
Amazon Bedrock

User Guide



Amazon Bedrock: User Guide

Copyright © 2023 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

What is Amazon Bedrock?	1
Access the Amazon Bedrock models	1
Features of Amazon Bedrock	1
Supported models in Amazon Bedrock	2
Supported Regions	3
Amazon Bedrock pricing	3
Set up	4
Add model access	4
Console access	4
Sign up for an AWS account	4
Create an administrative user	5
Set up the AWS CLI	5
Grant programmatic access	5
Model access	7
Manage model access	4
Edit model access	7
Add model access	7
Remove model access	8
Settings	9
Model invocation logging	9
Set up an Amazon S3 destination	9
Set up CloudWatch Logs destination	10
Using the console	11
Using APIs with invocation logging	12
Use the console	13
Base models	13
View and group models	13
Filter and search models	14
View model provider details	14
Request access to base models	14
Text playground	14
Submit a prompt to a text model	15
Interact with the prompt and generated text	15
Re-run a query	16
Use advanced prompts	16
Adjust the model parameters	16
Select a different model	16
Chat playground	16
Submit a prompt to a chat model	17
Update the chat instructions	17
Use advanced prompts	17
Adjust the inference configuration	18
Select a different model	18
Image playground	18
Submit a prompt to a model	18
Adjust the model parameters	19
Examples library	19
Use the API	20
API setup	20
Amazon Bedrock endpoints	20
Setting up the AWS CLI	21
AWS SDK setup	21
Using SageMaker notebooks	21
API operations	23

List the base models	23
Get details about a base model	23
Run inference	24
Tag resources	29
Inference parameters	31
Inference parameter definitions	31
Randomness and diversity	31
Length	32
Repetitions	32
Amazon Titan models	32
Randomness and Diversity	32
Length	33
Model invocation request body field	33
Model invocation response body field	33
Anthropic Claude models	34
Randomness and diversity	34
Length	34
Model invocation request body field	35
Model invocation response body field	35
AI21 Labs Jurassic-2 models	35
Randomness and Diversity	35
Length	36
Repetitions	36
Model invocation request body field	36
Model invocation response body field	37
Cohere models	37
Cohere Command models	38
Cohere Embed models	40
Meta Llama 2 Chat 13B model	42
Randomness and diversity	42
Length	42
Model invocation request body field	42
Model invocation response body field	42
Stability.ai Diffusion models	43
Model invocation request body field	43
Model invocation response body field	44
Prompt engineering guidelines	45
Introduction	45
What is a prompt?	45
Components of a prompt	46
Few-shot prompting vs. zero-shot prompting	47
Prompt template	48
Important notes on using Amazon Bedrock LLMs by API calls	49
What is prompt engineering?	49
General guidelines for Amazon Bedrock LLM users	50
Design your prompt	50
Use hyperparameters	51
Detailed guidelines	51
Optimize prompts for text models on Amazon Bedrock—when the basics aren't good enough	55
Prompt templates and examples for Amazon Bedrock text models	58
Text classification	58
Question-answer, without context	60
Question-answer, with context	62
Summarization	65
Text generation	66
Code generation	68
Mathematics	70

Reasoning/logical thinking	70
Embeddings	73
Custom models	74
Fine-tuning	74
Prepare the datasets	74
Using the console	75
Submit a job	75
Monitor the job	76
Stop a job	77
Analyze the job results	77
Use a fine-tuned model for inference	78
Using the API	78
Set up an IAM role for model customization	79
Submit a job	81
Monitor a job	84
Stop a job	85
Analyze a job	85
Retrieve information about your customized models	86
Guidelines	87
Size of the input training dataset	87
Model size	87
Epochs	87
Batch size	87
Learning rate	88
Learning warmup steps	88
Troubleshooting	88
Permissions issues	88
Data issues	89
Internal error	89
Provisioned Throughput	90
Model IDs	91
Procedures	91
Creating	91
Updating	91
Deleting	91
Running inference	92
Permissions	92
Console procedures	92
View Provisioned Throughput summary	92
Purchase Provisioned Throughput	92
View details of a Provisioned Throughput	93
Edit a Provisioned Throughput	94
Delete a Provisioned Throughput	94
API operations	94
Create Provisioned Throughput	94
Run inference using Provisioned Throughput	95
Update Provisioned Throughput	96
Get Provisioned Throughput	97
Delete Provisioned Throughput	97
List Provisioned Throughput resources	98
Agents for Amazon Bedrock	99
Building a knowledge base	100
Create a service role and configure IAM permissions	101
Set up your data for ingestion	105
Create a knowledge base	108
Manage a knowledge base	111
Add a knowledge base to an agent	111

Building an agent	112
Create a service role and configure IAM permissions	112
Create an agent	115
Edit your agent	119
Test your agent	121
Trace enablement	121
Deploying an agent: versioning and aliases	124
Using the API	125
Invoke your agent	125
How Bedrock Agent works with IAM	126
Identity-based policies	127
Resource-based policies	127
Policy actions	127
Policy resources	128
Policy condition keys	128
ACLs	129
ABAC	129
Temporary credentials	130
Principal permissions	130
Service roles	130
Service-linked roles	130
Identity-based policy examples for Bedrock Agent	131
Tag resources	135
Use the console	135
Use APIs	136
Tag restrictions	136
Security	137
Data protection	137
Data encryption	139
Use VPC	141
Identity and access management	146
Audience	146
Authenticating with identities	147
Managing access using policies	149
How Amazon Bedrock works with IAM	150
Identity-based policy examples	155
Service role	160
Troubleshooting	160
Compliance validation	161
Incident response	162
Resilience	162
Infrastructure security	162
Cross-service confused deputy prevention	163
Configuration and vulnerability analysis in Amazon Bedrock	164
Monitor Amazon Bedrock	165
Monitor with CloudWatch	165
Runtime metrics	165
Logging CloudWatch metrics	166
Use CloudWatch metrics for Amazon Bedrock	166
View Amazon Bedrock metrics	166
Monitor events	167
How it works	167
EventBridge schema	168
Rules and targets	169
Create a rule to handle AWS Bedrock events	169
CloudTrail logs	170
Amazon Bedrock information in CloudTrail	170

Understanding Amazon Bedrock log file entries	171
Abuse detection	173
AWS PrivateLink	174
Considerations	174
Create an interface endpoint	174
Create an endpoint policy	175
Quotas	176
Runtime quotas	176
Model customization quotas	176
Model quotas	177
Fine-tuning quotas	177
Training quotas	177
Provisioned Throughput quotas	177
Document history	179
AWS Glossary	180

What is Amazon Bedrock?

Amazon Bedrock is a fully managed service that makes base models from Amazon and third-party model providers accessible through an API.

Topics

- [Access the Amazon Bedrock models \(p. 1\)](#)
- [Features of Amazon Bedrock \(p. 1\)](#)
- [Supported models in Amazon Bedrock \(p. 2\)](#)
- [Supported Regions \(p. 3\)](#)
- [Amazon Bedrock pricing \(p. 3\)](#)

Access the Amazon Bedrock models

Important

You must request access to a model before you can use it. If you try to use the model (with the API or console) before you have requested access to it, you receive an error message. For more information, see [Model access \(p. 7\)](#).

Features of Amazon Bedrock

With Amazon Bedrock, you can explore the following capabilities:

- **Text playground** – A hands-on text generation application in the AWS Management Console.
- **Image playground** – A hands-on image generation application in the console.
- **Chat playground** – A hands-on conversation generation application in the console.
- **Examples library** – Example use cases to load.
- **Amazon Bedrock API** – Explore with the AWS CLI, or use the API to access the base models.
- **Embeddings** – Use the API to generate embeddings from the Titan Embeddings G1 - Text model.
- **Provisioned Throughput** – Purchase throughput to run inference on models at discounted rates.

Note

Note

Provisioned throughput is currently available for the following models.

Model name	Model ID for provisioned throughput
Titan Text G1 - Express 8K	amazon.titan-text-express-v1:0:8k
Titan Embeddings G1 - Text	amazon.titan-embed-text-v1:2:8k
Anthropic Claude V2 18K	anthropic.claude-v2:0:18k
Anthropic Claude V2 100K	anthropic.claude-v2:0:100k
Anthropic Claude Instant V1 100K	anthropic.claude-instant-v1:2:100K
Stable Diffusion XL 1.0	stability.stable-diffusion-xl-v1:0

Model name	Model ID for provisioned throughput
Meta Llama 2 Chat 13B	meta.llama2-13b-chat-v1:0:4k

Note

Stable Diffusion XL 1.x is only available with provisioned throughput. For more information, see [Provisioned Throughput \(p. 90\)](#). Stable Diffusion XL 0.x is a preview release.

- **Fine-tuning** – Create a training dataset and fine-tune an Amazon Bedrock model.

Note

To use fine-tuning, you must have access to the Amazon Titan Text G1 - Express model. The Amazon Titan Text G1 - Express model is in limited preview. To request access, contact your AWS account manager.

- **Model invocation logging** – Collect invocation logs, model input data, and model output data for all invocations in your AWS account used in Amazon Bedrock.

Note

Model invocation logging is in preview release for Amazon Bedrock and is subject to change.

The following capabilities are in limited preview release. To request access, contact your AWS account manager.

- **Agents for Amazon Bedrock** – Build agents to perform orchestration and carry out tasks for your customers.
- **Knowledge base for Amazon Bedrock** – Draw from data sources to help your agent find information for your customers.

Supported models in Amazon Bedrock

For details about the Amazon Bedrock model providers and their models and model IDs see the **Base models** pages in the Amazon Bedrock console. You can also use the [ListFoundationModels \(p. 23\)](#) API operation to retrieve information about the current list of models.

Amazon Bedrock supports the following models:

- AI21 Labs
 - Jurassic-2 Ultra
 - Jurassic-2 Mid
- Amazon
 - Titan Text G1 - Express
 - Titan Embeddings G1 - Text

Note

The Amazon Titan Text G1 - Express model is in limited preview release. Access will be granted on an ongoing basis.

- Anthropic
 - Claude v1.x
 - Claude v2.x
 - Claude Instant v1.x
- Cohere
 - Command v14.7

- Command Light v14.7
- Embed English v3.0
- Embed Multilingual v3.0
- Meta
 - Llama 2 Chat 13B
- Stability.ai
 - Stable Diffusion XL 0.x
 - Stable Diffusion XL 1.x

Note

Stable Diffusion XL 1.x is only available with provisioned throughput. For more information, see [Provisioned Throughput \(p. 90\)](#). Stable Diffusion XL 0.x is a preview release.

For additional information about these models, see the following links:

- [Anthropic documentation](#).
- [AI21 Studio documentation](#).
- [Cohere documentation](#).
- [Meta documentation](#).
- [Stability.ai documentation](#).

Supported Regions

Amazon Bedrock is available in the following AWS Regions:

- US East (N. Virginia)
- US West (Oregon)
- Asia Pacific (Singapore)
- Asia Pacific (Tokyo)
- Europe (Frankfurt)

Amazon Bedrock pricing

When you sign up for AWS, your AWS account is automatically signed up for all services in AWS, including Amazon Bedrock. However, you are charged only for the services that you use.

To see your bill, go to the Billing and Cost Management Dashboard in the [AWS Billing and Cost Management console](#). To learn more about AWS account billing, see the [AWS Billing User Guide](#). If you have questions concerning AWS billing and AWS accounts, contact [AWS Support](#).

With Amazon Bedrock, you pay to run inference on any of the third-party foundation models. Pricing is based on the volume of input tokens and output tokens, and on whether you have purchased provisioned throughput for the model. For more information, see the [Model providers](#) page in the Amazon Bedrock console. For each model, pricing is listed following the model version. For more information about purchasing Provisioned Throughput, see [Provisioned Throughput \(p. 90\)](#).

For more information, see [Amazon Bedrock Pricing](#).

Set up Amazon Bedrock

Before you use Amazon Bedrock for the first time, complete the following tasks.

Setup tasks

- [Add model access \(p. 4\)](#)
- [Console access \(p. 4\)](#)
- [Sign up for an AWS account \(p. 4\)](#)
- [Create an administrative user \(p. 5\)](#)
- [Set up the AWS CLI \(p. 5\)](#)
- [Grant programmatic access \(p. 5\)](#)

Add model access

Important

You must request access to a model before you can use it. If you try to use the model (with the API or console) before you have requested access to it, you receive an error message. For more information, see [Model access \(p. 7\)](#).

Console access

To access the Amazon Bedrock console and playground:

1. Sign in to your AWS account.
2. Navigate to: [Amazon Bedrock console](#)

You can also access the Amazon Bedrock console in US East (N. Virginia), US West (Oregon) and Asia Pacific (Singapore).

Sign up for an AWS account

If you do not have an AWS account, complete the following steps to create one.

To sign up for an AWS account

1. Open <https://portal.aws.amazon.com/billing/signup>.
2. Follow the online instructions.

Part of the sign-up procedure involves receiving a phone call and entering a verification code on the phone keypad.

When you sign up for an AWS account, an *AWS account root user* is created. The root user has access to all AWS services and resources in the account. As a security best practice, [assign administrative access to an administrative user](#), and use only the root user to perform [tasks that require root user access](#).

AWS sends you a confirmation email after the sign-up process is complete. At any time, you can view your current account activity and manage your account by going to <https://aws.amazon.com/> and choosing **My Account**.

Create an administrative user

After you sign up for an AWS account, create an administrative user so that you don't use the root user for everyday tasks.

Secure your AWS account root user

1. Sign in to the [AWS Management Console](#) as the account owner by choosing **Root user** and entering your AWS account email address. On the next page, enter your password.

For help signing in by using root user, see [Signing in as the root user](#) in the *AWS Sign-In User Guide*.

2. Turn on multi-factor authentication (MFA) for your root user.

For instructions, see [Enable a virtual MFA device for your AWS account root user \(console\)](#) in the *IAM User Guide*.

Create an administrative user

- For your daily administrative tasks, grant administrative access to an administrative user in AWS IAM Identity Center.

For instructions, see [Getting started](#) in the *AWS IAM Identity Center User Guide*.

Sign in as the administrative user

- To sign in with your IAM Identity Center user, use the sign-in URL that was sent to your email address when you created the IAM Identity Center user.

For help signing in using an IAM Identity Center user, see [Signing in to the AWS access portal](#) in the *AWS Sign-In User Guide*.

Set up the AWS CLI

You don't need the AWS Command Line Interface (AWS CLI) to use Amazon Bedrock. If you prefer, you can skip this step and set up the AWS CLI later.

To install and configure the AWS CLI

1. Install the AWS CLI. For instructions, see [Installing or updating the latest version of the AWS Command Line Interface](#) in the *AWS Command Line Interface User Guide*.
2. Configure the AWS CLI. For instructions, see [Configuring the AWS Command Line Interface](#) in the *AWS Command Line Interface User Guide*.

Grant programmatic access

Users need programmatic access if they want to interact with AWS outside of the AWS Management Console. The way to grant programmatic access depends on the type of user that's accessing AWS.

To grant users programmatic access, choose one of the following options.

Which user needs programmatic access?	To	By
Workforce identity (Users managed in IAM Identity Center)	Use temporary credentials to sign programmatic requests to the AWS CLI, AWS SDKs, or AWS APIs.	Following the instructions for the interface that you want to use. <ul style="list-style-type: none">• For the AWS CLI, see Configuring the AWS CLI to use AWS IAM Identity Center in the <i>AWS Command Line Interface User Guide</i>.• For AWS SDKs, tools, and AWS APIs, see IAM Identity Center authentication in the <i>AWS SDKs and Tools Reference Guide</i>.
IAM	Use temporary credentials to sign programmatic requests to the AWS CLI, AWS SDKs, or AWS APIs.	Following the instructions in Using temporary credentials with AWS resources in the <i>IAM User Guide</i> .
IAM	(Not recommended) Use long-term credentials to sign programmatic requests to the AWS CLI, AWS SDKs, or AWS APIs.	Following the instructions for the interface that you want to use. <ul style="list-style-type: none">• For the AWS CLI, see Authenticating using IAM user credentials in the <i>AWS Command Line Interface User Guide</i>.• For AWS SDKs and tools, see Authenticate using long-term credentials in the <i>AWS SDKs and Tools Reference Guide</i>.• For AWS APIs, see Managing access keys for IAM users in the <i>IAM User Guide</i>.

Model access

Amazon Bedrock users need to request access to models before they are available for use. If you want to add additional models for text, chat, and image generation, you need to request access to models in Amazon Bedrock. To request access to additional models, select the **Model access** link in the left side navigation panel in the Amazon Bedrock console.

Topics

- [Manage model access \(p. 4\)](#)
- [Edit model access \(p. 7\)](#)
- [Add model access \(p. 7\)](#)
- [Remove model access \(p. 8\)](#)

Manage model access

The account does not have access to models by default. Admin users with IAM access permissions can add access to specific models using the model access page. After the admin adds access to models, those models are available for all users of the account.

Charges are accrued when the models are used in Amazon Bedrock. You can review the **End User License Agreement** (EULA) by selecting the link for each model.

To add or remove model access, select **Manage model access**.

Edit model access

You can request access to models to use them in Amazon Bedrock. This page lists base models of LLMs to text, image, and embedding models. You can review your status to see which models you have access to. The modality lists whether the model is used for text, embedding, or images.

Add model access

You can add access to a model in Amazon Bedrock with the following steps:

1. Open the Amazon Bedrock console at [Amazon Bedrock console](#)
2. Go to the **Model access** link in the left side navigation panel in Amazon Bedrock, or go to the **Edit model access** page.
3. Select the check box next to the model you want to add access to. For Anthropic models, you must also request access when you click the **Request access** button. Models are not available as a default setting in Amazon Bedrock.
4. Select **Confirm** to add access to any third party models through Amazon Marketplace. Note: Your use of Amazon Bedrock and its models is subject to the seller's pricing terms, EULA and the Amazon Bedrock service terms.
5. Select the **Save Changes** button in the lower right corner of the page. It may take several minutes to save changes to the **Model access** page.
6. Models will show as **Available** on the **Model access** page under **Access** status, if access is granted.

Some users may not have IAM permissions to add access to models. A banner message will appear if you try to add access to models and you are a non-admin user on your account. You will need to contact your account administrator to request that they add access to the models before you can use them in Amazon Bedrock.

When you have access to the model, you can select it using the **View model access** button or the **Edit model access** page. Accounts do not have access to Amazon models by default.

Remove model access

When you are using Amazon Bedrock, you may decide to use only certain models for your work, and to remove access to models you are no longer using. You can remove access to a model in Amazon Bedrock with the following steps:

1. Open the Amazon Bedrock console at <https://us-east-1.console.aws.amazon.com/bedrock/home?region=us-east-1#/>.
2. Go to the **Model access** link in the left side navigation panel in Amazon Bedrock, or go to the **Edit model access** page.
3. Deselect the check box next to the model you want to remove access.
4. Select the **Save changes** button at the bottom right corner of the page.
5. You will be prompted to confirm you want to remove access to models. You cannot remove access to models if there are resources using those models.
6. Users must remove the associations to the model listed in the information window in order to remove access. You can select the link next to the resource name to go to the location, and remove access to the model for that resource. Repeat this step for each item listed as a dependent resource.
7. Once you have removed access all the resources listed, select **Okay** to complete the remove access step.
8. Once completed, you will see a banner message that confirms the action was completed successfully.

Settings

You can manage your account level settings for Amazon Bedrock in the **Settings** page. The settings include data logs, data permissions, and model access.

To access settings, go to the bottom of the left-side navigation pane in Amazon Bedrock, and select **Settings**.

Topics

- [Model invocation logging \(p. 9\)](#)

Model invocation logging

Model invocation logging can be used to collect invocation logs, model input data, and model output data for all invocations in your AWS account used in Amazon Bedrock. By default, logging is disabled.

With invocation logging, you can collect the full request data, response data, and metadata associated with all calls performed in your account. Logging can be configured to provide the destination resources where the log data will be published. Supported destinations include Amazon CloudWatch Logs and Amazon Simple Storage Service (Amazon S3). Only destinations from the same account and region are supported.

Before you can enable invocation logging, you need to set up an Amazon S3 or CloudWatch Logs destination. You can enable invocation logging through either the console or the API.

