



A generative AI powered self-service assistant for contact centers

Generative AI in the contact center

Many organizations operate contact centers to service requests from their customers

- Staffing for live agents in contact centers can be a significant contributor to an organization's operating expense
- In many cases, live agents spend time responding to customer issues that can easily be solved by referencing available information on the company's website or in knowledge base documents

Generative AI can help reduce operating expenses, while streamlining and improving the customer experience

- LLMs on Amazon Bedrock, including Anthropic's Claude Haiku model, can use information in **Amazon Bedrock Knowledge Bases** to answer questions within a few seconds, making it a viable solution for both text messaging interactions as well as voice calls
- AWS customers using Amazon Connect can **easily integrate with LLMs on Amazon Bedrock via Amazon Lex**

AWS Generative AI Innovation Center solution

- The Contact Center Generative AI Assistant Solution provides an easy to deploy, easy to operate solution that can reduce call volumes requiring a live agent, by **answering customer questions based on content available in a knowledge base**
- The solution includes **automated testing, in-depth conversation analytics, hallucination detection and prevention, and monitoring**

Solution benefits

- Reduction in operating expense, by deflecting inbound calls from human agents to automated agents
 - Costs for Bedrock-powered automated agents range from **\$0.03 and \$0.06 per call**
- Improved and streamlined customer experience
 - Customers can get **questions answered and issues resolved immediately**, without having to wait for a live agent
 - The generative AI model can be instructed to interact with customers in a **consistent, friendly, and informative manner that is always on-brand**, and stays within specified guardrails



Building a Generative AI Contact Center Solution for DoorDash Using Amazon Bedrock

Challenges

DoorDash receives hundreds of thousands of calls to its support center each day and wanted to help Dashers get answers to routine questions efficiently while freeing up live agents to handle more complex issues.

Solutions

Using Anthropic's Claude Haiku on Amazon Bedrock, DoorDash achieved a response latency of 2.5 seconds or less. By handling common inquiries with generative AI, DoorDash has improved self-service workflows and reduced issue resolution speeds.

Results

- 50x increase in testing capacity
- 50% reduction in response latency
- 2.5 seconds or less response latency
- 50% reduction in development time



Using AWS, we've built a solution that gives Dashers reliable access to the information they need, when they need it.

Chaitanya Hari

Contact Center Product Lead, DoorDash



CUSTOMER PROFILE



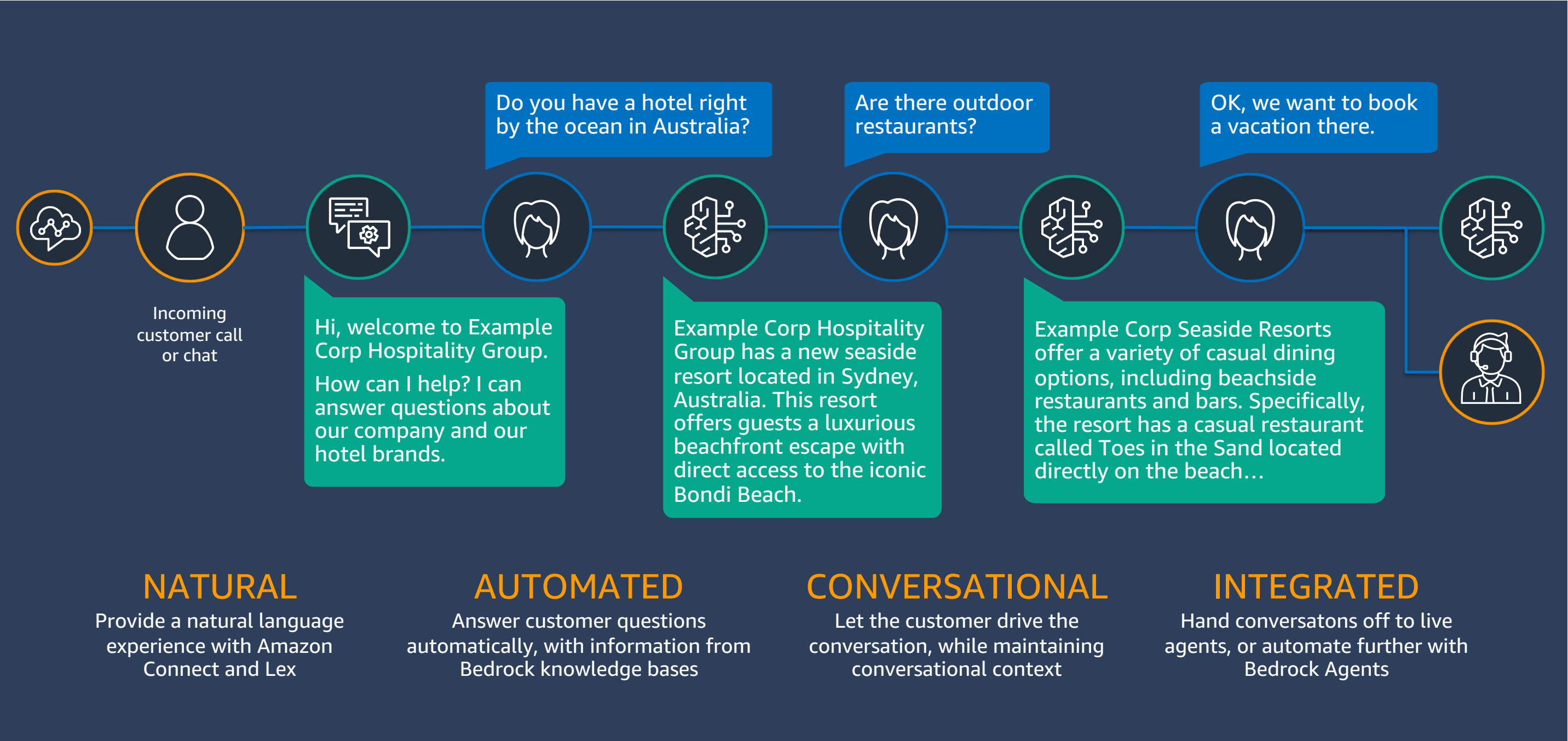
INDUSTRY
Technology

REGION
United States

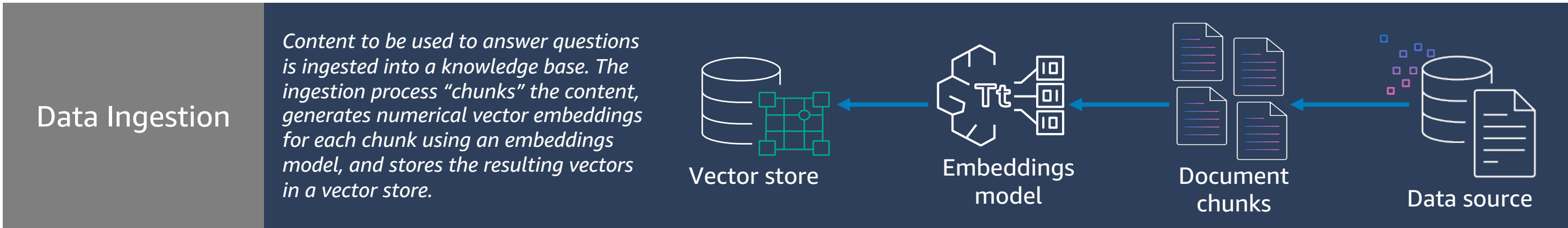
DoorDash is a local commerce platform dedicated to helping Merchants thrive in the convenience economy, giving consumers access to more of their communities, and providing work that empowers.

