Implementation Guide

Generative AI Application Builder on AWS



Copyright © 2025 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Generative AI Application Builder on AWS: Implementation Guide

Copyright © 2025 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon's trademarks and trade dress may not be used in connection with any product or service that is not Amazon's, in any manner that is likely to cause confusion among customers, or in any manner that disparages or discredits Amazon. All other trademarks not owned by Amazon are the property of their respective owners, who may or may not be affiliated with, connected to, or sponsored by Amazon.

Table of Contents

| Solution overview | |
|---|---|
| Features and benefits | 3 |
| Use cases | 4 |
| Concepts and definitions | 4 |
| Architecture overview | 6 |
| Architecture diagrams | 6 |
| Deployment dashboard | 6 |
| Text use case | |
| Agent use case 1 | 1 |
| AWS Well-Architected design considerations1 | |
| Operational excellence 1 | 4 |
| Security 1 | 4 |
| Reliability 1 | |
| Performance efficiency 1 | 5 |
| Cost optimization 1 | 5 |
| Sustainability | 5 |
| Architecture details 1 | |
| AWS services in this solution | |
| Deployment dashboard 1 | 9 |
| API Gateway custom authorizers 1 | |
| Text use case 2 | |
| Streaming support | 0 |
| How the Generative AI Application Builder on AWS solution works | 1 |
| Plan your deployment 2 | 5 |
| Supported AWS Regions | 5 |
| Cost 2 | 6 |
| Sample costs for running the Deployment dashboard | 8 |
| Sample costs for a text-based proof of concept 2 | 8 |
| Sample costs for a highly scalable generative AI query engine | |
| Costs for adding a knowledge base 3 | |
| Incremental cost of enabling Amazon VPC for a use case | 4 |
| Cost implications when using Provisioned Throughput | 5 |
| Cost for using cross-region inference 3 | |
| Sample costs for an agent-based proof of concept 3 | 5 |

| Security | 39 |
|---|----|
| Using foundation models on Amazon Bedrock | 39 |
| IAM roles | 39 |
| CloudWatch Logs | 39 |
| VPC | 39 |
| Let the solution build an Amazon VPC for you | 40 |
| Managing your own Amazon VPC | 40 |
| Amazon CloudFront | 41 |
| Quotas | 42 |
| Quotas for AWS services in this solution | 42 |
| Deploy the solution | 43 |
| Deployment process overview | 43 |
| AWS CloudFormation template | 44 |
| Step 1: Launch the Deployment dashboard stack | 44 |
| Step 2: Deploy a use case | 49 |
| Step 3: Deploy a use case using the Deployment dashboard wizard | 50 |
| Step 3a: Deploy a Text use case | 50 |
| Step 4: Post-deployment configuration | 56 |
| Amazon S3 bucket versioning, lifecycle policies, and cross-Region replication | 56 |
| Amazon DynamoDB backups | 56 |
| Amazon CloudWatch dashboard and alarms | 56 |
| Amazon CloudWatch Logs | 57 |
| Custom web domains with TLS v1.2 or higher certificates | 57 |
| Scaling with Amazon Kendra | 57 |
| Setting up SSO using Idp federation | 58 |
| Customizing login screen | 59 |
| Additional security considerations | 59 |
| Deploying a standalone Text use case | 59 |
| Deploying a standalone Agent use case | 70 |
| Supplying a DynamoDB chat configuration | 77 |
| Monitor the solution with Service Catalog AppRegistry | 79 |
| Activate CloudWatch Application Insights | 79 |
| Confirm cost tags associated with the solution | 81 |
| Activate cost allocation tags associated with the solution | 82 |
| AWS Cost Explorer | 83 |
| Update the solution | 84 |

| Step 1: Update Deployment dashboard | 84 |
|---|-----|
| Step 2: Migrate use case configurations (Only updates from versions below 2.0.0) | 85 |
| Step 3: Update use cases | 86 |
| Troubleshooting | 87 |
| Problem: Deploying a VPC-enabled configuration, with Create a VPC for me, fails | 87 |
| Resolution | 87 |
| Problem: Use case stack can't be deleted in CloudFormation after the Deployment dashboard | |
| stack gets deleted | 88 |
| Resolution | 88 |
| Problem: Use case UI does not reflect changes in settings | 89 |
| Resolution | 89 |
| Contact Support | 89 |
| Create case | 89 |
| How can we help? | 89 |
| Additional information | 90 |
| Help us resolve your case faster | 90 |
| Solve now or contact us | 90 |
| Uninstall the solution | 91 |
| Using the AWS Management Console | 91 |
| Using AWS Command Line Interface | 91 |
| Manual uninstall steps | 91 |
| Deleting the Amazon S3 buckets | 91 |
| Deleting the Amazon Kendra indexes | 92 |
| Deleting the CloudWatch Logs | 92 |
| Use the solution | 94 |
| Accessing the UI | 94 |
| How to update a deployment | 94 |
| How to clone a deployment | 95 |
| How to delete a deployment | 95 |
| Configuring a Large Language Model (LLM) | 95 |
| Using Amazon SageMaker AI as an LLM Provider | 96 |
| Creating a SageMaker AI endpoint | 96 |
| Advanced LLM Settings 1 | 100 |
| Amazon Bedrock Guardrails 1 | |
| Provisioned Throughput for Amazon Bedrock 1 | 101 |
| Model parameters 1 | 102 |

| Tips for managing model token limits | 103 |
|---|-----|
| Configuring a knowledge base | 103 |
| Advanced knowledge base settings | 104 |
| Knowledge base filtering | 104 |
| RAG with Role Based Access Control with Amazon Kendra | 105 |
| Configuring your prompts | 107 |
| Using the deployed Text use case | 109 |
| Chat window | 110 |
| Chat input box | 110 |
| Settings | 110 |
| Clear conversation | 110 |
| Accessing and analyzing user collected feedback | 111 |
| Custom Feedback Mappings | 113 |
| Analyzing feedback data | 115 |
| Viewing operation metrics for a deployment | 117 |
| Access CloudWatch Logs insights | 117 |
| Developer guide | 120 |
| Source code | 120 |
| Integration guide | 120 |
| Expanding supported LLMs | 120 |
| Expanding supported knowledge bases and conversation memory types | 123 |
| Building and deploying the code changes | 124 |
| Customization guide | 124 |
| Managing Cognito user pool | 124 |
| API reference | 125 |
| Deployment dashboard | 125 |
| Shared Use Case APIs | 129 |
| Text use case | 130 |
| Agent use case | 135 |
| Reference | 138 |
| Supported LLM providers | 138 |
| Anonymized data collection | 139 |
| Contributors | 141 |
| Revisions | 143 |
| Notices | 144 |

This solution facilitates the development, rapid experimentation, and deployment of generative artificial intelligence (AI) applications

Generative AI Application Builder on AWS facilitates the development, rapid experimentation, and deployment of generative artificial intelligence (AI) applications without requiring deep experience in AI. This AWS Solution accelerates development and streamlines experimentation by helping you:

- Ingest your business-specific data and documents
- Evaluate and compare the performance of large language models (LLMs)
- Run multi-step tasks and workflows with AI agents
- Rapidly build extensible applications, and deploy those applications with an enterprise-grade architecture

Generative AI Application Builder on AWS includes integrations with:

- LLMs available on Amazon Bedrock
- LLMs that you have deployed on <u>Amazon SageMaker AI</u>
- Amazon Bedrock Knowledge Bases for Retrieval-Augmented Generation (RAG)
- Amazon Bedrock Guardrails to implement safeguards and reduce hallucinations
- <u>Amazon Bedrock Agents</u> to build agentic workflows that can carry out task orchestrations and completion

Additionally, this solution enables connections to your choice of model by using LangChain connectors. These connectors are available in an <u>AWS Lambda</u> function that deploys with the solution. You can start with the no-code deployment wizard to build generative AI applications for conversational search, AI- generated chatbots, text generation, and text summarization.

This implementation guide provides an overview of the Generative AI Application Builder on AWS solution, its reference architecture and components, considerations for planning the deployment, and configuration steps for deploying the solution to the Amazon Web Services (AWS) Cloud.

This guide is intended for solution architects, business decision makers, DevOps engineers, data scientists, and cloud professionals who want to implement Generative AI Application Builder on AWS in their environment.

Use this navigation table to quickly find answers to these questions:

| If you want to | Read |
|---|-----------------------|
| Know the cost for running this solution. | Cost |
| The estimated cost for running this solution varies based on the components you deploy and the number of queries. | |
| The cost to run the Deployment dashboard with default parameters and 100 active users in the US East (N.Virginia) Region for one month is approximately \$20.12 USD per month. | |
| The cost for a Text use case deployed without RAG for 1 business user performing 100 queries per day with the LLM is approximately \$12.39 USD per month. | |
| The cost for a RAG-enabled use case with an Amazon Kendra index supporting 8,000 interactions per day is approximately \$204.26 USD per month, plus the cost of the knowledge base. | |
| Understand the security considerations for this solution. | Security |
| Know how to plan for quotas for this solution. | Quotas |
| Know which AWS Regions support this solution. | Supported AWS Regions |

| If you want to | Read |
|--|-----------------------------|
| View or download the AWS CloudForm ation template included in this solution to automatically deploy the infrastructure resources (the "stack") for this solution. | AWS CloudFormation template |
| Access the source code and optionally use the AWS Cloud Development Kit (AWS CDK) to deploy the solution. | <u>GitHub repository</u> |

Features and benefits

The Generative AI Application Builder on AWS solution provides the following features:

Rapid experimentation

This solution allows users to experiment quickly by removing the heavy lifting required to deploy multiple instances with different configurations and compare outputs and performance. Experiment with multiple configurations of various LLMs, prompt engineering, enterprise knowledge bases, guardrails, AI agents, and other parameters.

Choice and configurability

With pre-built connectors to a variety of LLMs, such as models available through Amazon Bedrock, this solution gives you the flexibility to deploy the model of your choice, as well as the AWS and leading FM services you prefer. You can also enable Amazon Bedrock Agents to fulfill various tasks and workflows.

Production-ready

Built with AWS Well-Architected design principles, this solution offers enterprise-grade security and scalability with high availability and low latency, ensuring seamless integration into your applications with high performance standards.

Extensible modular architecture

Extend this solution's functionality by integrating your existing projects or natively connecting additional AWS services. Because this is an open-source application, you can use the included LangChain orchestration layer or Lambda functions to connect with the services of your choice.

Integration with Service Catalog AppRegistry and Application Manager, a capability of AWS Systems Manager

This solution includes a <u>Service Catalog AppRegistry</u> resource to register the solution's CloudFormation template and its underlying resources as an application in both AWS Service Catalog AppRegistry and <u>AWS Systems Manager Application Manager</u>. With this integration, you can centrally manage the solution's resources.

Use cases

Question answering over enterprise data

LLMs and other foundation models have been pre-trained on a large corpus of data enabling them to perform well at many natural language processing (NLP) tasks. But most foundation models and LLMs are static and have been pre-trained, limiting their ability to accurately answer questions on topics which are either new, specialized, or proprietary. Using prompt-based learning, you can leverage the powerful NLP and text generation features of an LLM to provide richer customer experiences over your enterprise data.

Rapid generative AI prototyping

Out of the box, the solution comes bundled with various model providers and use cases. With an easy to use deployment wizard, customers can deploy pre-built use cases to enable the rapid experimentation of different generative AI prototypes and workloads.

Multi LLM comparison and experimentation

LLMs perform differently, and given your application's specific needs, you may find that one LLM suits your application better than another. This may be for reasons related to performance, accuracy, cost, creativity, or many other factors. This solution lets you quickly deploy multiple use cases enabling you to experiment with and compare different configurations until you've found what meets your needs.

Concepts and definitions

This section describes key concepts and defines terminology specific to this solution:

admin user

Within the context of this guide, the admin user is the one responsible for managing the content contained within the deployment. This user gets access to the Deployment dashboard UI and is primarily responsible for curating the business user experience. This is our primary target customer.

business user

Within the context of this guide, the business user represents the individuals who the use case has been deployed for. They are the consumers of the knowledge base and the customer responsible for evaluating and experimenting with the LLMs.

Deployment dashboard

The Deployment dashboard is a web interface that serves as a management console for admin users to view, manage, and create their *use cases*. This dashboard enables customers to rapidly experiment, iterate, and productionize various AI/ML workloads leveraging LLMs.

DevOps user

Within the context of this guide, the DevOps user is the one responsible for deploying the solution within the AWS account and for managing the infrastructure, updating the solution, monitoring performance, and maintaining the overall health and lifecycle of the solution.

use case

Use cases are isolated applications from the overall solution which integrate with LLMs to enable richer customer experiences by enabling the addition of a natural language interface into new or existing applications. Use cases are deployable through the Deployment dashboard or on their own.

🚯 Note

For a general reference of AWS terms, see the AWS Glossary.

Architecture overview

This section provides two reference implementation architecture diagrams for the components deployed with this solution.

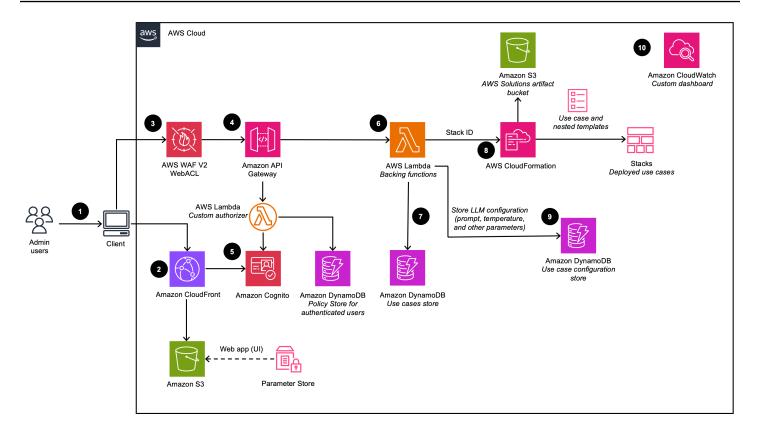
Architecture diagrams

To support multiple use cases and business needs, this solution provides two AWS CloudFormation templates:

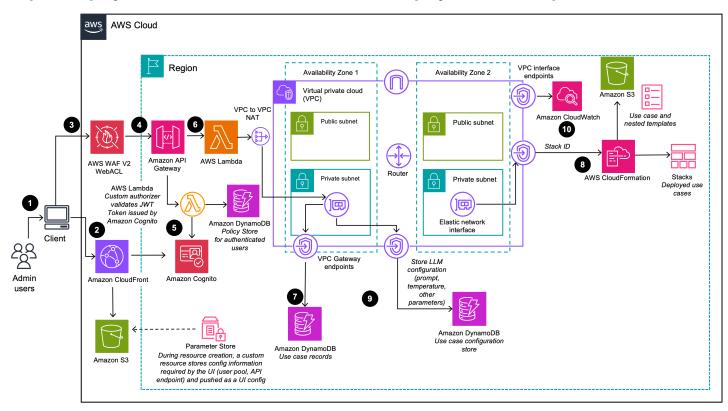
- Deployment dashboard The Deployment dashboard is a web interface that serves as a management console for admin users to view, manage, and create their use cases. This dashboard enables customers to rapidly experiment, iterate, and productionize various AI/ML workloads leveraging LLMs.
- 2. **Text use case** The Text use case enables users to experience a natural language interface using generative AI. This use case can be integrated into new or existing applications, and is deployable through the Deployment dashboard or independently through a provided URL.

Deployment dashboard

Depicts Deployment dashboard architecture (when deployed with VPC option disabled)



Depicts Deployment dashboard architecture (when deployed with VPC option enabled)



i Note

AWS CloudFormation resources are created from AWS Cloud Development Kit (AWS CDK) constructs.

The high-level process flow for the solution components deployed with the AWS CloudFormation template is as follows:

- 1. Admin users log in to the Deployment Dashboard user interface (UI).
- <u>Amazon CloudFront</u> delivers the web UI, which is hosted in an <u>Amazon Simple Storage Service</u> (Amazon S3) bucket.
- <u>AWS WAF</u> protects the APIs from attacks. This solution configures a set of rules called a web access control list (web ACL) that allows, blocks, or counts web requests based on configurable, user defined web security rules and conditions.
- 4. The web UI leverages a set of REST APIs that are exposed using Amazon API Gateway.
- 5. Amazon Cognito authenticates users and backs both the CloudFront web UI and API Gateway.
- <u>AWS Lambda</u> provides the business logic for the REST endpoints. This *backing* Lambda function manages and creates the necessary resources to perform use case deployments using <u>AWS</u> <u>CloudFormation</u>.
- 7. <u>Amazon DynamoDB</u> stores the list of deployments.
- 8. When a new use case is created by the admin user, the *backing* Lambda function initiates a CloudFormation stack creation event for the requested use case.
- 9. All of the LLM configuration options provided by the admin user in the deployment wizard are saved in DynamoDB. The deployment uses this DynamoDB table to configure the LLM at runtime.
- 10Using <u>Amazon CloudWatch</u>, this solution collects operational metrics from various services to generate custom dashboards that allow you to monitor the solution's performance and operational health.

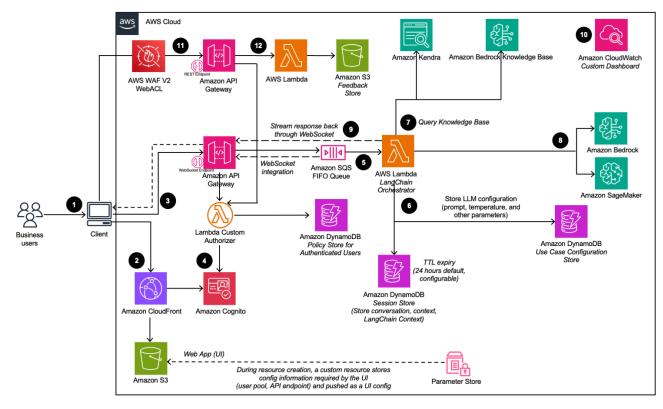
1 Note

• If you choose to deploy this solution in an Amazon VPC, the data will be routed within your private network.

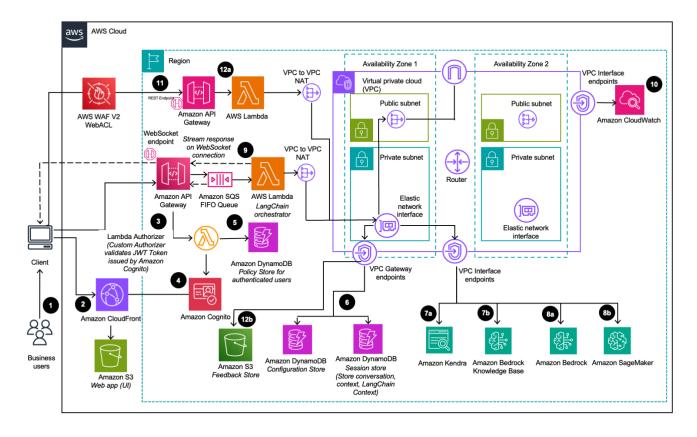
 Although the Deployment dashboard can be launched in most AWS Regions, the deployed use cases have certain restrictions based on service availability. See <u>Supported</u> AWS Regions for more details.

Text use case





Depicts Text use case architecture (when deployed with VPC option enabled)



The high-level process flow for the solution components deployed with the AWS CloudFormation template is as follows:

- 1. Admin users deploy the use case using the Deployment Dashboard. <u>Business users</u> log in to the use case UI.
- 2. CloudFront delivers the web UI which is hosted in an S3 bucket.
- 3. The web UI leverages a WebSocket integration built using API Gateway. The API Gateway is backed by a custom <u>Lambda authorizer</u> function, which returns the appropriate <u>AWS Identity</u> <u>and Access Management</u> (IAM) policy based on the Amazon Cognito group the authenticating user belongs to. The policy is stored in DynamoDB.
- 4. Amazon Cognito authenticates users and backs both the CloudFront web UI and API Gateway.
- 5. Incoming requests from the business user are passed from API Gateway to an <u>Amazon SQS</u> <u>queue</u> and then to the *LangChain Orchestrator*. The *LangChain Orchestrator* is a collection of Lambda functions and layers that provide the business logic for fulfilling requests coming from the business user. The queue enables the asynchronous operation of the API Gateway to Lambda integration. The queue passes connection information to the Lambda functions which will then post results directly back to the API Gateway websocket connection to support long running inference calls.

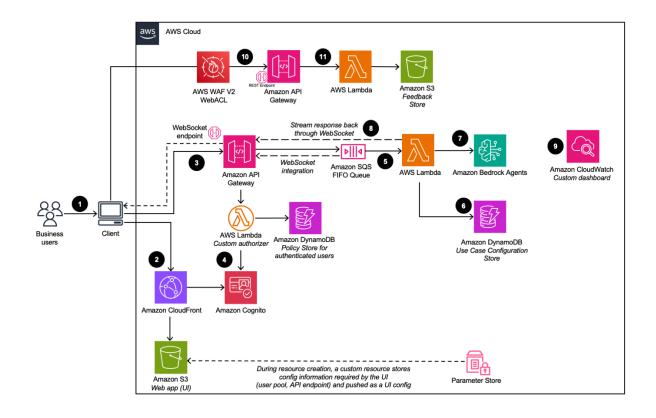
- 6. The *LangChain Orchestrator* uses Amazon DynamoDB to get the configured LLM options and necessary session information (such as the chat history).
- If the deployment has a knowledge base enabled, then the LangChain Orchestrator leverages <u>Amazon Kendra</u> or <u>Knowledge Bases for Amazon Bedrock</u> to run a search query to retrieve document excerpts.
- 8. Using the chat history, query, and context from the knowledge base, the *LangChain Orchestrator* creates the final prompt and sends the request to the LLM hosted on <u>Amazon Bedrock</u> or <u>Amazon SageMaker AI</u>.
- 9. When the response comes back from the LLM, the *LangChain Orchestrator* streams the response back through the API Gateway WebSocket to be consumed by the client application.
- 10Using Amazon CloudWatch, this solution collects operational metrics from various services to generate custom dashboards that allow you to monitor the deployment's performance and operational health.
- 11If feedback collection is enabled, a REST API endpoint, leveraging Amazon API Gateway is made available for the collection of user feedback.
- 12.The feedback backing lambda, augments the submitted feedback with additional use case specific metadata (e.g. model used) and stores the data in Amazon S3 for later analysis and reporting by the DevOps users.

🚯 Note

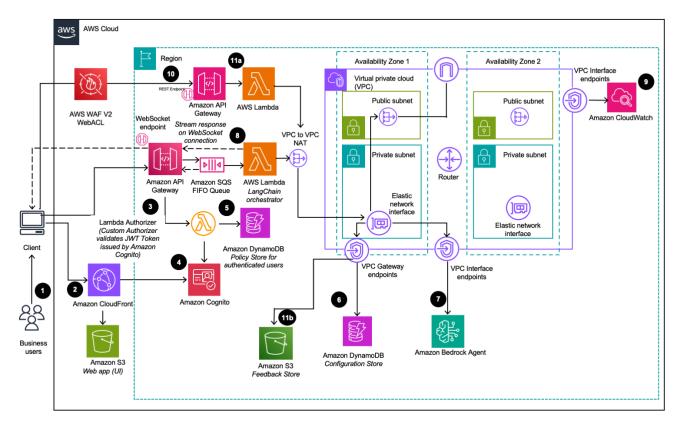
If you choose to deploy this solution in an Amazon VPC, the data will be routed to your private network.