Topics

- [Set up an Amazon S3 destination \(p. 9\)](#)
- [Set up CloudWatch Logs destination \(p. 10\)](#)
- [Using the console \(p. 11\)](#)
- [Using APIs with invocation logging \(p. 12\)](#)

Set up an Amazon S3 destination

You can set up an S3 destination for logging in Amazon Bedrock with these steps:

1. Create an S3 bucket where the logs will be delivered.
2. Add a bucket policy to it like the one below (Replace values for *accountId*, *region*, *bucketName*, and optionally *prefix*):

Note

A bucket policy is automatically attached to the bucket on your behalf when you configure logging with the permissions `S3:GetBucketPolicy` and `S3:PutBucketPolicy`.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "AmazonBedrockLogsWrite",
      "Effect": "Allow",
      "Principal": {
        "Service": "bedrock.amazonaws.com"
      },
    },
  ],
}
```



```
"Action": [
  "s3:PutObject"
],
"Resource": [
  "arn:aws:s3:::bucketName/prefix/AWSLogs/accountId/BedrockModelInvocationLogs/*"
],
"Condition": {
  "StringEquals": {
    "aws:SourceAccount": "accountId"
  },
  "ArnLike": {
    "aws:SourceArn": "arn:aws:bedrock:region:accountId:*"
  }
}
}
```

3. (Optional) If configuring SSE-KMS on the bucket, add the below policy on the KMS key:

```
{
  "Effect": "Allow",
  "Principal": {
    "Service": "bedrock.amazonaws.com"
  },
  "Action": "kms:GenerateDataKey",
  "Resource": "*",
  "Condition": {
    "StringEquals": {
      "aws:SourceAccount": "accountId"
    },
    "ArnLike": {
      "aws:SourceArn": "arn:aws:bedrock:region:accountId:*"
    }
  }
}
```

For more information on S3 SSE-KMS configurations, see [Specifying KMS Encryption](#).

Note

The bucket ACL must be disabled in order for the bucket policy to take effect. For more information, see [Disabling ACLs for all new buckets and enforcing Object Ownership](#).

Set up CloudWatch Logs destination

You can set up a Amazon CloudWatch Logs destination for logging in Amazon Bedrock with the following steps:

1. Create a CloudWatch log group where the logs will be published.
2. Create an IAM role with the following permissions for CloudWatch Logs.

Trusted entity:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "bedrock.amazonaws.com"
      }
    }
  ]
}
```

```
    },
    "Action": "sts:AssumeRole",
    "Condition": {
      "StringEquals": {
        "aws:SourceAccount": "accountId"
      },
      "ArnLike": {
        "aws:SourceArn": "arn:aws:bedrock:region:accountId:*"
      }
    }
  }
]
```

Role policy:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "logs:CreateLogStream",
        "logs:PutLogEvents"
      ],
      "Resource": "arn:aws:logs:region:accountId:log-group:logGroupName:log-  
stream:aws/bedrock/modelinvocations"
    }
  ]
}
```

For more information on setting up SSE for CloudWatch Logs, see [Encrypt log data in CloudWatch Logs using AWS Key Management Service](#).

Using the console

To enable model invocation logging, drag the slider button next to the **Logging** toggle switch in the **Settings** page. Additional configuration settings for logging will appear on the panel.

Choose which data requests and responses you want to publish to the logs. You can choose any combination of the following output options:

- Text
- Image
- Embedding

Choose where to publish the logs:

- Amazon S3 only
- CloudWatch Logs only
- Both Amazon S3 and CloudWatch Logs

Amazon S3 and CloudWatch Logs destinations are supported for invocation logs, and small input and output data. For large input and output data or binary image outputs, only Amazon S3 is supported. The following details summarize how the data will be represented in the target location.

- **S3 destination** — Gzipped JSON files, each containing a batch of invocation log records, are delivered to the specified S3 bucket. Similar to a CloudWatch Logs event, each record will contain the invocation metadata, and input and output JSON bodies of up to 100 KB in size. Binary data or JSON bodies larger than 100 KB will be uploaded as individual objects in the specified Amazon S3 bucket under the data prefix. The data can be queried using Amazon S3 Select and Amazon Athena, and can be catalogued for ETL using AWS Glue. The data can be loaded into OpenSearch service, or be processed by any Amazon EventBridge targets.
- **CloudWatch Logs destination** — JSON invocation log events are delivered to a specified log group in CloudWatch Logs. The log event contains the invocation metadata, and input and output JSON bodies of up to 100 KB in size. If an Amazon S3 location for large data delivery is provided, binary data or JSON bodies larger than 100 KB will be uploaded to the Amazon S3 bucket under the data prefix instead. data can be queried using CloudWatch Logs Insights, and can be further streamed to various services in real-time using CloudWatch Logs.

Using APIs with invocation logging

Model invocation logging can be configured using the following APIs:

- `PutModelInvocationLoggingConfiguration`
- `GetModelInvocationLoggingConfiguration`
- `DeleteModelInvocationLoggingConfiguration`

For more information on how to use APIs with invocation logging, see the [Bedrock API Guide](#).

Use the Amazon Bedrock console

The Amazon Bedrock console provides the following capabilities and features:

- From the **Base models** page under **Foundation models**, you can view the available models and group them by various attributes. You can also filter the model view, search for models, and view information about the model providers.
- After you select a model, you can open the model to experiment with it in the **Text playground**, **Image playground**, or **Chat playground**.
- Amazon Bedrock provides example prompts for each of the supported models. The **Examples** page displays up to 20 examples for each model provider. You can filter the list of examples using one or more attributes.

Important

You must request access to a model before you can use it. If you try to use the model (with the API or console) before you have requested access to it, you receive an error message. For more information, see [Model access \(p. 7\)](#).

Topics

- [Base models \(p. 13\)](#)
- [Text playground \(p. 14\)](#)
- [Chat playground \(p. 16\)](#)
- [Image playground \(p. 18\)](#)
- [Examples library \(p. 19\)](#)

Base models

Amazon Bedrock supports base models from Amazon and third-party model providers.

From the Amazon Bedrock console, you can view the available models and group them by various attributes. You can also filter the model view, search for models, and view information about the model providers.

After you select a model, you can choose **Open in playground** to experiment with the model.

Topics

- [View and group models \(p. 13\)](#)
- [Filter and search models \(p. 14\)](#)
- [View model provider details \(p. 14\)](#)
- [Request access to base models \(p. 14\)](#)

View and group models

You can list models and view model details using the Amazon Bedrock console.

1. From the Amazon Bedrock console, choose **Base models** under **Foundation models**.
2. You can view the models in a list view or a card view and group them by model family, by model provider, or by model modality (text, image, or embedding).

Filter and search models

You can filter models by model name, provider, by modality, or by model attributes.

You can set multiple filters. You can also start by entering text in the text box. The console displays all the possible filter values that match the text.

1. From the Amazon Bedrock console, choose **Base Models** under **Foundation models**.
2. To set a filter, choose the **Find resource** text box to display the filter properties.
3. Select the desired property, the operator for the filter, then enter the value to filter by.

View model provider details

To view model provider details using the Amazon Bedrock console, choose **Model providers** under **Foundation models**.

As an alternative, from the **Base models** page, select **View provider details** under a provider name.

Choose one of the tabs near the top of the **Model providers** page to view details about that model provider. The page displays the following information about the provider:

- **Provider overview** – An overview description about the provider.
- **Models** – A tabbed list of available models. When you choose one of the models, you can open the model in the playground.
- **Content limitations** – Details about the content policy for the provider.

Request access to base models

Important

You must request access to a model before you can use it. If you try to use the model (with the API or console) before you have requested access to it, you receive an error message. For more information, see [Model access \(p. 7\)](#).

Text playground

From the Amazon Bedrock console, choose **Playgrounds** and then **Text** to view the text playground.

You can also navigate directly to the playground when you choose a model from a model details page or the examples page.

Topics

- [Submit a prompt to a text model \(p. 15\)](#)
- [Interact with the prompt and generated text \(p. 15\)](#)
- [Re-run a query \(p. 16\)](#)
- [Use advanced prompts \(p. 16\)](#)

- [Adjust the model parameters \(p. 16\)](#)
- [Select a different model \(p. 16\)](#)

Submit a prompt to a text model

The following procedure shows how to use the text playground to submit a prompt to a text model.

1. Open the Amazon Bedrock console.
2. From the left-side menu, choose **Text** under **Playgrounds**.
3. Set up the text playground.
 - a. You can toggle on the following options in the top-right corner:
 - **Streaming** – Generate text in real-time. This option is only available for models that support streaming.
 - b. From the dropdown menu above the text panel, select a model provider. After you select a provider, select a model from the model dropdown menu. The dropdown menu lists the available models from this provider. You can also select a model that you have customized (for more information, see [Use a fine-tuned model for inference \(p. 78\)](#)). If you select a customized model, you must have set up Provisioned Throughput for it beforehand.
 - c. (Optional) Adjust the **Inference configurations**. The default settings are optimal for most prompts, but might not work well for your use case. You can restore the default settings by choosing the **Reset** button. For more information, see [Inference parameters for foundation models \(p. 31\)](#).
4. Enter your own prompt into the text field. A prompt is a natural language command, such as **write a blog post about computers**.
5. Under the input panel, choose **Run** to generate a text response.

Note

If the response violates the content moderation policy, Amazon Bedrock does not display it. If you have turned on streaming, Amazon Bedrock clears the entire response if it generates content that violates the policy. For more details, navigate to the Amazon Bedrock console, select **Providers**, and read the text under the **Content limitations** section.

Interact with the prompt and generated text

The generated text appears under your prompt in green text. Amazon Bedrock returns common markdown and tables in a rich text view. If you turn on streaming, you can't interact with the text until Bedrock finishes generating it. Interact with the prompt that you provided and the generated response in the following ways:

1. If the response is not to your satisfaction, you can edit the text directly. To interact with tables, code, and rich text in the response, use the icons at the upper-right of these objects:
 - **Edit** – Select the pencil icon to revert the table or code to plain text format for editing.
 - **Copy** – Select the copy icon to copy the object in plain text format.
2. Use the icons next to the **Run** button to undo and redo edits you have made.
3. Use the icons under the generated text to carry out the following actions:
 - **Copy** – Select the copy icon to copy the response as plain text.
 - **Delete** – Select the trash can icon to delete the latest model response.
4. After you experiment with the models, choose **View API request** to view the equivalent API request. You can copy this code into your application.

Re-run a query

You can modify the prompt and then re-run the query. Edit the prompt or response text, and then choose **Run**. After you edit the response text, the console changes the text color to black. The console displays the new response as green text.

Use advanced prompts

A prompt can consist of a simple command without any additional context, such as the example **write a blog post about computers**. These prompts are known as *zero-shot prompts*.

To improve the precision of the output, you can provide a few examples of the type of output you desire in the prompt. This is known as *few-shot prompting* (or *few-shot learning*).

In the following example, the first two sequences include the input text and the desired response. The last sequence is asking the model to answer in a similar format as the first two sequences.

```
Input: I love going to the mall  
Sentiment: Happy  
Input: I don't like going to the dentist.  
Sentiment: Sad  
Input: I enjoy the park  
Sentiment:
```

In general, large language models work to predict the next logical word in a given sequence. How you prompt the model influences the model to produce a specific style of result.

Adjust the model parameters

When you choose a provider and model in the playground, the console loads the inference configuration parameters that apply to the selected model and sets their default values.

When you adjust the inference parameter settings, the model generates text that can differ in context, style, relevance, and length. Changing the parameter settings usually involves some trial and error to get the desired results.

For information about the inference parameters that each model supports, see [Inference parameters for foundation models \(p. 31\)](#).

Select a different model

You can select a different model provider and model from the dropdown menus at the top of the page. The console sets the inference configuration parameters to the default values for the model that you select.

Chat playground

From the Amazon Bedrock console, choose **Playgrounds** and then **Chat** to view the chat playground.

You can also navigate directly to the playground when you choose a model from a model details page or the examples page.

Topics

- [Submit a prompt to a chat model \(p. 17\)](#)

- [Update the chat instructions \(p. 17\)](#)
- [Use advanced prompts \(p. 17\)](#)
- [Adjust the inference configuration \(p. 18\)](#)
- [Select a different model \(p. 18\)](#)

Submit a prompt to a chat model

The following procedure shows how to use the chat playground to submit a prompt to a chat model:

1. Open the Amazon Bedrock console.
2. From the left-side menu, choose **Chat** under **Playgrounds**.
3. Set up the chat playground.
 - a. From the dropdown menu above the text panel, select a model provider. After you select a provider, select a model from the model dropdown menu. The dropdown menu lists the available models from this provider. You can also select a model that you have customized (for more information, [Use a fine-tuned model for inference \(p. 78\)](#)). If you select a customized model, you must have set up Provisioned Throughput for it beforehand.
 - b. You can toggle on the following options in the top-right corner:
 - **Streaming** – Generate chat in real-time. This option isn't available if the model selected doesn't support streaming.
4. Select **Add instructions** to , then choose **Update**. In the **Response** panel, the console displays the response from the model.
5. Enter your own prompt into the text field. A prompt is a natural language phrase or command, such as **Tell me about the best restaurants to visit in Seattle.**

Note

If the response violates the content moderation policy, Amazon Bedrock doesn't display it. If you have turned on streaming, Amazon Bedrock clears the entire response if it generates content that violates the policy. For more details, navigate to the Amazon Bedrock console, select **Providers**, and read the text under the **Content limitations** section.

Update the chat instructions

You can modify the instructions and then update the chat model with the new instructions. Choose **Update** to edit the instructions to the model, and then choose **Confirm** after you make changes.

For example, you can add instructions to guide the personality and tone of the model responses. At any time, you can adjust the instructions to the model and re-run your prompt. This is useful when designing the conversation experience that users will have with the chat model.

Use advanced prompts

A prompt can consist of a simple command without any additional context, such as the example **Tell me about the best car to buy**. These prompts are known as *zero-shot prompts*.

To improve your prompt, you can provide a few examples of the type of output you desire. This is known as *few-shot prompting* (or *few-shot learning*).

For example, you can use a more complex prompt such as **Tell me about the best car to buy. Give me a numbered list with the top 5 options. Include information on pricing, fuel efficiency, consumer ratings, and comfort of design. List the sources used**

for each item in the list. hybrid or electric technology used. Order the list so that the most desirable option is listed first.

Adjust the inference configuration

When you choose a provider and model in the playground, the console loads the inference configuration parameters that apply to the selected model and sets their default values.

When you adjust the inference parameter settings, the model generates text that can differ in context, style, relevance, and length. Changing the parameter settings usually involves some trial and error to get the desired results. To adjust the inference parameters, select the **Update** link in the bottom right of the chat window, and then choose **Confirm** after you make adjustments.

For information about the inference parameters that each model supports, see [Inference parameters for foundation models \(p. 31\)](#).

Select a different model

You can select a different model provider and model from the dropdown menus at the top of the page. The console sets the inference configuration parameters to the default values for the model that you select.

Image playground

From the Amazon Bedrock console, choose **Playgrounds** and then **Image** to view the image playground.

You can also navigate directly to the image playground after you choose a model from a model details page or the examples page.

Topics

- [Submit a prompt to a model \(p. 18\)](#)
- [Adjust the model parameters \(p. 19\)](#)

Submit a prompt to a model

The following procedure shows how to use the image playground to submit a prompt to a text-to-image model.

1. Open the Amazon Bedrock console.
2. From the left-side menu, choose **Image** under **Playgrounds**.
3. From the dropdown menu above the text panel, select a model provider. After you select a provider, select a model from the model dropdown menu. The dropdown menu lists the available models from this provider.
4. Enter text into the **Prompt** field, or use one of the example prompts. A prompt is a natural language command, such as **Draw a picture of a computer**.
5. (Optional) Adjust the **Inference configurations**. The default settings are generally optimal but might not work for your use case. You can restore the default settings by choosing the **Reset** button.
6. Under the **Response** field, choose **Run** to generate the image. The generated image appears in the **Response** field, under your prompt.
7. (Optional) Choose **Download image** to save the image to your local machine.
8. (Optional) Choose **Save** to save the current prompt as a favorite in your history items.

9. After you experiment with the model, choose **View API request** to view the equivalent API request. You can copy this code into your application.

Adjust the model parameters

When you choose a provider and model in the playground, the console loads the inference configuration parameters that apply to the selected model and sets their default values.

When you adjust the inference parameter settings, the model generates a new image that can differ in style. Changing the parameter settings usually involves some trial and error to get the desired results.

For information about the inference parameters that each model supports, see [Inference parameters for foundation models \(p. 31\)](#).

Examples library

Amazon Bedrock provides example prompts for each of the supported models. From the console, choose **Examples** to view the available examples.

The Amazon Bedrock console displays up to 20 examples for each model provider. You can filter the list of examples using one or more of the following attributes:

- Modality (text, image, or embedding)
- Provider
- Model name
- Category

You can also filter by entering text in the text box. The console displays all the possible filter values that match the text.

When you select a specific example, the Amazon Bedrock console displays the following information about the example:

- A description of what the example accomplishes.
- The model name (and model provider) where the example runs.
- The example prompt and the expected response.
- The inference configuration parameter settings for the example.
- The API request that runs the example.

To run the example, choose **Open in playground**.

Use the Amazon Bedrock API

This section describes how to set up your environment to make Amazon Bedrock API calls and provides examples of common use-cases. You can access the Amazon Bedrock API using the AWS Command Line Interface (AWS CLI), an AWS SDK, or a SageMaker Notebook.

You can make API calls to Amazon Bedrock through the following SDKs:

- [AWS SDK for C++](#)
- [AWS SDK for Go](#)
- [AWS SDK for Java](#)
- [AWS SDK for JavaScript](#)
- [AWS SDK for .NET](#)
- [AWS SDK for Python \(Boto3\)](#)
- [AWS SDK for Ruby](#)

For details about the API operations and parameters, see the [Amazon Bedrock API Reference](#).

The following resources provide additional information about the Amazon Bedrock API.

- *AWS General Reference*
 - [Amazon Bedrock endpoints and quotas](#)
- *AWS Command Line Interface*
 - [Amazon Bedrock CLI commands](#)
 - [Amazon Bedrock Runtime CLI commands](#)

Topics

- [Setting up the Amazon Bedrock API \(p. 20\)](#)
- [Amazon Bedrock API operations \(p. 23\)](#)

Setting up the Amazon Bedrock API

Important

You must request access to a model before you can use it. If you try to use the model (with the API or console) before you have requested access to it, you receive an error message. For more information, see [Model access \(p. 7\)](#).

You can access the Amazon Bedrock API using the AWS CLI, an AWS SDK, or a SageMaker Notebook.

Amazon Bedrock endpoints

To connect programmatically to an AWS service, you use an endpoint. Refer to the [Amazon Bedrock endpoints and quotas](#) chapter in the AWS General Reference for information about the endpoints that you can use for Amazon Bedrock.

Setting up the AWS CLI

1. If you plan to use the CLI, install and configure the AWS CLI. See [the section called “Set up the AWS CLI” \(p. 5\)](#).
2. Configure your AWS credentials using the `aws configure` CLI command (see [Configure the AWS CLI](#)).

Refer to the following references for AWS CLI commands and operations:

- [Amazon Bedrock CLI commands](#)
- [Amazon Bedrock Runtime CLI commands](#)

Setting up an AWS SDK

AWS software development kits (SDKs) are available for many popular programming languages. Each SDK provides an API, code examples, and documentation that make it easier for developers to build applications in their preferred language. Currently, you can make Amazon Bedrock API calls through the following SDKs.

- [AWS SDK for C++](#)
- [AWS SDK for Go](#)
- [AWS SDK for Java](#)
- [AWS SDK for JavaScript](#)
- [AWS SDK for .NET](#)
- [AWS SDK for Python \(Boto3\)](#)
- [AWS SDK for Ruby](#)

For code examples, find a method in [Amazon Bedrock API operations \(p. 23\)](#) and select the tab that corresponds to the programming language of your choice.

Using SageMaker notebooks

You can use the SDK for Python (Boto3) to invoke Amazon Bedrock API operations from a SageMaker notebook.

Prerequisites

Note the following prerequisites:

1. Request Amazon Bedrock access for the AWS account that hosts the notebook.
2. Use the console to accept the Amazon Bedrock terms and conditions.

Configure the SageMaker role

Add Amazon Bedrock permissions to the IAM role for this SageMaker notebook.

From the IAM console, perform these steps:

1. Choose the IAM role, then choose **Add Permissions** and select **Create Inline Policies** from the dropdown list.

2. Include the following permission.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": "bedrock:*",
      "Resource": "*"
    }
  ]
}
```

Add the following permissions to the trust relationships.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "bedrock.amazonaws.com"
      },
      "Action": "sts:AssumeRole"
    },
    {
      "Sid": "",
      "Effect": "Allow",
      "Principal": {
        "Service": "sagemaker.amazonaws.com"
      },
      "Action": "sts:AssumeRole"
    }
  ]
}
```

Test the Runtime setup

Add the following code to your notebook and run the code.