The AWS case study is at <https://aws.amazon.com/solutions/case-studies/doorDash-bedrock-case-study>

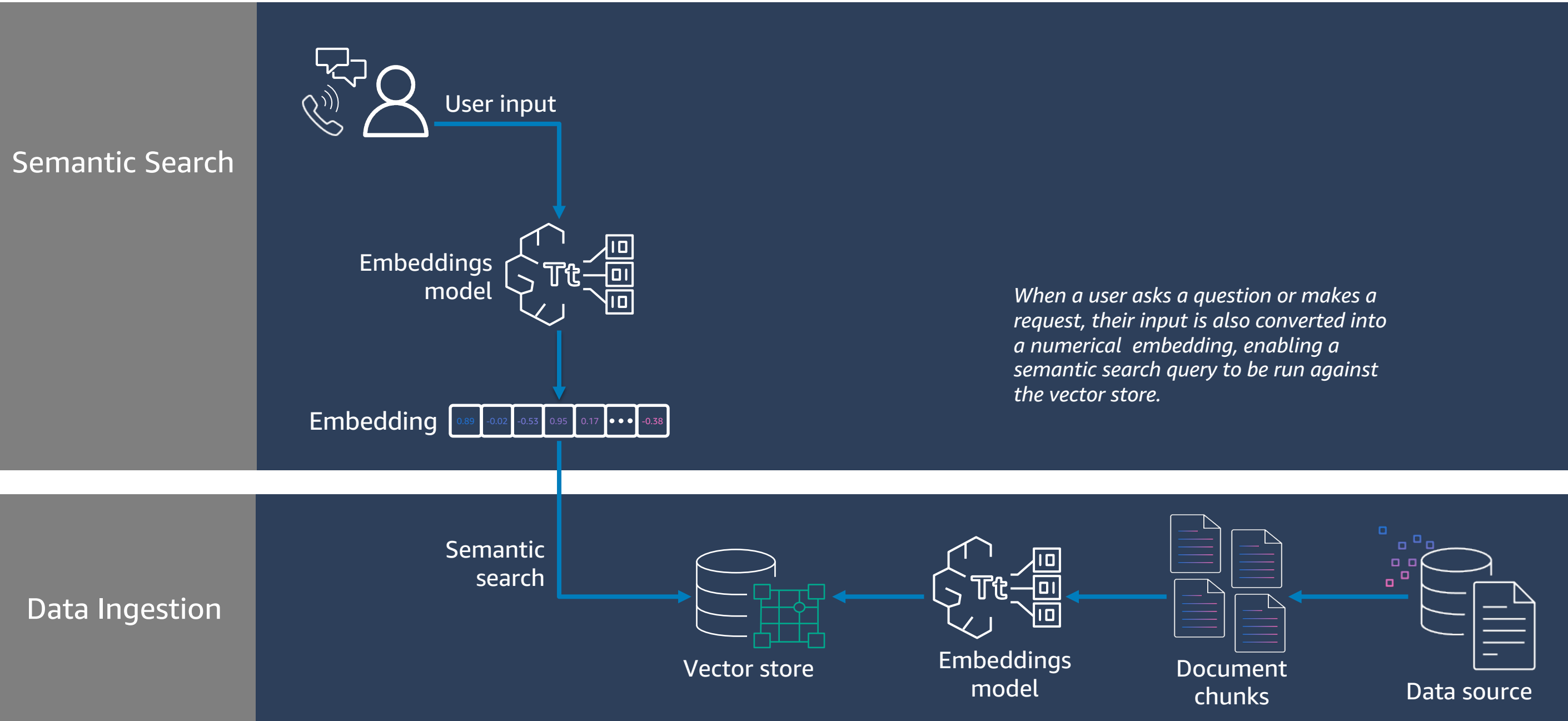
Customer experience (using a fictional hotel chain example)



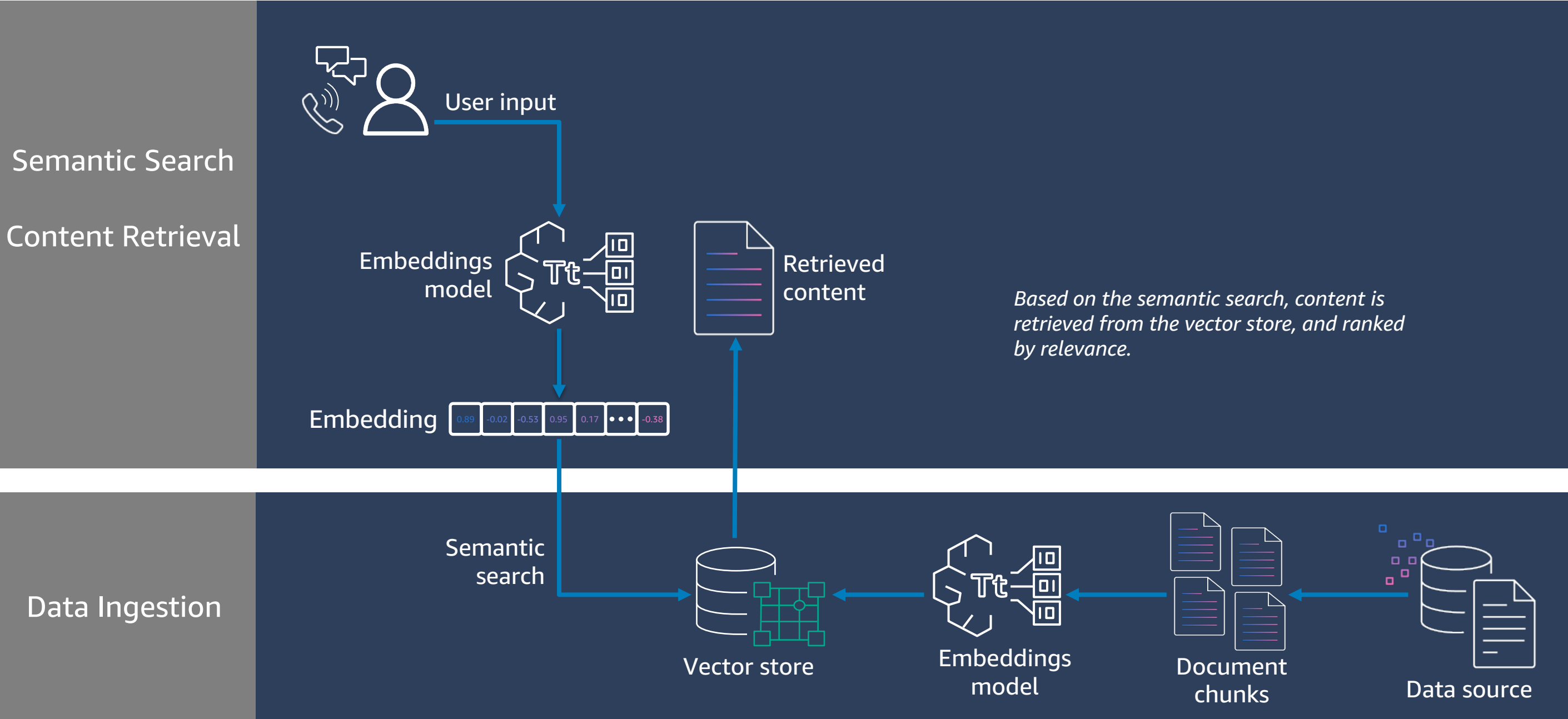
Retrieval Augmented Generation (or, in-context learning)