Agent use case

Depicts Agent use case architecture (when deployed with VPC option disabled)



Depicts Agent use case architecture (when deployed with VPC option enabled)



The high-level process flow for the solution components deployed with the AWS CloudFormation template is as follows:

- 1. Admin users deploy the use case using the Deployment Dashboard. <u>Business users</u> sign in to the use case UI.
- 2. CloudFront delivers the web UI which is hosted in an S3 bucket.
- 3. The web UI leverages a WebSocket integration built using API Gateway. The API Gateway is backed by a custom Lambda authorizer function, which returns the appropriate <u>AWS Identity</u> <u>and Access Management</u>(IAM) policy based on the Amazon Cognito group the authenticating user belongs to. The policy is stored in DynamoDB.
- 4. Amazon Cognito authenticates users and backs both the CloudFront web UI and API Gateway.
- 5. Incoming requests from the business user are passed from API Gateway to an <u>Amazon SQS</u> <u>queue</u> and then to the AWS Lambda function. The queue enables the asynchronous operation of the API Gateway to Lambda integration. The queue passes connection information to the Lambda function which will then post results directly back to the API Gateway websocket connection to support long running inference calls.
- 6. The AWS Lambda function uses Amazon DynamoDB to get the use case configurations as needed
- 7. Using the user input and any relevant use case configurations, the AWS Lambda function builds and sends a request payload to the configured <u>Amazon Bedrock Agent</u> to fulfill the user intent.
- 8. When the response comes back from the Amazon Bedrock Agent, the Lambda function streams the response back through the API Gateway WebSocket to be consumed by the client application.
- 9. Using Amazon CloudWatch, this solution collects operational metrics from various services to generate custom dashboards that allow you to monitor the deployment's performance and operational health.
- 10If feedback collection is enabled, a REST API endpoint, leveraging Amazon API Gateway is made available for the collection of user feedback.
- 11. The feedback backing lambda, augments the submitted feedback with additional use case specific metadata and stores the data in Amazon S3 for later analysis and reporting by the DevOps users.

🚯 Note

If you choose to deploy this solution in an Amazon VPC, data will be routed within your private network.

AWS Well-Architected design considerations

This solution was designed with best practices from the <u>AWS Well-Architected Framework</u> which helps customers design and operate reliable, secure, efficient, and cost-effective workloads in the cloud.

This section describes how the design principles and best practices of the Well-Architected Framework were applied when building this solution.

Operational excellence

This section describes how we architected this solution using the principles and best practices of the <u>operational excellence pillar</u>.

- We built the solution as infrastructure-as-code using Amazon CloudFormation.
- Lambda functions push custom metrics to CloudWatch and a custom CloudWatch dashboard to monitor the health of the solution.
- The solution components are highly modularized, providing the flexibility to choose which components to deploy.

Security

This section describes how we architected this solution using the principles and best practices of the <u>security pillar</u>.

- The Deployment dashboard and all use cases are authenticated and authorized with Amazon Cognito.
- All inter-service communications use AWS IAM roles.
- All solution roles follows least-privilege access; meaning, only the minimum permissions required are granted.
- All data storage including S3 buckets, DynamoDB, and Amazon Kendra have encryption at rest.

Reliability

This section describes how we architected this solution using the principles and best practices of the <u>reliability pillar</u>.

- Architecture based on serverless paradigm.
- We built the architecture for on-demand, horizontal scalability, and automatic recovery from failure of underlying infrastructure.
- The architecture includes buffering and throttling requests to not overwhelm underlying endpoints.

Performance efficiency

This section describes how we architected this solution using the principles and best practices of the <u>performance efficiency pillar</u>.

- The solution uses DynamoDB, a fully managed serverless NoSQL database with on-demand scaling.
- The solution uses Amazon S3 for object storage and to host a website (through CloudFront) to provide low cost, scalable, with 11 9s durability.

Cost optimization

This section describes how we architected this solution using the principles and best practices of the <u>cost optimization pillar</u>.

• Where possible, we built the solution to use serverless architecture; so you only pay for what you use.

Sustainability

This section describes how we architected this solution using the principles and best practices of the <u>sustainability pillar</u>.

• The solution's modular, componentized architecture provides the flexibility to customize resources to be provisioned for individual use cases.

- The architecture uses serverless compute and storage, which optimizes resource utilization.
- As a cloud-based solution, this solution benefits from shared resources, networking, power cooling, and physical facilities.

Architecture details

This section describes the components and AWS services that make up this solution and the architecture details on how these components work together.

AWS services in this solution

| AWS service | Description |
|------------------------|---|
| Amazon API Gateway | Core . This service provides the REST APIs for the Deployment dashboard and the WebSocket API for the use case. |
| AWS CloudFormation | Core . This solution is distributed as a CloudFormation template, and CloudForm ation deploys the AWS resources for the solution. |
| Amazon CloudFront | Core . CloudFront serves the web content hosted in Amazon S3. |
| Amazon Cognito | Core . This service handles user management and authentication for the API. |
| <u>Amazon DynamoDB</u> | Core . DynamoDB stores deployment information and configuration details for the Deployment dashboard. It stores chat history and conversation IDs in the Text use case to enable conversation history and query disambiguation. |
| AWS Lambda | Core . The solution uses Lambda functions to: |
| | * Back the REST and WebSocket API endpoints * Handle the core logic of each use case orchestrator * Implement custom resources during CloudFormation deployment |

| AWS service | Description |
|----------------------------|---|
| Amazon S3 | Core . Amazon S3 hosts the static web content. |
| <u>Amazon CloudWatch</u> | Supporting . This solution publishes logs from solution resources to <u>CloudWatch Logs</u> , and publishes metrics to <u>CloudWatch metrics</u> . The solution also creates a <u>CloudWatch dashboard</u> to view this data. |
| <u>AWS Systems Manager</u> | Supporting . Systems Manager provides application-level resource monitoring and visualization of resource operations and cost data. Also used to store configuration data in Parameter Store. |
| AWS WAF | Supporting. AWS WAF is deployed in front of the API Gateway deployment to protect it. |
| <u>Amazon Bedrock</u> | Optional . The solution leverages Amazon Bedrock to access foundation or customized models, Amazon Bedrock Agents, and Amazon Bedrock Knowledge Bases. Amazon Bedrock is the recommended integration to keep your data from leaving the AWS network. |
| <u>Amazon Kendra</u> | Optional . In the Text use case, admin users can optionally decide to connect an Amazon Kendra index to use as a knowledge base for the conversation with the LLM. This can be used to inject new information into the LLM giving it the ability to use that information in its responses. |

| Generative AI Application Builder on AWS | Implementation Guide |
|--|--|
| AWS service | Description |
| <u>Amazon SageMaker AI</u> | Optional . The solution can integrate with an Amazon SageMaker AI inference endpoint to access FMs that are hosted within your AWS account and Region and is a preferred integration to keep your data from leaving the AWS network. |
| | Solution in the solution in the same Region where the inference endpoint is available. |
| <u>Amazon Virtual Private Cloud</u> | Optional . The solution provides the option to deploy components with a VPC-enabled configuration. While deploying the solution with a VPC-enabled configuration, you have the option to let the solution create a VPC for you, or use an existing VPC that exists in the same account and Region where the solution will be deployed (Bring Your Own VPC). If the solution creates the VPC, it creates the necessary network components that includes, subnets, security groups and its rules, route tables, network ACLs, NAT Gateways, Internet Gateways, VPC endpoints, and its policies. |

Deployment dashboard

API Gateway custom authorizers

Beneath the surface, Lambda custom authorizers for API Gateway are used for all API calls (both RESTful and WebSocket based) to validate if a given user has permission to perform an action based on the group(s) they belong to. This custom authorizer is backed by a DynamoDB table containing the policies for each group. On invocation of an API, API Gateway invokes the custom authorizer Lambda function, which decodes the provided Amazon Cognito access token to determine which user groups the user belongs to. The policy table is then queried by group name to return the relevant policy for that group.

On every new use case deployment, the admin policy is updated to store a new statement allowing the **execute-api:Invoke** action on that use case's API. When use cases are deleted, the corresponding statement is removed from the policy.

For the groups created for an individual use case, only a single statement is present in the policy, allowing the **execute-api:Invoke** action on only that use case's API.

Due to this structure, any user belonging to a use case's group can access that use case's API. A single user can also be manually added to multiple groups to allow that user to use multiple use cases.

<u> M</u>arning

You can also manually edit the policies for a given group in the policy table if you want to grant access to a new use case to an existing group of users. The use case group is deleted when the use case is deleted (even if you have made manual edits), so proceed with caution when deleting a use case.

In the case where a use case stack is deployed standalone (without the use of the Deployment dashboard), an <u>Amazon Cognito user pool</u> is created for that deployment containing a single user with access to that use case's API. This user pool belongs only to this use case and is not shared across other standalone deployments.

Text use case

Streaming support

In a chat application, latency is an important metric to enable a responsive user experience. The potential for LLM inferences to take from seconds to minutes, provides challenges in how to best serve content to customers. For this reason, several LLM providers allow streaming responses back to the caller. Instead of waiting for the entire inference to complete before returning a response, each token can be returned when it is available.

To support the use of this feature, the Text use case has been designed to use a WebSocket API to back the chat experience. This WebSocket is deployed through API Gateway. The use of a WebSocket API enables a connection to be created at the beginning of a chat session and for responses to be streamed through that socket. This allows frontend applications to provide a better user experience.

1 Note

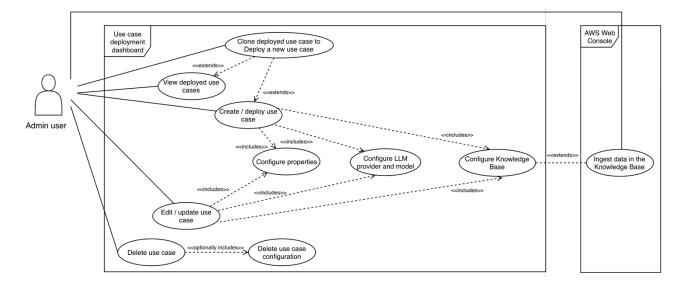
Even if a model provides streaming support, this does not necessarily mean that the solution will be able to stream responses back through the WebSocket API. There is a need for the solution to enable custom logic to support streaming for each model provider. If streaming is available, admin users will be able to enable/disable this feature at deployment time.

How the Generative AI Application Builder on AWS solution works

The admin user primarily interfaces with the Deployment dashboard to view, create, and manage new and existing use case deployments. Through this dashboard, the admin user has access to the following actions:

- View list of deployments
- Create new deployments
- Edit existing deployments
- Clone a deployment's configuration to create a new deployment
- Delete a deployment (deprovision the resources through a CloudFormation delete)
- · Permanently delete the configuration details of a deployment

Depicts Use case diagram for the admin user of the Deployment dashboard



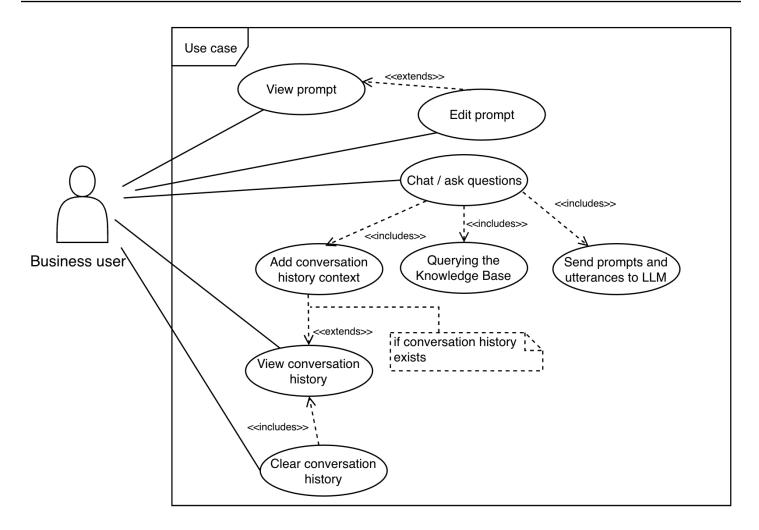
1 Note

The admin user might not have direct access to the AWS console. In that case, the admin user must work with the DevOps user to support actions such as ingesting data into a Kendra knowledge base.

For the Text use case, the business user gets access to a user interface enabling them to chat with the LLM. The specifics of this configuration are controlled by the deployment settings configured by the admin user. In the Text use case, the business user has access to the following actions:

- Send messages through the chat interface
- View conversation history
- Clear the conversation history
- View prompt
- Edit prompt

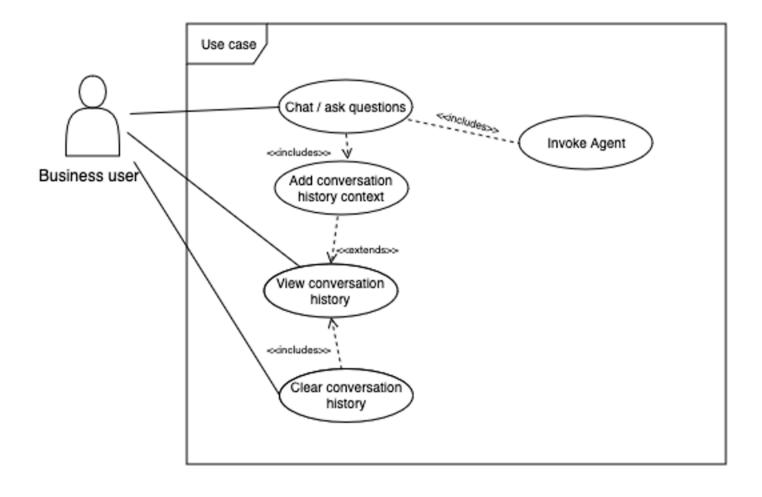
Depicts Use case diagram for the business user of the Text use case



With the Agent use case, the business user can access a UI for chatting with the configured Amazon Bedrock Agent. The admin user can configure these specifics in the deployment settings. In the Agent use case, the business user has access to the following actions:

- Send messages through the chat interface
- View conversation history
- Clear the conversation history

Depicts Use case diagram for the business user of the Agent use case



Plan your deployment

This section describes the <u>cost</u>, <u>security</u>, <u>Region</u>, and <u>quota</u> considerations for planning your deployment.

🛕 Important

This solution leverages Amazon Bedrock as the primary service for accessing AI-generated models. You must first request access to models before they are available for use within the solution. For details, refer to <u>Model access</u> in the *Amazon Bedrock User Guide*.

Supported AWS Regions

🛕 Important

This solution optionally uses the Amazon Bedrock and Amazon Kendra services, which are not currently available in all AWS Regions. You must launch this solution in an AWS Region where these services are available. For the most current availability of AWS services by Region, see the AWS Regional Services List.

Generative AI Application Builder on AWS is supported in the following AWS Regions:

| Region name | |
|-------------------------------|--------------------|
| US East (Ohio) | Canada (Central) |
| US East (N. Virginia) | Europe (Frankfurt) |
| US West (Northern California) | Europe (Ireland) |
| US West (Oregon) | Europe (London) |
| Asia Pacific (Mumbai) | Europe (Milan) |
| Asia Pacific (Seoul) | Europe (Paris) |

| Region name | |
|--------------------------|---------------------------|
| Asia Pacific (Singapore) | Europe (Stockholm) |
| Asia Pacific (Sydney) | Middle East (Bahrain) |
| Asia Pacific (Tokyo) | South America (São Paulo) |

Note

If using a foundation model accessed outside of AWS in your deployments, check with the model provider which Regions their APIs are available in. If their APIs are only available in certain Regions, you might experience instability in the form of high latency or even time outs. It's also important to check with your organization's legal and compliance teams to evaluate the considerations of data crossing regional boundaries.

Cost

With this AWS Solution, you pay only for the resources you use and there are no minimum fees or setup charges. Users pay for the dashboard used to launch Generative AI use cases and, and for any use cases that are deployed. The cost of deployed use cases depends on the configurations. Example configurations:

- 1. A simple Deployment dashboard which costs approximately \$20 USD per month.
- A simple production-ready chatbot use case deployed with default settings running in US East (N. Virginia), powered by Amazon Bedrock without access to documents, which also costs around \$200 USD per month.
- 3. A scaled system in an Amazon VPC use case that supports 8,000 queries per day over tens of thousands of documents, which costs around \$1,400 USD per month. The cost of the use case will vary depending on the configuration, such as Text use cases with different model providers, with or without Retrieval Augmented Generation (RAG) enabled, and so on.

| Workload description | Estimated cost (USD/month) |
|---|----------------------------|
| Sample cost for Deployment dashboard | \$20/month |
| Sample costs for a text-based proof of concept (includes Deployment dashboard and 1 Text use case, ~100 interactions per day) | \$40/month |
| Sample costs for a highly scalable generative AI query engine (Includes Deployment dashboard, 1 Text use case, and an Amazon Kendra Index for RAG up to 100K documents with ~8K queries per day, with <u>VPC enabled</u> | \$1,400/month |
| Sample costs for an agent-based proof of concept (Includes Deployment dashboard, 1 Agent use case with Amazon Bedrock Knowledge Bases and Amazon Bedrock Guardrails enabled, ~100 interactions per day) | \$840/month |

🔥 Important

These examples are only intended to help you estimate the costs for your specific workloads. The use of different LLMs, configurations, or AWS services can change your costs (example, serverless/on-demand billing vs. provisioned/time-billed). To manage costs, we recommend <u>creating a budget</u> through <u>AWS Cost Explorer</u>. Prices are subject to change. For full details, refer to the pricing webpage for each AWS service used in this solution.

Sample costs for running the Deployment dashboard

The following table provides the cost breakdown for a Deployment dashboard with default parameters and 100 active users in the US East (N. Virginia) Region for one month, which will cost about \$20/month.

| AWS service | Dimensions | Cost [USD] |
|--|--|------------|
| API Gateway, DynamoDB, CloudFront, Amazon S3, Lambda, Systems Manager Parameter Store | 5,000 512 KB REST API calls per month without caching enabled | \$1.97 |
| Amazon Cognito | 100 active users per month with advanced security features enabled and no users signing in through SAML or OIDC federation | \$5.55 |
| AWS WAF | 10,000 web requests across 1 web ACL and 7 defined rules without any rule groups | \$12.60 |
| Total Deployment dashboard cost | | \$20.12 |

Sample costs for a text-based proof of concept

A Deployment dashboard can have many use cases deployed at a given time. The following table shows the cost breakdown of a use case deployed without RAG for 1 business user performing 100 queries per day with the LLM. Queries are sent as a text message on the WebSocket and the response is streamed back as tokens with the assumption that streaming is enabled. With an Amazon Bedrock Titan Text Express model, the cost of running this use case is about \$15/month.

| AWS service | Dimensions | Cost [USD] |
|--|--|------------|
| API Gateway (WebSocket), CloudFront, Lambda, Amazon S3, AWS Systems Manager Parameter Store | 100 chat interactions per day. Average message size 32 KB per message and 5 minutes per connection. | \$0.61 |
| CloudWatch | 1.5 GB CloudWatch logs with verbose mode on for experimentation | \$7.23 |
| Amazon DynamoDB | Conversation history table, 1 GB storage | \$3.05 |
| | LLM configuration table, 1 GB storage | |
| Subtotal of the use case costs (not including LLMs) | | \$10.89 |
| Amazon Bedrock (Titan Text Express) | Assumptions for 100 interacti ons per day: | \$1.50 |
| | * Monthly cost for 190K input tokens per day = \$0.04 × 30 * Monthly cost for 16K output tokens per day = \$0.01 × 30 | |
| Total application cost with Amazon Bedrock (Titan Text Express) | \$10.89 (Use Case cost) + \$1.50 (Amazon Bedrock cost) | \$12.39 |

🚯 Note

The costs of inference calls made to services outside the AWS network are not included in these estimates. Refer to the pricing guide of your LLM provider if you're not using an AWS model provider.

Pricing guides for AWS services can be found at: <u>Amazon Bedrock pricing</u> and <u>Amazon</u> SageMaker AI pricing.

Sample costs for a highly scalable generative AI query engine

The following table provides the cost breakdown of a RAG-enabled use case with a Kendra index supporting 8000 interactions/day. With Amazon Bedrock's Titan Text Express model as the LLM, this use case costs about \$1200/month

| AWS service | Dimensions | Cost [USD] |
|--|---|------------|
| API Gateway (WebSocket) | 8000 chat interactions per day. Average message size 32 KB per message and 5 minutes per connection. | \$38.89 |
| CloudFront | 240,000 requests per month with 100 GB data transferred out to the internet and 1 GB data transferred out to the origin | \$8.76 |
| Amazon Bedrock (Titan Text Express) | Assumptions: Input tokens = promptTem plate (400) + context (400)+ chatHistory (1080) + query Input tokens (20)= 1,900 Output tokens = 160 (average) With 8,000 transactions a day, Daily Input Tokens cost (1,900 x 8,000 = 15,200,000 tokens x 0.0002/1000 price per token) | \$114.30 |

| AWS service | Dimensions | Cost [USD] |
|---------------------|--|--|
| | Daily Output Tokens cost (160 x 8,000 = 1,280,000 tokens x 0.0006/1000 price per token) | |
| | Monthly cost ((\$3.04 + \$0.77) x 30) | |
| CloudWatch | 24 metrics using 5 GB data ingested for logs and 1 dashboard | \$9.72 |
| DynamoDB | DynamoDB table to keep track of conversation history with each record up to 1 KB data, 8,000 read and writes per day | \$11.70 |
| Lambda | Container size - 128 MB, 512 MB ephemeral storage, 2 Lambda functions used for authorization Container size - 256 MB, 512 MB ephemeral storage, 5 requests per second with 20 seconds average compute time | \$20.89 |
| Total use case cost | | \$204.26/month + knowledge base cost (see below) |

i Note

The costs of API calls made to any services outside of the AWS network are not included in these estimates. See the pricing guide of your LLM provider if not using Amazon Bedrock.

Costs for adding a knowledge base

Knowledge base costs will vary based on the type of knowledge base used, and (in the case of Bedrock) the backing vector store used by the knowledge base. Provisioning and managing the knowledge bases is outside of the scope of the solution.

Amazon Kendra

The solution can provision a Kendra index for you, or you can bring your own. The cost for running a configuration suited to the above highly scalable generative AI query engine is as follows:

| AWS service | Dimensions | Cost [USD] |
|---------------|---|------------|
| Amazon Kendra | 0-8,000 queries a day and up to 100,000 documents with Amazon Kendra Enterpris e Edition with 0-50 data sources | \$1,008.00 |

Note

You can share the Amazon Kendra index between use cases, but this can drive up the number of queries per index. If this falls outside the Amazon Kendra Enterprise edition, additional charges will apply.