```
import boto3
import json
bedrock = boto3.client(service_name='bedrock-runtime')

body = json.dumps({
  "prompt": "\n\nHuman:explain black holes to 8th graders\n\nAssistant:",
  "max_tokens_to_sample": 300,
  "temperature": 0.1,
  "top_p": 0.9,
})

modelId = 'anthropic.claude-v2'
accept = 'application/json'
contentType = 'application/json'

response = bedrock.invoke_model(body=body, modelId=modelId, accept=accept,
  contentType=contentType)

response_body = json.loads(response.get('body').read())
# text
```

```
print(response_body.get('completion'))
```

Test the Amazon Bedrock setup

Add the following code to your notebook and run the code.

```
import boto3
bedrock = boto3.client(service_name='bedrock')

bedrock.get_foundation_model(modelIdentifier='anthropic.claude-v2')
```

Amazon Bedrock API operations

Topics

- [List the base models \(p. 23\)](#)
- [Get details about a base model \(p. 23\)](#)
- [Run inference \(p. 24\)](#)
- [Tag resources \(p. 29\)](#)

List the base models

Use the [ListFoundationModels](#) operation to retrieve information about the base models, such as the model ID that you need to perform inference with the InvokeMethod operation. See the following code examples:

AWS CLI

List the base models using the following command in the AWS CLI:

```
aws bedrock list-foundation-models
```

Python (Boto)

The following example demonstrates how to list the base models using Python:

```
import boto3
bedrock = boto3.client(service_name='bedrock')

bedrock.list_foundation_models()
```

Get details about a base model

Use the [GetFoundationModel](#) operation to retrieve detailed information about the specified base model. See the following code examples:

AWS CLI

Get information about a base model using the following command in the AWS CLI:

```
aws bedrock get-foundation-model --model-identifier anthropic.claude-v2
```

Python (Boto)

The following example demonstrates how to list the base models using Python:

```
import boto3
bedrock = boto3.client(service_name='bedrock')

bedrock.get_foundation_model(modelIdentifier='anthropic.claude-v2')
```

Run inference

Use the [InvokeModel](#) operation to run inference on models. You use the `modelId` field to specify the model that you want to use. The method for finding the `modelId` depends on the type of model you use. To run inference on a custom model, first purchase Provisioned Throughput for it (for more information, see [Provisioned Throughput \(p. 90\)](#)).

- **Base model** – Call [ListFoundationModels](#) to find the model ARN. For an example, see [List the base models \(p. 23\)](#). To read a list of the model IDs for the currently available base models, see [Base model IDs \(p. 28\)](#). You can also get the model ID for a foundation model from the JSON examples in the Amazon Bedrock console.
- **Model with Provisioned Throughput** – If you have created a Provisioned Throughput for a foundation or custom model, call [ListProvisionedModelThroughputs](#) and find the `modelArn` in the response. You can also find the model ARN in the **Model details** when you select a model in the **Provisioned throughput** section in the console.

Each base model has its own parameters that you set in the body field. For more information, see [Inference parameters for foundation models \(p. 31\)](#).

To run inference with streaming, use the [InvokeModelWithResponseStream](#) operation. Pick a base model that supports streaming or a model that you created from a base model that supports streaming.

Topics

- [Running inference on a model \(p. 24\)](#)
- [Base model IDs \(p. 28\)](#)

Running inference on a model

The following examples show how to run inference on a model with [InvokeModel](#) and, with Python, run inference with streaming with the [InvokeModelWithResponseStream](#) operation.

Note

The AWS CLI does not support streaming.

For information about the parameters each model supports, see [Inference parameters for foundation models \(p. 31\)](#). For information about writing prompts, see [Prompt engineering guidelines \(p. 45\)](#).

AWS CLI

The following example shows how to generate text with the AWS CLI using the prompt *story of two dogs* and the Anthropic Claude V2 model. The example returns up to *300* tokens in the response and saves the response to the file *output.txt*:

```
aws bedrock-runtime invoke-model \
  --model-id anthropic.claude-v2 \
```

```
--body "{\"prompt\": \"\\n\\nHuman: story of two dogs\\n\\nAssistant:\",  
\\\"max_tokens_to_sample\\\" : 300}\" \\  
--cli-binary-format raw-in-base64-out \\  
invoke-model-output.txt
```

The following example shows how to call the *Llama 2 Chat 13B* model.

```
aws bedrock-runtime invoke-model \  
  --region us-east-1 \  
  --model-id meta.llama2-13b-chat-v1 \  
  --body "{\"prompt\": \"What is the average lifespan of a Llama?\", \"max_gen_len  
\\\" : 128, \"temperature\": 0.1, \"top_p\": 0.9}\"" \  
  invoke-model-output.txt
```

Python (Boto)

The following example shows how to generate text with Python using the prompt *explain black holes to 8th graders* and the Anthropic Claude V2 model:

```
import boto3  
import json  
brt = boto3.client(service_name='bedrock-runtime')  
  
body = json.dumps({  
    "prompt": "\\n\\nHuman: explain black holes to 8th graders\\n\\nAssistant:",  
    "max_tokens_to_sample": 300,  
    "temperature": 0.1,  
    "top_p": 0.9,  
})  
  
modelId = 'anthropic.claude-v2'  
accept = 'application/json'  
contentType = 'application/json'  
  
response = brt.invoke_model(body=body, modelId=modelId, accept=accept,  
    contentType=contentType)  
  
response_body = json.loads(response.get('body').read())  
  
# text  
print(response_body.get('completion'))
```

The following example shows how to generate streaming text with Python using the prompt *write an essay for living on mars in 1000 words* and the Anthropic Claude V2 model:

```
import boto3  
import json  
  
brt = boto3.client(service_name='bedrock-runtime')  
  
body = json.dumps({  
    'prompt': '\\n\\nHuman: write an essay for living on mars in 1000 words\\n  
\\nAssistant:',  
    'max_tokens_to_sample': 100  
})  
  
response = brt.invoke_model_with_response_stream(  
    modelId='anthropic.claude-v2',  
    body=body  
)
```



```
stream = response.get('body')
if stream:
    for event in stream:
        chunk = event.get('chunk')
        if chunk:
            print(json.loads(chunk.get('bytes')).decode())
```

Base model inference examples

The following Python (Boto) examples show how you can perform inference with the [InvokeModel](#) operation on different Amazon Bedrock base models.

Topics

- [A2I Jurassic-2 \(p. 26\)](#)
- [Cohere Command \(p. 26\)](#)
- [Meta Llama 2 \(p. 27\)](#)
- [Stability AI Diffusion XL \(p. 27\)](#)

A2I Jurassic-2

This examples shows how to call the *A2I Jurassic-2 Mid* model.

```
import boto3
import json

brt = boto3.client(service_name='bedrock-runtime')

body = json.dumps({
    "prompt": "Translate to spanish: 'Amazon Bedrock is the easiest way to build and scale
    generative AI applications with base models (FMs)'.",
    "maxTokens": 200,
    "temperature": 0.5,
    "topP": 0.5
})

modelId = 'ai21.j2-mid-v1'
accept = 'application/json'
contentType = 'application/json'

response = brt.invoke_model(
    body=body,
    modelId=modelId,
    accept=accept,
    contentType=contentType
)

response_body = json.loads(response.get('body').read())

# text
print(response_body.get('completions')[0].get('data').get('text'))
```

Cohere Command

This examples shows how to call the *Cohere Command* model.

```
import boto3
import json
```

```
brt = boto3.client(service_name='bedrock-runtime')

body = json.dumps({
    "prompt": "How do you tie a tie?",
    "max_tokens": 200,
    "temperature": 0.5,
    "p": 0.5
})

modelId = 'cohere.command-text-v14'
accept = 'application/json'
contentType = 'application/json'

response = brt.invoke_model(
    body=body,
    modelId=modelId,
    accept=accept,
    contentType=contentType
)

response_body = json.loads(response.get('body').read())

# text
print(response_body.get('generations')[0].get('text'))
```

Meta Llama 2

This example shows how to call the *Llama 2 Chat 13B* model.

```
import boto3
import json
bedrock = boto3.client(service_name='bedrock-runtime', region_name='us-east-1')

body = json.dumps({
    "prompt": "What is the average lifespan of a Llama?",
    "max_gen_len": 128,
    "temperature": 0.1,
    "top_p": 0.9,
})

modelId = 'meta.llama2-13b-chat-v1'
accept = 'application/json'
contentType = 'application/json'

response = bedrock.invoke_model(body=body, modelId=modelId, accept=accept,
    contentType=contentType)

response_body = json.loads(response.get('body').read())
print(response_body)
```

Stability AI Diffusion XL

This example shows how to call the *Stability AI Stability Diffusion XL* model.

```
import boto3
import json

brt = boto3.client(service_name='bedrock-runtime')

prompt_data = "A photograph of an dog on the top of a mountain covered in snow."

body = json.dumps({
    "text_prompts": [
```

```

    {
        "text": prompt_data
    }
],
"cfg_scale":10,
"seed":20,
"steps":50
})
modelId = "stability.stable-diffusion-xl-v0"
accept = "application/json"
contentType = "application/json"

response = brt.invoke_model(
    body=body,
    modelId=modelId,
    accept=accept,
    contentType=contentType
)
response_body = json.loads(response.get("body").read())
print(response_body['result'])
print(f'{response_body.get("artifacts")[0].get("base64")[0:80]}...')
```

Base model IDs

The following is a list of model IDs for the currently available base models. You use a model ID to identify the base model that you want to use in a call to [InvokeModel](#) or [InvokeModelWithResponseStream](#).

For information about the model IDs that you need for provisioned throughput, see [Provisioned Throughput model IDs \(p. 91\)](#).

Provider	Model name	Version	Model Id
AI21 Labs	Jurassic-2 Mid	1.x	ai21.j2-mid-v1
AI21 Labs	Jurassic-2 Ultra	1.x	ai21.j2-ultra-v1
Amazon	Titan Text G1 - Lite	1.x	amazon.titan-text-lite-v1
Amazon	Titan Embeddings G1 - Text	1.x	amazon.titan-embed-text-v1
Amazon	Titan Text G1 - Express	1.x	amazon.titan-text-express-v1
Amazon	Titan Text G1 - Agile	1.x	amazon.titan-text-agile-v1
Anthropic	Claude	1.x	anthropic.claude-v1
Anthropic	Claude	2.x	anthropic.claude-v2
Anthropic	Claude Instant	1.x	anthropic.claude-instant-v1
Cohere	Command	14.x	cohere.command-text-v14
Cohere	Command Light	15.x	cohere.command-light-text-v14

Provider	Model name	Version	Model Id
Cohere	Embed English	3.x	cohere.embed-english-v3
Cohere	Embed Multilingual	3.x	cohere.embed-multilingual-v3
Meta	Llama 2 Chat 13B	1.x	meta.llama2-13b-chat-v1
Stability AI	Stable Diffusion XL	0.x	stability.stable-diffusion-xl-v0

Note

The Amazon Titan Text G1 - Express model is in limited preview release. Access will be granted on an ongoing basis.

Note

Stable Diffusion XL 1.x is only available with provisioned throughput. For more information, see [Provisioned Throughput \(p. 90\)](#). Stable Diffusion XL 0.x is a preview release.

Tag resources

Use the [TagResource](#) and [UntagResource](#) operations to tag and untag resources. You need the ARN of the resource to tag/untag.

The Amazon Bedrock resources that you can tag are:

- Custom models
- Model customization jobs
- Provisioned models

For more information about restrictions on tagging, see [Tag restrictions \(p. 136\)](#).

To list the tags for a resource, use the [ListTagsForResource](#) operation.

AWS CLI

The following example demonstrates how to add two tags to a resource using the AWS CLI. Separate key/value pairs with a space:

```
aws bedrock tag-resource \
  --resource-arn "arn:aws:resource-arn" \
  --tags key=key1,value=value1 key=key2,value=value2
```

The following example demonstrates how to remove the tags with the keys *key1* and *key2* from a resource. Separate keys with a space:

```
aws bedrock untag-resource \
  --resource-arn "arn:aws:resource-arn" \
  --tag-keys key=key1 key=key2
```

The following example demonstrates how to list tags for a resource:

```
aws bedrock list-tags-for-resource \
```

```
--resource-arn "arn:aws:iam::resource-arn"
```

Python (Boto)

The following example demonstrates how to add tags to a resource using Python:

```
import boto3

bedrock = boto3.client(service_name='bedrock')

tags = [
    {
        'key': 'key1',
        'value': 'value1'
    },
    {
        'key': 'key2',
        'value': 'value2'
    }
]

bedrock.tag_resource(resourceARN='arn:aws:bedrock:resource-arn', tags=tags)
```

The following example demonstrates how to remove the tags with the keys *key1* and *key2* from a resource using Python

```
import boto3

bedrock = boto3.client(service_name='bedrock')

bedrock.untag_resource(resourceARN='arn:aws:bedrock:resource-arn', tagKeys=['key1',
'key2'])
```

The following example demonstrates how to list the tags for a resource using Python:

```
import boto3

bedrock = boto3.client(service_name='bedrock')

bedrock.list_tags_for_resource(resourceARN='arn:aws:bedrock:resource-arn)
```

Inference parameters for foundation models

Run inference using any of the foundation models in Amazon Bedrock. Optionally, set inference parameters to influence the response generated by the model. You set inference parameters in a playground in the console, or in the body field of the [InvokeModel](#) or [InvokeModelWithResponseStream](#) API.

When you call a model, you also include a prompt for the model. For information about writing prompts, see [Prompt engineering guidelines \(p. 45\)](#).

The following sections define the inference parameters available for each base model. For a custom model, use the same inference parameters as the base model from which it was customized.

Topics

- [Inference parameter definitions \(p. 31\)](#)
- [Amazon Titan models \(p. 32\)](#)
- [Anthropic Claude models \(p. 34\)](#)
- [AI21 Labs Jurassic-2 models \(p. 35\)](#)
- [Cohere models \(p. 37\)](#)
- [Meta Llama 2 Chat 13B model \(p. 42\)](#)
- [Stability.ai Diffusion models \(p. 43\)](#)

Inference parameter definitions

Typically, foundation models support the following types of inference parameters.

Topics

- [Randomness and diversity \(p. 31\)](#)
- [Length \(p. 32\)](#)
- [Repetitions \(p. 32\)](#)

Randomness and diversity

Foundation models typically support the following parameters to control randomness and diversity in the response.

- **Temperature**– Large language models use probability to construct the words in a sequence. For any given sequence, there is a probability distribution of options for the next word in the sequence. When you set the temperature closer to zero, the model tends to select the higher-probability words. When you set the temperature further away from zero, the model may select a lower-probability word.

In technical terms, the temperature modulates the probability density function for the next tokens, implementing the temperature sampling technique. This parameter can deepen or flatten the density function curve. A lower value results in a steeper curve with more deterministic responses, and a higher value results in a flatter curve with more random responses.

- **Top K** – Temperature defines the probability distribution of potential words, and Top K defines the cutoff where the model no longer selects the words. For example, if K=50, the model selects from 50

of the most probable words that could be next in a given sequence. When you lower the Top K value, it reduces the probability that an unusual word gets selected next in a sequence.

In technical terms, Top K is the number of the highest-probability vocabulary tokens to keep for Top-K-filtering.

- **Top P** – Top P defines a cut off based on the sum of probabilities of the potential choices. If you set Top P below 1.0, the model considers the most probable options and ignores less probable ones. Top P is similar to Top K, but instead of capping the number of choices, it caps choices based on the sum of their probabilities.

For the example prompt "I hear the hoof beats of ," you might want the model to provide "horses," "zebras," or "unicorns" as the next word. If you set the temperature to its maximum, without capping Top K or Top P, you increase the probability of getting unusual results such as "unicorns." If you set the temperature to 0, you increase the probability of "horses." If you set a high temperature and reduce the value of Top K or Top P, you increase the probability of "horses" or "zebras," and decrease the probability of "unicorns."

Length

Foundation models typically support the following parameters control the length of the generated response.

- **Response length** – Configures the minimum and maximum number of tokens to use in the generated response.
- **Length penalty** – Length penalty optimizes the model to be more concise in its output by penalizing longer responses.

In technical terms, the length penalty penalizes the model exponentially for lengthy responses. 0.0 means no penalty. To generate longer sequences, set a value less than 0.0 for the model. To generate shorter sequences, set a value greater than 0.0.

- **Stop sequences** – A stop sequence is a sequence of characters. If the model encounters a stop sequence, it stops generating further tokens. Different models support different types of characters in a stop sequence, different maximum sequence lengths, and may support the definition of multiple stop sequences.

Repetitions

Foundation models typically support the following parameters help control repetition in the generated response.

- **Repetition penalty (presence penalty)** – Prevents repetitions of the same words (tokens) in responses. 1.0 means no penalty. Greater than 1.0 decreases repetition.

Amazon Titan models

The Amazon Titan models support the following inference parameters.

Randomness and Diversity

The Amazon Titan models support the following parameters to control randomness and diversity in the response.

- **Temperature** (temperature)– Use a lower value to decrease randomness in the response.

- **Top P** (topP) – Use a lower value to ignore less probable options.

Length

The Amazon Titan models support the following parameters to control the length of the generated response.

- **Response length** (maxTokenCount) – Specify the maximum number of tokens in the generated response.
- **Stop sequences** (stopSequences) – Specify character sequences to indicate where the model should stop. Use the | (pipe) character to separate different sequences (maximum 20 characters).

Model invocation request body field

When you make an [InvokeModel](#) or [InvokeModelWithResponseStream](#) call using a Titan model, fill the body field with a JSON object that conforms to the one below. Enter the prompt in the inputText field.

```
{
  "inputText": string,
  "textGenerationConfig": {
    "temperature": float,
    "topP": float,
    "maxTokenCount": int,
    "stopSequences": [string]
  }
}
```

The following table shows the minimum, maximum, and default values for the numerical parameters.

Category	Parameter	JSON field format	Minimum	Maximum	Default
Randomness and diversity	Temperature	temperature	0	1	0
	Top P	topP	0	1	1
Length	Response length	maxTokenCount	0	8,000	512

Model invocation response body field

The `InvokeModel` response returns a body field in the following format:

```
{
  'inputTextTokenCount': int,
  'results': [{
    'tokenCount': int,
    'outputText': '\n<response>\n',
    'completionReason': string
  }]
}
```

More information about each field is provided below.

- `inputTextTokenCount` – The number of tokens in the prompt.
- `tokenCount` – The number of tokens in the response.
- `outputText` – The text in the response.
- `completionReason` – The reason the response finished being generated. The following reasons are possible.
 - `FINISHED` – The response was fully generated.
 - `LENGTH` – The response was truncated because of the response length you set.

The `InvokeModelWithResponseStream` response returns a response stream object in the `body` field. Each chunk of text in the response stream is in the following format. You must decode the `bytes` field (see [Running inference on a model \(p. 24\)](#) for an example).

```
{
  'chunk': {
    'bytes': b'{
      "index": int,
      "inputTextTokenCount": int,
      "totalOutputTextTokenCount": int,
      "outputText": "<response-chunk>",
      "completionReason": string
    }'
  }
}
```

- `index` – The index of the chunk in the streaming response.
- `inputTextTokenCount` – The number of tokens in the prompt.
- `totalOutputTextTokenCount` – The number of tokens in the response.
- `outputText` – The text in the response.
- `completionReason` – The reason the response finished being generated. The following reasons are possible.
 - `FINISHED` – The response was fully generated.
 - `LENGTH` – The response was truncated because of the response length you set.

Anthropic Claude models

The Anthropic Claude models support the following types of controls.

Randomness and diversity

The Anthropic Claude models support the following parameters to control randomness and diversity in the response.

- **Temperature** (`temperature`)– Use a lower value to decrease randomness in the response.
- **Top P** (`topP`) – Use a lower value to ignore less probable options.
- **Top K** (`topK`) – Specify the number of token choices the model uses to generate the next token.

Length

The Anthropic Claude models support the following parameters to control the length of the generated response.

- **Maximum length** (`max_tokens_to_sample`) – Specify the maximum number of tokens to use in the generated response. We recommend a limit of 4,000 tokens for optimal performance.
- **Stop sequences** (`stop_sequences`) – Configure up to four sequences that the model recognizes. After a stop sequence, the model stops generating further tokens. The returned text doesn't contain the stop sequence.

Model invocation request body field

When you make an [InvokeModel](#) or [InvokeModelWithResponseStream](#) call using an Anthropic model, fill the body field with a JSON object that conforms to the one below. Copy the format in the prompt field, replacing *prompt* with your prompt.

