. The prompt chain converts user inputs into structured metadata, retrieves relevant logs for coffee roasts, and generates a personalized roast recommendation for each customer.

Users in multiple AWS Regions report inconsistent roast recommendations for identical inputs, slow inference during the retrieval step, and unsafe recommendations such as brewing at excessively high temperatures. The company must improve the stability of outputs for repeated inputs. The company must also improve app performance and the safety of the app's outputs. The updated solution must ensure 99.5% output consistency for identical inputs and achieve inference latency of less than 1 second. The solution must also block unsafe or hallucinated recommendations by using validated safety controls.

Which solution will meet these requirements?

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

Answer: A

NEW QUESTION 54

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 59

A company is planning to deploy multiple generative AI (GenAI) applications to five independent business units that operate in multiple countries in Europe and the Americas.

Each application uses Amazon Bedrock Retrieval Augmented Generation (RAG) patterns with business unit-specific knowledge bases that store terabytes of unstructured data.

The company must establish well-architected, standardized components for security controls, observability practices, and deployment patterns across all the GenAI applications. The components must be reusable, versioned, and governed consistently.

Which solution will meet these requirements?

- A. Configure Amazon API Gateway REST API endpoints for the GenAI application
- B. Deploy common security, observability, and RAG patterns based on the AWS Well-Architected Generative AI Lens in standardized AWS CloudFormation template
- C. Use CloudFormation Guard after deployment to validate policy compliance in each business unit.
- D. Create standardized AWS CloudFormation templates to implement security, observability, and RAG patterns based on the AWS Well-Architected Generative AI Lens
- E. Establish a centralized repository for version control
- F. Integrate a CI/CD pipeline with CloudFormation Guard to enforce consistent and repeatable deployments across business units.
- G. Use AWS Service Catalog to define standardized portfolios and versioned products for each business unit
- H. Use the portfolios to enforce security, observability, and RAG patterns based on the AWS Well-Architected Generative AI Lens
- I. Require business units to use the Service Catalog console to deploy resources.
- J. Document security controls, observability requirements, and RAG patterns based on the AWS Well-Architected Generative AI Lens in a shared design document
- K. Use Amazon Macie to enforce deployment
- L. Delegate implementation responsibility to each business unit.

Answer: B

NEW QUESTION 62

A finance company is developing an AI assistant to help clients plan investments and manage their portfolios. The company identifies several high-risk conversation patterns such as requests for specific stock recommendations or guaranteed returns. High-risk conversation patterns could lead to regulatory violations if the company cannot implement appropriate controls.

The company must ensure that the AI assistant does not provide inappropriate financial advice, generate content about competitors, or make claims that are not factually grounded in the company's approved financial guidance. The company wants to use Amazon Bedrock Guardrails to implement a solution.

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

- A. Add the high-risk conversation patterns to a denied topics guardrail.
- B. Configure a content filter guardrail to filter prompts that contain the high-risk conversation patterns.
- C. Configure a content filter guardrail to filter prompts that contain competitor names.
- D. Add the names of competitors as custom word filter
- E. Set the input and output actions to block.
- F. Set a low grounding score threshold.
- G. Set a high grounding score threshold.