Amazon Bedrock Knowledge Bases

The solution does not manage or provision any resources related to Amazon Bedrock Knowledge Bases. Amazon Bedrock does not incur cost for using the knowledge base feature itself, however you will be charged for the usage of the embedding model used by your use case on each query. Additionally, the backing vector store for your knowledge base (for example, an index in <u>Amazon</u> <u>OpenSearch Service</u>, or a database inside Amazon Relational Database Service) will have an associated cost which cannot be provided or calculated here.

For the above highly scalable generative AI query engine scenario, the costs incurred by this service for calling the Amazon Bedrock embeddings model are as follows:

| AWS service | Dimensions | Cost [USD] |
|--|---|------------|
| Amazon Bedrock (Amazon Titan Text Embeddings) | 8,000 queries a day with 1,900 input tokens per query = 15,200,000 tokens = \$0.30 USD per day. Daily cost x 30 days = \$9.00 USD monthly cost | \$9.00 |
| Amazon OpenSearch Service (Serverless) Sample Usage | Basic serverless configura tion with 4 x OpenSearch Compute Unit (OCU) (billable minimum) = \$23.04 USD per day Daily cost x 30 days = \$691.20 USD [NOTE] ==== This provides a rough estimate, as some workloads will require more OCUs, while customers with existing provisioned OpenSearch resources will incur less cost here. ==== | \$691.20 |
| Total additional cost | | \$ 700.20 |

Incremental cost of enabling Amazon VPC for a use case

The following table provides the cost breakdown of enabling Amazon VPC for a use case deployed in two AZs.

| AWS service | Dimensions | Cost [USD] |
|------------------------------------|--|------------|
| Amazon NAT Gateway | Assumption: 2 AZ deploymen t, with a NAT Gateway in each AZ. 100 GB of data processed through NAT Gateway 730 hours, 100 GB data processed per month | \$74.70 |
| AWS PrivateLink (VPC Endpoints) | Assumptions: 2 AZ deployment, with 1 private subnet in each AZ and 1 VPC Endpoint with 2 elastic network interfaces (ENIs). 6 VPC endpoints, 2 ENIs per VPC endpoint, 730 hours with 1,024 GB data processed in a month | \$97.84 |
| Public IPv4 address | Assumption: 2 AZ deploymen t, 1 public subnet in each AZ with a NAT Gateway in each public subnet. Each NAT Gateway configured with 1 active public IPv4. 2 active public IPv4 address x 730 hours in a month x \$0.005 hourly charge = \$7.3 USD | \$7.30 |
| Additional cost | | \$179.93 |

| AWS service | Dimensions | Cost [USD] |
|------------------|------------|------------|
| (for Amazon VPC) | | |

Cost implications when using Provisioned Throughput

Provisioned throughput costs will vary based on the type of model you've provisioned and your commitment period as well as Model Units selected for the commitment period. There is an additional cost associated with using Provisioned Throughput. As an example, when using Anthropic Claude Instant or Claude 2.x models or Amazon Titan Text Express, your prices per hour would look like:

| Anthropic models | Price per hour per model with no commitment | Price per hour per model unit for 1- month commitment | Price per hour per model unit for 6- month commitment |
|------------------------------|---|---|---|
| Claude Instant | \$44.00 | \$39.60 | \$22.00 |
| Claude 2.0/2.1 | \$70.00 | \$63.00 | \$35.00 |
| Amazon Titan Text Express | \$20.50 | \$18.40 | \$14.80 |

For more information and most up-to-date pricing, you can refer **Bedrock Pricing**.

Cost for using cross-region inference

There is no additional cost for routing or data transfer for using <u>cross-region inference</u>. You pay the same price per token for models as in your source or primary Region.

Sample costs for an agent-based proof of concept

When you use Amazon Bedrock Agents, you're charged based on the components comprising the agent, such as the backing model and knowledge base (if RAG is enabled), along with additional capabilities that you add. The following table shows the cost breakdown of an Agent use case configured with an on-demand Claude 3.5 Sonnet model, Amazon Bedrock Knowledge Bases, and Amazon Bedrock Guardrails.

Similar to the <u>cost for adding Amazon Bedrock Knowledge Bases</u>, this solution doesn't manage or provision resources related to Amazon Bedrock Agents. The solution also doesn't incur cost for using Amazon Bedrock Knowledge Bases, but does incur cost for:

- Using the embedding model for each query that is sent to it
- The backing vector store for your knowledge base (for example, an index in Amazon OpenSearch Service, or a database inside Amazon RDS)

The following table assumes 100 interactions per day with 1,900 input tokens and 160 output tokens per query.

Note

For this sample Agent use case, if there were an action group configured to use an external API, those costs would be additional. They are outside the scope of the calculations in this table.

| AWS service | Dimensions | Cost [USD] |
|---|--|------------|
| API Gateway (WebSocke t), CloudFront, Lambda, Amazon S3, Systems Manager Parameter Store | 100 chat interactions per day, average message size 32 KB per message, 5 minutes per connection | \$0.61 |
| CloudWatch | 1.5 GB CloudWatch Logs with verbose mode on for experimentation | \$7.23 |
| DynamoDB | LLM configuration table for 1KB record size and 1 GB storage | \$0.25 |
| Subtotal of costs (not including LLMs) | | \$8.09 |

| AWS service | Dimensions | Cost [USD] |
|---|--|------------|
| Anthropic Claude 3.5 Sonnet | * Daily cost for 190K input tokens per day (0.003/1,000 tokens) = \$0.57 + | \$24.30 |
| | Daily cost × 30 days = \$17.10 * Daily cost for 16K output tokens per day (0.015/1,000 tokens) = \$0.24 + | |
| | Daily cost × 30 days = \$7.20 | |
| Amazon Bedrock (Amazon Titan Text Embeddings v2) for Amazon Bedrock Knowledge | Daily cost for 190K input tokens per day (0.00002/ 1000 tokens) = 0.004 | \$0.12 |
| Bases | Daily cost × 30 days = \$0.12 | |
| Amazon OpenSearch Service (Serverless) sample usage | Basic serverless configura tion with 4 × OpenSearch Compute Unit (OCU) (billable minimum) = \$23.04 per day | \$691.20 |
| | Daily cost × 30 days = \$691.20 | |

| AWS service | Dimensions | Cost [USD] |
|---|--|------------|
| Amazon Bedrock Guardrails | 190K tokens is roughly equivalent of 760K (190,000 × 4) characters and 3,800 text units (760K characters / 200) Consider a guardrail configure d with content filters, personally identifiable information (PII) filter, sensitive information filter (regular expression) and word filters Daily content filter cost (0.75/1000 text units) + PII filter cost (\$0.1/1000 text units) + sensitive information filters = \$2.85 + \$0.38 + \$0 + \$0 Monthly cost = Daily cost × 30 days = \$96.90 | \$96.90 |
| Total application cost for an agent backed by Anthropic Claude 3.5 Sonnet | \$8.09 (use case cost) + \$812.52 (other agent configurations) | \$820.61 |

🚯 Note

Refer to the pricing guide of your LLM provider if you're not using an AWS model provider. Pricing guides for AWS services can be found at: <u>Amazon Bedrock pricing</u> and <u>Amazon</u> <u>SageMaker AI pricing</u>.

Security

When you build systems on AWS infrastructure, security responsibilities are shared between you and AWS. This <u>shared responsibility model</u> reduces your operational burden because AWS operates, manages, and controls the components including the host operating system, virtualization layer, and physical security of the facilities in which the services operate. For more information about AWS security, visit AWS Cloud Security.

Using foundation models on Amazon Bedrock

Amazon Bedrock hosts a collection of models from Amazon Titan models to other leading foundation models (FMs). When using Amazon Bedrock, all models are hosted within the AWS infrastructure. This means that when using Amazon Bedrock as the LLM provider, all of your inference requests will remain within the AWS network and network traffic will not leave your Region.

🚯 Note

All foundation models (FMs) available through Amazon Bedrock are hosted directly on AWS infrastructure managed and owned by AWS. Model providers do not have access to customer data such as prompts and continuations, or Amazon Bedrock service logs. For additional information about Amazon Bedrock's security posture, refer to <u>Data protection</u> in Amazon Bedrock in the *Amazon Bedrock User Guide*.

IAM roles

IAM roles allow customers to assign granular access policies and permissions to services and users on the AWS Cloud. This solution creates IAM roles that grant the solution's Lambda functions access to create Regional resources.

CloudWatch Logs

If verbose logs are enabled, depending on the data and prompt used, sensitive information may be logged.

VPC

The solution provides two options for Amazon VPC configuration:

- 1. Let the solution build an Amazon VPC for you.
- 2. Managing and bringing your own Amazon VPC for use within the solution.

Let the solution build an Amazon VPC for you

If you select the option to let the solution build an Amazon VPC, it will deploy as a 2-AZ architecture by default with a CIDR range 10.10.0.0/20. You have the option to use <u>Amazon VPC IP</u> <u>Address Manager (IPAM)</u>, with 1 public subnet and 1 private subnet in each AZ. The solution creates NAT Gateways in each of the public subnets, and configures Lambda functions to create the <u>ENIs</u> in the private subnets. Additionally, this configuration creates route tables and its entries, security groups and its rules, network ACLs, VPC endpoints (gateway and interface endpoints).

Managing your own Amazon VPC

When deploying the solution with an Amazon VPC, you have the option to use an existing Amazon VPC in your AWS account and Region. We recommended that you make your VPC available in at least two availability zones to ensure high availability. Your VPC must also have the following VPC endpoints and their associated IAM policies for your VPC and route table configurations.

For a Deployment dashboard Amazon VPC

- 1. Gateway endpoint for DynamoDB.
- 2. Gateway endpoint for S3.
- 3. Interface endpoint for CloudWatch.
- 4. Interface endpoint for AWS CloudFormation.

For a use case Amazon VPC

- 1. Gateway endpoint for DynamoDB.
- 2. Gateway endpoint for S3.
- 3. Interface endpoint for CloudWatch.
- 4. Interface endpoint for Systems Manager Parameter Store.

Note

The solution only requires com.amazonaws.region.ssm.

- 5. <u>Interface endpoint for Amazon Bedrock (bedrock-runtime, agent-runtime, bedrock-agent-runtime)</u>.
- 6. Optional: If the deployment will use Amazon Kendra as a knowledge base, then an <u>interface</u> <u>endpoint for Amazon Kendra</u> is needed.
- 7. Optional: if the deployment will use any LLM under Amazon Bedrock, then an <u>interface endpoint</u> for Amazon Bedrock is needed.

1 Note

The solution only requires com.amazonaws.region.bedrock-runtime.

8. Optional: If the deployment will use Amazon SageMaker AI for the LLM, then an <u>interface</u> endpoint for Amazon SageMaker AI is needed.

🚯 Note

The solution will not delete or modify the VPC configuration when using the **Bring your own VPC deployment** option. However, it will delete any VPCs that are created by the solution in the **Create a VPC for me** option. For this reason, you must be careful when sharing a solution-managed VPC across stacks/deployments. For example, deployment A uses **Create a VPC for me** option. Deployment B uses **Bring**

my own VPC using the VPC created by deployment A. If deployment A is deleted before deployment B, then deployment B will no longer work because the VPC has been deleted. Also because deployment B is using the ENIs created by the Lambda functions, deleting deployment A might have errors and retention of residual resources.

Amazon CloudFront

This solution deploys a web console <u>hosted</u> in an Amazon S3 bucket. To help reduce latency and improve security, this solution includes a CloudFront distribution with an origin access identity, which is a CloudFront user that provides public access to the solution's website bucket contents. For more information, see <u>Restricting Access to Amazon S3 Content by Using an Origin Access</u> Identity in the *Amazon CloudFront Developer Guide*.

🚯 Note

CloudFront has an account-level soft quota limit of 20 response header policies. This solution creates custom response header policies for security purposes. If you have more than 20 deployments of the Generative AI Application Builder on AWS or its use cases, new deployments may fail due to hitting the quota limit.

To resolve this issue, you can request a quota increase for the **Response Header Policies** quota in the AWS Service Quotas console by following these steps:

- 1. Open the AWS Service Quotas console.
- 2. In the navigation pane, select AWS services.
- 3. Search for and select Amazon CloudFront.
- 4. Scroll to the **Response Header Policies** quota and choose **Request quota increase**.
- 5. Follow the prompts to request an increase in the quota limit for your AWS account.

By increasing the **Response Header Policies** quota, you can ensure that new deployments of the Generative AI Application Builder on AWS or its use cases do not fail due to the quota limit.

Quotas

Service quotas, also referred to as limits, are the maximum number of service resources or operations for your AWS account.

Quotas for AWS services in this solution

Make sure you have sufficient quota for each of the <u>services implemented in this solution</u>. For more information, refer to <u>AWS service quotas</u>.

Use the following links to go to the page for that service. To view the service quotas for all AWS services in the documentation without switching pages, view the information in the <u>Service</u> <u>endpoints and quotas</u> page in the PDF instead.

Deploy the solution

This solution uses <u>AWS CloudFormation templates and stacks</u> to automate its deployment. The CloudFormation template specifies the AWS resources included in this solution and their properties. The CloudFormation stack provisions the resources that are described in the template.

Deployment process overview

Before you launch the solution, review the <u>cost</u>, <u>architecture</u>, <u>security</u>, and other considerations discussed in this guide.

🔥 Important

If you plan to use Amazon Bedrock, you must request access to models before they are available for use. Refer to <u>Model access</u> in the *Amazon Bedrock User Guide* for more details.

Time to deploy: Approximately 10 minutes

- Step 1: Launch the Deployment dashboard stack
- Step 2: Deploy a use case
- Step 3: Deploy a use case using the Deployment dashboard wizard
- Step 4: Post-deployment configuration

Optionally, you can deploy the use cases separately from the solution, if you prefer not to have the Deployment dashboard UI or APIs.

- Deploying a standalone Text use case
- Deploying a standalone Agent use case

You can also supply a DynamoDB chat configuration.

🔥 Important

This solution includes an option to send anonymized operational metrics to AWS. We use this data to better understand how customers use this solution and related services and

products. AWS owns the data gathered though this survey. Data collection is subject to the AWS Privacy Policy.

To opt out of this feature, download the template, modify the AWS CloudFormation mapping section, and then use the AWS CloudFormation console to upload your updated template and deploy the solution. For more information, see the <u>Anonymized data</u> <u>collection</u> section of this guide.

AWS CloudFormation template

You can download the CloudFormation template for this solution before deploying it.

View template

generative-ai-application-builder-on-aws.template - Use this template to launch the solution and all associated components. The default configuration deploys the core and supporting solutions found in the <u>AWS services in this solution</u> section, but you can customize the template to meet your specific needs.

🚯 Note

AWS CloudFormation resources are created from AWS Cloud Development Kit (AWS CDK) constructs.

This AWS CloudFormation template deploys Generative AI Application Builder on AWS in the AWS Cloud.

Step 1: Launch the Deployment dashboard stack

Follow the step-by-step instructions in this section to configure and deploy the solution into your account.

Time to deploy: Approximately 10 minutes

 Sign in to the <u>AWS Management Console</u> and select the button to launch the generativeai-application-builder-on-aws.template CloudFormation template.

Launch solution

2. The template launches in the US East (N. Virginia) Region by default. To launch the solution in a different AWS Region, use the Region selector in the console navigation bar.

i Note

This solution uses Amazon Kendra and Amazon Bedrock, which are not currently available in all AWS Regions. If using these features, you must launch this solution in an AWS Region where these services are available. For the most current availability by Region, see the AWS Regional Services List.

- 3. On the **Create stack** page, verify that the correct template URL is in the **Amazon S3 URL** text box and choose **Next**.
- 4. On the **Specify stack details** page, assign a name to your solution stack. For information about naming character limitations, see <u>IAM and STS Limits</u> in the AWS Identity and Access Management User Guide.
- 5. Under **Parameters**, review the parameters for this solution template and modify them as necessary. This solution uses the following default values.

| Parameter | Default | Description |
|------------------|--------------------|--|
| Admin User Email | <_Requires input_> | The email address of the admin user who will have access to the Deploymen t dashboard. An Amazon Cognito user will be created with permissions to deploy and manage use cases. |
| VpcEnabled | No | Should the Deployment dashboard be deployed within a VPC |

| Parameter | Default | Description |
|--------------|------------------|--|
| CreateNewVpc | No | Only available, if VpcEnable d is Yes. If the value is Yes, the stack will create the VPC and deploy the solution within the created VPC. If VpcEnabled is Yes and CreateNewVpc is No, then you must provide an existing VPC configuration (ExistingV pcld , ExistingPrivateSub netIds , ExistingSecurityGr oupIds , VpcAzs). |
| IPAMPoolId | (Optional input) | You can configure IPAM and provide the created id as input to assign the IP address range that the deployment of this stack should use. For details regarding IPAM, see <u>How</u> <u>IPAM works.</u> |

| Parameter | Default | Description |
|--------------------------|------------------|---|
| DeployUI | Yes | You have the option to deploy the Deployment dashboard without the web user interface (and the AWS resources required for the web deployment). In which case, the solution will deploy all infrastructure including REST API endpoints. This option is useful to integrate your own web interface with the Deployment dashboard APIs. |
| ExistingVpcId | (Optional input) | Required only if you want to deploy the solution in an existing VPC that you have created. |
| ExistingPrivateSubnetIds | (Optional input) | Required only if you want to deploy the solution in an existing VPC that you have created. The Lambda functions will be deployed in this subnet. |
| ExistingSecurityGroupIds | (Optional input) | Required only if you want to deploy the solution in an existing VPC that you have created. Ensure that security groups have the permissio ns for an outbound TCP connection. |

| Parameter | Default | Description |
|-----------------------------------|------------------|---|
| VpcAzs | (Optional input) | Required only if you want to deploy the solution in an existing VPC that you have created. |
| CognitoDomainPrefix | (Optional input) | Required only if you want to deploy the solution in an existing Amazon Cognito user pool that you created. If you don't provide a value, the solution generates it. |
| ExistingCognitoUserPoolId | (Optional input) | Required only if you want to deploy the solution in an existing Amazon Cognito user pool that you created. |
| ExistingCognitoUse rPoolClient | (Optional input) | Required only if you want to deploy the solution in an existing Amazon Cognito user pool that you created. If you don't provide a value, the solution creates a user pool client. This parameter can only be provided if you provide an ExistingC ognitoUserPoolId value. |

6. Choose Next.

- 7. On the **Configure stack options** page, choose **Next.**
- 8. On the **Review and create** page, review and confirm the settings. Select the box acknowledging that the template will create AWS Identity and Access Management (IAM) resources.
- 9. Choose **Submit** to deploy the stack.

You can view the status of the stack in the AWS CloudFormation console in the **Status** column. You should receive a CREATE_COMPLETE status in approximately 10 minutes.

Step 2: Deploy a use case

<u> Important</u>

Once the stack has been successfully deployed, a sign-up email is sent to the configured admin user email. Using those credentials, the admin user can sign in to the Deployment dashboard to use the web application.

Note

The DevOps user with access to the AWS Management Console must provide the admin user with the CloudFront URL of the Deployment dashboard UI when the stack completes. The URL can be found in the **Outputs** tab of the CloudFormation stack.

- 1. Sign in to the Deployment dashboard as an admin user.
- 2. On the application landing page, choose **Deploy new use case**.

This launches the deployment wizard, which walks you through building the use case.

Depicts Deployment dashboard landing page - fresh deployment

| Generative AI Application | n Buik | ler on AWS | | | | | | | | Sign out |
|--------------------------------------|--------|------------|---------------------|-------------------|------------------|--------------------------------------|----------------------|-----------------------|--------------|----------|
| Generative Al Application Builder | < | | Deployments (0 | D) Info | | | C View details Edit | Clone Delete Deploy | new use case | |
| Deployments | | | Q Find deployments | | | | | | < 1 > © | |
| Deploy New Use Case | | | Deployment Stack ID | マ Use Case Name | ▼ Application St | atus 🗢 🕴 Date Created | 🔹 🔻 🛛 Model Provider | ▼ Application Acces | ss ⊽ | |
| | | | | | | No matches We can't find a match. | | | | |
| | | | | | | | | | | |
| | | | | | | | | | < 1 > | |

Note

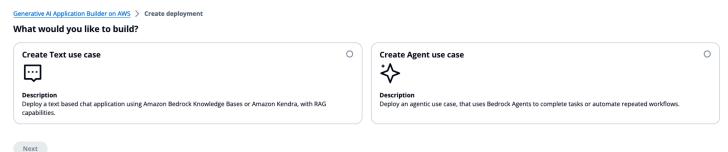
If you need to add additional users to your deployment, refer to the <u>Managing Cognito user</u> pool for more details.

Step 3: Deploy a use case using the Deployment dashboard wizard

In the Deployment dashboard wizard, you must choose between the following:

- Text use case Deploys a chat application, with optional RAG capabilities
- <u>Agent use case</u> Uses Amazon Bedrock Agents to complete tasks or automate repeated workflows

Shows two options: Create Text use case or Create Agent use case.



Step 3a: Deploy a Text use case

This section provides instructions for deploying a Text use case.

Select use case

When you choose **Create Text use case**, the UI opens the **Select use case** screen. Provide the following information:

- Use case name.
- Optional email address for the default user of the use case to be added to the Amazon Cognito user pool for the use case, and to be given permissions to interact with it.
- Whether you want to deploy a UI with this use case. If you don't want to deploy a UI with the use case, you can use the deployed API endpoints for use with your application.

By default, the Text use case creates and configures an Amazon Cognito user pool for you when the solution deploys the Deployment dashboard. The solution authenticates new use cases with a newly created client in the same user pool. However, you can provide an existing user pool ID and client ID in this step if you want to use your own Amazon Cognito user pool and client with the use case.

A Important

Admin users have access to all deployed use cases when the Amazon Cognito user pool is created via the deployment wizard. If you provide your own user pool during the deployment, you must ensure that the admin has the permissions to access the deployed use cases.

You will also need to update the Allowed callback URLs and Allowed sign-out URLs in your App clients in Cognito. To do this:

- 1. Navigate to the Cognito console
- 2. Choose User Pools.
- 3. Choose your user pool.
- 4. Choose App Clients on the left menu.
- 5. Choose the app client you want to modify.
- 6. Choose the Login pages tab.
- 7. Choose Edit and add your URLs.
- 8. Choose Save changes.

Additionally, if you need to add more users to a use case, refer to the <u>Managing Cognito</u> user pool section.

Select network configuration

This wizard step allows you to deploy the use case with a pre-existing or new <u>Amazon Virtual</u> <u>Private Cloud</u> (Amazon VPC). If selecting pre-existing VPC, you are required to provide a VPC ID, up to 16 subnet Ids and up to 5 security group IDs to use with this VPC. If you're not using a preexisting VPC, these settings will be configured for you.