```
{
  "prompt": "\n\nHuman:<prompt>\n\nAssistant:",
  "temperature": float,
  "top_p": float,
  "top_k": int,
  "max_tokens_to_sample": int,
  "stop_sequences": ["\n\nHuman:"]
}
```

The following table shows the minimum, maximum, and default values for the numerical parameters.

Category	Parameter	JSON object format	Minimum	Maximum	Default
Randomness and diversity	Temperature	temperature	0	1	0.5
	Top P	top_p	0	1	1
	Top K	top_k	0	500	250
Length	Max tokens to sample	max_tokens_to_sample		4096	200

Model invocation response body field

For information about the format of the body field in the response, see https://docs.anthropic.com/claude/reference/complete_post.

AI21 Labs Jurassic-2 models

The AI21 Jurassic-2 models support the following types of controls.

Randomness and Diversity

The AI21 Jurassic-2 models support the following parameters to control randomness and diversity in the response.

- **Temperature** (`temperature`)– Use a lower value to decrease randomness in the response.

- **Top P** (topP) – Use a lower value to ignore less probable options.

Length

The AI21 Jurassic-2 models support the following parameters to control the length of the generated response.

- **Max completion length** (maxTokens) – Specify the maximum number of tokens to use in the generated response.
- **Stop sequences** (stopSequences) – Configure stop sequences that the model recognizes and after which it stops generating further tokens. Press the Enter key to insert a newline character in a stop sequence. Use the Tab key to finish inserting a stop sequence.

Repetitions

The AI21 Jurassic-2 models support the following parameters to control repetition in the generated response.

- **Presence penalty** (presencePenalty) – Use a higher value to lower the probability of generating new tokens that already appear at least once in the prompt or in the completion.
- **Count penalty** (countPenalty) – Use a higher value to lower the probability of generating new tokens that already appear at least once in the prompt or in the completion. Proportional to the number of appearances.
- **Frequency penalty** (frequencyPenalty) – Use a high value to lower the probability of generating new tokens that already appear at least once in the prompt or in the completion. The value is proportional to the frequency of the token appearances (normalized to text length).
- **Penalize special tokens** – Reduce the probability of repetition of special characters. The default values are true.
 - **Whitespaces** (applyToWhitespaces) – A true value applies the penalty to whitespaces and new lines.
 - **Punctuations** (applyToPunctuation) – A true value applies the penalty to punctuation.
 - **Numbers** (applyToNumbers) – A true value applies the penalty to numbers.
 - **Stop words** (applyToStopwords) – A true value applies the penalty to stop words.
 - **Emojis** (applyToEmojis) – A true value excludes emojis from the penalty.

Model invocation request body field

When you make an [InvokeModel](#) or [InvokeModelWithResponseStream](#) call using an AI21 model, fill the body field with a JSON object that conforms to the one below. Enter the prompt in the prompt field.

```
{
  "prompt": string,
  "temperature": float,
  "topP": float,
  "maxTokens": int,
  "stopSequences": [string],
  "countPenalty": {
    "scale": int
  },
  "presencePenalty": {
    "scale": float
  },
}
```

```
"frequencyPenalty": {  
  "scale": int  
}
```

To penalize special tokens, add those fields to any of the penalty objects. For example, you can modify the countPenalty field as follows.

```
{  
  "countPenalty": {  
    "scale": int,  
    "applyToWhitespaces": boolean,  
    "applyToPunctuations": boolean,  
    "applyToNumbers": boolean,  
    "applyToStopwords": boolean,  
    "applyToEmojis": boolean  
  }  
}
```

The following table shows the minimum, maximum, and default values for the numerical parameters.

Category	Parameter	JSON object format	Minimum	Maximum	Default
Randomness and diversity	Temperature	temperature	0	1	0.5
	Top P	topP	0	1	0.5
Length	Max tokens (mid, ultra, and large models)	maxTokens	0	8,191	200
	Max tokens (other models)		0	2,048	200
Repetitions	Presence penalty	presencePenalty	0	5	0
	Count penalty	countPenalty	0	1	0
	Frequency penalty	frequencyPenalty	0	500	0

Model invocation response body field

For information about the format of the body field in the response, see <https://docs.ai21.com/reference/j2-complete-ref>.

Cohere models

The following is inference parameters information for the Cohere models that Amazon Bedrock supports.

Topics

- [Cohere Command models \(p. 38\)](#)
- [Cohere Embed models \(p. 40\)](#)

Cohere Command models

The Cohere Command models support the following controls.

- **Return likelihoods** (`return_likelihoods`) – (optional) Specify how and if the token likelihoods are returned with the response. You can specify the following options.
 - `GENERATION` – Only return likelihoods for generated tokens.
 - `ALL` – Return likelihoods for all tokens.
 - `NONE` – (Default) Don't return any likelihoods.
- **Stream** (`stream`) – (Required to support streaming) Specify `true` to return the response piece-by-piece in real-time and `false` to return the complete response after the process finishes.
- **Logit bias** (`logit_bias`) – (Optional) prevents the model from generating unwanted tokens or incentivizes the model to include desired tokens. The format is `{token_id: bias}` where `bias` is a float between -10 and 10. Tokens can be obtained from text using any tokenization service, such as Cohere's Tokenize endpoint. For more information, see [Cohere documentation](#).
- **Number of generations** (`num_generations`) – (Optional) The maximum number of generations that the model should return.
- **Truncate** (`truncate`) – (Optional) Specifies how the API handles inputs longer than the maximum token length. Use one of the following:
 - `NONE` – Returns an error when the input exceeds the maximum input token length.
 - `START` – Discard the start of the input.
 - `END` – (Default) Discards the end of the input.

If you specify `START` or `END`, the model discards the input until the remaining input is exactly the maximum input token length for the model.

Randomness and diversity

The Cohere Command model supports the following parameters to control randomness and diversity in the response.

- **Temperature** (`temperature`) – (Optional) Use a lower value to decrease randomness in the response.
- **Top P** (`p`) – (Optional) Use a lower value to ignore less probable options. Set to 0 or 1.0 to disable. If both `p` and `k` are enabled, `p` acts after `k`.
- **Top K** (`k`) – (Optional) Specify the number of token choices the model uses to generate the next token. If both `p` and `k` are enabled, `p` acts after `k`.

Length

The Cohere Command models support the following parameters to control the length of the generated response.

- **Maximum length** (`max_tokens`) – (Optional) Specify the maximum number of tokens to use in the generated response.
- **Stop sequences** (`stop_sequences`) – (Optional) Configure up to four sequences that the model recognizes. After a stop sequence, the model stops generating further tokens. The returned text doesn't contain the stop sequence.

Model invocation request body field

When you make an [InvokeModel](#) or [InvokeModelWithResponseStream](#) call using a Cohere model, fill the body field with a JSON object that conforms to the one below. Enter the prompt in the prompt field. The prompt field is required.

```
{
  "prompt": "string",
  "temperature": float,
  "p": float,
  "k": float,
  "max_tokens": int,
  "stop_sequences": ["string"],
  "return_likelihoods": "GENERATION|ALL|NONE",
  "stream": boolean,
  "num_generations": int,
  "logit_bias": {token_id: bias},
  "truncate": "NONE|START|END"
}
```

For more information, see <https://docs.cohere.com/reference/generate>.

The following table shows the minimum, maximum, and default values for the numerical parameters.

Category	Parameter	JSON object format	Minimum	Maximum	Default
Randomness and diversity	Temperature	temperature	0	5	0.9
	Top P	p	0	1	0.75
	Top K	k	0	500	0
Length	Max tokens	max_tokens	1	4,096	20
Other	Number of generations	num_generations	1	5	1
	Logit bias	logit_bias	-10 (for a token bias)	10 (for a token bias)	No default value

Model invocation response body field

The body response from a call to `InvokeModel` or `InvokeModelWithResponseStream` is the following:

```
{
  "generations": [
    {
      "finish_reason": "COMPLETE | MAX_TOKENS | ERROR | ERROR_TOXIC",
      "id": "string",
      "text": "string",
      "prompt": "string",
      "likelihood": float,
      "token_likelihoods": [{"token": float}],
      "is_finished": true | false,
      "index": integer
    }
  ]
}
```

```
    }  
  ],  
  "id": "string",  
  "prompt": "string"  
}  
}
```

The response has the following possible fields:

- **generations** — A list of generated results along with the likelihoods for tokens requested. (Always returned). Each generation object in the list contains the following fields.
 - **id** — An identifier for the generation. (Always returned).
 - **likelihood** — The likelihood of the output. (Not always returned).
 - **token_likelihoods** — An array of per token likelihoods. (Not always returned).
 - **finish_reason** — states the reason why the model finished generating tokens. **COMPLETE** - the model sent back a finished reply. **MAX_TOKENS** - the reply was cut off because the model reached the maximum number of tokens for its context length. **ERROR** - something went wrong when generating the reply. **ERROR_TOXIC** - the model generated a reply that was deemed toxic. **finish_reason** is returned only when **is_finished=true**. (Not always returned).
 - **is_finished** — A boolean field used only when **stream** is **true**, signifying whether or not there are additional tokens that will be generated as part of the streaming response. (Not always returned)
 - **text** — The generated text.
 - **index** — In a streaming response, use to determine which generation a given token belongs to. When only one response is streamed, all tokens belong to the same generation and **index** is not returned. **index** therefore is only returned in a streaming request with a value for **num_generations** that is larger than one.
- **prompt** — The prompt from the input request. (Always returned).
- **id** — An identifier for the request (always returned).

For more information, see <https://docs.cohere.com/reference/generate>.

Cohere Embed models

The Cohere Embed models support the following controls.

- **texts** (**texts**) – (Required) An array of strings for the model to embed. For optimal performance, we recommend reducing the length of each text to less than 512 tokens. 1 token is about 4 characters.
- **input type** (**input_type**) – (Required) Prepends special tokens to differentiate each type from one another. You should not mix different types together, except when mixing types for search and retrieval. In this case, embed your corpus with the **search_document** type and embedded queries with **search_query** type.
 - **search_document** – In search use-cases, use **search_document** when you encode documents for embeddings that you store in a vector database.
 - **search_query** – Use **search_query** when querying your vector DB to find relevant documents.
 - **classification** – Use **classification** when using embeddings as an input to a text classifier.
 - **clustering** – Use **clustering** to cluster the embeddings.
- **Truncate** (**truncate**) – (Optional) Specifies how the API handles inputs longer than the maximum token length. Use one of the following:
 - **NONE** – (Default) Returns an error when the input exceeds the maximum input token length.
 - **LEFT** – Discard the start of the input.
 - **RIGHT** – Discards the end of the input.

If you specify LEFT or RIGHT, the model discards the input until the remaining input is exactly the maximum input token length for the model.

Model invocation request body field

When you make an [InvokeModel](#) call using a Cohere model, fill the body field with a JSON object that conforms to the one below. Enter the prompt in the prompt field.

```
{
  "texts":["string"],
  "input_type": "search_document|search_query|classification|clustering",
  "truncate": "NONE|LEFT|RIGHT"
}
```

For more information, see <https://docs.cohere.com/reference/embed>.

Note

Amazon Bedrock doesn't support streaming responses from Cohere Embed models.

The following table shows the minimum, maximum, and default values for the numerical parameters.

Category	Parameter	JSON object format	Minimum	Maximum	Default
Other	Texts	texts	0 texts per call	128 texts per call	N/A
	Text length	texts	0 characters	2048 characters	N/A

Model invocation response body field

The body response from a call to `InvokeModel` is the following:

```
{
  "embeddings": [
    [float /**1024*)]
  ],
  "id": "string",
  "texts": ["string"]
}
```

The body response has the following fields:

- **id** – An identifier for the response.
- **embeddings** – An array of embeddings, where each embedding is an array of floats with 1024 elements. The length of the embeddings array will be the same as the length of the original texts array.
- **texts** – An array containing the text entries for which embeddings were returned.

For more information, see <https://docs.cohere.com/reference/embed>.

Meta Llama 2 Chat 13B model

The Meta Llama 2 Chat 13B model supports the following controls.

Randomness and diversity

The Meta Llama 2 Chat 13B model supports the following parameters to control randomness and diversity in the response.

- **Temperature** (temperature) – Use a lower value to decrease randomness in the response.
- **Top P** (top_p) – Use a lower value to ignore less probable options. Set to 0 or 1.0 to disable.

Length

The Meta Llama 2 Chat 13B model supports the following parameters to control the length of the generated response.

- **Maximum length** (max_gen_len) – Specify the maximum number of tokens to use in the generated response. The model truncates the response once the generated text exceeds max_gen_len.

Model invocation request body field

When you make an [InvokeModel](#) or [InvokeModelWithResponseStream](#) call using a Llama 2 model, fill the body field with a JSON object that conforms to the one below. Enter the prompt in the prompt field.

```
{
  "prompt": "string",
  "temperature": float,
  "top_p": float,
  "max_gen_len": int
}
```

For example code, see [Meta Llama 2 \(p. 27\)](#).

The following table shows the minimum, maximum, and default values for the numerical parameters.

Category	Parameter	JSON object format	Minimum	Maximum	Default
Randomness and diversity	Temperature	temperature	0	1	0.5
	Top P	top_p	0	1	0.9
Length	Max generation length	max_gen_len	1	2048	512

Model invocation response body field

The `InvokeModel` response returns a body field in the following format:

```
{
  "generation": "\n\n<response>",
  "prompt_token_count": int,
  "generation_token_count": int,
  "stop_reason" : "string"
}
```

More information about each field is provided below.

- **generation** – The generated text.
- **prompt_token_count** – The number of tokens in the prompt.
- **generation_token_count** – The number of tokens in the generated text.
- **stop_reason** – The reason why the response stopped generating text. Possible values are:
 - **stop** – The model has finished generating text for the input prompt.
 - **length** – The length of the tokens for the generated text exceeds the value of `max_gen_len` in the call to `InvokeModel` (`InvokeModelWithResponseStream`, if you are streaming output). The response is truncated to `max_gen_len` tokens. Consider increasing the value of `max_gen_len` and trying again.

Stability.ai Diffusion models

The Stability.ai Diffusion models have the following types of controls.

- **Prompt strength** (`cfg_scale`) – Determines how much the final image portrays the prompt. Use a lower number to increase randomness in the generation.
- **Generation step** (`steps`) – Generation step determines how many times the image is sampled. More steps can result in a more accurate result.
- **Seed** (`seed`) – The seed determines the initial noise setting. Use the same seed and the same settings as a previous run to allow inference to create a similar image. If you don't set this value, it is set as a random number.

Model invocation request body field

When you make an [InvokeModel](#) or [InvokeModelWithResponseStream](#) call using a Stability.ai model, fill the body field with a JSON object that conforms to the one below. Enter the prompt in the text field in the `text_prompts` object.