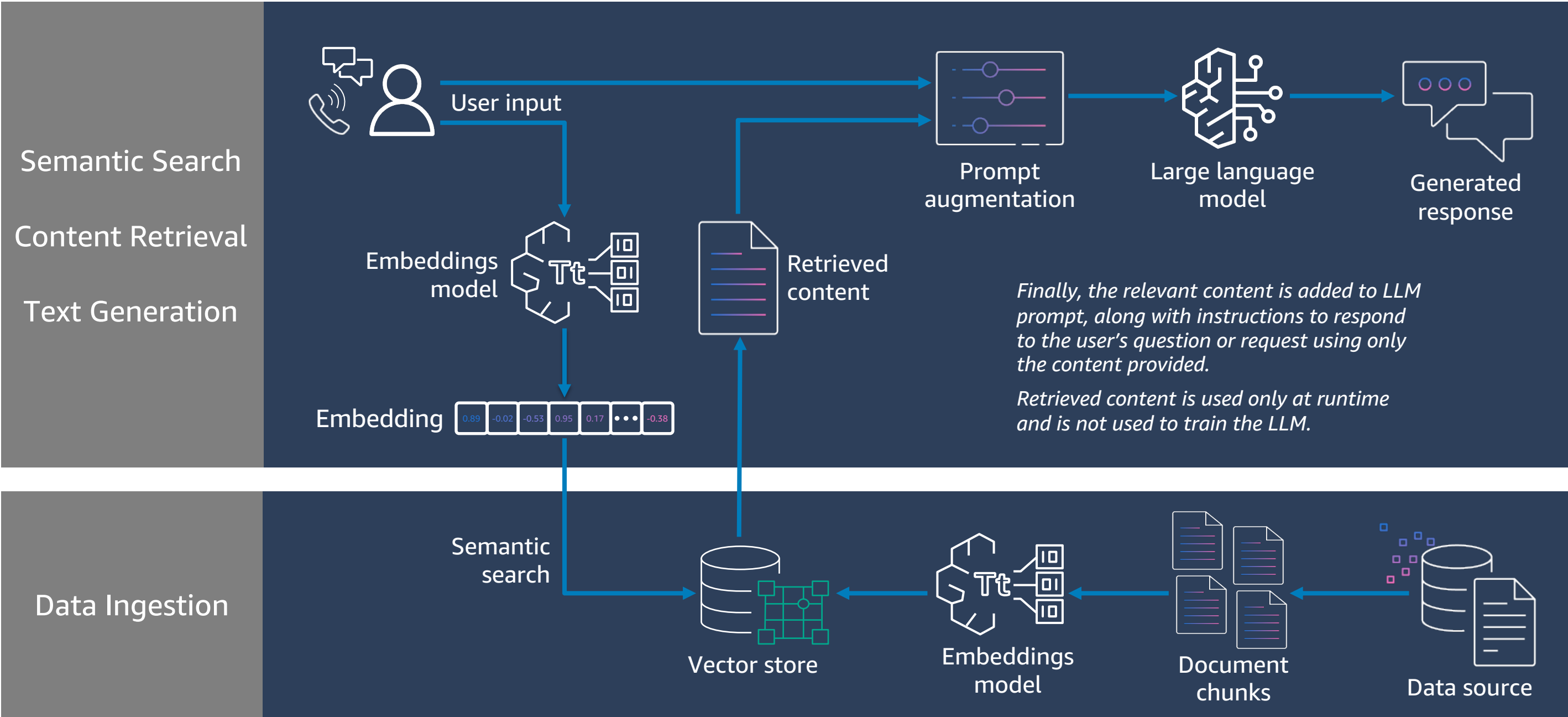
Retrieval Augmented Generation (or, in-context learning)



Retrieval Augmented Generation (or, in-context learning)



Retrieval Augmented Generation (or, in-context learning)



Amazon Bedrock Knowledge Bases

Fully-managed native support for retrieval augmented generation (RAG)