Select model

In the **Select model** step, you can choose your model provider from the dropdown menu. There are two options: **Bedrock** and **SageMaker**.

If you select **SageMaker**, you can create a SageMaker AI model endpoint in the SageMaker AI console and provide the input schema that the model expects and output JSONPath for the LLM response. You can refer to the <u>Using Amazon SageMaker AI as an LLM Provider</u> section and <u>SageMaker AI payload examples</u> provided in the solution's GitHub repository.

If you select Amazon Bedrock, you will be presented with four options:

- **Quick Start Models** Get started quickly with a collection of models with different price/ performance characteristics. Recommended for building your first apps. This option allows you to select a model name from the provided list.
- Other Foundation Models Access the full range of foundation models with different capabilities and specializations. This option allows you to enter the model ID for your desired Bedrock on-demand foundation model.
- Inference Profiles Inference profiles leverage Bedrock's cross-region inference to increase throughput and improve resiliency by routing your requests across multiple AWS Regions during peak utilization bursts. This option allows you to enter the ID of the inference profile you want to use.
- **Provisioned Models** Dedicated throughput capacity for production workloads requiring consistent performance. This option allows you to enter the ARN of the provisioned/custom model to use from Amazon Bedrock.

Model selection step also allows you to choose your advanced model settings. Refer to <u>Advanced</u> <u>LLM settings</u> for details on configuring Amazon Bedrock Guardrails, provisioned throughput for Amazon Bedrock, and additional model parameters.

Cross-region inference

Cross-region inference helps Amazon Bedrock users to seamlessly manage unplanned traffic bursts by using compute across different AWS Regions. To use cross-region inference, you need the *inference profile*. An inference profile is an abstraction over an on-demand pool of resources from a configured set of AWS Regions. It can route your inference request, originating from your source Region, to another Region configured in that pool. This allows traffic distribution across multiple AWS Regions. This helps enable higher throughput and enhanced resilience during periods of peak demands.

Inference profiles are named after the model and Regions that they support. You must call an inference profile from one of the Regions that it includes. For example, as shown in the following

table, the inference profile ID us.anthropic.claude-3-haiku-20240307-v1:0 allows distribution of traffic over us-east-1 and us-west-2 Regions of the model you choose. Certain models are only available with an inference profile in a particular Region.

| Inference profile | Inference profile ID | Regions included |
|-----------------------------|---|---|
| US Anthropic Claude 3 Haiku | us.anthropic.claud e-3-haiku-20240307- v1:0 | US East (N. Virginia) (us- east-1) US West (Oregon) (us- west-2) |

If you want to use an inference profile ID instead of a model ID, then you must identify the appropriate inference profile ID. See <u>Supported Regions and models for inference profiles</u> in the *Amazon Bedrock User Guide* for more information. In the <u>Amazon Bedrock console</u>, the cross-region inference option in the left navigation menu provides these inference profile IDs.

After you identify the inference profile ID to use, you can use this during the **Select model** stage by performing the following steps:

- 1. Select Amazon Bedrock as the model provider.
- 2. Select the Inference Profiles radio button option.
- 3. Enter your inference profile ID in the text box that appears.

Refer to <u>Improve resilience with cross-region inference</u> in the *Amazon Bedrock User Guide* for more details on inference profiles.

Select knowledge base

If you're looking to deploy a non-Retrieval Augmented Generation (RAG) use case, you can skip this step.

However, if you wish to enable RAG as a part of your deployment, you can now provide either a pre-configured *Amazon Kendra Index Id* or an *Amazon Bedrock Knowledge Base ID*. You can also create a new Amazon Kendra Index for use with the solution. The solution currently supports Amazon Kendra and Amazon Bedrock Knowledge Bases as knowledge bases for your RAG-based use case deployment.

Refer to the <u>Configuring a Knowledge Base</u> section for guidelines on ingesting data into the knowledge base for use with your RAG-based deployment.

Advanced RAG configurations

The wizard allows you to select advanced options for use with your RAG deployment such as the **number of documents to retrieve** each time a query is sent to your knowledge base, a **static text response** from the LLM when no documents are found in the knowledge base, whether you wish to **display document sources** with your LLM response for sanity checks, etc. You can additionally also configure knowledge base specific configurations for Amazon Kendra such as <u>Role-based Access</u> <u>Control (RBAC)</u>, or <u>Override Search Type</u> when using Amazon OpenSearch Serverless with Amazon Bedrock Knowledge Bases. Refer to the <u>Advanced Knowledge Base settings</u> section for more details on these advanced settings.

Note

Your knowledge base must be in the same account and Region as the deployed Deployment dashboard and use case stacks.

Select prompts and token limits

In this step, you can configure your prompt for use with the LLM. Prompts may require placeholders such as {input}, {history} and {context}. These placeholders instruct the LLM on where to draw user input, conversation history, and information retrieved from the knowledge base from.

- For Bedrock model provider, the system prompt must be provided which has no restrictions for a non-RAG use-case. The disambiguation prompt for Bedrock model provider however, requires a minimum of two placeholders - {input} and {history}
- For SageMaker model provider, system and disambiguation prompts, both require a minimum of two placeholders - {input} and {history}.
- For RAG use cases, for each model provider, the {context} placeholder is additionally required.

For more information, see <u>Configuring your prompts</u>. You can also refer to the <u>Tips for managing</u> <u>model token limits</u> section while selecting token limit sizes for your prompts.

Review and deploy

After this step, review the settings you selected and choose **Deploy Use Case**. The new use case then deploys and becomes visible in your Deployment dashboard view to manage further.

Step 3b: Deploy an Agent use case

The Agent use case provides a powerful and secure mechanism for invoking Amazon Bedrock Agents within your use cases. This feature allows developers to seamlessly integrate the capabilities of AI-powered autonomous agents that can orchestrate and execute multi-step tasks across various foundation models, data sources, software applications, and user conversations while maintaining robust security measures.

Prerequisites

Before creating an Amazon Bedrock agent, ensure that you have the following:

- 1. The AWS account where Generative AI Application Builder on AWS is deployed, with an access to the Amazon Bedrock console.
- 2. Appropriate IAM permissions to create and manage Amazon Bedrock Agents.

Creating an Amazon Bedrock Agent

Refer to the <u>Create and configure agent manually</u> in the *Amazon Bedrock User Guide* for detailed instructions on creating an agent. You can configure options such as:

- Instructions (prompts) for your agent
- Knowledge base, which is used to look up additional information based on user's input
- Agent's memory to allow agents to remember information across multiple sessions (for a maximum of 30 days)

After you successfully create an Amazon Bedrock agent, you can proceed to the Generative AI Application Builder on AWS Agent use case wizard flow. To do so, choose **Deploy a new use case** on the Deployment dashboard and select **Create Agent Use Case**. Follow the wizard and use the following steps to configure the use case.

Select use case

This step is the same as the Text use case described previously.

Select network configuration

This step is the same as the Text use case described previously

Select agent

In this step, you must provide the **Agent ID** and **Alias ID** of the Amazon Bedrock agent that you created.

Step 4: Post-deployment configuration

This section provides recommendations for configuring the solution after deployment.

Amazon S3 bucket versioning, lifecycle policies, and cross-Region replication

This solution doesn't enforce lifecycle configurations on the buckets it creates. We recommend the following:

- Setting lifecycle configurations for production deployments. For details, see <u>Setting lifecycle</u> <u>configuration on a bucket</u> in the *Amazon Simple Storage Service User Guide*.
- Enabling <u>versioning</u> and <u>cross-Region replication</u> for Amazon S3 buckets based on the use case for which the solution is deployed.

Amazon DynamoDB backups

This solution uses DynamoDB for several purposes (see <u>AWS services in this solution</u>). The solution doesn't enable backups for the tables it creates. We recommend creating a backup of this feature for production deployments. See <u>Backing up a DynamoDB table</u> and <u>Using AWS Backup for</u> <u>DynamoDB</u> for details.

Amazon CloudWatch dashboard and alarms

The solution deploys a custom dashboard in CloudWatch to render charts from custom published metrics and AWS service metrics. We recommend creating CloudWatch <u>alarms</u> and adding notifications based on the use case for which the solution is deployed.

Amazon CloudWatch Logs

Lambda logs are configured to never expire and API Gateway logs are configured with a 10-year expiry. You can update the expiry of the respective log groups to align with your enterprise's record retention policy.

Custom web domains with TLS v1.2 or higher certificates

The solution deploys a web UI and Edge Optimized API Gateway using CloudFront. CloudFront's domain doesn't enforce TLS v1.2 or higher certificates. We recommend creating a custom domain using <u>Amazon Route 53</u>, creating a certificate using <u>AWS Certificate Manager</u>, or using an existing certificate if your organization has one.

For additional details, refer to the <u>Amazon Route 53 Developer Guide</u> and <u>Choosing a minimum</u> TLS version for a custom domain in API Gateway.

Scaling with Amazon Kendra

This solution provides the ability to use Amazon Kendra to perform NLP-powered intelligent search across the ingested documents. You can increase the capacity of Amazon Kendra using the following CloudFormation parameters for larger workloads:

| Parameter | Default | Description |
|--|---------|---|
| <u>Amazon Kendra additional</u> <u>query capacity</u> | 0 | The amount of extra query capacity for an index and <u>GetQuerySuggestions</u> capacity. An additional capacity unit for an index provides approximately 8,000 queries per day. |
| <u>Amazon Kendra additional</u> <u>storage capacity</u> | 0 | The amount of extra storage capacity for an index. A single capacity unit provides 30 GB of storage space or 100,000 documents, whichever reaches first. |

| Parameter | Default | Description |
|------------------------------|-----------|--|
| <u>Amazon Kendra edition</u> | Developer | Amazon Kendra provides Developer and Enterprise Editions to create indexes. For more information about the differences between Amazon Kendra Editions, see <u>Amazon</u> <u>Kendra pricing</u> . |

To modify the values of these CloudFormation parameters, select the appropriate values at the time of stack deployment. For more information on query and storage capacity units, see <u>Adjusting</u> <u>capacity</u>.

🚯 Note

If the Text use case is not deployed with RAG enabled, then an Amazon Kendra index is not used or created.

Setting up SSO using Idp federation

This solution allows integration with external identity providers that support SAML or OIDC based identity federation. When the solution deploys, it creates an Amazon Cognito user pool and individual app client integration for the Deployment dashboard and individual use cases. Based on the external Idp, follow the steps provided in the <u>Configuring identity providers for your user</u> <u>pool</u> section of the *Amazon Cognito Developer Guide* and choose the app client integration for the Deployment dashboard or use case you would like to setup SSO with.

To pass the user group information to knowledge base or vector stores in a RAG based architecture, you will need to map user groups from the external Idp to Amazon Cognito user groups. The solution provides an initial scaffolding Lambda function trigger to be mapped with the pre token generation phase. The Lambda function has the group_mapping.json file which must be updated to provide the group mappings. Refer to <u>Customizing user pool workflows with Lambda triggers</u> for Lambda triggers supported by Amazon Cognito.

Customizing login screen

This solution uses <u>Amazon Cognito hosted UI</u> to render the login page. To customize the built-in sign-in page, refer to <u>Customizing the built-in sign-in and sign-up webpages</u> in the *Amazon Cognito Developer Guide*.

Additional security considerations

Based on the use case for which you deploy the solution, review the following security recommendations:

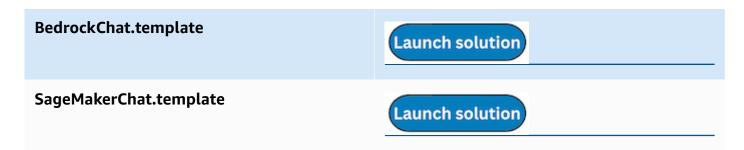
- Customer managed AWS KMS encryption keys The solution uses AWS managed AWS KMS keys by default, since these are available at no additional cost. Review your use case to determine if you should update the solution to use <u>customer managed AWS KMS keys</u>.
- API Gateway throttling rules The solution deploys with default throttling rules on API Gateway. Based on your use case and expected transaction volumes, we recommend that you configure throttling for the APIs. For details, see <u>Throttle API requests for better throughput</u> in the Amazon API Gateway Developer Guide.
- Enabling AWS CloudTrail As a recommended security practice, consider enabling <u>AWS</u> <u>CloudTrail</u> in the AWS account where the solution is deployed to log API calls in the AWS account. For details, see the AWS CloudTrail User Guide.
- **Drift detection** We recommend configuring drift detection on CloudFormation stacks to identify and be notified of unintentional or malicious changes to the deployed solution stack. For details, see <u>Implementing an alarm to automatically detect drift in AWS CloudFormation stacks</u>.
- Cognito JSON Web Tokens (JWTs) The solution uses Amazon Cognito-issued JWTs to authenticate with the REST API endpoints. We configured the solution with a five-minute expiry for <u>ID tokens</u> and <u>access tokens</u>. When a user logs out, their ability to generate new tokens is revoked (<u>refresh token</u> is revoked). However, until the expiry of the current token, any requests to the API endpoint will be successfully authenticated, since they have a valid token. Review the security considerations for your use case and adjust the token validity period.

Deploying a standalone Text use case

Follow the step-by-step instructions in this section to configure and deploy the solution into your account.

Time to deploy: Approximately 10-30 minutes

1. Sign in to the <u>AWS Management Console</u> and select the button to launch the CloudFront template that you want to deploy.



2. The template launches in the US East (N. Virginia) Region by default. To launch the solution in a different AWS Region, use the Region selector in the console navigation bar.

Note: This solution uses Amazon Kendra and Amazon Bedrock, which are not currently available in all AWS Regions. If using these features, you must launch this solution in an AWS Region where these services are available. For the most current availability by Region, see the <u>AWS</u> Regional Services List.

- 3. On the Create stack *page, verify that the correct template URL is in the *Amazon S3 URL *text box and choose *Next.
- 4. On the *Specify stack details *page, assign a name to your solution stack. For information about naming character limitations, see <u>IAM and STS Limits</u> in the AWS Identity and Access Management User Guide.
- 5. Under **Parameters**, review the parameters for this solution template and modify them as necessary. This solution uses the following default values.

| UseCaseUUID | <_Requires input_> | 36 character long UUIDv4 to identify this deployed use case within an application. |
|------------------------|--------------------|--|
| UseCaseConfigRecordKey | <_Requires input_> | Key corresponding of the record containing configura tions required by the chat provider Lambda at runtime. The record in the table must have a key attribute matching this value, and a config attribute containin |

| | | g the desired configuration. This record will be populated by the deployment platform if in use. For standalone deployments of this use case, a manually created entry in the table defined in UseCaseConfigTableName is required. |
|------------------------|--------------------|---|
| UseCaseConfigTableName | <_Requires input_> | The stack will read the configuration from the table with this name at the key UseCaseConfigRecordKey |

| ExistingRestApild | (Optional input) | Existing API Gateway REST API ID to use. If not provided, a new API Gateway REST API will be created. Typically provided when deploying from the Deployment dashboard. |
|-------------------|------------------|---|
| | | Note: Using Existing APIs can help reduce resource duplication and simplify management of APIs when you need to deploy multiple standalone use cases. When supplying existing APIs for a standalone use case, you are responsible for ensuring that the API is configured with the required route(s) with expected models. A required pre-configured / details route (fetches use case details during chat) and optionally, a /feedback route (if FeedbackEnabled is set to Yes to enable collection of feedback for LLM chat responses) must be configured. Additionally, ExistingApiRootResourceld , ExistingCognitoUserPoolld |
| | | and ExistingCognitoGro upPolicyTableName must also be provided. |

| ExistingApiRootResourceId | (Optional input) | Existing API Gateway REST API Root Resource ID to use. REST API Root Resource ID can be obtained from the AWS console by selecting the root resource (/) in the "Resources" section of the API. The Resource ID will then be displayed in the Resource details panel. You can alternatively run a describe API call on your REST API to find the Root Resource ID. |
|--------------------------------|-----------------------------|---|
| FeedbackEnabled | No | If set to No, the deployed use case stack will not have access to the feedback feature. |
| ExistingModelInfoT ableName | (Optional input) | DynamoDB table name for the table which contains model info and defaults. Used by the deploymen t platform. If omitted, a new table will be created to house model defaults. |
| DefaultUserEmail | placeholder@exampl e.com | Email of the default user for this use case. An Amazon Cognito user for this email is created to access the use case. |

| ExistingCognitoUserPoolId | (Optional input) | UserPoolId of an existing Amazon Cognito user pool which this use case will be authenticated with. Typically provided when deploying from the Deployment dashboard, but can be omitted when deploying this use case stack standalone. |
|---|------------------|--|
| CognitoDomainPrefix | (Optional input) | Enter a value if you want to provide a domain for the Cognito User Pool Client. If you don't provide a value, the deployment will generate one. |
| ExistingCognitoUse rPoolClient | (Optional input) | Provide a User Pool Client (App Client) to use an existing one. If you don't provide a User Pool Client, a new one will be created. This parameter can only be provided if an existing User Pool Id is provided. |
| ExistingCognitoGro upPolicyTableName | (Optional input) | Name of the DynamoDB table containing user group policies. This is used by the custom authorizer on the use case's API. Typically, you can provide an input when deploying from the deployment platform, but can be omitted when deploying this use case stack standalone. |

| RAGEnabled | true | If set to true, the deployed use case stack uses the provided Amazon Kendra index created to provide RAG functionality. If set to false, the user interacts directly with the LLM. |
|-----------------------|------------------|--|
| KnowledgeBaseType | Bedrock | Knowledge base type to be used for RAG. Only set if RAGEnabled is true. Can be Bedrock or Kendra. Note: Only relevant if RAGEnabled is true. |
| ExistingKendraIndexId | (Optional input) | Index ID of an existing Kendra index to be used for the use case. If none is provided and Knowledge BaseType is Kendra, a new index will be created for you. Note: Only relevant if RAGEnabled is true and KnowledgeBaseType is Kendra. |
| NewKendraIndexName | (Optional input) | Name for the new Kendra index to be created for this use case. Only applies if ExistingKendraIndexId is not supplied. Note: Only relevant if RAGEnabled is true and KnowledgeBaseType is Kendra. |

| NewKendraQueryCapa cityUnits | 0 | Additional query capacity units for the new Amazon Kendra index to be created for this use case. Only applies if ExistingK endraIndexId is not supplied, see <u>CapacityU</u> <u>nitsConfiguration</u> . Note: Only relevant if RAGEnabled is true and KnowledgeBaseType is Kendra. |
|-----------------------------------|---|--|
| NewKendraStorageCa pacityUnits | 0 | Additional storage capacity units for the new Amazon Kendra index to be created for this use case. Only applies if ExistingK endraIndexId is not supplied, see <u>CapacityU</u> <u>nitsConfiguration</u> . Note: Only relevant if RAGEnabled is true and KnowledgeBaseType is Kendra. |

| NewKendraIndexEdition | (Optional input) | The edition of Amazon Kendra to use for the new Amazon Kendra index to be created for this use case. Only applies if ExistingKendraIndexId is not supplied, see <u>Amazon</u> <u>Kendra Editions</u> . Note: Only relevant if RAGEnabled is true and KnowledgeBaseType is Kendra. |
|------------------------|------------------|---|
| BedrockKnowledgeBaseId | (Optional input) | ID of the bedrock knowledge base to use in a RAG use case. Cannot be provided if ExistingKendraIndexId or NewKendraIndexName are provided. Note: Only relevant if RAGEnabled is true and KnowledgeBaseType is Bedrock. |
| VpcEnabled | No | Should the stacks resources be deployed within a VPC. |
| CreateNewVpc | No | Select Yes, if you want the solution to create a new VPC for you and be used for this use case. Note: Only relevant if VpcEnabled is Yes. |

| IPAMPoolId | (Optional input) | If you want to assign the CIDR range using Amazon VPC IP Address Manager, provide the IPAM pool Id to use. Note: Only relevant if VpcEnabled is Yes and CreateNewVpc is No. |
|--------------------------|------------------|--|
| ExistingVpcId | (Optional input) | VPC ID of an existing VPC to be used for the use case. Note: Only relevant if VpcEnabled is Yes and CreateNewVpc is No. |
| ExistingPrivateSubnetIds | (Optional input) | Comma separated list of subnet IDs of existing private subnets to be used to deploy the Lambda function. Note: Only relevant if VpcEnabled is Yes and CreateNewVpc is No. |
| ExistingSecurityGroupIds | (Optional input) | Comma separated list of security groups of the existing VPC to be used for configuring Lambda functions. Note: Only relevant if VpcEnabled is Yes and CreateNewVpc is No. |

| VpcAzs | (Optional input) | Comma separated list of AZs in which subnets of the VPCs are created Note: Only relevant if VpcEnabled is Yes and CreateNewVpc is No. |
|---------------------|------------------|--|
| UseInferenceProfile | No | If the model configured is Bedrock, you can indicate if you are using Bedrock Inference Profile. This will ensure that the required IAM policies will be configured during stack deployment. For more details, refer to the following <u>https://docs.aws.</u> <u>amazon.com/bedrock/latest</u> <u>/userguide/cross-region-i</u> <u>nference.html</u> |
| DeployUI | Yes | Select the option to deploy the frontend UI for this deployment. Selecting No, will only create the infrastru cture to host the APIs, the authentication for the APIs, and backend processing. |

- 6. Choose Next.
- 7. On the **Configure stack options** page, choose **Next**.
- 8. On the **Review** page, review and confirm the settings. Select the box acknowledging that the template will create AWS Identity and Access Management (IAM) resources.
- 9. Choose **Create stack** to deploy the stack.

You can view the status of the stack in the AWS CloudFormation console in the **Status** column. You should receive a CREATE_COMPLETE status in approximately 10-30 minutes.

Deploying a standalone Agent use case

Follow the step-by-step instructions in this section to configure and deploy the solution into your account.

Time to deploy: Approximately 10-30 minutes

1. Sign in to the <u>AWS Management Console</u> and select the button to launch the CloudFront template.

BedrockAgent.template

Launch solution

2. The template launches in the US East (N. Virginia) Region by default. To launch the solution in a different AWS Region, use the Region selector in the console navigation bar.

Note

This solution uses Amazon Bedrock, which is not currently available in all AWS Regions. If you're using these features, you must launch this solution in an AWS Region where these services are available. For the most current availability by Region, see the <u>AWS Regional</u> <u>Services List</u>.