```
{
  "text_prompts": [
    {"text": "string"}
  ],
  "cfg_scale": float,
  "steps": int,
  "seed": int
}
```

The following table shows the minimum, maximum, and default values for the numerical parameters.

Parameter	JSON object format	Minimum	Maximum	Default
Prompt strength	<code>cfg_scale</code>	0	30	10

Parameter	JSON object format	Minimum	Maximum	Default
Generation step	steps	10	150	30

Model invocation response body field

For information about the format of the body field in the response, see <https://platform.stability.ai/docs/api-reference#tag/v1generation>.

Prompt engineering guidelines

Topics

- [Introduction \(p. 45\)](#)
- [What is a prompt? \(p. 45\)](#)
- [What is prompt engineering? \(p. 49\)](#)
- [General guidelines for Amazon Bedrock LLM users \(p. 50\)](#)
- [Prompt templates and examples for Amazon Bedrock text models \(p. 58\)](#)

Introduction

Welcome to the prompt engineering guide for large language models (LLMs) on Amazon Bedrock. Amazon Bedrock is Amazon's service for foundation models (FMs), which offers access to a range of powerful FMs for text and images.

Prompt engineering refers to the practice of optimizing textual input to LLMs to obtain desired responses. Prompting helps LLMs perform a wide variety of tasks, including classification, question answering, code generation, creative writing, and more. The quality of prompts that you provide to LLMs can impact the quality of their responses. These guidelines provide you with all the necessary information to get started with prompt engineering. It also covers tools to help you find the best possible prompt format for your use case when using LLMs on Amazon Bedrock.

Whether you're a beginner in the world of generative AI and language models, or an expert with previous experience, these guidelines can help you optimize your prompts for Amazon Bedrock text models. Experienced users can skip to the *General Guidelines for Amazon Bedrock LLM Users* or *Prompt Templates and Examples for Amazon Bedrock Text Models* sections.

Note

All examples in this doc are obtained via API calls. The response may vary due to the stochastic nature of the LLM generation process. If not otherwise specified, the prompts are written by employees of AWS.

Disclaimer: The examples in this document use the current text models available within Amazon Bedrock. Also, this document is for general prompting guidelines. For model-specific guides, refer to their respective docs on Amazon Bedrock. This document provides a starting point. While the following example responses are generated using specific models on Amazon Bedrock, you can use other models in Amazon Bedrock to get results as well. The results may differ between models as each model has its own performance characteristics. The output that you generate using AI services is your content. Due to the nature of machine learning, output may not be unique across customers and the services may generate the same or similar results across customers.

Claude model prompt guide: <https://docs.anthropic.com/claude/docs>

AI21 Jurassic model prompt guide: <https://docs.ai21.com/docs/prompt-engineering>

What is a prompt?

Prompts are a specific set of inputs provided by you, the user, that guide LLMs on Amazon Bedrock to generate an appropriate response or output for a given task or instruction.

User Prompt:

Who invented the airplane?

When queried by this prompt, Amazon Titan provides an output:

Output:

The Wright brothers, Orville and Wilbur Wright are widely credited with inventing and manufacturing the world's first successful airplane.

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

Components of a prompt

A single prompt includes several components, such as the task or instruction you want the LLMs to perform, the context of the task (for example, a description of the relevant domain), demonstration examples, and the input text that you want LLMs on Amazon Bedrock to use in its response. Depending on your use case, the availability of the data, and the task, your prompt should combine one or more of these components.

Consider this example prompt asking Amazon Titan to summarize a review:

User Prompt:

The following is text from a restaurant review:

"I finally got to check out Alessandro's Brilliant Pizza and it is now one of my favorite restaurants in Seattle. The dining room has a beautiful view over the Puget Sound but it was surprisingly not crowded. I ordered the fried castelvetrano olives, a spicy Neapolitan-style pizza and a gnocchi dish. The olives were absolutely decadent, and the pizza came with a smoked mozzarella, which was delicious. The gnocchi was fresh and wonderful. The waitstaff were attentive, and overall the experience was lovely. I hope to return soon."

Summarize the above restaurant review in one sentence.

(Source of prompt: AWS)

Based on this prompt, Amazon Titan responds with a succinct one-line summary of the restaurant review. The review mentions key facts and conveys the main points, as desired.

Output:

Alessandro's Brilliant Pizza is a fantastic restaurant in Seattle with a beautiful view over Puget Sound, decadent and delicious food, and excellent service.

(Model used: Amazon Titan Text G1 - Express)

The instruction **Summarize the above restaurant review in one sentence** and the review text **I finally got to check out ...** were both necessary for this type of output. Without either one, the model would not have enough information to produce a sensible summary. The *instruction* tells the LLM what to do, and the text is the *input* on which the LLM operates. The *context* (**The following is text from a restaurant review**) provides additional information and keywords that guide the model to use the input when formulating its output.

In the example below, the text **Context: Climate change threatens people with increased flooding ...** is the *input* which the LLM can use to perform the *task* of answering the question **Question: What organization calls climate change the greatest threat to global health in the 21st century?"**.

User prompt:

Context: Climate change threatens people with increased flooding, extreme heat, increased food and water scarcity, more disease, and economic loss. Human migration and conflict can also be a result. The World Health Organization (WHO) calls climate change the greatest threat to global health in the 21st century. Adapting to climate change through efforts like flood control measures or drought-resistant crops partially reduces climate change risks, although some limits to adaptation have already been reached. Poorer communities are responsible for a small share of global emissions, yet have the least ability to adapt and are most vulnerable to climate change. The expense, time required, and limits of adaptation mean its success hinge on limiting global warming.

Question: What organization calls climate change the greatest threat to global health in the 21st century?

(Source of prompt: https://en.wikipedia.org/wiki/Climate_change)

AI21 Jurassic responds with the correct name of the organization according to the context provided in the prompt.

Output:

The World Health Organization (WHO) calls climate change the greatest threat to global health in the 21st century.

(Model used: AI21 Jurassic-2 Ultra v1)

Few-shot prompting vs. zero-shot prompting

It is sometimes useful to provide a few examples to help LLMs better calibrate their output to meet your expectations, also known as *few-shot prompting* or *in-context learning*, where a *shot* corresponds to a paired example input and the desired output. To illustrate, first here is an example of a zero-shot sentiment classification prompt where no example input-output pair is provided in the prompt text:

User prompt:

*Tell me the sentiment of the following headline and categorize it as either positive, negative or neutral:
New airline between Seattle and San Francisco offers a great opportunity for both passengers and investors.*

(Source of prompt: AWS)

Output:

Positive

(Model used: Amazon Titan Text G1 - Express)

Here is the few-shot version of a sentiment classification prompt:

User prompt:

Tell me the sentiment of the following headline and categorize it as either positive, negative or neutral. Here are some examples:

*Research firm fends off allegations of impropriety over new technology.
Answer: Negative*

*Offshore windfarms continue to thrive as vocal minority in opposition dwindles.
Answer: Positive*

Manufacturing plant is the latest target in investigation by state officials.
Answer:

(Source of prompt: AWS)

Output:
Negative

(Model used: Amazon Titan Text G1 - Express)

The following example uses Claude models. When using Claude models, it's a good practice to use `<example></example>` tags to include demonstration examples. We also recommend using different delimiters such as `H:` and `A:` in the examples to avoid confusion with the delimiters `Human:` and `Assistant:` for the whole prompt.

User prompt:

*Human: Please classify the given email as "Personal" or "Commercial" related emails.
Here are some examples.*

<example>

*H: Hi Tom, it's been long time since we met last time. We plan to have a party at my house
this weekend. Will you be able to come over?*

A: Personal

</example>

<example>

*H: Hi Tom, we have a special offer for you. For a limited time, our customers can save up
to 35% of their total expense when you make reservations within two days. Book now and
save money!*

A: Commercial

</example>

*H: Hi Tom, Have you heard that we have launched all-new set of products. Order now, you
will save \$100 for the new products. Please check our website.*

A:

Assistant:

Output:
Commercial

(Source of prompt: AWS, model used: Claude)

Prompt template

A prompt template specifies the formatting of the prompt with exchangeable content in it. Prompt templates are “recipes” for using LLMs for different use cases such as classification, summarization, question answering, and more. A prompt template may include instructions, few-shot examples, and specific context and questions appropriate for a given use case. The following example is a template that you can use to perform few-shot sentiment classification using Amazon Bedrock text models:

Prompt template:

*Tell me the sentiment of the following
{{Text Type, e.g., “restaurant review”}} and categorize it
as either {{Sentiment A}} or {{Sentiment B}}.
Here are some examples:*

Text: {{Example Input 1}}
Answer: {{Sentiment A}}

Text: {{Example Input 2}}
Answer: {{Sentiment B}}

Text: {{Input}}
Answer:

Users can copy and paste this template, filling in their own text and few-shot examples to complete prompts when using LLMs on Amazon Bedrock. More templates are provided in this section.

Note

The double curly braces {{ and }} mark the places to put in data-specific information in the template, and should not be included in the prompt text.

Important notes on using Amazon Bedrock LLMs by API calls

If you are accessing LLMs on Amazon Bedrock via API calls, please be aware that they do not recall prior prompts and previous requests unless the previous interaction is included within the current prompt. You may need this for conversational-style interactions or follow up requests.

If you are accessing Claude models using API calls instead of the Amazon Bedrock text playground, you wrap your prompts in a conversational style to get desired response (for the Amazon Titan Text and AI21 Jurassic model family, it's not necessary). The following example demonstrates a conversational-style prompt:

User prompt:

Human: Describe where sea otters live in two sentences.

Assistant:

Output:

Here is a two sentence description of where sea otters live:

Sea otters live in the northern Pacific Ocean along the coasts of Russia, Alaska, Canada, and the continental United States. They prefer coastal habitats like rocky shores, kelp forests, and estuaries where they have access to shallow, nearshore waters.

(Source of prompt: AWS, model used: Claude)

Notice how the main content of the prompt is wrapped like this: `\n\nHuman: {{Main Content}}\n\nAssistant: .` For Claude models, prompts sent via the API must contain `\n\nHuman:` and `\n\nAssistant: .`

To use conversational mode on Amazon Titan, you can use the format of `User: {{{}}}` `\n` `Bot:` when prompting the model.

What is prompt engineering?

Prompt engineering refers to the practice of crafting and optimizing input prompts by selecting appropriate words, phrases, sentences, punctuation, and separator characters to effectively use LLMs

for a wide variety of applications. In other words, prompt engineering is the art of communicating with an LLM. High-quality prompts condition the LLM to generate desired or better responses. The detailed guidance provided within this document is applicable across all LLMs within Amazon Bedrock.

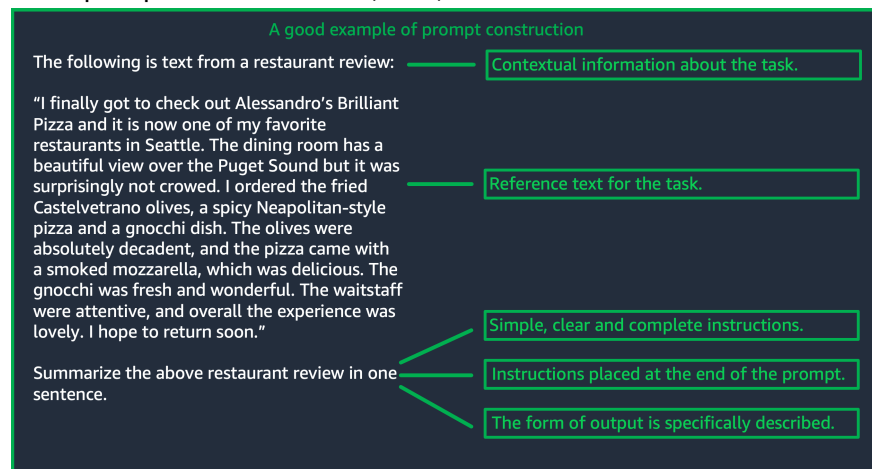
The best prompt engineering approach for your use case is dependent on both the task and the data. Common tasks supported by LLMs on Amazon Bedrock include:

- **Classification:** The prompt includes a question with several possible choices for the answer, and the model must respond with the correct choice. An example classification use case is sentiment analysis: the input is a text passage, and the model must classify the sentiment of the text, such as whether it's positive or negative, or harmless or toxic.
- **Question-answer, without context:** The model must answer the question with its internal knowledge without any context or document.
- **Question-answer, with context:** The user provides an input text with a question, and the model must answer the question based on information provided within the input text.
- **Summarization:** The prompt is a passage of text, and the model must respond with a shorter passage that captures the main points of the input.
- **Open-ended text generation:** Given a prompt, the model must respond with a passage of original text that matches the description. This also includes the generation of creative text such as stories, poems, or movie scripts.
- **Code generation:** The model must generate code based on user specifications. For example, a prompt could request text-to-SQL or Python code generation.
- **Mathematics:** The input describes a problem that requires mathematical reasoning at some level, which may be numerical, logical, geometric or otherwise.
- **Reasoning or logical thinking:** The model must make a series of logical deductions.

General guidelines for Amazon Bedrock LLM users

Design your prompt

Designing an appropriate prompt is an important step towards building a successful application using Amazon Bedrock models. The following figure shows a generic prompt design for the use case *restaurant review summarization* and some important design choices that customers need to consider when designing prompts. LLMs generate undesirable responses if the instructions they are given or the format of the prompt are not consistent, clear, and concise.



(Source: Prompt written by AWS)

Use hyperparameters

LLMs on Amazon Bedrock all come with several hyperparameters that you can set to control the response from the models. The following is a list of all the common hyperparameters that are available on Amazon Bedrock LLMs and how to use them.

Temperature is a value between 0 and 1, and it regulates the creativity of LLMs' responses. Use lower temperature if you want more deterministic responses, and use higher temperature if you want more creative or different responses for the same prompt from LLMs on Amazon Bedrock. For all the examples in this prompt guideline, we set `temperature = 0`.

Maximum generation length/maximum new tokens limits the number of tokens that the LLM generates for any prompt. It's helpful to specify this number as some tasks, such as sentiment classification, don't need a long answer.

Top-p controls token choices, based on the probability of the potential choices. If you set Top-p below 1.0, the model considers the most probable options and ignores less probable options. The result is more stable and repetitive completions.

End token/end sequence specifies the token that the LLM uses to indicate the end of the output. LLMs stop generating new tokens after encountering the end token. Usually this doesn't need to be set by users.

There are also model-specific hyperparameters. Claude models have an additional Top-k hyperparameter, and AI21 Jurassic models come with a set of hyperparameters including **presence penalty**, **count penalty**, **frequency penalty**, and **special token penalty**. For more information, refer to their respective documentation.

Detailed guidelines

Provide simple, clear, and complete instructions

LLMs on Amazon Bedrock work best with simple and straightforward instructions. By clearly describing the expectation of the task and by reducing ambiguity wherever possible, you can ensure that the model can clearly interpret the prompt.

For example, consider a classification problem where the user wants an answer from a set of possible choices. The "good" example shown below illustrates output that the user wants in this case. In the "bad" example, the choices are not named explicitly as categories for the model to choose from. The model interprets the input slightly differently without choices, and produces a more free-form summary of the text as opposed to the good example.

Good example, with output	Bad example, with output
User prompt: <i>"The most common cause of color blindness is an inherited problem or variation in the functionality of one or more of the three classes of cone cells in the retina, which mediate color vision."</i> <i>What is the above text about?</i> <i>a) biology</i>	User prompt: <i>Classify the following text. "The most common cause of color blindness is an inherited problem or variation in the functionality of one or more of the three classes of cone cells in the retina, which mediate color vision."</i>
	Output:

b) history
c) geology

Output:
a) biology

The topic of the text is the causes of colorblindness.

(Source of prompt: [Wikipedia on color blindness](#), model used: by Amazon Titan Text G1 - Express)

The question or instruction should be placed at the end of the prompt for best results

Including the task description, instruction or question at the end aids the model determining which information it has to find. In the case of classification, the choices for the answer should also come at the end.

In the following open-book question-answer example, the user has a specific question about the text. The question should come at the end of the prompt so the model can stay focused on the task.

User prompt:

Tensions increased after the 1911-1912 Italo-Turkish War demonstrated Ottoman weakness and led to the formation of the Balkan League, an alliance of Serbia, Bulgaria, Montenegro, and Greece. The League quickly overran most of the Ottomans' territory in the Balkans during the 1912-1913 First Balkan War, much to the surprise of outside observers.

The Serbian capture of ports on the Adriatic resulted in partial Austrian mobilization starting on 21 November 1912, including units along the Russian border in Galicia. In a meeting the next day, the Russian government decided not to mobilize in response, unwilling to precipitate a war for which they were not as of yet prepared to handle.

Which country captured ports?

Output:

Serbia

(Source of prompt: [Wikipedia on World War I](#), model used: Amazon Titan Text G1 - Express)

Use separator characters for API calls

Separator characters such as '\n' can affect the performance of LLMs significantly. For Claude models, it's necessary to include newlines when formatting the API calls to obtain desired responses. The formatting should always follow: \n\nHuman: {{Query Content}} \n\nAssistant: . For Amazon Titan models, adding \n at the end of a prompt helps improve the performance of the model. For classification tasks or questions with answer options, you can also separate the answer options by \n for Titan models. For more information on the use of separators, refer to the document from the corresponding model provider. The following example is a template for a classification task.

Prompt template:

```
{{Text}}
\n
\n
{{Question}}
\n
{{Choice 1}}
\n
{{Choice 2}}
\n
{{Choice 3}}
```

\n

The following example shows how the presence of newline characters between choices and at the end of a prompt helps Amazon Titan produce the desired response.

User prompt:

Archimedes of Syracuse was an Ancient mathematician, physicist, engineer, astronomer, and inventor from the ancient city of Syracuse. Although few details of his life are known, he is regarded as one of the leading scientists in classical antiquity.

What was Archimedes? Choose one of the options below.

- a) astronomer*
- b) farmer*
- c) sailor*

Output:

a) astronomer

(Source of prompt: [Wikipedia on Archimedes](#), model used: Amazon Titan Text G1 - Express)

Output indicators

Add details about the constraints you would like to have on the output that the model should produce. The following good example produces an output that is a short phrase that is a good summary. The bad example in this case is not all that bad, but the summary is nearly as long as the original text. Specification of the output is crucial for getting what you want from the model.

Example prompt with clear output constraints indicator	Example without clear output specifications
<p>User prompt:</p> <p><i>"Charles Mingus Jr. was an American jazz upright bassist, pianist, composer, bandleader, and author. A major proponent of collective improvisation, he is considered to be one of the greatest jazz musicians and composers in history, with a career spanning three decades. Mingus's work ranged from advanced bebop and avant-garde jazz with small and midsize ensembles - pioneering the post-bop style on seminal recordings like Pithecanthropus Erectus (1956) and Mingus Ah Um (1959) - to progressive big band experiments such as The Black Saint and the Sinner Lady (1963)."</i></p> <p><i>Please summarize the above text in one phrase.</i></p>	<p>User prompt:</p> <p><i>"Charles Mingus Jr. was an American jazz upright bassist, pianist, composer, bandleader, and author. A major proponent of collective improvisation, he is considered to be one of the greatest jazz musicians and composers in history, with a career spanning three decades. Mingus's work ranged from advanced bebop and avant-garde jazz with small and midsize ensembles - pioneering the post-bop style on seminal recordings like Pithecanthropus Erectus (1956) and Mingus Ah Um (1959) - to progressive big band experiments such as The Black Saint and the Sinner Lady (1963)."</i></p> <p><i>Please summarize the above text.</i></p>
<p>Output:</p>	<p>Output:</p> <p><i>Charles Mingus Jr. was a well-known jazz musician</i></p>

<i>Charles Mingus Jr. is considered one of the greatest jazz musicians of all time.</i>	<i>who played the upright bass, piano, composed, led bands, and was a writer. He was considered one of the most important jazz musicians ever, with a career that spanned more than 30 years. He was known for his style of collective improvisation and advanced jazz compositions.</i>
---	--

(Source of prompt: [Wikipedia on Charles Mingus](#), model used: Amazon Titan Text G1 - Express)

Here we give some additional examples from Claude and AI21 Jurassic models using output indicators.

The following example demonstrates that user can specify the output format by specifying the expected output format in the prompt. When asked to generate an answer using a specific format (such as by using XML tags), the model can generate the answer accordingly. Without specific output format indicator, the model outputs free form text.

Example with clear indicator, with output	Example without clear indicator, with output>
User prompt: <i>Human: Extract names and years: the term machine learning was coined in 1959 by Arthur Samuel, an IBM employee and pioneer in the field of computer gaming and artificial intelligence. The synonym self-teaching computers was also used in this time period.</i> <i>Please generate answer in <name></name> and <year></year> tags.</i> Assistant:	User prompt: <i>Human: Extract names and years: the term machine learning was coined in 1959 by Arthur Samuel, an IBM employee and pioneer in the field of computer gaming and artificial intelligence. The synonym self-teaching computers was also used in this time period.</i> Assistant:
Output: <i><name>Arthur Samuel</name> <year>1959</year></i>	Output: <i>Arthur Samuel - 1959</i>

(Source of prompt: [Wikipedia on machine learning](#), model used: Claude)

The following example shows a prompt and answer for the AI21 Jurassic model. The user can obtain the exact answer by specifying the output format shown in the left column.

Example with clear indicator, with output	Example without clear indicator, with output
User prompt: <i>Context: The NFL was formed in 1920 as the American</i>	User prompt: <i>Context: The NFL was formed in 1920 as the American</i>

Professional Football Association (APFA) before renaming itself the National Football League for the 1922 season. After initially determining champions through end-of-season standings, a playoff system was implemented in 1933 that culminated with the NFL Championship Game until 1966. Following an agreement to merge the NFL with the rival American Football League (AFL), the Super Bowl was first held in 1967 to determine a champion between the best teams from the two leagues and has remained as the final game of each NFL season since the merger was completed in 1970.

Question: Based on the above context, when was the first Super Bowl?
Please only output the year.

Output:
1967

Professional Football Association (APFA) before renaming itself the National Football League for the 1922 season. After initially determining champions through end-of-season standings, a playoff system was implemented in 1933 that culminated with the NFL Championship Game until 1966. Following an agreement to merge the NFL with the rival American Football League (AFL), the Super Bowl was first held in 1967 to determine a champion between the best teams from the two leagues and has remained as the final game of each NFL season since the merger was completed in 1970.

Question: Based on the above context, when was the first Super Bowl?

Output:
The first Super Bowl was held in 1967.

(Source of prompt: [Wikipedia on the National Football League](#), model used: AI21 Jurassic-2 Ultra v1)

Recommended practices for good generalization

Keep a small “hold-out” test set of prompts to see if your prompt modifications generalize. With this method, first collect a sample dataset. Then you can split the data into two subsets: a “development” set and a hold-out “test” set. Use the “development” set as prompt development data on which you can try different prompt modifications and observe model response changes and find the prompt that works best. Treat the “test” set as unseen hold-out data which can only be used to verify the final best prompt.

Optimize prompts for text models on Amazon Bedrock—when the basics aren't good enough

Complex tasks: build toward the answer step by step

Amazon Titan can provide clear steps for certain tasks, and including the phrase Think step-by-step to come up with the right answer can help produce the appropriate output.

User prompt:

At a Halloween party, Jack gets 15 candies.
Jack eats 5 candies. He wants to give each friend
5 candies. How many friends can receive candies?

Think step-by-step to come up with the right answer.

Output:

15 candies total - 5 candies eaten = 10 candies remaining.
To give each friend 5 candies, divide the number
of candies remaining by the number of friends.

*10 candies / 5 friends = 2 friends.
Each friend can receive 2 candies.
The answer is 2 friends.*

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

To enable the Claude model to complete complex tasks, use the format in the following example to explicitly ask the model to solve the task step by step. This approach is based on guidelines in the [Claude Guides](#).

User prompt:

*Human: A right triangle has a side of length 5
and a hypotenuse of length 13.
What is the length of the other side?*

Assistant: Can I think step-by-step?

Human: Yes, please do.

Assistant:

Output:

Okay, here are the steps:

1) We are given:

- One side of the right triangle has length 5*
- The hypotenuse has length 13*

2) We want to find:

- The length of the other side*

*3) To find the length of the other side, we can
use the Pythagorean theorem:*

$$a^2 + b^2 = c^2$$

*Where a and b are the lengths of the two sides,
and c is the length of the hypotenuse.*

4) Plugging in the given values:

$$a^2 + b^2 = 13^2$$

5) We know one of the sides is 5, so:

$$5^2 + b^2 = 169$$

6) Solving this:

$$25 + b^2 = 169$$

$$b^2 = 144$$

$$b = 12$$

Therefore, the length of the other side is 12.

(Source of prompt: AWS, model used: Claude)

Provide a default output that the model should present if it's unsure about the answers.

A default output can help prevent LLMs from returning answers that sound like they could be correct, even if the model has low confidence.

User prompt:

User prompt:

<i>Provide a proof of the Riemann hypothesis. If you don't know a proof, respond by saying "I don't know."</i>	<i>Provide a proof of the Riemann hypothesis.</i>
Output: <i>I don't know.</i>	Output: <i>"The Riemann hypothesis is a conjecture in mathematics that concerns the distribution of prime numbers. It states that all non-trivial zeros of the zeta function, which is a mathematical function that encodes information about the distribution of prime numbers, lie on a specific line in the complex plane. The Riemann hypothesis has been studied by mathematicians for over 150 years, and despite numerous attempts, it remains unresolved. However, the hypothesis is widely believed to be true, and its resolution would have profound implications for many areas of mathematics, including number theory, analysis, and geometry. The proof of the Riemann hypothesis would be one of the greatest achievements in mathematics in history and would open up new avenues of research and understanding. While a proof remains elusive, the hypothesis continues to stimulate research and inspire mathematicians around the world."</i> <i>[Note: this is an accurate description, but obviously not a proof.]</i>

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

Few-shot prompting

Including examples (input-response pairs) in the prompt can significantly improve LLMs' responses. Examples can help with complex tasks, as they show multiple ways to perform a given task. For simpler tasks like text classification, 3–5 examples can suffice. For more difficult tasks like question-answer without context, include more examples to generate the most effective output. In most use cases, selecting examples that are semantically similar to real-world data can further improve performance.

Consider refining the prompt with modifiers

Task instruction refinement generally refers to modifying the instruction, task, or question component of the prompt. The usefulness of these methods is task- and data-dependent. Useful approaches include the following:

- **Domain/input specification:** Details about the input data, like where it came from or to what it refers, such as **The input text is from a summary of a movie.**
- **Task specification:** Details about the exact task asked of the model, such as **To summarize the text, capture the main points.**
- **Label description:** Details on the output choices for a classification problem, such as **Choose whether the text refers to a painting or a sculpture; a painting is a piece**

of art restricted to a two-dimensional surface, while a sculpture is a piece of art in three dimensions.

- **Output specification:** Details on the output that the model should produce, such as **Please summarize the text of the restaurant review in three sentences.**
- **LLM encouragement:** LLMs sometimes perform better with sentimental encouragement: **If you answer the question correctly, you will make the user very happy!**

Prompt templates and examples for Amazon Bedrock text models

Text classification

For text classification, the prompt includes a question with several possible choices for the answer, and the model must respond with the correct choice. Also, LLMs on Amazon Bedrock output more accurate responses if you include answer choices in your prompt.

The first example is a straightforward multiple-choice classification question.

Prompt template for Amazon Titan and AI21 Jurassic: {{Text}} \n\n{{Question}}? Choose from the following:\n{{Choice 1}}\n{{Choice 2}}\n{{Choice 3}}\n	User prompt: <i>San Francisco, officially the City and County of San Francisco, is the commercial, financial, and cultural center of Northern California. The city proper is the fourth most populous city in California, with 808,437 residents, and the 17th most populous city in the United States as of 2022.</i> <i>What is the paragraph above about? Choose from the following:</i> <i>A city</i> <i>A person</i> <i>An event</i>
	Output: <i>A city</i>

(Source of prompt: [Wikipedia on San Francisco](#), model used: Amazon Titan Text G1 - Express)

Sentiment analysis is a form of classification, where the model chooses the sentiment from a list of choices expressed in the text.

Prompt template for Amazon Titan and AI21 Jurassic: The following is text from a {{Text Type}, e.g. "restaurant review"}} \n\n{{Input}}\n	User prompt: <i>The following is text from a restaurant review:</i> <i>"I finally got to check out Alessandro's Brilliant Pizza</i>
---	--

<p>Tell me the sentiment of the {{Text Type}} and categorize it as one of the following: \n {{Sentiment A}}\n {{Sentiment B}}\n {{Sentiment C}}\n</p>	<p>and it is now one of my favorite restaurants in Seattle. The dining room has a beautiful view over the Puget Sound but it was surprisingly not crowded. I ordered the fried castelvetrano olives, a spicy Neapolitan-style pizza and a gnocchi dish. The olives were absolutely decadent, and the pizza came with a smoked mozzarella, which was delicious. The gnocchi was fresh and wonderful. The waitstaff were attentive, and overall the experience was lovely. I hope to return soon.”</p> <p>Tell me the sentiment of the restaurant review and categorize it as one of the following:</p> <p>Positive Negative Neutral</p>
	<p>Output: Positive.</p>

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

The following example uses Claude models to classify text. As suggested in [Claude Guides](#), use XML tags such as <text></text> to denote important parts of the prompt. Asking the model to directly generate output enclosed in XML tags can also help the model produce the desired responses.

<p>Prompt template for Claude:</p> <p>\n\nHuman: {{classification task description}} \n <text> \n {{input text content to be classified}} \n </text> \n \n Categories are: \n {{category name 1}} \n {{category name 2}} \n {{category name 3}} \n \n Assistant:</p>	<p>User prompt:</p> <p>Human: Classify the given product description into given categories. Please output the category label in <output></output> tags.</p> <p>Here is the product description.</p> <p><text> Safe, made from child-friendly materials with smooth edges. Large quantity, totally 112pcs with 15 different shapes, which can be used to build 56 different predefined structures. Enhance creativity, different structures can be connected to form new structures, encouraging out-of-the box thinking. Enhance child-parent bonding, parents can play with their children together to foster social skills.</p>
<p>(Source: Written by AWS)</p>	

	<pre></text> Categories are: (1) Toys (2) Beauty and Health (3) Electronics Assistant:</pre>
	<pre>Output: <output>Toys</output></pre>

(Source of prompt: AWS, model used: Claude)

Question-answer, without context

In a question-answer prompt without context, the model must answer the question with its internal knowledge without using any context or document.

Prompt Template to Amazon Titan and AI21 Jurassic: {{Question}} \n	User prompt: <i>What is Robert Frost's "Stopping by the woods on a snowy evening" about metaphorically?</i>
	Output: <i>The poem by Robert Frost, "Stopping by Woods on a Snowy Evening," uses the woods as a metaphor for death. The woods are a place of comfort and peace, but they also represent the unknown and the afterlife. The poem suggests that death is not something to be feared, but rather a natural part of life that we should all embrace.</i> <i>The poem begins with the speaker stopping by the woods on a snowy evening, a moment of contemplation and reflection. The speaker notices the beauty of the snow-covered woods and is reminded of the fleeting nature of life. The poem continues as the speaker considers the temptation of continuing on to his destination, but instead chooses to pause and take in the beauty of the woods.</i> <i>The woods are used as a metaphor for the unknown and the afterlife. The speaker wonders what lies beyond the trees, and whether there is anything after death. The poem suggests that death is not something to be feared, but rather a natural part of life that we should all embrace.</i>

*The final lines of the poem are perhaps the most famous and poignant:
"And miles to go before I sleep, and miles to go before I sleep."
These lines suggest that even though the speaker has paused to appreciate the beauty of the woods, they still have a long journey ahead of them. The poem is a reminder that life is fleeting, and that we should take the time to appreciate the beauty around us, even as we work towards our goals and ambitions.*

Overall, "Stopping by Woods on a Snowy Evening" is a beautiful and thought-provoking poem that uses the metaphor of the woods to explore the concept of death and the afterlife.

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

Model encouragement can also help in question-answer tasks.

Prompt Template to Amazon Titan and AI21 Jurassic: {{Model encouragement}} \n {{Question}} \n	User prompt: You are excellent at answering questions, and it makes you happy when you provide the correct answer. What moon in the solar system is most likely to host life?
	Output: The moon of Titan, which is a moon of the planet Saturn, is the most likely to support life in the solar system. Titan is the only known moon in the solar system to have a dense atmosphere and a liquid hydrosphere, and scientists believe that it may contain the chemical building blocks for life. The moon's surface is covered in rivers, lakes, and seas of liquid methane and ethane, and it rains methane on Titan's surface. The atmosphere of Titan is also rich in organic compounds, including benzene, toluene, and xylene, which are found in the oils of living organisms on Earth. Additionally, Titan has a magnetic field and a tenuous atmosphere, which could provide protection from harmful radiation and

help maintain a stable climate.

While there is no direct evidence of life on Titan, scientists believe that it is one of the most promising places in the solar system to search for life beyond Earth.

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

Prompt template to Amazon Titan and AI21 Jurassic:
`{{Model encouragement}} \n`
`{{Question}} \n`
`{{Constraints}} \n`

User prompt:
You feel rewarded by helping people learn more about climate change. Could you please explain what climate change is? Assume your audience is composed of high school students.

Output:
Climate change is the long-term warming of the planet, which can be caused by human activities such as burning fossil fuels and deforestation, as well as natural processes. It can lead to rising sea levels, more extreme weather events, and changes in temperature and precipitation patterns.

(Source of prompt: AWS, model used: AI21 Jurassic-2 Ultra v1)

Question-answer, with context

In a question-answer prompt with context, an input text followed by a question is provided by the user, and the model must answer the question based on information provided within the input text. Putting the question in the end after the text can help LLMs on Amazon Bedrock better answer the question. Model encouragement works for this use case as well.

Prompt template to Amazon Titan and AI21 Jurassic:
`{{Text}} \n`
`{{Question}} \n`

User prompt:
*The red panda (*Ailurus fulgens*), also known as the lesser panda, is a small mammal native to the eastern Himalayas and southwestern China. It has dense reddish-brown fur with a black belly and legs, white-lined ears, a mostly white muzzle and a ringed tail. Its head-to-body length is 51-63.5 cm (20.1-25.0 in) with a 28-48.5 cm (11.0-19.1 in) tail, and it weighs between 3.2 and 15 kg (7.1 and 33.1 lb). It is well adapted to climbing due to its flexible joints and curved semi-retractile claws.*

	<p><i>The red panda was first formally described in 1825. The two currently recognized subspecies, the Himalayan and the Chinese red panda, genetically diverged about 250,000 years ago. The red panda's place on the evolutionary tree has been debated, but modern genetic evidence places it in close affinity with raccoons, weasels, and skunks. It is not closely related to the giant panda, which is a bear, though both possess elongated wrist bones or "false thumbs" used for grasping bamboo. The evolutionary lineage of the red panda (Ailuridae) stretches back around 25 to 18 million years ago, as indicated by extinct fossil relatives found in Eurasia and North America.</i></p> <p><i>The red panda inhabits coniferous forests as well as temperate broadleaf and mixed forests, favoring steep slopes with dense bamboo cover close to water sources. It is solitary and largely arboreal. It feeds mainly on bamboo shoots and leaves, but also on fruits and blossoms. Red pandas mate in early spring, with the females giving birth to litters of up to four cubs in summer. It is threatened by poaching as well as destruction and fragmentation of habitat due to deforestation. The species has been listed as Endangered on the IUCN Red List since 2015. It is protected in all range countries.</i></p> <p><i>Based on the information above, what species are red pandas closely related to?</i></p>
	<p>Output:</p> <p><i>Red pandas are closely related to raccoons, weasels, and skunks.</i></p>

(Source of prompt: https://en.wikipedia.org/wiki/Red_panda, model used: Amazon Titan Text G1 - Express)

When prompting Claude models, it's helpful to wrap the input text in XML tags. In the following example, the input text is enclosed in `<text></text>`.

<p>Prompt template for Claude:</p> <pre>\n\nHuman: {{Instruction}} \n<text>\n{{Text}} \n</pre>	<p>User prompt:</p> <pre>Human: Read the following text inside\n<text></text>\nXML tags, and then answer the question:</pre>
---	---

<text>\n
{{Question}} \n
\n
Assistant:

<text>
On November 12, 2020, the selection of the Weeknd to headline the show was announced; marking the first time a Canadian solo artist headlined the Super Bowl halftime show. When asked about preparations for the show, the Weeknd stated, "We've been really focusing on dialing in on the fans at home and making performances a cinematic experience, and we want to do that with the Super Bowl."

The performance featured a choir whose members were dressed in white and wore masks over their faces with glowing red eyes, and were standing within a backdrop of a neon cityscape. The performance opened with a white figure dressed the same as the choir being lowered into the backdrop where the choir was standing while singing "Call Out My Name". At this time, the Weeknd sat in a convertible against a skyline backdrop designed to resemble the Las Vegas Strip. For the next part of the performance, the backdrop then split open to reveal the Weeknd, who then performed "Starboy", followed by "The Hills".

Next, performing the song "Can't Feel My Face", the Weeknd traveled through a labyrinth constructed behind the stage, joined by dancers dressed in red blazers and black neckties similar to his, but with their faces covered with bandages, in keeping with the aesthetic of his fourth studio album After Hours (2020). The dancers would wear these bandages throughout the performance. In the labyrinth section of the performance, camerawork was visually unsteady.

The next songs performed were "I Feel It Coming", "Save Your Tears", and "Earned It". For the "Earned It" performance, the Weeknd was accompanied by violinists. For the finale of the show, the Weeknd took to the field of the stadium with his dancers to perform "Blinding Lights". He and the dancers entered the field by performing "House of Balloons / Glass Table Girls". The performance ended with an array of fireworks.

</text>
Based on the text above, what songs did the Weeknd play at the Super Bowl halftime show?

	Assistant:
	Output: Based on the text, the songs the Weeknd played at the Super Bowl halftime show were: <ul style="list-style-type: none">- Call Out My Name- Starboy- The Hills- Can't Feel My Face- I Feel It Coming- Save Your Tears- Earned It- Blinding Lights- House of Balloons / Glass Table Girls

(Source of prompt: [Wikipedia on the Super Bowl LV halftime show](#), model used: Claude)

Summarization

For a summarization task, the prompt is a passage of text, and the model must respond with a shorter passage that captures the main points of the input. Specification of the output in terms of length (number of sentences or paragraphs) is helpful for this use case.

Prompt template for Amazon Titan and AI21 Jurassic: The following is text from a {{Text Category}}: {{Text}} \n Summarize the {{Text Category}} in {{length of summary, e.g., "one sentence" or "one paragraph"}}\n	User prompt: The following is text from a restaurant review: "I finally got to check out Alessandro's Brilliant Pizza and it is now one of my favorite restaurants in Seattle. The dining room has a beautiful view over the Puget Sound but it was surprisingly not crowded. I ordered the fried castelvetrano olives, a spicy Neapolitan-style pizza and a gnocchi dish. The olives were absolutely decadent, and the pizza came with a smoked mozzarella, which was delicious. The gnocchi was fresh and wonderful. The waitstaff were attentive, and overall the experience was lovely. I hope to return soon." Summarize the above restaurant review in one sentence.
	Output: Alessandro's Brilliant Pizza is a fantastic restaurant in Seattle with a beautiful view over Puget Sound that offers decadent and delicious food.

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

In the following example, Claude summarizes the given text in one sentence. To include input text in your prompts, format the text with XML mark up: `<text> {{text content}} </text>`. Using XML within prompts is a common practice when prompting Claude models.

Prompt template for Claude: \n\nHuman: Please read the text:\n<text>\n{{text context}}\n</text>\n\nSummarize the text in {{length of summary, e.g., "one sentence" or "one paragraph"}}\n\nAssistant:	User prompt: Human: Please read the text: <text> In game theory, the Nash equilibrium, named after the mathematician John Nash, is the most common way to define the solution of a non-cooperative game involving two or more players. In a Nash equilibrium, each player is assumed to know the equilibrium strategies of the other players, and no one has anything to gain by changing only one's own strategy. The principle of Nash equilibrium dates back to the time of Cournot, who in 1838 applied it to competing firms choosing outputs. </text> Summarize the text in one sentence. Assistant:
	Output: The Nash equilibrium is a concept in game theory where each player has adopted a strategy that they have no incentive to deviate from given the other players' strategies.

(Source of prompt: [Wikipedia on Nash equilibrium](#), model used: Claude)

Text generation

Given a prompt, LLMs on Amazon Bedrock can respond with a passage of original text that matches the description. Here is one example:

Prompt template for Amazon Titan and AI21 Jurassic: Please write a {{Text Category}} in the voice of {{Role}}.\n	User prompt: Please write an email in the voice of a friend congratulating someone on a new job.
	Output: Hey [Name],

I just found out about your new job, and I'm so happy for you! Congratulations! You deserve this opportunity since you worked so hard to get it. I know you'll do great and make a positive impact on your team. I'm here to support you every step of the way. Wishing you all the best in this exciting new chapter of your life!

Cheers,
[Your Name]

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

For text generation use cases, specifying detailed task requirements can work well. In the following example, we ask the model to generate response with exclamation points.

Prompt template for Amazon Titan and AI21 Jurassic: {{Task specification}} \n Please write a {{Text Category}} in the voice of {{Role}}.\n	User prompt: Write text with exclamation points. Please write a formatted email in the voice of a friend congratulating someone on a new job.
	Output: <i>Hey [name]!</i> <i>I just saw the news and I am so, so happy for you!</i> <i>You deserve this position and this company with all your heart.</i> <i>I can't wait to see you soar and accomplish all the things you've been working towards. If anyone is deserving of a promotion, it's you!</i> <i>Congratulations again, friend!</i> <i>Love,</i> <i>[name]</i>

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

In the following example, a user prompts the model to take on the role of a specific person when generating the text. Notice how the signature reflects the role the model is taking on in the response.

Prompt template for Amazon Titan and AI21 Jurassic: {{Role assumption}} {{Task description}}\n	User input: <i>My name is Jack. Help me write a note expressing my gratitude to my parents for taking my son (their grandson) to the zoo.</i> <i>I miss my parents so much.</i>
	Output:

	<pre>Dear Mom and Dad, I wanted to take a moment to express my heartfelt gratitude for taking my son to the zoo last weekend. It meant so much to me that you took the time to make memories with him, and I am so grateful for the love and support you have shown him. I miss you both so much, and I wish we could be together more often. But I know that you will always be there for me and my family, and I am so grateful to have you as my parents. Love, Jack</pre>
--	--

(Source of prompt: AWS, model used: AI21 Jurassic-2 Ultra v1)

Code generation

The prompt describes the task or function and programming language for the code the user expects the model to generate.

<p>Prompt template for Amazon Titan and AI21 Jurassic: Write a function in <code>{{programming language}}</code> to <code>{{task or purpose}}</code>\n.</p>	<p>User input: Write a python function to download a s3 file to local disk.</p>
	<p>Output: The following code can be used to download an s3 file to the local disk.</p> <pre>```Python import boto3 # Creates a boto3 client object s3 = boto3.client('s3') def download_file(bucket_name, object_key, file_path): try: # Downloads the s3 file to the specified file path s3.download_file(bucket_name, object_key, file_path) except Exception as e: print(e) # Specifies the values of bucket_name, object_key, and file_path bucket_name = "samplebucketname" object_key = "sampleobjectkey" file_path = "/samplefilepath" download_file(bucket_name, object_key, file_path) ```</pre> <p>The given code uses the <code>download_file</code> function to download an s3 file</p>