Fully managed support for end-to-end RAG workflow



Securely connect LLMs and agents to data sources



Automatically converts text documents into embedding vectors



Stores embeddings in your vector database



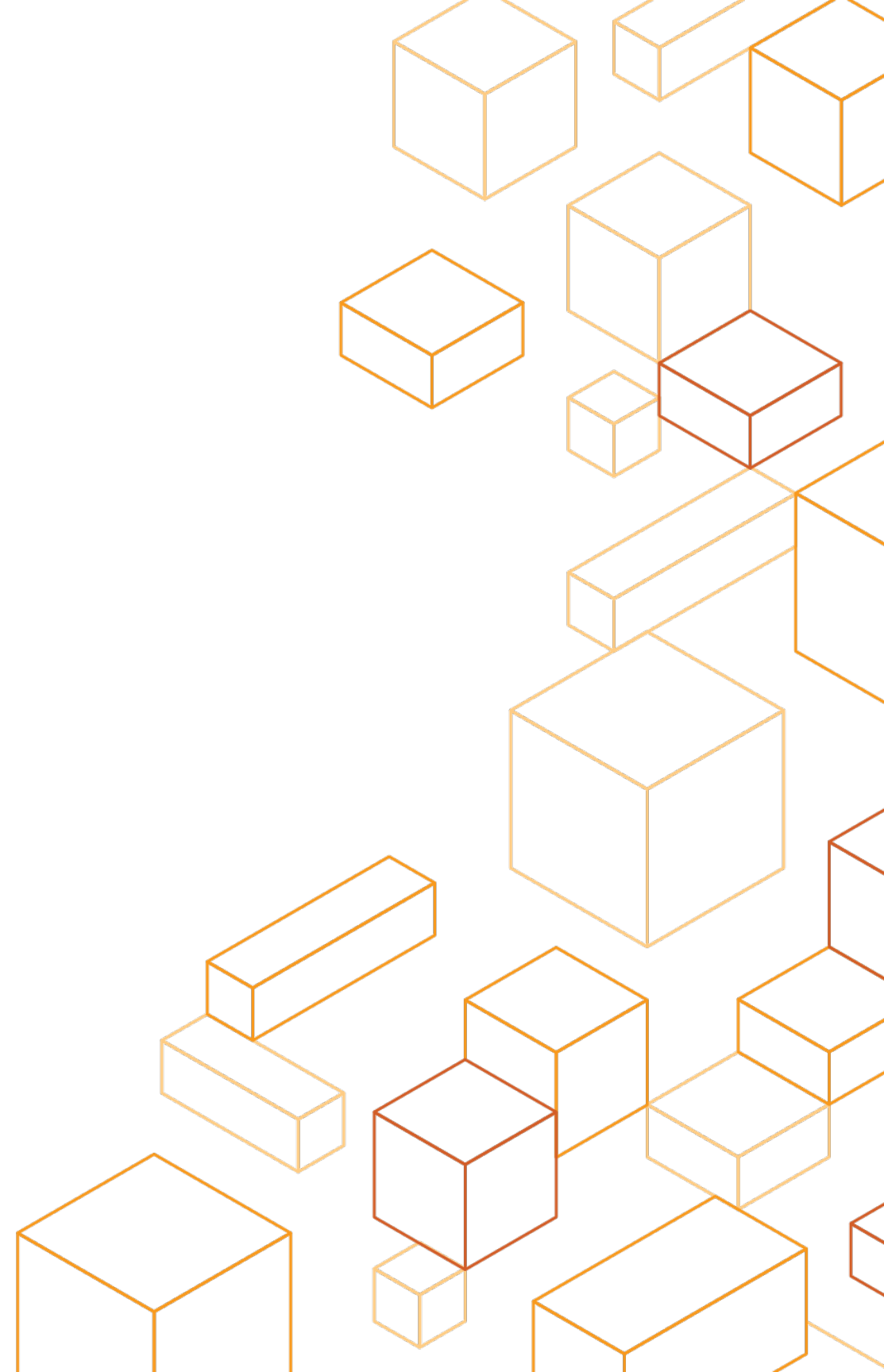
Retrieves content based on semantic search and augments LLM prompts