- 3. On the **Create stack** page, verify that the correct template URL is in the **Amazon S3 URL** text box and choose **Next**.
- 4. On the **Specify stack details** page, assign a name to your solution stack. For information about naming character limitations, see {https---docs-aws-amazon-com-https---docs-aws-amazon-com-IAM-latest-UserGuide-reference-iam-limits-html}[IAM and AWS STS quotas] in the AWS *Identity and Access Management User Guide*.
- 5. Under **Parameters**, review the parameters for this solution template and modify them as necessary. This solution uses the following default values.

| Parameter | Default entry | Description |
|-------------|--------------------|--|
| UseCaseUUID | <_Requires input_> | 36 character long UUIDv4 to identify this deployed use case within an application. |

| Parameter | Default entry | Description |
|------------------------|---------------------------------|--|
| UseCaseConfigRecordKey | <requires input=""></requires> | Key corresponding to the record that contains configurations required by the chat provider Lambda function at runtime. The record in the table must have a key attribute matching this value, and a config attribute containing the desired configuration. This record will be populated by the deployment platform if it's in use. For standalon e deployments of this use case, a manually created entry in the table defined in UseCaseConfigTableName is required. |
| UseCaseConfigTableName | <requires input="">`</requires> | The stack will read the use case configuration from the table provided here and using the record key defined in UseCaseConfigRecordKey . |
| DefaultUserEmail | placeholder@exampl e.com | Email of the default user for this use case. The solution creates an Amazon Cognito user for this email to access the use case. |

| Parameter | Default entry | Description |
|-------------------|------------------|--|
| ExistingRestApild | (Optional input) | Existing API Gateway REST API ID to use. If not provided, a new API Gateway REST API will be created. Typically provided when deploying from the Deployment dashboard. Note: Using Existing APIs can help reduce resource duplication and simplify management of APIs when you need to deploy multiple standalone use cases. When supplying existing APIs for a standalone use case, you are responsible for ensuring that the API is configured with the required route(s) with expected models. A required pre-configured // details route (fetches use case details during chat) and optionally, a /feedback route (if FeedbackEnabled is set to Yes to enable collection of feedback for LLM chat responses) must be configured. Additionally, ExistingApiRootResourceId, ExistingCognitoUserPoolId and ExistingCognitoGro upPolicyTableName must also be provided. |

| Parameter | Default entry | Description |
|---------------------------|------------------|--|
| ExistingApiRootResourceId | (Optional input) | Existing API Gateway REST API Root Resource ID to use. REST API Root Resource ID can be obtained from the AWS console by selecting the root resource (/) in the "Resources" section of the API.The Resource ID will then be displayed in the Resource details panel. You can alternatively run a describe API call on your REST API to find the Root Resource ID. |
| FeedbackEnabled | No | If set to No, the deployed use case stack will not have access to the feedback feature. |
| CognitoDomainPrefix | (Optional input) | Enter a value if you want to provide a domain for the Amazon Cognito user pool client. If you don't provide a value, the solution generates one. |

| Parameter | Default entry | Description |
|---|------------------|--|
| ExistingCognitoUserPoolId | (Optional input) | UserPoolId of an existing Amazon Cognito user pool that you want to authentic ate this use case with. NOTE: You typically provide this ID when deploying from the Deployment dashboard , but you can omit it when deploying this use case stack standalone. |
| ExistingCognitoUse rPoolClient | (Optional input) | Provide a user pool client (app client) to use an existing one. If you don't provide a user pool client, the solution creates one. You can only provide this parameter if you provided an ExistingC ognitoUserPoolId. |
| ExistingCognitoGro upPolicyTableName | (Optional input) | Name of the DynamoDB table containing user group policies. This is used by the custom authorizer on the use case's API. NOTE: You typically provide this name when deploying from the Deployment dashboard, but you can omit it when deploying this use case stack standalone. |
| VpcEnabled | No | Whether the stacks resources be deployed within a VPC. |

| Parameter | Default entry | Description |
|--------------------------|------------------|--|
| CreateNewVpc | No | Select Yes if you want the solution to create a new VPC for you and to use it for this use case. NOTE: This parameter is only relevant if VpcEnabled is Yes. |
| IPAMPoolId | (Optional input) | If you want to assign the CIDR range using IPAM, provide the IPAM pool ID to use. NOTE: This parameter is only relevant if VpcEnabled is Yes and CreateNewVpc is No. |
| ExistingVpcId | (Optional input) | VPC ID of an existing VPC to be used for the use case. NOTE: This parameter is only relevant if VpcEnabled is Yes and CreateNewVpc is No. |
| ExistingPrivateSubnetIds | (Optional input) | Comma separated list of subnet IDs of existing private subnets to be used to deploy the Lambda function. NOTE: This parameter is only relevant if VpcEnabled is Yes and CreateNewVpc is No. |

| Parameter | Default entry | Description |
|--------------------------|--------------------------------|--|
| ExistingSecurityGroupIds | (Optional input) | Comma-separated list of security groups of the existing VPC to be used for configuring Lambda functions. NOTE: This parameter is only relevant if VpcEnabled is Yes and CreateNewVpc is No. |
| VpcAzs | (Optional input) | Comma separated list of AZs in which subnets of the VPCs are created Note: Only relevant if VpcEnabled is Yes and CreateNewVpc is No. |
| BedrockAgentId | <requires input=""></requires> | The ID of the Amazon Bedrock Agent to be used. |
| BedrockAgentAliasId | <requires input=""></requires> | The alias ID of the Amazon Bedrock Agent to be used. |
| DeployUI | Yes | Select the option to deploy the frontend chat UI for this deployment. Selecting No results in creating the infrastructure to host the APIs, the authentication for the APIs, and backend processing without the chat UI. |

6. Choose Next.

7. On the **Configure stack options** page, choose **Next**.

- 8. On the **Review** page, review and confirm the settings. Select the box acknowledging that the template will create IAM resources.
- 9. Choose **Create stack** to deploy the stack.

You can view the status of the stack in the AWS CloudFormation console in the **Status** column. You should receive a CREATE_COMPLETE status in approximately 10-30 minutes.

Supplying a DynamoDB chat configuration

When deploying a use case, **UseCaseConfigRecordKey** and **UseCaseConfigTableName** are required CloudFormation parameters which are normally populated by the Deployment dashboard. The deployment dashboards stack handles the creation and configuration of this table, while calls to the deployment API trigger population of the parameters.

When performing a standalone deployment, you must do the following:

- 1. Create a DynamoDB table with a hash key of **key**.
- 2. Create a record in the table containing the configuration for the use case as a record of the format: {key: some_use_case_key, config: {your_configuration}.
- 3. Pass the chosen UseCaseConfigTableName and UseCaseConfigRecordKey (some_use_case_key in this example) parameters to the use case stack when deploying.

To create a suitable configuration for a standalone deployment, you can create a required use case from the Deployment dashboard, and copy record from the configuration table. Otherwise, you can craft your own configuration based on the following example for a Bedrock deployment:

```
{
   "UseCaseName": "SampleUseCase",
   "ConversationMemoryParams": {
   "ConversationMemoryType": "DynamoDB",
   "HumanPrefix": "H",
   "AiPrefix": "A",
   "ChatHistoryLength": 20
   },
   "KnowledgeBaseParams": {
   "KnowledgeBaseType": "Bedrock",
   "NumberOfDocs": 2,
   "ScoreThreshold": 0,
}
```

```
"ReturnSourceDocs": false,
 "BedrockKnowledgeBaseParams": {
 "BedrockKnowledgeBaseId": "SOME_ID",
 "OverrideSearchType": null
}
 },
 "LlmParams": {
 "ModelProvider": "Bedrock",
 "BedrockLlmParams": { "ModelId": "anthropic.claude-v2" },
 "PromptParams": {
 "PromptTemplate": "some prompt",
 "MaxPromptTemplateLength": 187500,
 "MaxInputTextLength": 187500,
 "UserPromptEditingEnabled": true,
 "DisambiguationEnabled": true,
 "DisambiguationPromptTemplate": "some prompt"
 },
 "ModelParams": {},
 "Temperature": 1,
 "RAGEnabled": true,
 "Streaming": true,
"Verbose": false
}
}
```

Monitor the solution with Service Catalog AppRegistry

The solution includes a Service Catalog AppRegistry resource to register the CloudFormation template and underlying resources as an application in both Service Catalog AppRegistry and Systems Manager Application Manager.

Systems Manager Application Manager gives you an application-level view into this solution and its resources so that you can:

- Monitor its resources, costs for the deployed resources across stacks and AWS accounts, and logs associated with this solution from a central location.
- View operations data for the resources of this solution in the context of an application. For example, deployment status, CloudWatch alarms, resource configurations, and operational issues.

The following figure depicts an example of the application view for the solution stack in Application Manager.

Depicts solution stack in Application Manager

| Components (2) | AWS-Systems-Manager-Application-Manager | |
|---|---|--|
| Name Alarms | Application information Vlew in AppRegistry 🖾 | |
| AWS-Systems-Manager-Application-Manager AWS-Systems-Manager-A | Application type Name Application monitoring AWS-AppRegistry AWS-Systems-Manager-Application-Manager Application monitoring | |
| | Description Service Catalog application to track and manage all your resources for the solution | |
| | Overview Resources Instances Compliance Monitoring Opsitems Logs Runbooks Cost | |
| | Insights and Alarms Info View all Cost View all Monitor your application health with Amazon CloudWatch. View resource costs per application using AWS Cost Explorer. View all | |
| | Cost (USD) | |

Activate CloudWatch Application Insights

1. Sign in to the Systems Manager console.

- 2. In the navigation pane, choose **Application Manager**.
- 3. In **Applications**, search for the application name for this solution and select it.

The application name will have **App Registry** in the **Application Source** column, and will have a combination of the solution name, Region, account ID, or stack name.

- 4. In the **Components** tree, choose the application stack you want to activate.
- 5. In the Monitoring tab, in Application Insights, select Auto-configure Application Insights.

Application Insights dashboard showing no detected problems and option to auto-configure.

| Overview Resources Provisioning Compliance | e Monitoring Opsitems Logs Runbooks Cost |
|--|--|
| Application Insights (0) Info Problems detected by severity | View Ignored Problems Actions Add an application |
| Q Find problems | Last 7 days 🔻 🖸 < 1 > 💿 |
| Problem su Vertical Status Vertical Severity | v ⊽ Source ⊽ Start time ⊽ Insights |
| Advance | d monitoring is not enabled |
| | (SLR) is created in your account. The SLR is predefined by CloudWatch Application ne service requires to monitor AWS services on your behalf. |
| | figure Application Insights |

Monitoring for your applications is now activated and the following status box appears:

Application Insights dashboard showing successful monitoring activation message.

| opplication Insights (0) Info | View Ig | gnored Problems Actions V | Add an application |
|--------------------------------|---------|---------------------------|--------------------|
| Q Find problems | | Last 7 days 🔻 | C < 1 > @ |
| Problem su V Status | | ▽ Start time | |

Confirm cost tags associated with the solution

After you activate cost allocation tags associated with the solution, you must confirm the cost allocation tags to see the costs for this solution. To confirm cost allocation tags:

- 1. Sign in to the Systems Manager console.
- 2. In the navigation pane, choose **Application Manager**.
- 3. In **Applications**, choose the application name for this solution and select it.

The application name will have **App Registry** in the **Application Source** column, and will have a combination of the solution name, Region, account ID, or stack name.

4. In the Overview tab, in Cost, select Add user tag.

Screenshot depicting the Application Cost add user tag screen

| Cost View resource costs per application using AWS Cost Explorer. | View all |
|---|---------------|
| To enable cost tracking, add the "AppManagerCFNStackKey" user tag to your Clstack. Adding the user tag will require redeployment of the stack. Add user tag | loudFormation |

5. On the Add user tag page, enter confirm, then select Add user tag.

The activation process can take up to 24 hours to complete and the tag data to appear.

Activate cost allocation tags associated with the solution

After you activate Cost Explorer, you must activate the cost allocation tags associated with this solution to see the costs for this solution. The cost allocation tags can only be activated from the management account for the organization. To activate cost allocation tags:

- 1. Sign in to the AWS Billing and Cost Management and Cost Management console.
- 2. In the navigation pane, select **Cost Allocation Tags**.
- 3. On the **Cost allocation tags** page, filter for the AppManagerCFNStackKey tag, then select the tag from the results shown.
- 4. Choose Activate.

AWS Cost Explorer

You can see the overview of the costs associated with the application and application components within the Application Manager console through integration with AWS Cost Explorer, which must be first activated. Cost Explorer helps you manage costs by providing a view of your AWS resource costs and usage over time. To activate Cost Explorer for the solution:

- 1. Sign in to the <u>AWS Cost Management console</u>.
- 2. In the navigation pane, select **Cost Explorer** to view the solution's costs and usage over time.

Update the solution

If you have previously deployed the solution, follow this procedure to update the solution's CloudFormation stack to get the latest features and enhancements. There are three parts to the upgrade process:

- <u>Step 1: Update Deployment dashboard</u>
- Step 2: Migrate use case configurations
- Step 3: Update use cases

🚺 Note

- In v2.0.0, integration with Anthropic and Hugging Face was deprecated in favor of Amazon Bedrock and Amazon SageMaker AI. You can deploy models available through Hugging Face through SageMaker JumpStart. Refer to <u>Use Hugging Face with Amazon</u> <u>SageMaker AI</u> for more details.
- 2. Ensure you test the update process in a non-production environment before running these steps.

Step 1: Update Deployment dashboard

- Sign in to the <u>CloudFormation console</u>, select your existing CloudFormation stack, and select Update.
- 2. Select Replace current template.
- 3. Under Specify template:
 - a. Select Amazon S3 URL.
 - b. Copy the latest CloudFormation template link.
 - c. Paste the link in the Amazon S3 URL box.
 - d. Verify that the correct template URL shows in the **Amazon S3 URL** text box, and choose **Next**. Choose **Next** again.
- 4. Under **Parameters**, review the parameters for the template and modify them as necessary. For details about the parameters, see Step 1: Launch the Deployment dashboard stack.

5. Choose Next.

- 6. On the Configure stack options page, choose Next.
- 7. On the **Review** page, review and confirm the settings. Check the box acknowledging that the template will create IAM resources.
- 8. Choose View change set and verify the changes.
- 9. Choose Update stack to deploy the stack.

You can view the status of the stack in the AWS CloudFormation console in the **Status** column. You should receive an UPDATE_COMPLETE status in approximately 10 minutes.

If the existing Solution version was prior to v2.0.0, updating will create a web UI stack (which replaces the amplify-ui implementation of the login screen with a Cognito hosted UI) and a new CloudFront URL, which can be obtained from the Output section of the CloudFormation console once the stack status is UPDATE_COMPLETE.

Note

Existing use cases created using versions prior to v2.0.0 will NOT be displayed until you complete the steps outlined below.

Step 2: Migrate use case configurations (Only updates from versions below 2.0.0)

The schema for storing and the AWS service to store use case configuration has changed in version 2.0.0. Follow the steps described in <u>GAAB v2 Migration User Guide</u> using the <u>gaab_v2_migration.py</u> script. After you run the script, you can access the Deployment dashboard to view the deployed use cases.

Note

You must follow the steps below to complete migrating the use cases.

Step 3: Update use cases

You can edit the deployed use cases with new features available in the latest versions of GAAB. See Use the solution for information about how to use the features in this solution.

To update use cases to the latest version, you must complete the `Edit`use case steps in the Deployment dashboard (although you might not make any changes). This action triggers a CloudFormation stack update with the latest template version.

🚯 Note

Use cases created with 1.x or 2.x versions of the solution might not work with later versions. Hence, we recommend cloning existing use cases created with versions prior to v3.0.0 through the Deployment dashboard. Then, gradually migrate and replace with new use cases created using v3.0.0 or later.

Troubleshooting

This section provides troubleshooting instructions for deploying and using the solution.

If these instructions don't address your issue, <u>Contact Support</u> provides instructions for opening an Support case for this solution.

Problem: Deploying a VPC-enabled configuration, with Create a VPC for me, fails

The Deployment dashboard stack or the use case stack fails deployment because the CloudFormation was not able to provision VPC networking resources.

Resolution

Check the quota limits for VPCs, and Elastic IPs in your account. Default limits are 5 each for Elastic IPs and VPCs per AWS account, per AWS Region.

i Note

When the solution creates a VPC, a single VPC-enabled deployment (Deployment Dashboard or Use Case) is a 2-AZ deployment with 1 public and 1 private subnet in each AZ, each public subnet deploys 1 NAT Gateway. With 2 NAT Gateways, the deployment consumes 2 public IP addresses from the quota limit.

Some limits to be aware of (per account, per Region):

- Number of VPCs 5
- Number of public IP addresses 5
- Number of Gateway VPC Endpoints 20
- Number of Inteface VPC Endpoints 20

Problem: Use case stack can't be deleted in CloudFormation after the Deployment dashboard stack gets deleted

If the Deployment dashboard stack is deleted in CloudFormation before all of the use case stacks are deleted, the use cases can end up in a locked (unusable) state. This is due to an IAM role created by the Deployment dashboard stack no longer exists preventing modifications to the use case stack.

Resolution

🔥 Warning

Ensure you clean up any manually created roles immediately after usage. These are elevated permissions that users could exploit for role elevation.

Recreate the deleted IAM role to enable the deletion of the CloudFormation stacks:

- 1. Open the CloudFormation console and determine the role that is associated with your locked stack.
 - a. The role ARN can be found in the stack info section labeled IAM role.
 - b. The role name is what follows after :role/ in the IAM role ARN (for example, arn:aws:iam::<account-id>:role/<role-name>)
- 2. Create a new role in IAM with the same name as the deleted role.
 - a. Select **AWS service** as the trusted entity and select **CloudFormation** from the drop down.
 - b. Add the necessary permissions. If you're unsure about the required permissions, you can use the AWS managed **AdministratorAccess** policy.
 - c. Enter the role name exactly as obtained in Step 1.
- 3. Return to the CloudFormation console and delete the locked stacks.
- 4. Once all locked stacks have been successfully deleted, return to IAM and delete any roles created in Step 2.

Problem: Use case UI does not reflect changes in settings

When use cases are updated, the UI is deployed to CloudFront. However, because CloudFront caches deployments as well as the configuration file that dictates how some settings are shown to the user, these changes might not be reflected immediately.

Resolution

The CloudFront distribution can be invalidated to force the new configuration to be propagated to frontend users.

- 1. Open the CloudFormation console and determine the CloudFront distribution that is associated with your use case stack.
 - a. The use case stack should start with the same name you used when deploying the use case.
 - b. Locate the nested stack corresponding to the UI. The nested stack name should begin with **WebAppS3UINestedStackS3UINestedStackResource**.
 - c. Under the **Resources** tab, locate the resource of type **AWS::CloudFront::Distribution**, then select the physical ID. This will open the distribution in the CloudFront console.
- 2. Navigate to the **Invalidations** tab, then choose **Create Invalidation**, and input a path of /*. This will invalidate all paths.
- 3. In your own browser, delete any cookies and cached files related to the use case.

Contact Support

If you have <u>AWS Developer Support</u>, <u>AWS Business Support</u>, or <u>AWS Enterprise Support</u>, you can use the Support Center to get expert assistance with this solution. The following sections provide instructions.

Create case

- 1. Sign in to Support Center.
- 2. Choose Create case.

How can we help?

1. Choose Technical.

- 2. For Service, select Solutions.
- 3. For Category, select Other Solutions.
- 4. For **Severity**, select the option that best matches your use case.
- 5. When you enter the **Service**, **Category**, and **Severity**, the interface populates links to common troubleshooting questions. If you can't resolve your question with these links, choose **Next step: Additional information**.

Additional information

- 1. For **Subject**, enter text summarizing your question or issue.
- 2. For **Description**, describe the issue in detail.
- 3. Choose Attach files.
- 4. Attach the information that AWS Support needs to process the request.

Help us resolve your case faster

- 1. Enter the requested information.
- 2. Choose Next step: Solve now or contact us.

Solve now or contact us

- 1. Review the **Solve now** solutions.
- 2. If you can't resolve your issue with these solutions, choose **Contact us**, enter the requested information, and choose **Submit**.

Uninstall the solution

🚯 Note

Deployments created through the Deployment dashboard are not intended to be managed outside of the solution. Be sure to delete and clean up any deployments from within the Deployment dashboard, before deleting the stack in CloudFormation.

You can uninstall the Generative AI Application Builder on AWS solution from the AWS Management Console or by using the AWS Command Line Interface. You must manually delete the Amazon S3 buckets, Amazon Kendra indexes, or CloudWatch Logs created by this solution. AWS Solutions do not automatically delete Amazon S3 buckets, Amazon Kendra indexes, or CloudWatch Logs in case you have stored data to retain.

Using the AWS Management Console

- 1. Sign in to the <u>AWS CloudFormation console</u>.
- 2. On the **Stacks** page, select this solution's installation stack.
- 3. Choose Delete.

Using AWS Command Line Interface

Determine whether the AWS Command Line Interface (AWS CLI) is available in your environment. For installation instructions, see <u>What Is the AWS Command Line Interface</u> in the AWS CLI User Guide. After confirming that the AWS CLI is available, run the following command.

\$ aws cloudformation delete-stack --stack-name <installation-stack-name>

Manual uninstall steps

Deleting the Amazon S3 buckets

This solution is configured to retain the solution-created Amazon S3 bucket if you decide to delete the AWS CloudFormation stack to prevent accidental data loss. After uninstalling the solution, you

can manually delete this Amazon S3 bucket if you do not need to retain the data. Follow these steps to delete the Amazon S3 bucket.

- 1. Sign in to the Amazon S3 console.
- 2. In the navigation pane, select **Buckets**.
- 3. Locate the *<stack-name>* S3 buckets.
- 4. Select the S3 bucket and choose **Delete**.

To delete the S3 bucket using AWS CLI, run the following command. You won't need to empty the bucket first when using the --force option.

\$ aws s3 rb s3://<bucket-name> --force

Deleting the Amazon Kendra indexes

To prevent accidental data loss, this solution is configured to retain the solution-created Amazon Kendra indexes when the AWS CloudFormation stack has been deleted. After uninstalling the solution, you can manually delete the Amazon Kendra indexes that you no longer need to retain data for. Follow these steps to delete the Amazon Kendra index.

- 1. Sign in to the Amazon Kendra console.
- 2. In the navigation pane, select Indexes.
- 3. Locate and select the index you want to delete.
- 4. Choose **Delete** to delete the selected index.

To delete the Amazon Kendra index using AWS CLI, run the following command:

```
$ aws kendra delete-index --id<index-id>
```

Deleting the CloudWatch Logs

To prevent accidental data loss, we configured this solution to retain the CloudWatch Logs if you decide to delete the CloudFormation stack. After uninstalling the solution, you can manually delete the logs if you don't need to retain the data. Follow these steps to delete the CloudWatch Logs.

1. Sign in to the <u>Amazon CloudWatch console</u>.

- 2. In the navigation pane, select Log Groups.
- 3. Locate the log groups created by the solution.
- 4. Select one of the log groups.
- 5. Choose **Actions** and then choose **Delete**.

Repeat the steps until you have deleted all the solution log groups.

Use the solution

Accessing the UI

During the stack deployment process (for both the Deployment dashboard and use cases) an email is sent to the configured email address. The email contains the user's temporary credentials they can use to sign up and access the web interface.

🚯 Note

The DevOps user with access to the AWS Management Console must provide the admin user with the CloudFront URL of the Deployment dashboard UI when the stack completes.