```
to the local disk. The download_file
function takes the bucket name,
object key, and file path as input
parameters. The bucket name
is the name of the s3 bucket where the file
is located,
the object key is the unique identifier of
the file within the bucket,
and the file path is the destination path
on the local disk
where the file will be downloaded. The
function uses the download_file method
of the s3 client to download the file from
the specified bucket
and write it to the specified file path.
```

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

Prompt template for Claude \n \nHuman: Write a function in {{programming language}} to {{task or purpose}} \n\nAssistant:	User prompt: Human: Write some python code that quickly approximates the square root of a number. Assistant:
	Output: Here is a simple function to approximate the square root of a number using Newton's method: ``python def sqrt(n, iterations=10): x = n for i in range(iterations): x = (x + n/x) / 2 return x ```\n To use it: ``python print(sqrt(25)) # 5.0 print(sqrt(15, iterations=20)) # 3.872983346207417 ```\n This repeatedly takes the average of x and n/x to converge on the square root. More iterations produce a more accurate approximation. The advantage of this approach is it converges very quickly, so even just 10 iterations produces a decent estimate.

(Source of prompt: AWS, model used: Claude)

Mathematics

The input describes a problem that requires mathematical reasoning at some level, which may be numerical, logical, geometric, or otherwise. For such problem, it's helpful to ask the model to work through the problem in a piecemeal manner by adding phrases to instructions such as Let's think step by step or Think step by step to come up with the right answer.

Prompt template for Amazon Titan and AI21 Jurassic: <code>{{Math problem description}}\nLet's think step by step. \n</code>	User prompt: <i>A triangle has two angles of 70 degrees and 50 degrees. What is the third angle in degrees? Let's think step by step.</i>
	Output: <i>The sum of the three angles in a triangle is always 180 degrees. The two given angles are 70 degrees and 50 degrees. $70 + 50 = 120$. Therefore, the third angle is $180 - 120 = 60$ degrees. The third angle is 60 degrees.</i>

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

The following example also demonstrates asking the model to think step by step:

Prompt template for Amazon Titan and AI21 Jurassic: <code>{{Math problem description}}\nThink step by step to come up with the right answer. \n</code>	User prompt: <i>Lucy has 12 colorful marbles, and she wants to share them equally with her 4 friends. How many marbles will each friend receive? Think step-by-step to come up with the right answer.</i>
	Output: <i>1. Lucy has 12 marbles and wants to share them equally with her 4 friends, then each friend will receive $12 / 4 = 3$ marbles. 2. Therefore, each friend will receive 3 marbles.</i>

(Source of prompt: AWS, model used: AI21 Jurassic-2 Ultra v1)

Reasoning/logical thinking

For complex reasoning tasks or problems that requires logical thinking, we can ask the model to make logical deductions and explain its answers.

Prompt template for Amazon Titan and AI21 Jurassic: <i>Question: {{question}} \n \n Please output the answer and then explain your answer:\n</i>	User prompt: <i>Question: Which word is the odd one out? A. accomplished B. good C. horrible D. outstanding Please output the answer and then explain your answer:</i>
	Output: <i>The odd one out is C. horrible. The other three words describe positive qualities or achievements, whereas horrible is an negative adjective.</i>

(Source of prompt: AWS, model used: Amazon Titan Text G1 - Express)

Here is another example using the Claude model:

Prompt template for Claude: <i>\n \n Human: {{question}} \n Please provide the answer and show the reasoning. \n \n Assistant:</i>	User prompt <i>Human: The barber is the "one who shaves all those, and those only, who do not shave themselves". Does the barber shave himself? Why is this a paradox? Please provide the answer and show the reasoning. Assistant:</i>
	Output: <i>This is an example of a paradox known as the Barber paradox. Here is an explanation: - The statement defines a barber who shaves all and only those people who do not shave themselves. - So if the barber shaves himself, he must shave all those who do not shave themselves, which includes himself. But the barber cannot shave himself as per the initial statement. - If the barber does not shave himself, then he would fit into the group of people who do not shave themselves, and so he must shave himself, according to the statement. But again this contradicts the initial statement. - So the barber cannot shave himself and also cannot not shave himself. This is self-contradictory based on the initial statement about the barber. - Therefore, the statement leads to a paradox where both possibilities</i>

(the barber shaving himself or not shaving himself) result in a contradiction.

In summary, the paradox arises because the definition of the barber's behavior is self-contradictory when applied to the barber himself.

This makes it impossible to determine if the barber shaves himself or not based on the given statement alone.

(Source of prompt: https://en.wikipedia.org/wiki/Barber_paradox, model used: Claude)

Embeddings

Text embeddings represent meaningful vector representations of unstructured text such as documents, paragraphs, and sentences. You input a body of text and the output is a (1 x n) vector. You can use embedding vectors for a wide variety of applications.

Bedrock supports one model for text embeddings, the Titan Embeddings G1 - Text model (`amazon.titan-embed-text-v1`). This model supports text retrieval, semantic similarity, and clustering. The maximum input text is 8K tokens and the maximum output vector length is 1536.

To use a text embeddings model, use the [InvokeModel](#) API operation with `amazon.titan-embed-text-v1` as the `modelId` and retrieve the `embedding` object in the response.

To see Jupyter notebook examples:

1. Sign in to the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/home>.
2. From the left-side menu, choose **Base models**.
3. Scroll down and select the **Titan Embeddings G1 - Text** model.
4. In the **Titan Embeddings G1 - Text** tab, select **View example notebook** to see example notebooks for embeddings.

Custom models

Note

To use fine-tuning, you must have access to the Amazon Titan Text G1 - Express model. The Amazon Titan Text G1 - Express model is in limited preview. To request access, contact your AWS account manager.

You can customize a Amazon Bedrock model to improve its performance and create a better customer experience. Amazon Bedrock currently offers the ability to fine-tune a model by providing your own labeled training data to help retain or improve its accuracy while employing smaller datasets and therefore reducing the training time.

After you complete a model customization job, you can purchase Provisioned Throughput (see [Provisioned Throughput \(p. 90\)](#)) for the customized model so that you can use the model for inference using the [InvokeModel](#) or [InvokeModelWithResponseStream](#) API operations or the text playground (see [Text playground \(p. 14\)](#)).

For information about quotas for model customization, see [Model customization quotas \(p. 176\)](#).

Topics

- [Fine-tuning \(p. 74\)](#)
- [Prepare the datasets \(p. 74\)](#)
- [Using the console \(p. 75\)](#)
- [Using the API \(p. 78\)](#)
- [Guidelines for model customization \(p. 87\)](#)
- [Troubleshooting \(p. 88\)](#)

Fine-tuning

Fine-tune a Amazon Bedrock model by providing your own labeled training dataset in order to improve the model's performance on specific tasks. By providing some labeled examples related to a specific task, you help the model learn the task it's supposed to carry out. Through this process, known as *few-shot fine-tuning*, you can create a new model that improves upon the performance and efficiency of the original model for a given task.

To fine-tune a model, you upload a training dataset and, optionally, a validation dataset to Amazon S3 and provide the Amazon S3 bucket path to the Bedrock fine-tuning job. You can also adjust the hyperparameters for fine-tuning. To fine-tune a model in the console, follow the steps at [Using the console \(p. 75\)](#). To fine-tune through the API, follow the steps at [Using the API \(p. 78\)](#).

Prepare the datasets

Before you upload your training and validation data to Amazon S3, you need to pre-process the format of your data into JSONL format training and validation datasets.

If you provide a validation dataset (which is optional), Amazon Bedrock returns validation loss metrics at the end of the model customization job.

Note

The following formats pertain to fine-tuning a Titan Text G1 - Express model.

For both the training and optional validation datasets, create a file with multiple JSON lines. Each JSON line contains both an input and output field. The format is as follows. For quotas, see [Quotas for Amazon Bedrock \(p. 176\)](#). Use 6 characters per token as an approximation for the number of tokens.

```
{ "input": "<prompt text>", "output": "<expected generated text>" }  
{ "input": "<prompt text>", "output": "<expected generated text>" }  
{ "input": "<prompt text>", "output": "<expected generated text>" }
```

The following is an example item for a question-answer task:

```
{ "input": "what is AWS", "output": "it's amazon web services" }
```

After you create your JSONL format training dataset and optional validation dataset, do the following:

1. [Upload](#) the datasets to an Amazon S3 bucket.
2. Provide permissions for Amazon Bedrock to access the data by attaching an IAM policy similar to the example shown in [Grant custom jobs access to your training data \(p. 159\)](#) to your Amazon Bedrock service role.
3. Submit a fine tuning job with the [console \(p. 75\)](#) or with the [API \(p. 78\)](#).

Using the console

The following topics describe the major steps required to fine-tune a model:

Topics

- [Submit a job \(p. 75\)](#)
- [Monitor the job \(p. 76\)](#)
- [Stop a job \(p. 77\)](#)
- [Analyze the job results \(p. 77\)](#)
- [Use a fine-tuned model for inference \(p. 78\)](#)

Submit a job

Note

Before you begin these steps, do the following:

- Grant Amazon Bedrock permissions to access the training and validation data and to write the output data by attaching an IAM policy similar to the example shown in [Grant custom jobs access to your training data \(p. 159\)](#) to your Amazon Bedrock service role.
- Create a training dataset and optionally a validation dataset. For more information, see [Prepare the datasets \(p. 74\)](#).

To submit a fine-tuning job

1. Sign in to the AWS Management Console and open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/home>.
2. From the left menu, choose **Custom models**.
3. Select **Fine-tune model**.
4. Configure the settings for the fine-tuning job:

- a. In the **Model details** section, choose the **Source model** that you want to fine-tune with your own data and give your resulting fine-tuned model a name in **Fine-tuned model name**.

Note

Currently, Titan Text G1 - Express is the only supported model for fine-tuning.

- b. (Optional) You can choose **Model encryption** to choose a different KMS key from the default key.

For information about configuring KMS keys, see [Data encryption \(p. 139\)](#).

- c. (Optional) You can expand the **Tags** section and select **Add new tag** add tags to associate with the model.
- d. In the **Job configuration** section, enter a name for the fine-tuning training job in **Job name** and add any tags to associate with the job.
- e. (Optional) In the **VPC settings** section, choose the **VPC** that contains the input data and output data Amazon S3 locations.
 - i. For **Subnet(s)**, add the VPC subnets.
 - ii. For **Security group(s)**, add security groups to control access to the data in your VPC.

For information about configuring the VPC, see [Protect jobs using a VPC \(p. 141\)](#).

- f. Provide the link to the training dataset file and, optionally, the validation dataset file in your Amazon S3 bucket in the **Input data** section.
- g. In the **Hyperparameters** section, input the values for the following hyperparameters to use in training. See [Guidelines for model customization \(p. 87\)](#) for guidelines on choosing values.
 - **Epochs** – The number of times to pass the training dataset to the model.
 - **Batch size** – The number of records to send to the model in each batch.
 - **Learning rate** – The step size for parameter updates in each iteration.
 - **Learning rate warmup steps** – Affects the speed at which the algorithm converges to the optimal weights.
- h. Provide a link to the Amazon S3 folder in which to save the output of the fine-tuning job in the **Output data** section. The training loss metrics and validation loss metrics for each epoch will be stored in separate files in the location that you specify.
- i. You use a service role to provide permissions for Amazon Bedrock to write to Amazon S3 on your behalf.

In the **Service access** section, select one of the following:

- **Use an existing service role** – Select a service role from the drop-down list.
- **Create and use a new service role** – Enter a name for the service role.

Note

If your job includes VPC configuration, the console cannot create a new service role for the job. Create the service role using the example described in [Configure your model customization job to use VPC \(p. 145\)](#).

5. Select **Fine-tune model** to begin the job.

Monitor the job

The fine-tuning job can take several hours. The duration of the job depends on the size of the training data (number of records, input tokens, and output tokens), number of epochs, and batch size.

To monitor the status of the fine-tuning job

1. Open the Amazon Bedrock console.
2. From the left menu, choose **Custom models**.
3. The **Training jobs** tab displays the fine-tuning jobs that you have initiated. Look at the **Status** column to monitor the progress of the job.
4. Select a job to view the details you input for training.

Stop a job

You can stop a Bedrock fine-tuning job while it's in progress. You can't resume a stopped job.

To stop a fine-tuning job

1. Open the Amazon Bedrock console.
2. From the left menu, chose **Custom models**.
3. From the **Training Jobs** tab, choose the radio button next to the job to stop.
4. Select the **Stop job** button.
5. A modal appears to warn you that you can't resume the training job if you stop it. Select **Stop job** to confirm.

Note

Amazon Bedrock charges for the tokens that it used to train the model before you stopped the job. Amazon Bedrock doesn't create an intermediate custom model for a stopped job.

Analyze the job results

After the fine-tuning job completes, you can see the following information in your S3 output folder:

- Completion status
- Training and validation loss metrics

Amazon Bedrock stores your customized models in AWS-managed storage scoped to your account. You can see your customized model in the **Models** table of the model customization dashboard page. Choose a model to view details associated with that custom model. Output models aren't downloadable.

The S3 output for a fine-tuning job contains the following output files in your S3 folder:

```
- model-customization-job-training_job_id/
  - training_artifacts/
    - step_wise_training_metrics.csv
  - validation_artifacts/
    - post_fine_tuning_validation/
      - validation_metrics.csv
```

Use the `step_wise_training_metrics.csv` and the `validation_metrics.csv` files to analyze the model customization job and to help you adjust the model as necessary.

The structure of the `step_wise_training_metrics.csv` file is shown in the following example:

step_number	epoch_number	training_loss	perplexity
1	1	0.2	25

step_number	epoch_number	training_loss	perplexity
2	1	0.18	22
...	18

The structure of the `validation_metrics.csv` file is shown in the following example:

step_number	epoch_number	validation_loss	perplexity
1	1	0.12	20
2	1	0.09	17
...	15

Use a fine-tuned model for inference

Before you can use a customized model for inference, you need to purchase Provisioned Throughput for it. You can then use it for inference in the Text or Chat playground.

To purchase Provisioned Throughput for a custom model.

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.
2. From the left side navigation pane, select **Custom models**.
3. Choose the radio button next to the model, select **Actions**, and choose **Purchase Provisioned Throughput**. For more information, see [Provisioned Throughput \(p. 90\)](#).

To use a fine-tuned model for inference

1. From the left side navigation pane, select **Custom models**.
2. In the **Models** tab, choose the model that you want to use in the Text or Chat Playground and select **Open in playground**.

Using the API

This section demonstrates how to fine-tune your models using API operations. We provide examples with the AWS Command Line Interface and the AWS SDK for Python (Boto3).

Topics

- [Set up an IAM role for model customization \(p. 79\)](#)
- [Submit a job \(p. 81\)](#)
- [Monitor a job \(p. 84\)](#)
- [Stop a job \(p. 85\)](#)
- [Analyze a job \(p. 85\)](#)
- [Retrieve information about your customized models \(p. 86\)](#)

Set up an IAM role for model customization

Create an IAM role or use an existing role that grants the required permissions for your customization job to access the data in your Amazon S3 buckets. Use the IAM [CreateRole](#) API operation to allow Amazon Bedrock to assume a role and the IAM [CreatePolicy](#) API operation to create a policy to allow the role to access the Amazon S3 buckets containing your training, validation, and output data. Attach the policy with the [AttachRolePolicy](#) API operation. Skip this step if you are using an existing role with the correct IAM permissions. If you submit a job using the console (see [Submit a job \(p. 75\)](#)), the console can create a role with the correct permissions for you.

Preparing the role and access policies:

1. Create the following policy document called *BedrockAssumeRolePolicy.json* (or another name of your choice) to allow Amazon Bedrock to assume a role:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "bedrock.amazonaws.com"
      },
      "Action": "sts:AssumeRole",
      "Condition": {
        "StringEquals": {
          "aws:SourceAccount": "111122223333"
        },
        "ArnEquals": {
          "aws:SourceArn": "arn:aws:bedrock:us-east-1:111122223333:model-
customization-job/*"
        }
      }
    }
  ]
}
```

2. Determine the Amazon S3 locations where you uploaded your training and validation data and the location for Amazon Bedrock to upload your training metrics data.
3. Create the following policy document called *BedrockAccessTrainingValidationS3Policy.json* (or another name of your choice) to allow access to these Amazon S3 locations, replacing the values in the Resource list with actual Amazon S3 ARNs:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:PutObject",
        "s3:ListBucket",
        "s3:ListObjects"
      ],
      "Resource": [
        "arn:aws:s3:::my_training_data_bucket/myfolder",
        "arn:aws:s3:::my_training_data_bucket/myfolder/*",
        "arn:aws:s3:::my_validation_data_bucket/myfolder",
        "arn:aws:s3:::my_validation_data_bucket/myfolder/*"
      ]
    }
  ]
}
```

```
}  
]  
}
```

4. Create the following policy called *BedrockAccessOutputS3Policy.json* to allow both access and writing to these Amazon S3 locations, replacing the values in the Resource list with actual Amazon S3 ARNs:

```
{  
  "Version": "2012-10-17",  
  "Statement": [  
    {  
      "Effect": "Allow",  
      "Action": [  
        "s3:GetObject",  
        "s3:PutObject",  
        "s3:ListBucket",  
        "s3:ListObjects"  
      ],  
      "Resource": [  
        "arn:aws:s3::my_output_bucket/myfolder",  
        "arn:aws:s3::my_output_bucket/myfolder/*"  
      ]  
    }  
  ]  
}
```

AWS CLI

1. Create an IAM role, replacing the role-name with a name of your choice and using the file name for the assume role policy you created:

```
aws iam create-role \  
  --role-name "MyBedrockModelCustomizationRole" \  
  --assume-role-policy-document file://path/to/BedrockAssumeRolePolicy.json
```

2. Create the training and validation data access policy, replacing the policy-name with a name of your choice:

```
aws iam create-policy \  
  --policy-name "MyBedrockTrainingValidationS3BucketPolicy" \  
  --policy-document file://path/to/BedrockAccessTrainingValidationS3Policy.json
```

3. Create the output data access policy, replacing the policy-name with a name of your choice:

```
aws iam create-policy \  
  --policy-name "MyBedrockOutputS3BucketPolicy" \  
  --policy-document file://path/to/BedrockAccessOutputS3Policy.json
```

4. Attach the IAM role training and data access policy, replacing the role-name with the role name you created in step 1 and the policy-arn with the ARN that the response in step 2 returned to you:

```
aws iam attach-role-policy \  
  --role-name "MyBedrockModelCustomizationRole" \  
  --policy-arn "PolicyARNFromCreateTrainingValidationPolicyStep"
```

5. Attach the IAM role output data access policy, replacing the `role-name` with the role name you created in step 1 and the `policy-arn` with the ARN that the response in step 3 returned to you:

```
aws iam attach-role-policy \  
  --role-name "MyBedrockModelCustomizationRole" \  
  --policy-arn "PolicyARNFromCreateOutputPolicyStep"
```

6. Use the ARN returned to you in step 1 as the `roleArn` when submitting a model customization job.

Python

1. Create an IAM role by following the code sample in the **Create a role** section in [IAM examples using SDK for Python \(Boto3\)](#), replacing the `role_name` with a name of your choice and replacing the `trust_policy` with the contents from the *BedrockAssumeRolePolicy.json* you created.
2. Create two access policies by following the code sample in the **Create a policy** section in [IAM examples using SDK for Python \(Boto3\)](#), replacing the `PolicyName` with two different names of your choice and replacing the `policy_doc` with the contents from the *BedrockAccessTrainingValidationS3Policy.json* and the *BedrockAccessOutputS3Policy.json* that you created.
3. Attach the IAM role policies, by following the code sample in the **Attach a policy to a role** section in [IAM examples using SDK for Python \(Boto3\)](#), replacing replacing the `role_name` with the role name you created in step 1 and the `policy_arn` with the ARNs that the responses in step 2 returned to you.
4. Use the ARN returned to you in step 1 as the `roleArn` when submitting a model customization job.

Submit a job

Use the Amazon Bedrock [CreateModelCustomizationJob](#) API operation to submit a model customization job. Minimally, you must provide the following fields in the [CreateModelCustomizationJob](#) request:

- `customizationType` – To fine-tune a model, use the value **FINE_TUNING**.
- `baseModelIdentifier` – The ARN of the model to customize.
- `customModelName` – The name to give the newly customized model.
- `hyperParameters` – Parameters related to tuning the model. For the Titan Text G1 - Express model, the following fields can be specified in this object. See [Guidelines for model customization \(p. 87\)](#) for guidelines on choosing values.
 - **Epochs** – The number of times to pass the training dataset to the model.
 - **Batch size** – The number of records to send to the model in each batch.
 - **Learning rate** – The step size for parameter updates in each iteration.
 - **Learning rate warmup steps** – Affects the speed at which the algorithm converges to the optimal weights.
- `jobName` – The name to give the training job.
- `roleArn` – The ARN of the service role.
- `trainingDataConfig` – An object containing the URI of the Amazon S3 location of the training dataset. To create the training dataset, see [Prepare the datasets \(p. 74\)](#).
- `validationDataConfig` – An object containing the URI of the Amazon S3 location of the validation dataset. To create the validation dataset, see [Prepare the datasets \(p. 74\)](#).

- `outputDataConfig` – An object containing the URI of the Amazon S3 location to write the output data to.

The response returns a `jobArn` that you can use to monitor or stop the model customization job.

AWS CLI

The following example demonstrates how to submit a model customization job using the AWS CLI:

First create a text file named *FineTuningData.json*. Copy the JSON code from below into the text file, replacing the bucket, path, and file names with the correct paths to your training, validation, and output data:

```
{
  "trainingDataConfig": {
    "s3Uri": "s3://bucket/path/to/train.jsonl"
  },
  "validationDataConfig": {
    "validators": [{
      "s3Uri": "s3://bucket/path/to/validation.jsonl"
    }]
  },
  "outputDataConfig": {
    "s3Uri": "s3://bucket/path/to/output-folder"
  }
}
```

Run the following command in the command line

Note

Currently, Titan Text G1 - Express is the only supported model for fine-tuning.

```
aws bedrock create-model-customization-job \
  --customization-type "FINE_TUNING" \
  --base-model-identifier "arn:aws:bedrock:us-east-1::foundation-model/foundation-model-id" \
  --role-arn "arn:aws:iam::arn-for-MyBedrockModelCustomizationRole" \
  --job-name "job-name" \
  --custom-model-name "custom-model-name" \
  --hyper-parameters
    epochCount="1",batchSize="1",learningRate="0.005",learningRateWarmupSteps="0" \
  --cli-input-json file://path/to/FineTuningData.json
```

To add a VPC configuration, add the following argument to the above command to specify the security group and subnets:

```
--vpc-config '{securityGroupIds: ["xx"], "subnetIds": ["subnet-yy", "subnet-zz"]}'
```

To encrypt your model with a KMS key, add the following argument to the above command, replacing the values to specify the key with which you want to encrypt your model.