Answer: ADF

NEW QUESTION 67

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

Answer: B

NEW QUESTION 72

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

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

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

Answer: B

NEW QUESTION 75

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 77

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:

```
python
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 81

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 86

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 90

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

Which solution will meet these requirements?

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

Answer: C

NEW QUESTION 93

A hotel company wants to enhance a legacy Java-based property management system (PMS) by adding AI capabilities. The company wants to use Amazon Bedrock Knowledge Bases to provide staff with room availability information and hotel-specific details. The solution must maintain separate access controls for each hotel that the company manages. The solution must provide room availability information in near real time and must maintain consistent performance during peak usage periods.

Which solution will meet these requirements?

- A. Deploy a single Amazon Bedrock knowledge base that contains combined data for all hotels
- B. Configure AWS Lambda functions to synchronize data from each hotel's PMS database through direct API connection
- C. Implement AWS CloudTrail logging with hotel-specific filters to audit access logs for each hotel's data.
- D. Create an Amazon EventBridge rule for each hotel that is invoked by changes to the PMS database
- E. Configure the rule to send updates to a centralized Amazon Bedrock knowledge base in a management AWS account
- F. Configure resource-based policies to enforce hotel-specific access controls.
- G. Implement one Amazon Bedrock knowledge base for each hotel in a multi-account structure
- H. Use direct data ingestion to provide near real-time room availability information
- I. Schedule regular synchronization for less critical information.
- J. Build a centralized Amazon Bedrock Agents solution that uses multiple knowledge bases
- K. Implement AWS IAM Identity Center with hotel-specific permission sets to control staff access.

Answer: C

NEW QUESTION 94

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 96

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 environment
- B. Store conversation states in Amazon DynamoDB
- 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 roles
- E. Store conversation states in Amazon ElastiCache
- 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 API
- H. Store interaction history in application process memory in each Lambda function or EKS container
- I. Use environment variables to configure model parameters.
- J. Use the Amazon Bedrock Converse API and IAM roles for authentication
- K. Pass previous messages in the request messages array to maintain conversational context
- L. Use programming language-specific SDKs to establish consistent API interfaces.

Answer: D

NEW QUESTION 101

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 106

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

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

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

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

Answer: ABD

NEW QUESTION 111

An enterprise application uses an Amazon Bedrock foundation model (FM) to process and analyze 50 to 200 pages of technical documents. Users are experiencing inconsistent responses and receiving truncated outputs when processing documents that exceed the FM's context window limits.

Which solution will resolve this problem?

- A. Configure fixed-size chunking at 4,000 tokens for each chunk with 20% overlap
- B. Use application-level logic to link multiple chunks sequentially until the FM's maximum context window of 200,000 tokens is reached before making inference calls.
- C. Use hierarchical chunking with parent chunks of 8,000 tokens and child chunks of 2,000 tokens
- D. Use Amazon Bedrock Knowledge Bases built-in retrieval to automatically select relevant parent chunks based on query context
- E. Configure overlap tokens to maintain semantic continuity.
- F. Use semantic chunking with a breakpoint percentile threshold of 95% and a buffer size of 3 sentences
- G. Use the RetrieveAndGenerate API to dynamically select the most relevant chunks based on embedding similarity scores.

- H. Create a pre-processing AWS Lambda function that analyzes document token count by using the FM's tokenize
- I. Configure the Lambda function to split documents into equal segments that fit within 80% of the context window
- J. Configure the Lambda function to process each segment independently before aggregating the results.

Answer: C

NEW QUESTION 115

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

Which solution will meet these requirements?

- A. Use AWS Glue to run an ETL job on each uploaded file
- B. Configure the ETL job to use the AWS SDK to invoke the SageMaker AI model endpoint
- C. Use real-time inference with the endpoint to re-deploy the model after it is re-trained on the updated customer dataset.
- D. Create an AWS Lambda function and webhook handlers to generate an event when an employee uploads a new file
- E. Configure SageMaker Pipelines to re-deploy the model after it is re-trained on the updated customer dataset
- F. Use Amazon EventBridge to create an event bus
- G. Set the Lambda function event as the source and SageMaker Pipelines as the target.
- H. Create an AWS Step Functions Express workflow with AWS SDK integrations to retrieve the customer data from the S3 bucket when an employee uploads a new file to the S3 bucket
- I. Use a SageMaker Data Wrangler flow to export the data from the S3 bucket to SageMaker Autopilot
- J. Use the SageMaker Autopilot to re-deploy the model after it has been re-trained on the updated customer dataset.
- K. Create an AWS Step Functions Standard workflow
- L. Configure the first state to call an AWS Lambda function to respond when an employee uploads a new file to the S3 bucket
- M. Use a pipeline in SageMaker Pipelines to re-deploy the model after it has been re-trained on the updated customer dataset
- N. Use the next state in the workflow to run the pipeline when the first state receives a response.

Answer: D

NEW QUESTION 119

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 122