Provides source attribution



Solution Architecture

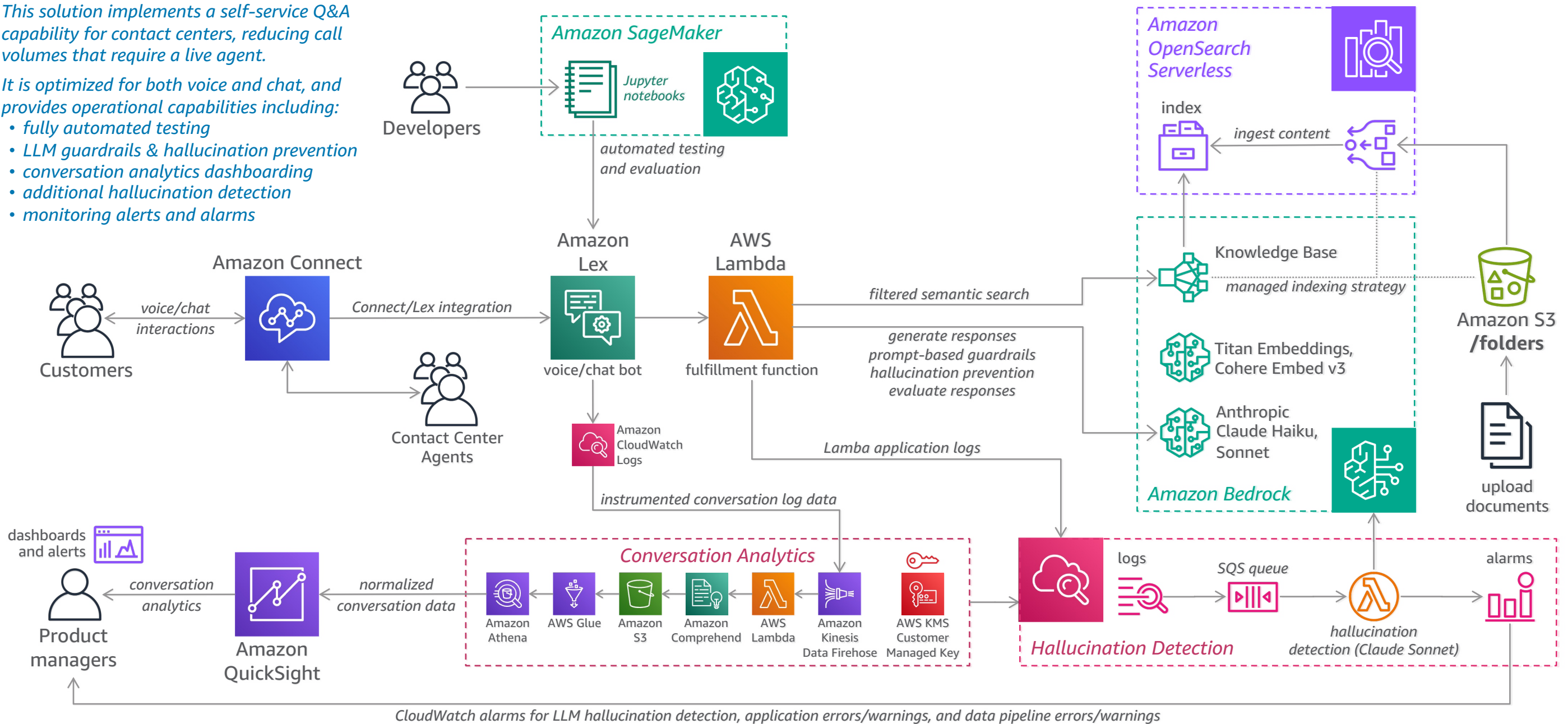


Contact center RAG solution architecture

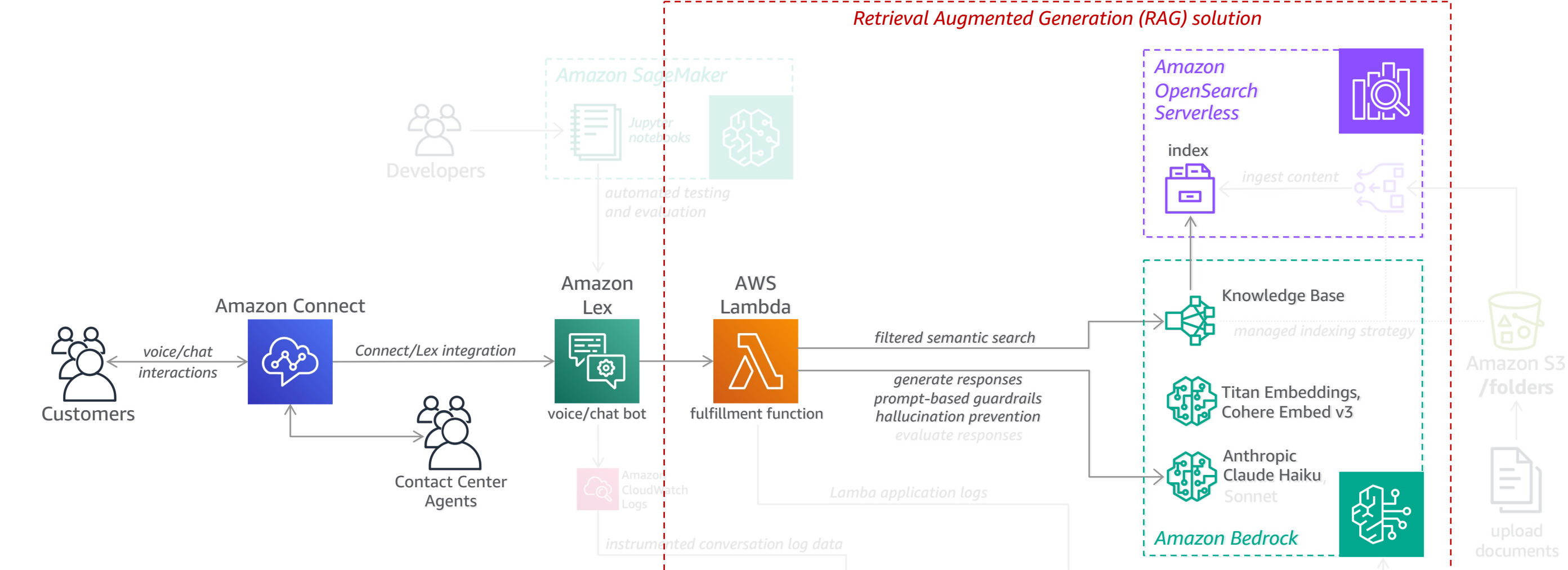
This solution implements a self-service Q&A capability for contact centers, reducing call volumes that require a live agent.

It is optimized for both voice and chat, and provides operational capabilities including:

- fully automated testing
- LLM guardrails & hallucination prevention
- conversation analytics dashboarding
- additional hallucination detection
- monitoring alerts and alarms



Core RAG solution

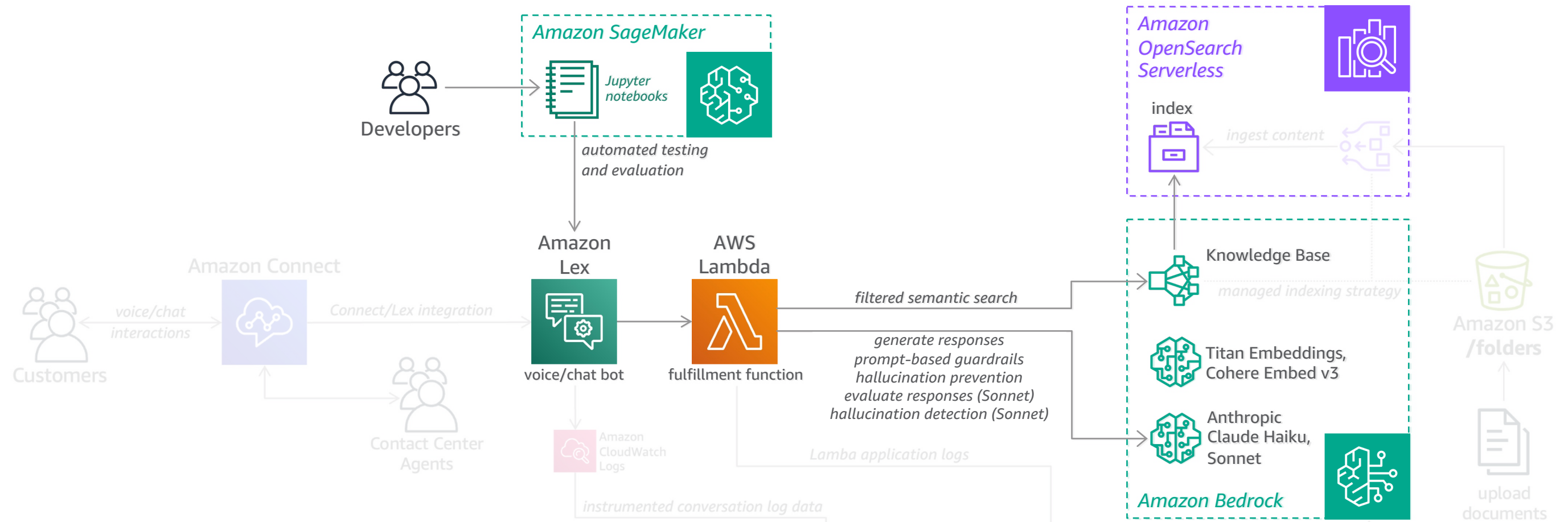


When a customer calls in or initiates a chat session, they're greeted by an Amazon Lex bot that asks them: "how can we help you today?" Based on the customer's response, Lex either routes the customer to the appropriate Amazon Connect live agent queue, or initiates the RAG solution to answer their questions.

The RAG solution uses Anthropic's Claude 3 Haiku LLM, which answers questions based on content available in an Amazon Bedrock Knowledge Base. The content resides in a single collection in Amazon OpenSearch Serverless. To zero in on the most relevant content, Knowledge Base uses a hybrid query that employs a metadata prefilter based on customizable attributes (such as customer type/persona, the specific customer intent/question type, and country), along with a semantic search based on the customer's specific question and conversation history (cosine similarity/k-nearest neighbors).

The solution makes it easy to use any LLM available on Amazon Bedrock, and any type of content repository supported by Bedrock Knowledge Bases. LLM prompts are provided for Claude Haiku and managed in an Amazon DynamoDB table, and can be customized as needed by customer persona, customer intent, country, etc. Guardrails are employed to redirect inappropriate or off-brand topics.

Automated testing



The automated testing component consists of a multi-threaded test script in a Jupyter notebook. The test script can execute 100s of test cases per minute, and can be incorporated into the test stage of a CI/CD deployment pipeline.