For the use cases, the admin user with access to the Deployment dashboard UI must provide the business user with the CloudFront URL of the use case UI when the deployment completes.

Once logged in, the user can interact with the solution UIs, either the Deployment dashboard in the case of admins, or the use case in the case of business users.

How to update a deployment

When on the Deployment dashboard home page (or the details page of a deployment) you can edit the configuration used by a deployment. You can only edit deployments that are in the CREATE_COMPLETE or UPDATE_COMPLETE statuses.

Except for the use case name, all other options are editable for a deployment. Just change the values you want to edit and redeploy.

Depending on the scope of edits made, the redeployment time will vary. It might take a few seconds if simple settings have changed (example, model parameters), to more than 30 minutes if larger infrastructure related options have changed (example, request to create the Amazon Kendra index for the Text use case RAG).

Once the edit has completed successfully, the application status will report an UPDATE_COMPLETE status. At this time, you can access the deployed UI through the CloudFront URL and interact with the modified deployment.

🚯 Note

It might be easier to run multiple deployments side-by-side if you want to compare different settings or LLMs. Use the **Clone** feature to quickly use an existing configuration to launch a new deployment.

How to clone a deployment

When on the Deployments dashboard home page (or the details page of a deployment) you can clone the configuration used by a deployment. Cloning a deployment launches the **Deploy new use case** wizard, but with most fields pre-filled with the same values.

This is a convenience operation to help you quickly duplicate deployments with changed settings, revive a deleted deployment, or compare multiple LLMs in otherwise identical deployments.

How to delete a deployment

When on the Deployments dashboard home page (or the details page of a deployment) you can delete it once you no longer need the deployment. Deleting a deployment invokes a CloudFormation stack delete operation and deprovisions the resources for the deployment.

By default, a deleted deployment still remains on the dashboard to enable the clone functionality. To completely remove a deployment from the dashboard so that it stops being tracked in the UI, choose **Permanently delete** on the delete confirmation window.

🔥 Important

Some resources are left behind during stack deletion and must be manually deleted. Refer to the <u>Manual uninstall</u> section for details on what resources are retained and how to clean them up.

Configuring a Large Language Model (LLM)

Which LLM is right for your use case depends on a large set of factors specific to your needs and the type of customer experience you want to curate. This solution does not look to be prescriptive, but rather aims to give you the necessary tools to evaluate what works best for your application.

The AI-generated space is evolving rapidly, so it is incumbent on you to keep up to date on the latest models, optimization techniques, and best practices to ensure you are building the right experiences for your customers.

🚺 Note

If you're working with non-public or sensitive data, then be sure to select an LLM option using AWS services (such as Amazon Bedrock or Amazon SageMaker AI). This improves the overall security posture of your deployment by keeping data within your Region and on the AWS network when compared to using an LLM hosted by a third-party provider.

Using Amazon SageMaker AI as an LLM Provider

As of v1.3.0, <u>Amazon SageMaker AI</u> is available as a model provider for Text use cases. This feature allows you to use a SageMaker AI inference endpoint already existing within the AWS account in the solution. Here are some ways to get started.

🛕 Important

The solution does not manage the lifecycle of your SageMaker AI endpoints. You are responsible for deleting the SageMaker AI endpoints once they are no longer needed to stop incurring additional charges.

Creating a SageMaker AI endpoint

You can use <u>Amazon SageMaker AI JumpStart</u> to quickly deploy an endpoint.

You can also use a text-generation based SageMaker AI endpoint and deploy using the base SageMaker AI service. Refer to the <u>SageMaker AI JumpStart documentation</u> for a step by step guide on <u>how to deploy a model</u> for inference.

🚯 Note

Foundation models/LLMs are typically quite large and can often require the use of large accelerated compute instances. Many of these larger instances might not be available by

default in your AWS account. Refer to the default <u>SageMaker AI quotas</u> and be sure to request a quota increase before deploying to avoid common deployment failures.

Use SageMaker AI endpoint to create a Text use case deployment

To deploy a new Text use case using a SageMaker AI endpoint for inference:

- 1. <u>Create a new use case</u> through the Deployment dashboard wizard and complete the forms until you reach the Models selection page.
- 2. On the Models page, select **SageMaker AI** as the model provider. This will generate a custom form requiring three key pieces of user input:
 - The name of the SageMaker AI endpoint you want to use. DevOps users can obtain this from the AWS console. Note that the endpoint must be in the same account and Region as the solution is deployed in.

Location of the endpoint name on the AWS console

| Amazon SageMaker > Endpoints > meta-textgeneration-llama-2- meta-textgeneration-llama-2-7b-f-2 | | Delete |
|---|-----------------------|-------------------|
| Endpoint summary | | |
| Name meta-textgeneration-llama-2-7b-f-2024-01-11-18-25-16-703 | Status ⊘ InService | Type Real-time |
| ARN | Creation time | Last updated |

- The schema of the input payload expected by the endpoint. To support the widest set of endpoints, admin users are required to tell the solution how their endpoint expects the input to be formatted. In the model selection wizard, provide the JSON schema for the solution to send to the endpoint. You can add placeholders to inject static and dynamic values into the request payload. The available options are:
 - Mandatory placeholders: \<\<prompt\>\> will be dynamically replaced with the full input (for example, history, context, and user input as per the prompt template) to be sent to the SageMaker AI endpoint at runtime.
 - Optional placeholders: \<\<temperature\>\> *,* as well as any parameters defined in advanced model parameters can be provided to the endpoint. Any string containing a placeholder enclosed in \<\< and \>\> (for example, \<\<max_new_tokens\>\>) will be replaced by the value of the advanced model parameter of the same name.

Example input schema - setting mandatory fields, prompt and temperature, along with a custom advanced parameter, max_new_tokens. Output path must be supplied as a valid JSONPath string

| <u>Select use case</u> | Select model Info | |
|--|--|--|
| Step 2 - optional Select network configuration | Model selection | |
| Step 3 Select model | Model provider Info Select the model provider you want to use. | |
| Step 4 - optional | SageMaker | |
| Select knowledge base | Sagemaker endpoint name - required Info Enter the name of the SageMaker inference endpoint in this AWS account to be us | ed. |
| O <u>Review and create</u> | meta-textgeneration-llama-2-7b-f-2024-01-11-18-25-16-703 | |
| | Note: The SageMaker endpoint name is case sensitive. | |
| | Input Payload Schema - <i>required</i> Provide the input schema that your endpoint expects. | Rendered Input Payload Rendered payload with the provided prompt and model parameters. |
| | <pre>1 { 2 "inputs": "<<pre>room parameters"; 3 "parameters"; { 4 "temperature": "<<temperature>>", 5 "max_new_tokens": "<<max_new_tokens>>" 6 } 7 }</max_new_tokens></temperature></pre></pre> | <pre>{ "inputs": "How many regions does AWS have?", "parameters": { "temperature": 1, "max_new_tokens": 1000 } }</pre> |
| | JSON Ln 5, Col 42 | |
| | Model Parameters section, wrapped with "< <key>>" to inject the values into the expected structure. Output path - <i>required</i> JSONPath expression that evaluates to the location of the generated text from the</key> | model's output response. |
| | \$[0].generated_text | |

3. The location of the LLMs generated string response within the output payload. This must be supplied as a JSONPath expression to indicate where the final text response shown to users is expected to be accessed from within the endpoint's return object and response.

Example of adding Advanced model parameters to use within SageMaker AI input schema (see Figure 2 for previous options/settings)

| ' Additional settings | | |
|---|---------------------------------------|---|
| lodel temperature his parameter regulates the randomness enerates creative responses. | s or creativity of the model's predic | tions. Use a temperature closer to 0 for analytical, deterministic or multiple choice queries. A higher temperature |
| 1 | | |
| in: 0, Max: 100. | | |
| erbose enabled, additional logs will be written t | o Amazon CloudWatch. | Streaming If enabled, the response from the model will be streamed |
| [INST] {history} {input} [/INST] | | |
| dvanced model param | | lease consult the model documentation to know what parameters the model accepts |
| | Value | Туре |
| ey | | |

(i) Note

SageMaker AI now supports hosting multiple models behind the same endpoint, and this is the default configuration when deploying an endpoint in the current version of SageMaker AI Studio (not Studio Classic).

If your endpoint is configured in this way, you will be required to add

InferenceComponentName to the advanced model parameters section, with a value corresponding to the name of the model you want to use.

Advanced LLM Settings

While using Amazon Bedrock, you can configure some advanced settings for your models such as Amazon Bedrock Guardrails, Provisioned Throughput for Amazon Bedrock, and additional model parameters.

Amazon Bedrock Guardrails

Amazon Bedrock Guardrails is a feature with Amazon Bedrock which evaluates user inputs and LLM responses based on user configured policies and provides an additional layer of safeguards, regardless of the underlying LLM that the user selects for a use case. A Guardrail consists of 2 policies to avoid content that falls into undesirable or harmful categories:

- 1. Denied topics to define a set of topics that are undesirable in the context of user's application, for example, investment advice in a financial application, and,
- 2. Content filters****which allows filtering input user prompts or model responses containing harmful content.

For usage in Generative AI Application Builder solution, a Guardrail must be configured in the *Amazon Bedrock* console using the *Create guardrail* wizard. Once created, you can add this Guardrail to your chat use case created through Generative AI Application Builder solution wizard in the **Additional settings** in the Model Selection step by supplying your Guardrail Identifier and Guardrail version.

Depicts Deployment wizard - enabling Amazon Bedrock Guardrails

| selection wider Info model provider you want to use. me* Info ame of the model from the model provider to use for this deployment. |
|--|
| nodel provider you want to use. |
| me* Info |
| |
| |
| ic.claude-3-sonnet-20240229-v1:0 |
| u like to use an on-demand model or a provisioned model? Info Irock supports Provisioned Throughput to support a higher rate of inputs and outputs processed by the model. Provisioned models have an unique ARN red to process queries. Provisioned throughput can be configured through the Bedrock console. mand oned |
| onal settings |
| nperature ter regulates the randomness or creativity of the model's predictions. Use a temperature closer to 0 for analytical, deterministic or multiple choice igher temperature generates creative responses. |
| |
| 1. |
| u like to enable guardrails? Info Identifier - <i>required</i> Info identifier of the Bedrock guardrail that you want to be applied to all LLM invocations. |
| ts012 |
| |

Provisioned Throughput for Amazon Bedrock

If enabled, additional logs will be written to Amazon CloudWatch

DRAFT

Each on-demand Amazon Bedrock model follows region-specific <u>account quota limit</u> for model inferencing. For example, Anthropic Claude 2.x on Bedrock currently allows for 500 requests and 500,000 tokens processed per minute in us-east-1 and us-west-2 regions. You may also want to use the solution with your fine-tuned or continued pre-trained models. For such instances, Amazon Bedrock allows <u>provisioned throughput</u> which allows running large consistent inference workloads for your base, fine-tuned or continued pre-trained models for use in production-grade applications.

Streaming

he response from the model will be streamed

Once Provisioned Throughput is purchased within the Amazon Bedrock console, a Model ARN is generated for usage. You can now supply this Model ARN in the Generative AI Application Builder wizard in the Model selection step. To do so, select Bedrock as the model provider and the base model name which was used to generate this provisioned Model ARN in Amazon Bedrock console.

Then, select '**Provisioned model'** when choosing between on-demand and provisioned models, and supply your Model ARN.

Depicts Deployment wizard - Enabling Provisioned Throughput for Amazon Bedrock

| Step 1 Select use case | Select model Info |
|---|--|
| Step 2 - optional Select network configuration | Model selection |
| Step 3 | Model provider Info Select the model provider you want to use. |
| Step 4 - optional | Bedrock |
| Select knowledge base Step 5 | Model name* Info Select the name of the model from the model provider to use for this deployment. |
| O Select prompt | anthropic.claude-3-sonnet-20240229-v1:0 |
| Step 6 Review and create | Would you like to use an on-demand model or a provisioned model? Info Amazon Bedrock supports Provisioned Throughput to support a higher rate of inputs and outputs processed by the model. Provisioned models have an unique ARN that is required to process queries. Provisioned throughput can be configured through the Bedrock console. On-Demand Provisioned Model ARN - required info ARN of the provisioned/custom model to use from Amazon Bedrock. arn:aws:bedrock:us-east-1:123456789012:provisioned-model/z8g9xzoxoxmw |
| | Additional settings |
| | Advanced model parameters Model parameters are passed to the model as they are inputted. Please consult the model documentation to know what parameters the model accepts |
| | Add new item |
| | Cancel Previous Next |

Note

Your guardrail and provisioned throughput must be in the same Region as the deployed Deployment Dashboard and use case stacks.

Model parameters

LLMs often accept a wide range of parameters specific to its implementation. Model providers often provide documentation outlining the set of supported parameters and their uses.

The solution passes model parameters directly through to the underlying model so it is important to ensure parameters are set correctly. Refer to the model provider's documentation for the latest information on supported parameters.

Tips for managing model token limits

Note: The solution does not directly attempt to manage token limits imposed by various LLMs. Test and ensure your prompt remains within the available limits enforced by the model provider.

To help control the size of prompts, try the following:

- 1. Familiarize yourself with the limits imposed by the model you want to use. These values can differ dramatically across models so it's important to know what your available budget is before getting started.
- 2. Craft your initial prompt with that budget in mind and consider how much you want to save for any dynamic elements of the prompt. For example, user input, chat history, document excerpts, and so on.
- 3. In the prompt configuration page, set a limit for **Size of trailing history** to limit the number of conversation turns included within the prompt.
- 4. Set document return limits in the Knowledge Base configuration wizard. You need to try and strike the right balance between providing the LLM with enough context to perform the task, but not so much as to exceed token limits or negatively affect latency.
- 5. Leave some buffer. Don't budget for the typical case, think about and experiment with the edge cases such as long input queries, large document excerpts, or long conversations.

Configuring a knowledge base

This section describes how to ingest data into the knowledge base you've selected for the solution. The solution currently supports Amazon Kendra and Amazon Bedrock Knowledge Bases as knowledge bases for your RAG-based use case deployment.

Amazon Kendra

If you're using Amazon Kendra as your knowledge base, refer to the <u>Amazon Kendra Developer</u> <u>Guide</u> for information on how to use various data source connectors to help you ingest data from a wide selection sources. Important: To prevent accidental data loss, the solution does not automatically delete the Kendra index (whether created by the solution or otherwise) when a deployment or stack is deleted. If you want to delete your knowledge base and stop incurring costs, see the <u>Manual uninstall</u> section for details on which resources are retained and how to clean them up.

Amazon Bedrock Knowledge Bases

Amazon Bedrock Knowledge Bases can be backed by a variety of different vector stores, each with the capability of indexing your data. To set up and populate your knowledge base, consult the <u>Amazon Bedrock User Guide</u>. Specifically, you will want to:

- First set up your data source
- Then <u>set up a vector index for your knowledge base in a supported vector store</u>. Note that this can be skipped if you use the "Quick create a new vector store" option in Bedrock console during knowledge base creation.
- Finally, you can create the knowledge base and sync your configured data sources.

Advanced knowledge base settings

Advanced Knowledge Base Settings such as Knowledge Base Filtering and RAG with Role Based Access Control are available for use with the solution. Knowledge Base Filtering can apply to either of the Knowledge Bases while RAG with Role Based Access Control is specifically available for Amazon Kendra.

Knowledge base filtering

The solution allows you to specify <u>Amazon Kendra attribute filters</u> or <u>Bedrock knowledge base</u> <u>retrieval filters</u> when deploying a use case in the Advanced RAG configurations section of the wizards knowledge base step. These filters define how data sources in the knowledge base are queried, such as search strategies, languages of the underlying document being queries, etc.

In both cases, a JSON object is used to specify the filter settings per the format specified in each services documentation (as linked above).

Example 1: Kendra AttributeFilter

```
{
    "EqualsTo": {
```

```
"Key": "_language_code",
"Value": {
  "StringValue": "es"
  }
 }
```

Example 2: Bedrock RetrievalFilter

```
{
  "equals": {
  "key": "language",
  "value": "es"
  }
}
```

RAG with Role Based Access Control with Amazon Kendra

<u>Role-based access control (RBAC)</u> allows controlling which users or groups can access certain documents in your Amazon Kendra index or see certain documents in their search results. To configure RBAC for your Amazon Kendra Index ID with your Generative AI Application Builder on AWS (GAAB) use case, follow these steps:

1. Configure Amazon Kendra Index

- 1. Ensure that you have an Amazon Kendra index created and at least one data source added to it.
- 2. Configure access control for your data source based on user groups. For an S3 data source, follow the <u>instructions available in the documentation</u> to set up access control lists (ACLs) using the same group names created in your Amazon Cognito User Pool. This ensures that users can only access the documents and search results they are authorized to view based on their group membership.

🚺 Note

Under User Access Control in the Kendra Index you created, leave Token-based user access control as No. When you enable Role Based Access Control in Step 2, Generative AI Application Builder on AWS extracts the appropriate claims from the user authentication token and creates an Attribute Filter.

2. Deploy RAG Use Case using GAAB Deployment Wizard

- 1. Follow the on-screen wizard instructions in the GAAB Deployment Wizard until you reach step 4 of the wizard to configure RAG.
- 2. In the **Select Knowledge Base** step of the deployment wizard, choose **Amazon Kendra** as the knowledge base type.
- 3. Specify whether you have an existing Amazon Kendra index or if you want to create a new one. If you have an existing index, provide the ID of your Amazon Kendra index that has been configured with access control lists (ACLs) based on user groups.
- 4. Enable the **Role Based Access Control** option. This option ensures that the search results returned from the Amazon Kendra index are filtered based on the user's role and group permissions.
- 5. Review and deploy the use case.

3. Configure Amazon Cognito

- 1. Locate the Amazon Cognito User Pool used by your GAAB deployment. This Amazon Cognito User Pool is typically created by the main deployment dashboard CloudFormation stack.
- 2. Create new users in the Amazon Cognito User Pool. When creating users, select the 'Send an email invitation' option so that users receive temporary login credentials via email. This allows new users to sign up and access the GAAB application.
- 3. Create user groups in the Amazon Cognito User Pool. Ensure that the group names exactly match the groups configured in your Amazon Kendra index ACLs. This is crucial for enabling RBAC, as the user's group membership will determine the search results they can access.
- 4. Assign users to the appropriate groups based on their roles and access permissions. Users must be added to both the group required for the Amazon Kendra index ACL, as well as the use case-specific group created during the GAAB deployment. This ensures that users have the necessary permissions to access the specific use case and the relevant search results.

By following these steps, you will have configured role-based access control (RBAC) for your GAAB deployment, ensuring that users can only access and interact with the information and features they are authorized for, based on their assigned user group and permissions.

1 Note

As of now, only Amazon Kendra supports RBAC for knowledge bases in the Generative AI Application Builder on AWS. For Amazon Bedrock Knowledge Base, RBAC is not supported, but you can use metadata filters to achieve some level of filtering. For more information, refer to the <u>Amazon Bedrock User Guide</u>.

Configuring your prompts

The Deployment dashboard wizard has a prompt configuration step which allows you to customize the prompt experience and template that will guide the interactions between users and the AI model. Properly configuring these settings is crucial for obtaining accurate and relevant responses from the AI assistant.

This section controls the overall experience and behavior of the AI prompt.

- Max prompt template length: This setting determines the maximum length (in characters) of the prompt template. A higher value allows for more context to be provided to the AI model, potentially leading to more accurate responses. However, excessively long prompts may also introduce noise and negatively impact performance. For Amazon Bedrock models, the default values for max prompt template length (in characters) is calculated using the underlying model token limits. If you edit and change a model name within Bedrock, 'Reset to default' button is highlighted and can be used to adopt the newly selected model's defaults. For Amazon SageMaker AI models, reasonable default values are provided, but it is recommended that you check your underlying model and choose these max prompt template length and input text lengths accordingly. Refer to the Tips on managing model token limits section for more information.
- Max input text length: This setting limits the maximum length (in characters) of the user's input text. Longer inputs may contain irrelevant information, increasing the risk of obtaining irrelevant or inaccurate responses from the AI model.
- User Prompt Editing: This option allows you to enable or disable the ability for users to modify the prompt template through the Chat UI. Disabling this feature can help maintain consistency and prevent unintended changes to the prompt.

Prompt template

This section allows you to define the actual prompt template that will be used by the AI model. The prompt template typically follows a structure that includes placeholders for various components, such as the user's input, reference passages, and chat history.

- **Prompt template**: This is the main text area where you can write or paste the desired prompt template. The template should be crafted to provide the necessary context and instructions to the AI model. It typically includes the following placeholders:
 - {input}: This placeholder is mandatory for Sagemaker AI deployments and will be substituted with the user's input or query.
 - {history}: This placeholder is mandatory for Sagemaker AI deployments and will be substituted with the chat history of the current conversation.
 - {context}: This placeholder is mandatory for RAG deployments and will be substituted with the document excerpts obtained from the configured knowledge base.
- **Rephrase Question?**: This option (available for RAG deployments only) determines whether the user's original input query should be rephrased or disambiguated before being passed to the AI model. Rephrasing the query can sometimes help the model better understand the user's intent, potentially leading to more accurate responses.

When configuring the prompt template and experience, it's essential to strike a balance between providing sufficient context and instructions to the AI model while avoiding excessively long or irrelevant information that may introduce noise or performance issues.

Advanced prompt settings

This section allows you to control how the conversation history is presented to the AI model.

- **Size of trailing history**: This setting determines the number of previous messages that should be included in the final prompt. Setting this value to zero would result in no history being injected into either the prompt template or the disambiguation prompt template. Please note: even when set to zero, a {history} placeholder is still required to exist in the prompt templates. At runtime, it will get replaced with an empty string.
 - Note: It is recommended to provide an even number for this value. Providing an odd number would result in only the AI response of a paired interaction being returned.
- Human Prefix: This is the prefix used to identify messages sent by the user in the conversation history.

• Al Prefix: This is the prefix used to identify messages returned by the Al model in the conversation history.

Disambiguation Prompt Configuration

This section allows you to configure the behavior and template for disambiguating user inputs before sending them to the configured knowledge base.

- Enable Disambiguation: This option determines whether user inputs should be disambiguated before sending to the knowledge base.
- **Disambiguation Prompt Template**: This is the prompt template used for disambiguating user inputs when connected to a knowledge base. The output generated from this prompt will be used as the query sent to the knowledge base. Disabling disambiguation would result in the user's raw query being sent to the knowledge base unchanged.

For example, with disambiguation enabled, a follow-up user query of "How much does it cost?" might be disambiguated to "How much does it cost renew my license plate?", leading to a better search query.

Using the deployed Text use case

The built-in UI for the Text use case is intended to enable business users to quickly explore and experiment with the deployment created by the admin user. Configuration changes made by the business user only take effect for their session. The business user must share these changes with the admin user who can update the base deployment with those changes for all to use.

The chat UI consists of the following components:

- Chat window
- Chat input box
- Settings
- Clear conversation

Chat window

Holds different turns of the conversation. Messages starting on the right are from the business user, and messages starting on the left are from the configured LLM. A small clipboard icon exists on all LLM responses to enable easy copying of responses.