A medical company uses Amazon Bedrock to power a clinical documentation summarization system. The system produces inconsistent summaries when handling complex clinical documents. The system performed well on simple clinical documents.

The company needs a solution that diagnoses inconsistencies, compares prompt performance against established metrics, and maintains historical records of prompt versions.

Which solution will meet these requirements?

- A. Create multiple prompt variants by using Prompt management in Amazon Bedrock
- B. Manually test the prompts with simple clinical documents
- C. Deploy the highest performing version by using the Amazon Bedrock console.
- D. Implement version control for prompts in a code repository with a test suite that contains complex clinical documents and quantifiable evaluation metrics
- E. Use an automated testing framework to compare prompt versions and document performance patterns.
- F. Deploy each new prompt version to separate Amazon Bedrock API endpoints
- G. Split production traffic between the endpoints
- H. Configure Amazon CloudWatch to capture response metrics and user feedback for automatic version selection.
- I. Create a custom prompt evaluation flow in Amazon Bedrock Flows that applies the same clinical document inputs to different prompt variants
- J. Use Amazon Comprehend Medical to analyze and score the factual accuracy of each version.

Answer: B

NEW QUESTION 125

A company is using Amazon Bedrock and Anthropic Claude 3 Haiku to develop an AI assistant. The AI assistant normally processes 10,000 requests each hour but experiences surges of up to 30,000 requests each hour during peak usage periods. The AI assistant must respond within 2 seconds while operating across multiple AWS Regions.

The company observes that during peak usage periods, the AI assistant experiences throughput bottlenecks that cause increased latency and occasional request timeouts. The company must resolve the performance issues.

Which solution will meet this requirement?

- A. Purchase provisioned throughput and sufficient model units (MUs) in a single Region
- B. Configure the application to retry failed requests with exponential backoff.
- C. Implement token batching to reduce API overhead
- D. Use cross-Region inference profiles to automatically distribute traffic across available Regions.
- E. Set up auto scaling AWS Lambda functions in each Region
- F. Implement client-side round-robin request distribution
- G. Purchase one model unit (MU) of provisioned throughput as a backup.
- H. Implement batch inference for all requests by using Amazon S3 buckets across multiple Regions
- I. Use Amazon SQS to set up an asynchronous retrieval process.

Answer: B

NEW QUESTION 130

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

- A. Use Amazon Bedrock Knowledge Bases with source attribution enabled
- B. Use the Anthropic Claude Messages API with RAG to set high-relevance thresholds for sourced documents
- C. Store reasoning and citations in Amazon S3 for auditing purposes.
- D. Use Amazon Bedrock with Anthropic Claude models and extended thinking
- E. Configure a 4,000-token thinking budget
- F. Store reasoning traces and citations in Amazon DynamoDB for auditing purposes.
- G. Configure Amazon SageMaker AI with a custom Anthropic Claude model
- H. Use the model's reasoning parameter and AWS Lambda to process responses
- I. Add source citations from a separate Amazon RDS database.
- J. Use Amazon Bedrock with Anthropic Claude models and chain-of-thought reasoning
- K. Configure custom retrieval tracking with the Amazon Bedrock Knowledge Bases API
- L. Use Amazon CloudWatch to monitor response latency metrics.

Answer: A

NEW QUESTION 134

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 queues
- K. Use Amazon ElastiCache to cache all human technician review decisions.

Answer: B

NEW QUESTION 137

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

Answer: D

NEW QUESTION 138

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 integration
- 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 integration
- 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 integration
- 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 140

A media company must use Amazon Bedrock to implement a robust governance process for AI-generated content. The company needs to manage hundreds of prompt templates. Multiple teams use the templates across multiple AWS Regions to generate content. The solution must provide version control with approval workflows that include notifications for pending reviews. The solution must also provide detailed audit trails that document prompt activities and consistent prompt parameterization to enforce quality standards. Which solution will meet these requirements?

- A. Configure Amazon Bedrock Studio prompt template
- B. Use Amazon CloudWatch dashboards to display prompt usage metrics
- C. Store approval status in Amazon DynamoDB
- D. Use AWS Lambda functions to enforce approvals.
- E. Use Amazon Bedrock Prompt Management to implement version control
- F. Configure AWS CloudTrail for audit logging
- G. Use AWS Identity and Access Management policies to control approval permissions
- H. Create parameterized prompt templates by specifying variables.
- I. Use AWS Step Functions to create an approval workflow
- J. Store prompts in Amazon S3. Use tags to implement version control
- K. Use Amazon EventBridge to send notifications.
- L. Deploy Amazon SageMaker Canvas with prompt templates stored in Amazon S3. Use AWS CloudFormation for version control
- M. Use AWS Config to enforce approval policies.

Answer: B

NEW QUESTION 141

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 policies
- B. Use Amazon Bedrock cross-Region inference to distribute the workload
- C. Use Amazon CloudWatch to log AI decision-making processes
- D. Manually track compliance certifications across Regions.
- E. Use SCPs with AWS Organizations to manage location-specific permissions
- 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 146

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 calls
- 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 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 151

.....

THANKS FOR TRYING THE DEMO OF OUR PRODUCT

Visit Our Site to Purchase the Full Set of Actual AIP-C01 Exam Questions With Answers.

We Also Provide Practice Exam Software That Simulates Real Exam Environment And Has Many Self-Assessment Features. Order the AIP-C01 Product From:

<https://www.2passeasy.com/dumps/AIP-C01/>

Money Back Guarantee

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