Sets of test cases are captured in an Excel workbook or .CSV file, and each test case can contain one or more steps (i.e., for multi-turn conversations, follow-on questions, etc). Each step of a test case includes a user input, a correct "ground truth" answer, and any configuration parameters (session attributes) needed for the test case. The test script sends the test case to the RAG solution, and then uses Anthropic Claude Sonnet to compare the response from the RAG solution to the ground truth answer. If the RAG solution's answer has the same semantic meaning as the ground truth answer, the test case passes. If the RAG solution's answer is semantically different than the ground truth answer, or is incomplete, the test case fails. Either way, a detailed explanation is provided by the LLM as to why the test case passed or failed.

In addition, the RAG solution's answer is evaluated by Claude Sonnet against the RAG content from the knowledge base, to detect any hallucinations (information provided in the answer that is not present in the knowledge base). Test results and hallucination detection results are saved to an output Excel workbook, and also to the conversation analytics dashboard.

Conversation analytics

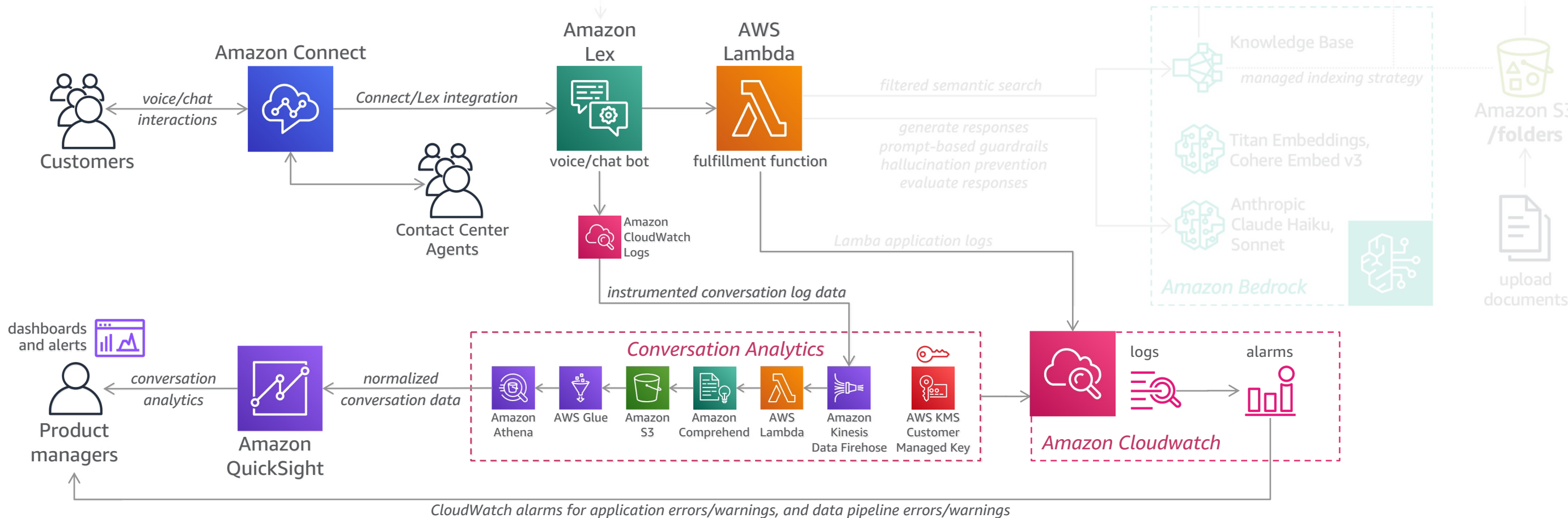
The conversation analytics subsystem consists of a data pipeline stack and an Amazon QuickSight dashboard, together with Amazon CloudWatch alarms.

The data pipeline extracts Lex conversation logs from Amazon CloudWatch Logs, “flattens” them into a tabular, name/value pair structure to simplify queries and analytics, and stores them in an S3 bucket. AWS Glue crawls the S3 bucket on a scheduled basis (defaulting to every 5 minutes) and updates a data catalog schema to make the data accessible via Amazon Athena via SQL queries.

For protecting sensitive data, the data pipeline can optionally redact PII data using Amazon Comprehend, and can apply an AWS Key Management Service customer-managed key (CMK) to limit access to both the CloudWatch Logs log group and the S3 bucket to specifically identified principals.

The resulting dataset includes hundreds of attributes for every conversation turn, including the user question, the RAG solution’s answer, Lex slot values and session attributes, the LLM model version used, the knowledge base identifier, retrieval and RAG latency, input/output token counts, and much more, providing fine-grained observability into the solution.

All testing data, including pass/fail results and explanations, as well as hallucination detection results and explanations, are also available for evaluating test runs.