Chat input box

Pinned to the bottom of the chat window is the chat input box. This is where business users can enter their messages to be sent to the LLM. Just above the input box is the connection status. If the connection is lost (for example, due to inactivity), a new connection is automatically created the next time a chat message is sent. This request is expected to take a little longer due to the additional WebSocket connection time.

Based on the specific configuration, there might be a maximum length enforced on the input. If this limit is exceeded, users receive an alert and the message is not sent.

Note: If using RAG with Amazon Kendra, the <u>Retrieve API</u> will truncate queries to 30 token words. If expecting longer user inputs, evaluate how this might affect search performance.

Settings

To enable business users to quickly experiment with different configurations, a settings panel is available, which enables on-the-fly editing of certain deployment configuration options

(example, prompt template). These changes can only be made at the start of a new session. Once a conversation is started, clearing the conversation re-enables the editing of the configuration settings.

Note: Admin users can choose to lock a deployment's settings. They can prevent live edits at deployment time through the wizard during the prompt step.

Clear conversation

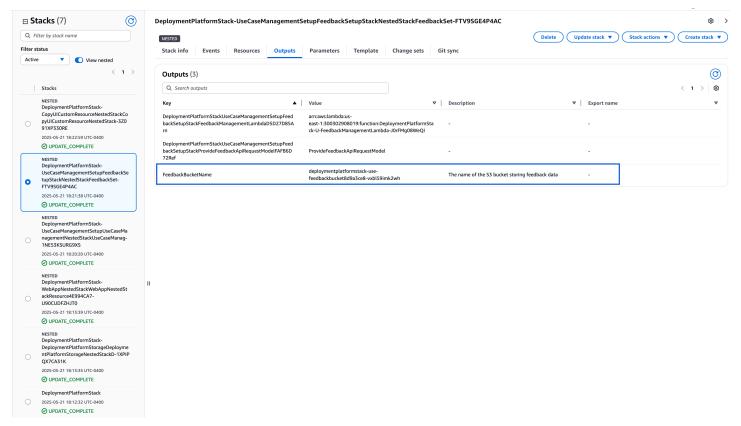
Over the course of the conversation, the solution maintains a chat history, which enables a conversational experience. This enables query disambiguation and follow-up questions. To reset a conversation and delete all chat history for this interaction, choose *Clear conversation *at the top of the chat window. Once the conversation has been cleared, a new session is created which re-enables editing of the settings.

Accessing and analyzing user collected feedback

As of v3.0.0, the Deployment Dashboard deploys a nested feedback stack which allows Text and Agent usecases deployed with the Dashboard to have the functionality of feedback collection for the responses that the LLM/Agent generates. Particularly, users can provide a positive or negative feedback along with an optional comment. If the user provides a negative feedback, they can further select one of these negative categories: 'Inaccurate', 'Incomplete or insufficient', 'Harmful' and/or 'Other'.

Once the user provides the feedback, the feedback is stored in an S3 bucket partitioned by Use Case ID, year and month. The Use Case ID can be found in the Deployment Dashboard and the Feedback S3 bucket can be found in the outputs of the feedback nested stack of the Deployment Dashboard stack:

Depicts Deployment stack - Finding Feedback Bucket Name



The user feedback is sent as an API request containing a minimal set of information:

```
{
    "useCaseRecordKey": "a1b2c3d4-e5f6g7h8",
    "conversationId": "12345678-1234-1234-1234-123456789012",
```

```
"messageId": "87654321-4321-4321-210987654321",
"rephrasedQuery": "What are the key features of the Generative AI Application Builder
on AWS?",
"sourceDocuments": [
"s3://bucket-name/document1.pdf",
"s3://bucket-name/document2.pdf"
],
"feedback": "positive",
"feedbackReason": [
"Incomplete or insufficient"
],
"comment": "The response was helpful but could include more details about important
features."
}
```

This payload is then processed by a lambda using the useCaseRecordKey which identifies the correct configuration of a usecase at the time of deployment. This configuration is used to get specific details for the feedback such as the ConversationTable's name (contains all the conversations and human and AI message sequences) which is further used to retrieve the the actual userInput and llmResponse. Additional details are also attached to this feedback record such as the agentId and agentAliasId for an Agent usecase, and modelProvider, bedrockModelId, etc. for a Text usecase using this configuration. For details on how to access this configuration, see <u>Custom Feedback Mappings</u> section below. Each incoming feedback request is stored as a JSON object and a sample feedback record can look like this for a Text usecase:

```
{
   "useCaseId": "12345678-1234-1234-1234-123456789012",
   "useCaseRecordKey": "c07a2e3b-2f31b1e0",
   "userId": "22345678-1234-1234-1234-123456789012",
   "conversationId": "dd51de5d-5af1-4ec6-91d2-aadf14352109",
   "messageId": "32345678-1234-1234-1234-123456789012",
   "userInput": "What are its key features?",
   "rephrasedQuery": "What are the key features of the Generative AI Application
Builder on AWS?",
   "llmResponse": "Generative AI Application Builder on AWS can help you build
production ready enterprise chatbots rapidly.",
   "feedback": "negative",
   "feedbackReason": [
      "Incomplete or insufficient"
  ],
   "comment": "The response was helpful but could include more details about important
 features.",
```

```
"timestamp": "2025-05-22T18:48:08.340Z",
    "feedbackId": "42345678-1234-1234-1234-123456789012",
    "useCaseType": "Text",
    "modelProvider": "Bedrock",
    "bedrockModelId": "amazon.nova-lite-v1:0",
    "ragEnabled": "false"
}
```

or like this for an Agent usecase:

```
{
   "useCaseId": "12345678-1234-1234-1234-123456789012",
   "useCaseRecordKey": "c07a2e3b-2f31b1e0",
   "userId": "22345678-1234-1234-1234-123456789012",
   "conversationId": "dd51de5d-5af1-4ec6-91d2-aadf14352109",
   "messageId": "32345678-1234-1234-1234-123456789012",
   "userInput": "What are its key features?",
   "llmResponse": "Generative AI Application Builder on AWS can help you build
 production ready enterprise chatbots rapidly.",
   "feedback": "negative",
   "feedbackReason": [
      "Incomplete or insufficient"
   ],
   "comment": "The response was helpful but could include more details about important
 features.",
   "timestamp": "2025-05-22T18:48:08.340Z",
   "feedbackId": "42345678-1234-1234-1234-123456789012",
   "useCaseType": "Agent",
   "agentId": "AHFXUJCAK1",
   "agentAliasId": "KSEDKOS0BL"
}
```

This feedback can then be used for further processing, analyzing and model re-training/feedback loops. You can also add custom mappings to enhance the feedback record being stored in the feedback lambda.

Custom Feedback Mappings

The Deployment Dashboard contains a LLMConfigTable which can be found in the stack outputs of the Deployment Dashboard stack with the key LLMConfigTableName. LLMConfigTable contains the configurations for each usecase based on the settings selected by the admin while deploying the usecase through the Deployment Dashboard wizard. Each usecase configuration

is identified by its useCaseRecordKey. Here is a sample usecase configuration record in the LLMConfigTable:

```
{
   "key": "2dd76cfa-bc1a14da",
   "config": {
      "ConversationMemoryParams": {
        . . .
      },
      "FeedbackParams": {
         "CustomMappings": {
            "NumberOfDocs": "$.KnowledgeBaseParams.NumberOfDocs",
            "ScoreThreshold": "$.KnowledgeBaseParams.ScoreThreshold"
         },
         "FeedbackEnabled": true
      },
      "IsInternalUser": "true",
      "KnowledgeBaseParams": {
         "KendraKnowledgeBaseParams": {
            "ExistingKendraIndexId": "d2831033-667f-4539-ab28-e6c7c7c5988b",
            "RoleBasedAccessControlEnabled": false
         },
         "KnowledgeBaseType": "Kendra",
         "NumberOfDocs": 5,
         "ReturnSourceDocs": false,
         "ScoreThreshold": 0.3
      },
      "LlmParams": {
         "BedrockLlmParams": {
            "BedrockInferenceType": "QUICK_START",
            "ModelId": "amazon.nova-lite-v1:0"
         },
         "ModelParams": {},
         "ModelProvider": "Bedrock",
         "PromptParams": {
            . . .
         },
         "RAGEnabled": true,
         "Streaming": false,
         "Temperature": 0.1,
         "Verbose": false
      },
      "UseCaseName": "test-rag-usecase",
```

}

```
"UseCaseType": "Text"
}
```

If feedback is enabled for a usecase, this configuration will contain a FeedbackParams object which allows a CustomMappings object inside it that can specify the JSONPaths for all the additional fields to be added to the feedback JSON record stored in the feedback S3 bucket. For example, for the above sample usecase configuration, the CustomMappings contains NumberOfDocs and ScoreThreshold JSONPaths additionally in the CustomMappings object which start with config as the root of the JSONPath. With this configuration, each JSON record stored in the feedback S3 bucket will start getting these 2 additional values apart from the fields that have already been provided.

Analyzing feedback data

The feedback data is stored in S3 as JSON objects. Here are some approaches to make this feedback data more accessible and actionable:

Using AWS Glue and Amazon Athena

<u>AWS Glue</u> and <u>Amazon Athena</u> provide a serverless way to catalog, query, and analyze your feedback data.

AWS Glue allows you to create an <u>AWS Glue crawler</u> that inspects the data in an S3 bucket, infers its schema, and records all the relevant metadata in a catalog. Post that, services like Amazon Athena can be used to query the data.

You can refer <u>AWS Athena Documentation</u> on steps for connecting the feedback S3 bucket with Amazon Athena using AWS Glue Data Catalog. You can also use some of Glue's more powerful features to perform Extract Transform & Load (ETL) jobs on this data and transform it into a format that suits your analytics or model re-training use cases. With Glue, you can perform operations such as filtering the records with certain feedback types, filling out any missing information, and you can also load this data into another storage location such as another S3 bucket or a different AWS data store.

🚯 Note

Depending upon your use case, consider scheduling the Glue crawler to run periodically (e.g., weekly) rather than nightly to optimize costs as feedback data can be sparse.

Using the solution's CloudWatch Dashboards

You also have access to a **CloudWatch Dashboard** packaged with the solution that can provide you trends for positive and negative feedback, negative feedback reason categories, etc. on a per use-case basis. You can find this dashboard using your usecase name in *Dashboards* inside the AWS CloudWatch console:

| test-d8b73596-us-east-1-Dashboard 🔻 🏠 | | (5)(2) 1h 3h 12h 1d 3d 1w Custor | □ Autosave: Off □ ■ UTC timezone ▼ C ▼ Actions ▼ Save + |
|---|---|--|--|
| Websocket Count Stats ① : | Websocket Connection Latency Stats O : | Cognito Sign-Ins & Sign-ups ① : | LangChain Stats : |
| Count | Milliseconds | Various units | Count 29.0 |
| 1,636 | 698 | 1.0 | 280 |
| 8 02:00 02:30 03:00 03:30 04:00 04:30 05:00 ConnectCount MessageCount ClientError ExecutionError | 187 02:00 02:30 03:60 03:30 04:00 04:30 05:00 ■ AverageIntegrationLatency | 0 02:00 02:30 03:00 03:30 04:00 04:30 05:00 SignifSuccessesCount AverageSignifSuccesses SignifpSuccessesCount AverageSignifpSuccesses | 27.0 0.00 0.230 0.500 0.400 0.430 0.500 0.200 0.230 0.500 0.430 0.500 0.430 0.500 I LangchainfailuresCount IncorrectinputFailuresCount IncorrectinputFailuresCount |
| LangChain LLM Query Processing Stats | Kendra Query Stats | Kendra Latency Stats : | Feedback Sentiment Count Stats |
| Various units | Count | Various units | Count |
| 28.00 | 1.00 No data available. Try adjusting the dashboard time range. | 1.00 No data available. Try adjusting the dashboard time range. | 9.0 • |
| 16.26 | 0.50 | 0.50 | 7.0 |
| 4.53 02:00 02:30 05:00 05:30 0.450 04:30 05:00 LangchainQueryProcessingTime AverageLangchainQueryProcessingTime | 0 0 22:0 02:30 05:00 05:30 04:00 04:30 05:00 0 EnndraQueriesCount | 0 02:00 02:30 05:00 05:30 04:00 04:30 05:00 KendraProcessingTime AverageKendraProcessingTime | 5.0 02:00 02:30 05:00 03:30 04:00 04:30 05:00 ProstiveFeedbackSubmittedCount |
| Feedback Error Stats : | Negative Feedback Reasons | | |
| Count | Count | | |
| 6.0 | 4.0 | | |
| 5.0 | | | |
| 4.0 02:00 02:30 03:00 03:30 04:00 04:30 05:00 | 1.0 02:00 02:30 03:00 03:30 04:00 04:30 05:00 | | |
| FeedbackRejectionCount FeedbackProcessingErrorCount FeedbackStorageErrorCount | InaccurateFeedbackCount IncompleteOrInsufficientFeedbackCount HarmfulFeedbackCount OtherNegativeFeedbackCount | | |

Depicts Usecase CloudWatch Dashboard

You can also build additional widgets in this Dashboard or create QuickSight dashboards.

Best practices for feedback data analysis

- Implement data lifecycle policies on your S3 bucket to archive older feedback data to lowercost storage tiers
- **Create separate analysis for each use case** to identify model-specific improvement opportunities
- Establish feedback thresholds that trigger alerts when negative feedback exceeds acceptable levels
- **Export critical insights** periodically for sharing with stakeholders and model improvement teams

Viewing operation metrics for a deployment

The Deployment dashboard and use case stacks each come with their own CloudWatch dashboard tracking various operational metrics of the solution. You can use these CloudWatch dashboards to help compare different deployments. To access the dashboards:

- 1. Navigate to the <u>CloudWatch console</u>.
- 2. Search for the pre-built dashboards either by looking up the stack name, or universally unique identifier (UUID).

For example, the Text use case comes with graphs tracking the number of WebSocket connections, the number of user sign ins and sign ups, the amount of time the LLM took to process a completion, and so on. Customers can use these graphs to compare various _quantitative _metrics of a deployment.

Example

It is difficult to compare the *qualitative* results of various models applied to different use cases. Use the <u>Clone feature</u> to spin up multiple deployments quickly so that you can compare the outputs side by side.

Access CloudWatch Logs insights

This solution logs error, warning, informational, and debugging messages for the Lambda functions. To choose the type of messages to log:

- 1. Locate the applicable function in the AWS Lambda console.
- 2. Add a **POWERTOOLS_LOG_LEVEL** environment variable.
- 3. Set the variable to the applicable type of message.

For further instructions, see <u>Create Lambda environment variables</u> in the AWS Lambda Developer Guide.

The following table lists the types of log levels that you can choose from.

| Level | Description |
|---------|--|
| ERROR | Logs include information on anything that causes an operation to fail. |
| WARNING | Logs include information on anything that could potentially cause inconsistencies in the function but might not necessarily cause the operation to fail. Logs also include ERROR messages. |
| INFO | Logs include high-level information about how the function is operating. Logs also include ERROR and WARNING messages. |
| DEBUG | Logs include information that might be helpful when debugging a problem with the function. Logs also include ERROR, WARNING, and INFO messages. |

Use the following procedure to add CloudWatch Logs insights to this solution.

- 1. Identify the relevant log groups:
 - a. Sign in to the <u>AWS CloudFormation console</u>.
 - b. Choose your target stack.
 - c. Select the **Resources** tab and search for your target Lambda functions.
 - d. Sign in to the AWS Lambda console and choose each of your target Lambda functions.
 - e. For each of your target Lambda functions, select the **Monitor** tab and choose **View CloudWatch Logs**.
 - f. Copy the names of the log groups that you want to extract insights from.
- 2. Navigate to the Amazon CloudWatch console.
- 3. On the navigation menu, under **Logs**, choose **Logs Insights**.
- 4. On the **Logs Insights** page, choose the **Logs** tab.
- 5. Search for log group names from step 1.
- 6. Copy one of the following example queries and paste it into the query field:

a. To identify all client exceptions:

b. To retrieve count of invocations by function name:

```
stats count(*) by function_name
```

c. To retrieve count of invocations over five-minute intervals:

```
stats count(*) as invocations by bin(5m)
```

d. To retrieve all <u>AWS X-Ray</u> trace IDs:

```
filter @message like "XRAY TraceId"
|parse @message "XRAY TraceId: * " as traceId|stats count(*) by traceId
```

e. To retrieve logs relating to a specific X-Ray Trace ID:

filter @message like "your-traceid-here"

f. To retrieve unauthorized WebSocket errors:

g. To retrieve count of metrics published:

```
filter @message like "CloudWatchMetrics"
|parse @message /"Metrics":\s*\[(?<metrics>.*?)\]/|stats count(*) as metric_count
    by metrics
```

Developer guide

This section provides the <u>source code</u> for the solution, an <u>integration guide</u>, a <u>customization guide</u>, and <u>API reference</u>.

Source code

Visit our <u>GitHub repository</u> to download the source files for this solution and to share your customizations with others.

The Generative AI Application Builder on AWS templates are generated using the <u>AWS Cloud</u> <u>Development Kit (AWS CDK)</u>. See the <u>README.md</u> file for additional information.

Integration guide

The entire solution is designed to be easily extensible. The orchestration layer of this solution is built using <u>LangChain</u>. You can add any model provider, knowledge base, or conversation memory type supported by LangChain (or a third party that provides LangChain connectors for these components) to this solution.

Expanding supported LLMs

To add another model provider, such as a custom LLM provider, you must update the following three components of the solution:

- 1. Create a new TextUseCase CDK stack, which deploys the chat application configured with your custom LLM provider:
 - a. Clone this solution's <u>GitHub repository</u>, and set up your build environment by following the instructions provided in the <u>README.md</u> file.
 - b. Copy (or create new) the source/infrastructure/lib/bedrock-chat-stack.ts file, paste it to the same directory, and rename it to custom-chat-stack.ts.
 - c. Rename the class in the file to a suitable one, such as CustomLLMChat.
 - d. You can choose to add a Secrets Manager secret to this stack, which stores your credentials for your custom LLM. You can retrieve these credentials during model invocation in the chat Lambda layer discussed in the next paragraph.

- 2. Build and attach a Lambda layer containing the Python library of the model provider to be added. For an Amazon Bedrock use case chat application, the langchain-aws Python library contains the custom connectors on top of the LangChain package to connect to the AWS model providers (Amazon Bedrock and SageMaker AI), knowledge bases (Amazon Kendra and Amazon Bedrock Knowledge Bases), and memory types (such as DynamoDB). Similarly, other model providers have their own connectors. This layer helps you attach this model provider's Python library so that you can use these connectors in the chat Lambda layer, which invokes the LLM (step 3). In this solution, a custom asset bundler is used to build Lambda layers, which are attached using CDK aspects. To create a new layer for the custom model provider library:
 - a. Navigate to the LambdaAspects class in the source/infrastructure/lib/utils/ lambda-aspects.ts file.
 - b. Follow the instructions on how to extend the functionality of the Lambda aspects class provided in the file (such as adding the getOrCreateLangchainLayer method). To use this new method (for example, getOrCreateCustomLLMLayer), also update the LLM_LIBRARY_LAYER_TYPES enum in the source/infrastructure/lib/utils/ constants.ts file.
- 3. Extend the chat Lambda function to implement a builder, client, and handler for the new provider.

The source/lambda/chat contains the LangChain connections for different LLMs along with the supporting classes to build these LLMs. These supporting classes follow Builder and Object Oriented design patterns to create the LLM.

Each handler (for example, bedrock_handler.py) first creates a *client*, checks the environment for required environment variables, and then calls a get_model method to get the LangChain LLM class. The generate method is then called to invoke the LLM and get its response. LangChain currently supports streaming functionality for Amazon Bedrock, but not SageMaker AI. Based on streaming or non-streaming functionality, appropriate WebSocket handler (WebsocketStreamingCallbackHandler or WebsocketHandler) is called to send the response back to the WebSocket connection using the post_to_connection method.

The clients/builder folder contains the classes which help build an LLM Builder using Builder pattern. First, a use_case_config is retrieved from a DynamoDB configurations store, which stores the details on what type of knowledge base, conversation memory, and model to construct. It also contains relevant model details such as model parameters and prompts. The Builder then helps in following the steps for creating a knowledge base, creating a conversation memory to maintain conversation context for LLM, setting the appropriate LangChain callbacks for streaming and non-streaming cases, and creating an LLM model based on the provided model configurations. The DynamoDB configuration is stored at the time of use case creation when you deploy a use case from the Deployment dashboard (or when it is provided by the users in standalone use case stack deployments without the Deployment dashboard).

The clients/factories subfolder helps set the appropriate conversation memory and knowledge base class, based on the LLM configruation. This enables easy extension to any other knowledge base or memory types that you want your implementation to support.

The shared subfolder contains specific implementations of knowledge base and conversation memory which are instantiated inside the factories by the builder. It also contains Amazon Kendra and Amazon Bedrock Knowledge Base retrievers called within LangChain to retrieve documents for the RAG use cases, along with callbacks, which are used by the LangChain LLM model.

The LangChain implementations use LangChain Expression Language (LCEL) to compose conversation chains together. RunnableWithMessageHistory class is used to maintain conversation history with custom LCEL chains, enabling functionalities such as returning source documents and using the rephrased (or disambiguated) question sent to the knowledge base to also be sent to the LLM.

To create your own implementation of a custom provider, you can:

- a. Copy the bedrock_handler.py file and create your custom handler (for example, custom_handler.py), which creates your custom client (for example, CustomProviderClient) (specified in the following step.)
- b. Copy bedrock_client.py in the clients folder. Rename it to custom_provider_client.py (or your specific model provider name, such as CustomProvider). Name the class within it appropriately, such as CustomProviderClient which inherits LLMChatClient.

You can use the methods provided by LLMChatClient or write your own implementations to override these.

The get_model method builds a CustomProviderBuilder (see the following step), and calls the construct_chat_model method that constructs the chat model using builder steps. This method acts as the *Director* in the builder pattern.

c. Copy clients/builders/bedrock_builder.py and rename it to custom_provider_builder.py and the class within it to CustomProviderBuilder that inherits LLMBuilder (llm_builder.py). You can use the methods provided by LLMBuilder or write your own implementations to override these. The builder steps are called in sequence inside the client's construct_chat_model method, such as set_model_defaults, set_knowledge_base, and set_conversation_memory.

The set_llm_model method would create the actual LLM model using all of the values that are set using the methods called before it. Specifically, you can create a RAG (CustomProviderRetrievalLLM) or non-RAG (CustomProviderLLM) LLM, based on the rag_enabled variable that is retrieved from the LLM configuration in DynamoDB.

This configuration is fetched in the retrieve_use_case_config method in the LLMChatClient class.

d. Implement your CustomProviderLLM or CustomProviderRetrievalLLM implementation in the llm_models subfolder based on whether you require RAG or non-RAG use case. Most functionalities to implement these models are provided in their BaseLangChainModel and RetrievalLLM classes respectively, for non-RAG and RAG use cases.

You can copy the llm_models/bedrock.py file and make the necessary changes to call the LangChain model that refer to your custom provider. For example, Amazon Bedrock uses a ChatBedrock class to create a chat model using LangChain.

The generate method generates the LLM response using the LangChain LCEL chains.

You can also use the get_clean_model_params method to sanitize the model parameters per LangChain or your model requirements.

Expanding supported knowledge bases and conversation memory types

To add your implementations of conversation memory or knowledge base, add the required implementations in the shared folder and then edit the factories and appropriate enumerations to create an instance of these classes.

When you supply the LLM configuration, which is stored inside the parameter store, the appropriate conversation memory and knowledge base will be created for your LLM. For example, when the ConversationMemoryType is specified as DynamoDB, an instance of DynamoDBChatMessageHistory (available inside shared_components/memory/ddb_enhanced_message_history.py) is created. When the KnowledgeBaseType

is specified as Amazon Kendra, an instance of KendraKnowledgeBase (available inside shared_components/knowledge/kendra_knowledge_base.py) is created.

Building and deploying the code changes

Build the program with the npm run build command. Once any errors are resolved, run cdk synth to generate the template files and all the Lambda assets.

- 1. You can use the 0/stage-assets.sh script to manually stage any generated assets to the staging bucket in your account.
- 2. Use the following command to deploy or update platform:

```
cdk deploy DeploymentPlatformStack --parameters AdminUserEmail='admin-
email@amazon.com'
```

Any additional AWS CloudFormation parameters should also be supplied along with the **AdminUserEmail** parameter.

Customization guide

Managing Cognito user pool

When the Deployment dashboard is deployed, an Amazon Cognito user pool along with an admin user are created to provide authentication for the application. This user pool is shared across the Deployment dashboard and all use cases. The admin user created on deployment of the dashboard is automatically granted access to all use cases deployed using the dashboard. This mechanism is provided via Amazon Cognito user pool groups.

When a use case is deployed from the dashboard, if an email is provided, a user will be created in the shared user pool, along with a user group named for the specific use case. The newly created user is then added to the group, granting the user access to the use case.

If you wish to add an additional user to a given use case, this can be achieved by creating a user in the Cognito user pool and adding them to the group(s) corresponding to the use case(s) you want the user to have access to. For a step-by-step guide, see <u>Creating a new user in the AWS</u> <u>Management Console</u>.

Similarly, if you want to create additional admin users, you must create a new user and add them to the **Admin** group in the user pool.

The user names are created by taking the portion of the provided email before the @, and appending the generated use case UUID (or -admin in the case of the admin user).

In the **Groups** tab, you can see that an **Admin** group and a group for each use case have been automatically created using the name of the use case (as provided in the wizard) and the use case UUID.

API reference

This section provides API references for the solution.

Deployment dashboard

| REST API | HTTP method | Functionality | Authorized callers |
|--------------------------------|-------------|--|--|
| /deployments | GET | Get all deployments. | Amazon Cognito authenticated JWT token |
| /deployments | POST | Creates a new use case deployment. | Amazon Cognito authenticated JWT token |
| /deployments/ {useCaseId} | GET | Gets deployment details for a single deployment. | Amazon Cognito authenticated JWT token |
| /deployments/ {useCaseId} | РАТСН | Updates a given deployment. | Amazon Cognito authenticated JWT token |
| /deployments/ {useCaseId} | DELETE | Deletes a given deployment. | Amazon Cognito authenticated JWT token |
| /model-info/ use-case-types | GET | Gets the available use case types for the deployment | Amazon Cognito authenticated JWT token |

| REST API | HTTP method | Functionality | Authorized callers |
|--|-------------|--|--|
| /model-info/ {useCaseType}/p roviders | GET | Gets the available model providers for the given use case type | Amazon Cognito authenticated JWT token |
| /model-info/ {useCaseType}/{ providerName} | GET | Gets the IDs of the models available for a given provider and use case type | Amazon Cognito authenticated JWT token |
| /model-info/ {useCaseType}/{ providerName}/ {modelId} | GET | Gets the info about the given model, including default parameters. | Amazon Cognito authenticated JWT token |

(i) Note

OpenAPI and Swagger files can also be exported from API Gateway for easier integration with the API. See Export a REST API from API Gateway.

POST and PATCH Payloads

See below for an example of a POST payload to the /deployments endpoint, which will create a new use case.

```
{
    "UseCaseName": "usecase1",
    "UseCaseDescription": "Description of the use case to be deployed. For display
    purposes", // optional
    "DefaultUserEmail": "email@example.com",
    "DeployUI": true, // optional
    "VpcParams": {
    "VpcEnabled": true,
    "CreateNewVpc": false,
    // provide these if not creating new vpc
    "ExistingVpcId": "vpc-id",
```

```
"ExistingPrivateSubnetIds": ["subnet-1", "subnet-2"],
 "ExistingSecurityGroupIds": ["sg-1", "sg-2"]
},
 "ConversationMemoryParams": {
 "ConversationMemoryType": "DynamoDB",
"HumanPrefix": "user", // optional
"AiPrefix": "ai", // optional
 "ChatHistoryLength": 10 // optional
},
 "KnowledgeBaseParams": {
 "KnowledgeBaseType": "Bedrock",
// one of the following based on selected provider
"BedrockKnowledgeBaseParams": {
 "BedrockKnowledgeBaseId": "my-bedrock-kb",
 "RetrievalFilter": {}, // optional
"OverrideSearchType": "HYBRID" // optional
},
"KendraKnowledgeBaseParams": {
"AttributeFilter": {}, // optional
"RoleBasedAccessControlEnabled": true, // optional
 "ExistingKendraIndexId": "12345678-abcd-1234-abcd-1234567890ab",
// provide the following in place of ExistingKendraIndexId if you want the solution to
deploy an index for you
"KendraIndexName": "index",
 "QueryCapacityUnits": 1, // optional
 "StorageCapacityUnits": 1, // optional
"KendraIndexEdition": "DEVELOPER" // optional
},
 "NoDocsFoundResponse": "Sorry, I couldn't find any relevant information for your
query.", // optional
 "NumberOfDocs": 3, // optional
"ScoreThreshold": 0.7, // optional
 "ReturnSourceDocs": true // optional
},
"LlmParams": {
"ModelProvider": "Bedrock | SAGEMAKER",
// one of the following based on selected provider
"BedrockLlmParams": {
"ModelId": "model-id", // use this for on demand models. Can't use with ModelArn
 "ModelArn": "model-arn", // use this for provisioned/custom models. Can't use with
ModelId,
"InferenceProfileId": "profile-id"
"GuardrailIdentifier": "arn:aws:bedrock:us-east-1:123456789012:guardrail/my-
guardrail", // optional
```

```
"GuardrailVersion": "1" // optional. Required if GuardrailIdentifier provided.
},
"SageMakerLlmParams": {
"EndpointName": "some-endpoint",
"ModelInputPayloadSchema": {},
"ModelOutputJSONPath": "$."
},
// optional. Passes on arbitrary params to the underlying LLM.
"ModelParams": {
"param1": {
"Value": "value1",
"Type": "string"
},
"param2": {
"Value": 1,
"Type": "integer"
}
},
// optional
"PromptParams": {
"PromptTemplate": "some template",
"UserPromptEditingEnabled": true,
"MaxPromptTemplateLength": 1000,
"MaxInputTextLength": 1000,
"DisambiguationPromptTemplate": "some disambiguation template",
"DisambiguationEnabled": true
},
"Temperature": 1.0, // optional
"Streaming": true, // optional
"RAGEnabled": true, // optional. Must be true if providing KnowledgeBaseParams above.
"Verbose": false // optional
},
"AgentParams": {
"AgentType": "Bedrock",
"BedrockAgentParams": {
"AgentId": "agent-id",
"AgentAliasId": "alias-id",
"EnableTrace": true
}
},
// optional
"AuthenticationParams": {
"AuthenticationProvider": "Cognito",
"CognitoParams": {
```

```
"ExistingUserPoolId": "user-pool-id",
"ExistingUserPoolClientId": "client-id" // optional. If not provided, the solution
will create a client for you in the provided pool
}
}
```

For updates, the structure is the same as above with some caveats:

- The use case name cannot be changed
- A use case can only change security groups and subnets once it has been deployed in a VPC. The VPC itself cannot be changed.
- If a Kendra index was created for you as a knowledge base, you cannot change the configuration of that index (for example, KendraIndexName, QueryCapacityUnits)

Shared Use Case APIs

The following REST API endpoints are available for both Text and Agent use cases:

| REST API | HTTP method | Functionality | Authorized callers |
|-----------|-------------|------------------------|--------------------|
| /details/ | GET | Gets configuration | Amazon Cognito |
| {useCaseC | | details for a specific | authenticated JWT |
| onfigKey} | | use case. | token |

| WebSocket API | Functionality | Authorized callers |
|---------------|--|--|
| /\$connect | Initiate WebSocket connectio n and authenticate user. | Amazon Cognito authentic ated JWT token |
| /\$disconnect | Endpoint called when a WebSocket connection has been disconnected. | Amazon Cognito authentic ated JWT token |

Use Case Details API

The details API endpoint retrieves information about a specific use case:

```
GET /details/{useCaseConfigKey}
```

This endpoint returns the configuration details for a specific use case, including model parameters, knowledge base settings, and other deployment information. It requires an Amazon Cognito authenticated JWT token for authorization.

Text use case

| WebSocket API | Functionality | Authorized callers |
|---------------|---|--|
| /sendMessage | Sends user's chat message to the WebSocket for processin g with the configured LLM experience. | Amazon Cognito authentic ated JWT token |

| REST API | HTTP method | Functionality | Authorized callers |
|---------------------------|-------------|--|--|
| /feedback/ {useCaseId} | POST | Submits user feedback for a specific use case. | Amazon Cognito authenticated JWT token |

sendMessage Payloads

If you're directly integrating with the /sendMessage API, you must adhere to the following request and response payload formats.

Request Payload

```
{
    "action": "sendMessage",
    "question": "the message to send to the api",
    "conversationId": "", // If not provided, a new conversation will be created, with the
    conversationId returned in the response. All subsequent messages in that conversation
    (where history is retained), should provide the conversationId there.
```

```
"promptTemplate": "", // Optional. Overrides the configured prompt
"authToken": "XXXX" // Optional. accessToken from cognito flow. Required for RAG with
RBAC
```

}

| Parameter Name | Туре | Description |
|----------------|-------------------|---|
| action | String | Currently we only support the "sendMessage" action on the WebSocket |
| question | String | The user input to send to the LLM |
| conversationId | String | A UUID identifying the conversation. If not provided, a new conversation will be created, with the conversat ionId returned in the response. All subsequent messages in that conversation (where you wish for history/c ontext to be retained), should provide the conversationId there. |
| promptTemplate | String [Optional] | Overrides the prompt template for this message. If empty or not provided, will default to the prompt set at deployment time. Must have the proper placehold ers specified for the given configuration (i.e. {history} and {input} for non-RAG Sagemaker AI deploymen ts, with the addition of |

| Parameter Name | Туре | Description |
|----------------|-------------------|---|
| | | {context} if using RAG for all deployments. |
| authToken | String [Optional] | accessToken as obtained from the cognito auth flow. This is required when invoking a chat websocket endpoint configured for RAG with Role Based Access Control (RBAC). The cognito:groups claim list in this JWT token is used to control access to documents in the Kendra index. This parameter is not required for non-RAG use cases. It is also not required for RAG use cases that has RBAC disabled. |

Response Payloads

Question Response

The WebSocket API will respond with 1 (if streaming is disabled) or many (if streaming is enabled) JSON objects structured as follows for each query.

```
{
  "data": "some data",
  "conversationId": "id",
}
```

| Parameter Name | Туре | Description |
|----------------|--------|---|
| data | String | A chunk of the response from the LLM if streaming is enabled, or the entire response. If using streaming |

| Parameter Name | Туре | Description |
|----------------|--------|---|
| | | , a response of this format with the data content being <i>END_CONVERSATION</i> will be sent to indicate the end of the response to a single question. |
| conversationId | String | The ID of the conversat ion this sourceDocument response belongs to. |

Source Document Response

If you have configured your RAG use case to return source documents, you will also receive the following payload at the end of every response for each source document used to create the response.

```
{
   "sourceDocument": {
    "excerpt": "some excerpt from the",
    "location": "s3://fake-bucket/test.txt",
    "score": 0.500,
   "document_title": null,
    "document_id": null,
    "additional_attributes": null
    },
    "conversationId": "some-id"
}
```

| Parameter Name | Туре | Description |
|----------------|--------|--|
| excerpt | String | An excerpt from the source document. |
| location | String | Location of the source document. This will depend on the data sources used and |

| Parameter Name | Туре | Description |
|-----------------------|---------------|--|
| | | type of knowledge base, but could be things like s3 URIs or websites. |
| score | Number String | The confidence that the document corresponds to the question asked. This will be a float from 0 to 1 for Bedrock, and a string (e.g. HIGH, LOW, etc.) for Kendra. |
| document_title | String | Title of the returned source document. Only available when using Kendra. |
| document_id | String | ID of the returned source document. Only available when using Kendra. |
| additional_attributes | String | This field will contain all additional attributes on the document as customized on your knowledge base at ingestion. |
| conversationId | String | The ID of the conversat ion this sourceDocument response belongs to. |

Feedback API Payload

Below is an example of a POST payload to the /feedback/{useCaseId} endpoint, which will submit user feedback for a specific use case:

```
"useCaseRecordKey": "12345678-12345678",
"conversationId": "12345678-1234-1234-1234-123456789012",
```

{

```
"messageId": "12345678-1234-1234-1234-123456789012",
"feedback": "positive",
"feedbackReason": ["accurate", "helpful"],
"comment": "This response was very helpful.",
"rephrasedQuery": "What are the key features of Amazon Bedrock?",
"sourceDocuments": [
    "s3://bucket-name/document1.pdf",
    "s3://bucket-name/document2.pdf"
]
}
```

Agent use case

| WebSocket API | Functionality | Authorized callers |
|---------------|---|--|
| /invokeAgent | Sends user's message to the WebSocket for processing with the configured agent. | Amazon Cognito authentic ated JWT token |

invokeAgent Payloads

If you're directly integrating with the /invokeAgent API, you must adhere to the following request and response payload formats.

Request payload

```
{
  "action": "invokeAgent",
  "inputText": "User query to the agent",
  "conversationId": "", // Optional. Empty conversationId implies a new conversation.
  When not provided, a new conversationId will be created and returned with the
  response. All subsequent messages in the same conversation should provide the same
  conversationId (i.e. chat memory/history is maintained).
  "authToken": "XXXX" // Optional. accessToken from cognito flow. If provided, it needs
  to be a valid JWT token associated with the user
}
```

| Parameter name | Туре | Description |
|----------------|------------------|--|
| action | String | We only support the invokeAgent action on the WebSocket. |
| inputText | String | The user input to send to the LLM. |
| conversationId | String[Optional] | A UUID that uniquely identifie s the conversation. If you don't provide this value, the solution creates a new conversation and returns the conversationId in the response. All subsequent messages in that conversat ion (where you want to retain history and context) provide the conversationId there. |
| authToken | String[Optional] | accessToken as obtained from the Amazon Cognito auth flow. This parameter is not required. If you provide it, the JWT token will be validated. This helps make it easier for this solution to be extended. |

Response payloads

Question response

The WebSocket API will respond with one (if streaming is disabled) or many (if streaming is enabled) JSON objects structured as follows for each query.

{

```
"data" "some data",
"conversationId": "id",
}
```

| Parameter name | Туре | Description |
|----------------|--------|---|
| data | String | The response from the agent invocation. |
| conversationId | String | The ID of the conversation. |

Reference

This section includes information about an optional feature for collecting unique metrics for this solution, pointers to related resources, and a list of builders who contributed to this solution.

Supported LLM providers

The solution can integrate with the following LLM providers:

- 1. Amazon Bedrock
 - Documentation: https://aws.amazon.com/bedrock/
 - Supported models:
 - Amazon
 - Nova Lite
 - Nova Micro
 - Nova Pro
 - AI21 Labs
 - Jamba 1.5 Mini
 - Jamba 1.5 Large
 - Anthropic
 - Claude v3 Haiku
 - Claude v3.5 Sonnet
 - Claude v3.7 Sonnet (through the use of inference profiles)
 - Cohere
 - Command R
 - Command R+
 - Deepseek
 - Deepseek-R1 (through the use of inference profiles)
 - Meta
 - Llama 3
 - Llama 3.2 (through the use of inference profiles)
 - Mistral AI

- Mistral 7B Instruct
- Mistral 8x7B Instruct
- Cross-region inference
 - Ability to use inference profiles defined in the same Region as the Deployment dashboard
- 2. Amazon SageMaker Al
 - Documentation: https://aws.amazon.com/sagemaker/
 - Supported models: Text to Text models

For the latest model parameters, best practices, and recommended uses, refer to the documentation from the model providers.

Anonymized data collection

This solution includes an option to send anonymized operational metrics to AWS. We use this data to better understand how customers use this solution and related services and products. When invoked, the following information is collected and sent to AWS:

- Solution ID The AWS solution identifier.
- **UUID** A randomly generated unique identifier for each Generative AI Application Builder on AWS deployment.
- Timestamp The timestamp when the data was collected.
- New Amazon Kendra Index Created Indicates whether a new Amazon Kendra index was created.
- Amazon Kendra Edition Specifies the Amazon Kendra edition selected.
- Usecase Type Indicates the selected use case type: Text or Agent.
- Model Provider The selected model provider, either Bedrock or SageMaker AI.
- **RAG Enabled** Indicates whether Retrieval-Augmented Generation (RAG) functionality is enabled.
- **Type of Knowledge Base** If RAG is enabled, specifies the type of knowledge base selected (Bedrock or Kendra).
- Streaming If the use case is Text, indicates whether the output is streamed.
- Verbose Indicates whether verbose logging is enabled.
- VPC Enabled Indicates whether a VPC is enabled.

- Create VPC Specifies whether a VPC was created by the solution or provided by the business user.
- Use Case Deployment Source Indicates whether the use case was created from the deployment dashboard or as a standalone use case.
- Model ID If Bedrock is the selected model provider, specifies the LLM model ID.
- Inference Profile ID If Bedrock is the selected model provider, specifies the Bedrock inference profile ID used when the model is deployed in a different region.
- **Guardrails Enabled** If Bedrock is the selected model provider, indicates whether Bedrock guardrails are enabled.
- Provisioned Model Enabled Whether a provisioned model is used.
- **Prompt Parameters** Configuration options related to prompt behavior, such as rephrasing, user editing, disambiguation, and input/prompt length limits.
- **Tracing Enabled** If the Agent use case is selected, indicates whether tracing is enabled for the Bedrock Agent.
- Client-Owned User Pool Indicates whether the Cognito user pool is client-owned.
- Feedback Enabled Indicates whether the user has opted to provide feedback (positive or negative) on LLM responses.
- Usage Counts Various metric counts collected from the solution's custom CloudWatch dashboard, which provides application usage analytics. Examples include LLM input and output token counts.

Note: Customer prompts and Model Arn are explicitly excluded and not collected.

AWS owns the data gathered through this survey. Data collection is subject to the <u>AWS Privacy</u> <u>Policy</u>. To opt out of this feature, complete the following steps before launching the AWS CloudFormation template.

- 1. Download the generative-ai-application-builder-on-aws.template <u>AWS</u> CloudFormation template to your local hard drive.
- 2. Open the AWS CloudFormation template with a text editor.
- 3. Modify the AWS CloudFormation template mapping section from:

```
Mappings:
Solution:
```

```
Data:
SendAnonymousUsageData: 'Yes'
```

to:

```
Mappings:
Solution:
Data:
SendAnonymousUsageData: 'No'
```

- 4. Sign in to the AWS CloudFormation console.
- 5. Choose Create stack.
- 6. On the Create stack page, specify template section, select Upload a template file.
- 7. Under **Upload a template file**, choose **Choose file** and select the edited template from your local drive.
- Choose Next and follow the steps in <u>Launch the stack</u> in the Deploy the solution section of this guide.

Note

By default, users can opt out of sending anonymous usage data for **deployment dashboards** only. However, **use cases** within the solution will continue to send anonymous metrics unless explicitly disabled.

To disable anonymous usage data for all use cases:

- 1. Visit the GitHub repository and download the source code for this solution.
- 2. In the file <u>use-case-stack.ts</u>, set the value of SendAnonymousUsageData to No.
- 3. Follow the instructions in the **README.md** to deploy your customized solution.

Contributors

- Tarek Abdunabi
- Majd Arbash
- Mukit Bin Momin
- Michael Connor

- Johny Duval
- Nihit Kasabwala
- Simon Krol
- Michael Lin
- Tim Mekari
- Ibrahim Mohamed
- Omar Radwan Mohsen
- James Nixon
- Jae Shim
- Ajay Swamy
- Reet Takkar
- Dimitri Tchikatilov
- Jason Wreath
- Kamyar Ziabari

Revisions

Publication date: October 2023 (last update: January 2025)

Check the <u>CHANGELOG.md</u> file in the GitHub repository to see all notable changes and updates to the software. The changelog provides a clear record of improvements and fixes for each version.

Notices

Customers are responsible for making their own independent assessment of the information in this document. This document: (a) is for informational purposes only, (b) represents AWS current product offerings and practices, which are subject to change without notice, and (c) does not create any commitments or assurances from AWS and its affiliates, suppliers or licensors. AWS products or services are provided "as is" without warranties, representations, or conditions of any kind, whether express or implied. AWS responsibilities and liabilities to its customers are controlled by AWS agreements, and this document is not part of, nor does it modify, any agreement between AWS and its customers.

Generative AI Application Builder on AWS is licensed under the terms of the <u>Apache License Version</u> <u>2.0</u>.

🔥 Important

Generative AI Application Builder on AWS allows you to build and deploy generative artificial intelligence applications on AWS by engaging the generative AI model of your choice, including third-party generative AI models that you can choose to use that AWS does not own or otherwise have any control over ("Third-Party Generative AI Models"). Your use of the Third-Party Generative AI Models is governed by the terms provided to you by the Third-Party Generative AI Model providers when you acquired your license to use them (for example, their terms of service, license agreement, acceptable use policy, and privacy policy).

You are responsible for ensuring that your use of the Third-Party Generative AI Models comply with the terms governing them, and any laws, rules, regulations, policies, or standards that apply to you.

You are also responsible for making your own independent assessment of the Third-Party Generative AI Models that you use, including their outputs and how Third-Party Generative AI Model providers use any data that might be transmitted to them based on your deployment configuration. AWS does not make any representations, warranties, or guarantees regarding the Third-Party Generative AI Models, which are "Third-Party Content" under your agreement with AWS. Generative AI Application Builder on AWS is offered to you as "AWS Content" under your agreement with AWS.