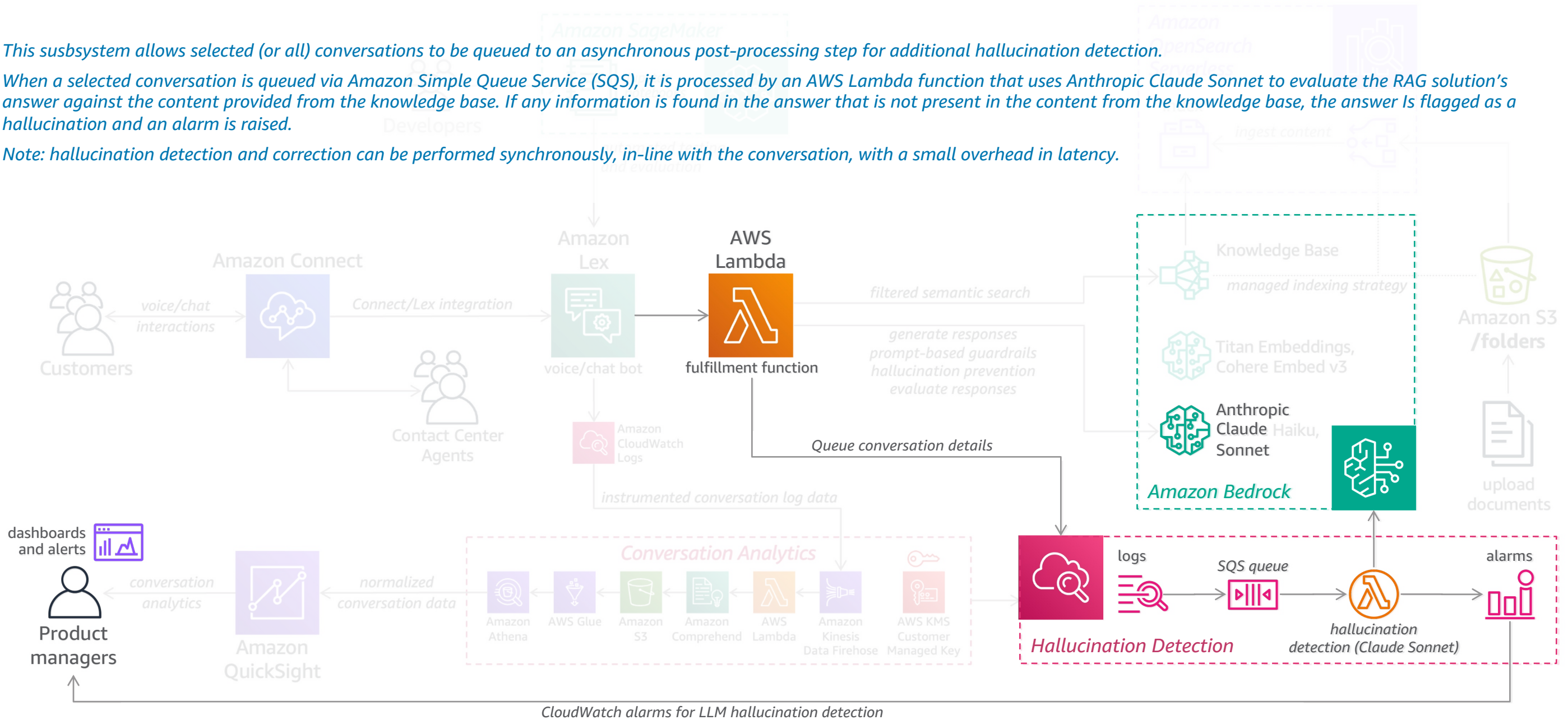


Hallucination detection

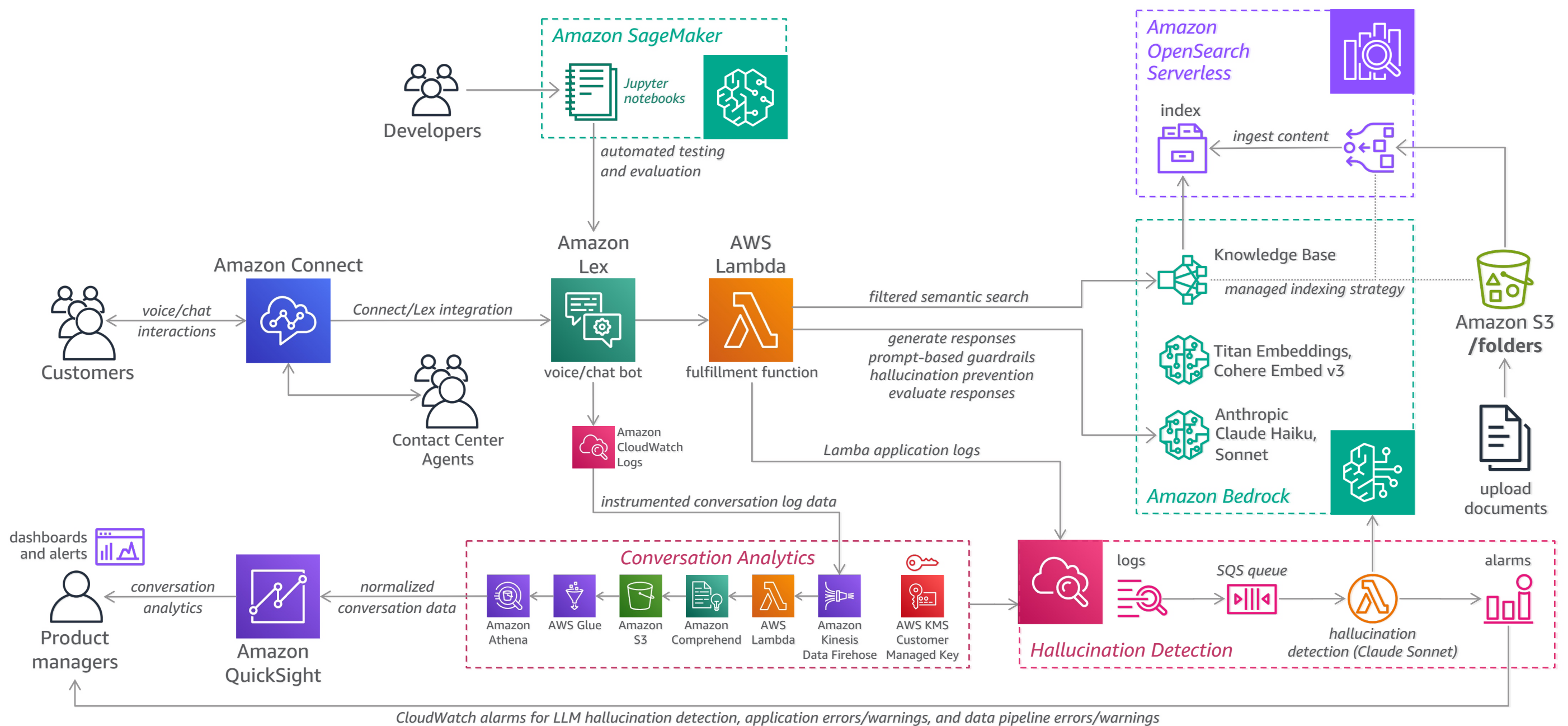
This subsystem allows selected (or all) conversations to be queued to an asynchronous post-processing step for additional hallucination detection.

When a selected conversation is queued via Amazon Simple Queue Service (SQS), it is processed by an AWS Lambda function that uses Anthropic Claude Sonnet to evaluate the RAG solution's answer against the content provided from the knowledge base. If any information is found in the answer that is not present in the content from the knowledge base, the answer is flagged as a hallucination and an alarm is raised.

Note: hallucination detection and correction can be performed synchronously, in-line with the conversation, with a small overhead in latency.

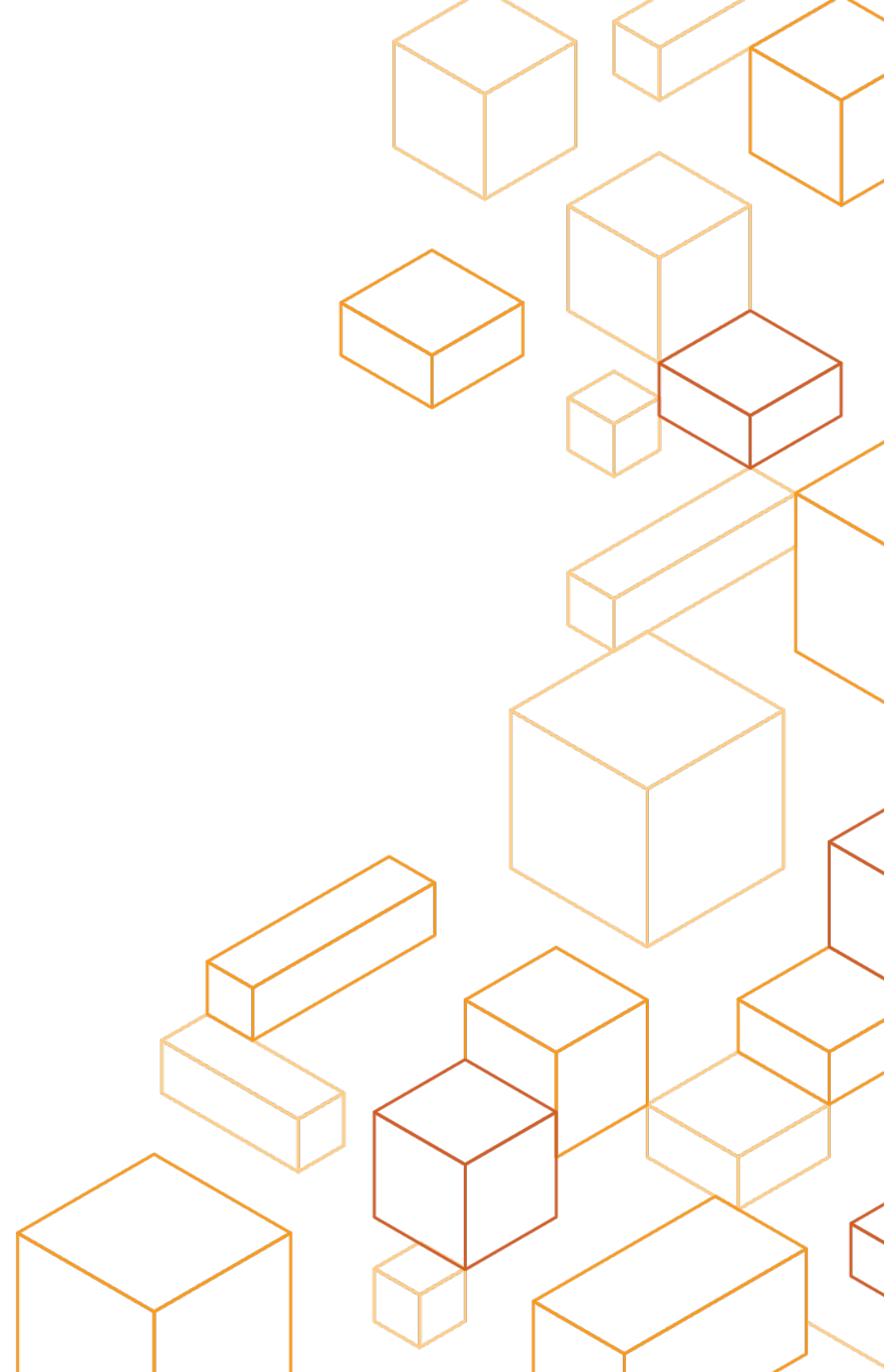


Full end-to-end contact center RAG solution architecture

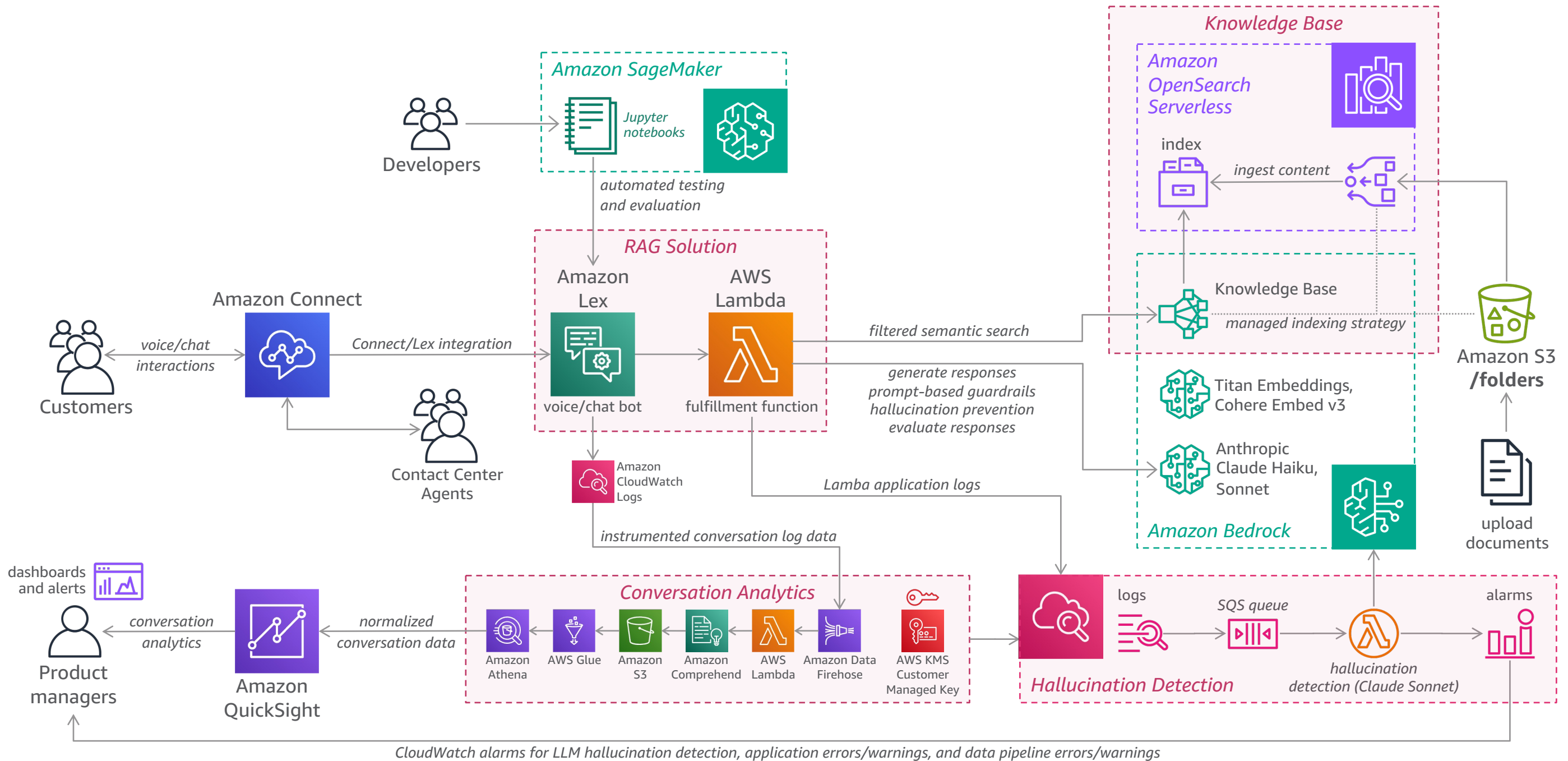




Appendix



CloudFormation stacks



Conversation analytics data pipeline – detail view