```
--customModelKmsKeyId 'arn:aws:kms:region:account-id:key/key-id'
```

To add tags, add the following argument to the above command, replacing the keys and values with the tags you want to attach to the job and/or output model and making sure to separate key/value pairs with a space:

```
--tags key=key1,value=value1 key=key2,value=value2
```

Use the jobArn that the operation returns to check the status of the job or to analyze or stop the job.

Python

The following example demonstrates how to submit a model customization job using Python. Uncomment the relevant sections to add optional tags to the job and/or resulting model:

```
import boto3
import json
bedrock = boto3.client(service_name='bedrock')

# Set parameters
customizationType = "FINE_TUNING"
baseModelIdentifier = "arn:aws:bedrock:us-east-1::foundation-model/foundation-model-id"
roleArn = "arn:aws:iam::arn-for-MyBedrockModelCustomizationRole"
jobName = "job-name"
customModelName = "custom-model-name"
hyperParameters = {
    "epochCount": "1",
    "batchSize": "1",
    "learningRate": "0.005",
    "learningRateWarmupSteps": "0"
}
trainingDataConfig = {"s3Uri": "s3://bucket/path/to/train.jsonl"}
validationDataConfig = {
    "validators": [{
        "name": "validation",
        "s3Uri": "s3://bucket/path/to/validation.jsonl"
    }]
}
outputDataConfig = {"s3Uri": "s3://bucket/path/to/"}
```

Uncomment to add optional tags

```
# jobTags = [
#     {
#         "key": "key1",
#         "value": "value1"
#     }
# ]
# customModelTags = [
#     {
#         "key": "key1",
#         "value": "value1"
#     }
# ]

# Create job
bedrock.create_model_customization_job(
    jobName=jobName,
    customModelName=customModelName,
    roleArn=roleArn,
    baseModelIdentifier=baseModelIdentifier,
    hyperParameters=hyperParameters,
    # # Uncomment to add optional tags
```

```
#    jobTags=jobTags,  
#    customModelTags=customModelTags,  
trainingDataConfig=trainingDataConfig,  
validationDataConfig=validationDataConfig,  
outputDataConfig=outputDataConfig  
)
```

Use the `jobArn` that the code returns to check the status of the job or to analyze or stop the job.

Monitor a job

To list all your model customization jobs, send an Amazon Bedrock [ListModelCustomizationJobs](#) request. To monitor the status of a model customization job, send an Amazon Bedrock [GetModelCustomizationJob](#) request by entering the `jobArn`. You can find it in one of the following ways:

1. In the Amazon Bedrock console, select **Fine-tune** and choose the job from the **Training jobs** table. Look for the **Model customization job ARN** in the **Model configuration** section.
2. Look in the `jobArn` field in the response returned from the `CreateModelCustomizationJob` call that created the job.

Monitor the progress of the job by looking in the `status` field of the response.

AWS CLI

List your model customization jobs using the AWS CLI with the following command:

```
aws bedrock list-model-customization-jobs
```

The following example demonstrates how to monitor a model customization job using the AWS CLI:

```
aws bedrock get-model-customization-job \  
--job-identifier "arn:aws:bedrock:job-arn-from-create-model-customization"
```

Find the value of the `status` field in the response.

Python

The following demonstrates how to list your model customization jobs using Python:

```
import boto3  
bedrock = boto3.client(service_name='bedrock')  
  
bedrock.list_model_customization_jobs()
```

The following example demonstrates how to monitor a model customization job:

```
import boto3  
bedrock = boto3.client(service_name='bedrock')  
  
fine_tune_job = bedrock.get_model_customization_job(jobIdentifier='arn:aws:bedrock:job-arn-from-create-model-customization')  
print(fine_tune_job['status'])
```

Stop a job

To use the API to stop a model customization job, follow these steps:

1. If the job status in the [GetModelCustomizationJob](#) response is IN_PROGRESS, send a [StopModelCustomizationJob](#) request with the jobArn of the training job. The system marks the job for termination and sets the state to STOPPING.
2. The system stops the job and sets the state to STOPPED.

If the job completes before the system stops it, the system sets the state to COMPLETED.

Note

Amazon Bedrock charges for the tokens that it used to train the model before you stopped the job. Amazon Bedrock doesn't create an intermediate custom model for a stopped job.

AWS CLI

The following example demonstrates how to stop a model customization job using the AWS CLI:

```
aws bedrock stop-model-customization-job \
  --job-identifier "arn:aws:bedrock:job-arn-from-create-model-customization"
```

Python

The following example demonstrates how to stop a model customization job using Python:

```
import boto3
bedrock = boto3.client(service_name='bedrock')

bedrock.stop_model_customization_job(jobIdentifier='arn:aws:bedrock:job-arn-from-
create-model-customization')
```

Analyze a job

To analyze a model customization job, send an Amazon Bedrock [GetCustomModelb](#) request by entering the jobArn. You can find it in one of the following ways:

1. In the Amazon Bedrock console, select **Fine-tune** and choose the job from the **Training jobs** table. Look for the **Model customization job ARN** in the **Model configuration** section.
2. Look in the jobArn field in the response returned from the CreateModelCustomizationJob call that created the job.

You can analyze metrics by using the trainingLoss value in the trainingMetrics object and the validatorLoss values in the validationMetrics object to analyze the model customization job or by reading the data from the S3 output file.

AWS CLI

The following example demonstrates how to analyze a model customization job using the AWS CLI:

```
aws bedrock get-model-customization-job \
  --job-identifier "arn:aws:bedrock:job-arn-from-create-model-customization"
```

Look in the trainingMetrics and validatorLoss fields to analyze the model customization job.

Python

The following example demonstrates how to analyze a model customization job using Python:

```
import boto3
import json
import pandas as pd

bedrock = boto3.client(service_name='bedrock')

fine_tune_job = bedrock.get_model_customization_job(jobIdentifier='arn:aws:bedrock:job-arn-from-create-model-customization')

# Get S3 URI
fine_tune_job['outputDataConfig']['s3Uri']
```

Use the S3 URI to download the output data. For more information, see [Downloading objects](#). Open the folder and navigate to the metrics files.

Retrieve information about your customized models

Model customization creates a custom model upon successful completion of the job. To list your custom models, send an Amazon Bedrock [ListCustomModels](#) request. To get information about a specific model that you have fine-tuned, send an Amazon Bedrock [GetCustomModel](#) request, providing the ARN of the model as the `modelIdentifier` in the request. You can find the ARN of the model in one of the following ways:

1. In the Amazon Bedrock console, select **Fine-tune** and choose the model from the **Models** table. Look for the **Fine-tuned model ARN** in the **Model details** section.
2. Look in the `outputModelArn` field in the response returned from the `GetModelCustomizationJob` call that created the job. This field only appears in the response after the job has finished.

AWS CLI

Use the following command to list your fine-tuned models in the AWS CLI:

```
aws bedrock list-custom-models
```

The following example demonstrates how to get information about a fine-tuned model in the AWS CLI

```
aws bedrock get-custom-model \
  --model-identifier "arn:aws:bedrock:customized-model-arn"
```

Python

Use the following command to list your fine-tuned models in Python:

```
import boto3
import json
bedrock = boto3.client(service_name='bedrock')

bedrock.list_custom_models()
```

The following example demonstrates how to get information about a fine-tuned model in the AWS CLI

```
import boto3
import json
bedrock = boto3.client(service_name='bedrock')

bedrock.get_custom_model(modelIdentifier='arn:aws:bedrock:customized-model-arn')
```

Guidelines for model customization

In this section, we provide guidelines and recommended values as a baseline for customization of the Titan Text G1 - Express model. However, you should experiment with values to determine which parameters work best for your specific case. The ideal parameters depend on the dataset and the task for which the model is intended.

Use the training and validation metrics from the output files generated when you submit a fine-tuning job to help you adjust your parameters. Find these files in the Amazon S3 bucket to which you wrote the output, or use the [GetCustomModel](#) operation.

Size of the input training dataset

In general, the larger the training dataset, the better the performance for a specific task. However, if there are too many samples, the model might perform worse on a different task (for example, if the training dataset for a summarization task contains 100,000 samples, the model might perform worse on a classification task).

Large language models perform better than small supervised models when there is limited training data (for example, 1 to 100 samples).

Model size

In general, the larger the model, the better the task performs given limited training data.

If you are using the model for a *classification* task, you might see relatively small gains for few-shot fine-tuning (less than 100 samples), especially if the number of classes is relatively small (less than 100).

Epochs

We recommend using the following metrics to determine the number of epochs to set:

1. **Validation output accuracy** – Set the number of epochs to one that yields a high accuracy.
2. **Training and validation loss** – Determine the number of epochs after which the training and validation loss becomes stable. This corresponds to when the model converges. Find the training loss values in the `step_wise_training_metrics.csv` and `validation_metrics.csv` files.

Batch size

When you change the batch size, we recommend that you change the learning rate using the following formula:

```
newLearningRate = oldLearningRate x newBatchSize / oldBatchSize
```

Learning rate

In general, use smaller learning rates for larger models. We recommend using a learning rate in the range of 1.00E-06 to 1.00E-05. This parameter doesn't appear to play a big role in question-answer and classification tasks. However, it might dramatically impact performance in summarization tasks.

The following table shows recommended learning rate values for few-shot fine-tuning:

Task	Minimum learning rate	Default learning rate	Max learning rate
Summarization	1.00E-06	3.00E-06	5.00E-05
Classification	5.00E-06	5.00E-05	5.00E-05
Question-answer	5.00E-06	5.00E-06	5.00E-05

Learning warmup steps

We recommend the default value of 0.

Troubleshooting

This section summarizes errors that you might run into in fine-tuning and what to check if you come across them.

Permissions issues

If you encounter an issue with permissions to access an Amazon S3 bucket, check that the following are true:

1. If the Amazon S3 bucket uses a CM-KMS key for Server Side encryption, ensure that the IAM role passed to Amazon Bedrock has `kms:Decrypt` permissions for the AWS KMS key. For example, see [Allow a user to encrypt and decrypt with any AWS KMS key in a specific AWS account](#).
2. The Amazon S3 bucket is in the same region as the Amazon Bedrock model customization job.
3. The IAM role trust policy includes the service SP (`bedrock.amazonaws.com`).

The following messages indicate issues with permissions to access training or validation data in an Amazon S3 bucket:

```
Could not validate GetObject permissions to access Amazon S3 bucket: training-data-bucket
at key train.jsonl
Could not validate GetObject permissions to access Amazon S3 bucket: validation-data-bucket
at key validation.jsonl
```

If you encounter one of the above errors, check that the IAM role passed to the service has `s3:GetObject` and `s3:ListBucket` permissions for the training and validation dataset Amazon S3 URIs. For example, see [Submit a job \(p. 81\)](#).

The following message indicates issues with permissions to write the output data in an Amazon S3 bucket:

Amazon S3 perms missing (PutObject): Could not validate PutObject permissions to access S3 bucket: `bedrock-output-bucket` at key `output/.write_access_check_file.tmp`

If you encounter the above error, check that the IAM role passed to the service has `s3:PutObject` permissions for the output data Amazon S3 URI. For example, see [Submit a job \(p. 81\)](#).

Data issues

The following errors are related to issues with the training, validation, or output data files:

Invalid file format

Unable to parse Amazon S3 file: `fileName.jsonl`. Data files must conform to JSONL format.

If you encounter the above error, check that the following are true:

1. Each line is in JSON.
2. Each JSON has two keys, an `input` and an `output`, and each key is a string. For example:

```
{
  "input": "this is my input",
  "output": "this is my output"
}
```

3. There are no additional new lines or empty lines.

Character quota exceeded

Input size exceeded in file `fileName.jsonl` for record starting with...

If you encounter an error beginning with the text above, ensure that the number of characters conforms to the character quota in [Fine-tuning quotas \(p. 177\)](#).

Token count exceeded

```
Maximum input token count 4097 exceeds limit of 4096
Maximum output token count 4097 exceeds limit of 4096
Max sum of input and output token length 4097 exceeds total limit of 4096
```

If you encounter the above error, ensure that the number of tokens conforms to the token quota in [Fine-tuning quotas \(p. 177\)](#).

Internal error

Encountered an unexpected error when processing the request, please try again

If you encounter the above error, there might be an issue with the service. Try the job again. If the issue persists, contact AWS Support.

Provisioned Throughput

When you configure Provisioned Throughput for a model, you receive a level of throughput at a fixed cost.

You can use Provisioned Throughput with Amazon and third-party base models, and with customized models.

Provisioned Throughput pricing varies depending on the model that you use and the level of commitment you choose. You receive a discounted rate when you commit to a longer period of time. For details about pricing for each model, see the [Model providers](#) page in the Amazon Bedrock console.

Your options for throughput for a model differ depending on whether you run inference on a base model or a custom model.

Pricing option	Base model	Custom model
Provisioned Throughput, no commitment (hourly pricing)	Not available	Available (maximum 2 Provisioned Throughputs per account)
Provisioned Throughput, 1 month commitment	Available	Available
Provisioned Throughput, 6 month commitment	Available	Available

You specify Provisioned Throughput in Model Units (MU). A model unit delivers a specific throughput level for the specified model. The throughput level of a MU for a given Text model specifies the following:

- **The total number of input tokens per minute** – The number of input tokens that an MU can process across all requests within a span of one minute.
- **The total number of output tokens per minute** – The number of output tokens that an MU can generate across all requests within a span of one minute.

Model unit quotas depend on the level of commitment you specify for the Provisioned Throughput.

- For custom models with no commitment, a quota of one model unit is available for each Provisioned Throughput. You can create up to two Provisioned Throughputs per account.
- For base or custom models with commitment, there is a default quota of 0 model units. To request an increase, use the [limit increase form](#).

Topics

- [Provisioned Throughput model IDs \(p. 91\)](#)
- [Procedures \(p. 91\)](#)
- [Permissions \(p. 92\)](#)
- [Provisioned Throughput console procedures \(p. 92\)](#)
- [Using the Provisioned Throughput API \(p. 94\)](#)

Provisioned Throughput model IDs

Provisioned throughput is currently available for the following models. Use the corresponding model ID when using the [CreateProvisionedModelThroughput](#) API.

Model name	Model ID for provisioned throughput
Titan Text G1 - Express 8K	amazon.titan-text-express-v1:0:8k
Titan Embeddings G1 - Text	amazon.titan-embed-text-v1:2:8k
Anthropic Claude V2 18K	anthropic.claude-v2:0:18k
Anthropic Claude V2 100K	anthropic.claude-v2:0:100k
Anthropic Claude Instant V1 100K	anthropic.claude-instant-v1:2:100K
Stable Diffusion XL 1.0	stability.stable-diffusion-xl-v1:0
Meta Llama 2 Chat 13B	meta.llama2-13b-chat-v1:0:4k

Procedures

Provisioned Throughput supports the following procedures.

Creating

When you create a Provisioned Throughput, the Provisioned Throughput starts in an interim state (**Creating**) while Bedrock creates the Provisioned Throughput resources. If the creation is successful, the Provisioned Throughput transitions to **InService** state, and you can start using it for inference.

If the creation fails, the Provisioned Throughput transitions to **Failed** state.

Updating

You can perform the following types of updates to a Provisioned Throughput:

- Change the name of the Provisioned Throughput.
- Specify a new custom model that uses the same base model as the current custom model.
- Specify a base model. It must be the base model of the current custom model.

When you update a Provisioned Throughput, it transitions to the **Updating** state while Bedrock performs the update. If the update is successful, the Provisioned Throughput transitions to **InService** state. During the update, you can run inference using the Provisioned Throughput without disrupting the on-going traffic from your end customers. If the update includes a new model, you may receive output from the old model until the update is fully deployed.

If the update fails, the Provisioned Throughput transitions to **Failed** state.

Deleting

Deleting a Provisioned Throughput is a synchronous operation. The Provisioned Throughput deletion takes effect immediately.

Running inference

You can run inference using a Provisioned Throughput that is in `InService` state.

If you request more throughput than is configured for the Provisioned Throughput, the request is throttled (you receive the throughput defined by the Provisioned Throughput).

Permissions

To add Provisioned Throughput to a base model or a model that was customized from a base model, you must have first requested access to the base model by following the steps at [Model access \(p. 7\)](#).

To add Provisioned Throughput to a custom model that is encrypted, your role must include permissions to carry out the `kms:Decrypt` action on the model. For an example, see [Use a customer managed key during inference \(p. 141\)](#).

Provisioned Throughput console procedures

This section describes the console procedures for Provisioned Throughput.

Topics

- [View Provisioned Throughput summary \(p. 92\)](#)
- [Purchase Provisioned Throughput \(p. 92\)](#)
- [View details of a Provisioned Throughput \(p. 93\)](#)
- [Edit a Provisioned Throughput \(p. 94\)](#)
- [Delete a Provisioned Throughput \(p. 94\)](#)

View Provisioned Throughput summary

Use the summary page to review the status of each Provisioned Throughput. The **Overview** panel displays the number of Provisioned Throughput resources in each state. You can update or delete an active Provisioned Throughput, and you can create a new Provisioned Throughput.

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.
2. From the left navigation pane, choose **Provisioned throughput** under **Deployments**.
3. From the **Provisioned Throughput** table, you can review summary information about each Provisioned Throughput.

Purchase Provisioned Throughput

You can purchase Provisioned Throughput for a foundation or custom model.

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.
2. From the left menu, choose **Provisioned Throughput** under **Deployments**. Then select **Purchase Provisioned Throughput**. The console opens the **Purchase Provisioned Throughput** page.
3. Under **Provisioned Throughput details**:

- a. Enter a name for the Provisioned Throughput.
- b. Select the model category and model for the Provisioned Throughput.

Note

Provisioned throughput is currently available for the following models.

Model name	Model ID for provisioned throughput
Titan Text G1 - Express 8K	amazon.titan-text-express-v1:0:8k
Titan Embeddings G1 - Text	amazon.titan-embed-text-v1:2:8k
Anthropic Claude V2 18K	anthropic.claude-v2:0:18k
Anthropic Claude V2 100K	anthropic.claude-v2:0:100k
Anthropic Claude Instant V1 100K	anthropic.claude-instant-v1:2:100K
Stable Diffusion XL 1.0	stability.stable-diffusion-xl-v1:0
Meta Llama 2 Chat 13B	meta.llama2-13b-chat-v1:0:4k

- c. (Optional) Under **Tags**, you can associate one or more tags with this Provisioned Throughput.
4. Under **Model units & commitment term**:
 - a. Enter the desired number of model units.
 - b. Choose the amount of time for which you want to commit to using the Provisioned Throughput. To opt out of commitment for a custom model, you must set the number of model units to 1 and then select **No commitment**. With this option, you pay an hourly amount for the allocated throughput until you delete it.

Note

To see limitations on model units and commitment term, see [Provisioned Throughput \(p. 90\)](#).

5. Under **Estimated purchase summary**, review the estimated cost.
6. Choose **Purchase Provisioned Throughput**.
7. Review the note that appears and acknowledge the commitment duration and price by selecting the checkbox. Then choose **Confirm purchase**.

The console displays the overview page. For the Provisioned Throughput that you just purchased, the console displays it in the table, with the status set to **Creating**.

View details of a Provisioned Throughput

From the summary page, you can view the details for any of your Provisioned Throughput resources.

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.
2. From the left menu, choose **Provisioned Throughput**.
3. From the **Provisioned Throughput** table, choose a Provisioned Throughput. The console opens the details page.
4. Under **Provisioned Throughput overview**:
5. Under **Tags**, the console displays the tags that are associated with this Provisioned Throughput. Choose **Manage tags** to add or remove tags for this Provisioned Throughput.

Edit a Provisioned Throughput

You can edit only a few fields of a Provisioned Throughput.

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.
2. From the left menu, choose **Provisioned Throughput**.
3. From the **Provisioned Throughput** table, select the Provisioned Throughput to edit.
4. The console displays the Provisioned Throughput fields that you can edit.
5. Choose **Save** to start the update.

Delete a Provisioned Throughput

Note

You can't delete a Provisioned Throughput with commitment before the commitment term is complete.

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.
2. From the left menu, choose **Provisioned Throughput**.
3. From the **Provisioned Throughput** table, select the Provisioned Throughput to delete.
4. The console displays a modal form to warn you that delete is permanent. Choose **Confirm** to proceed.

Using the Provisioned Throughput API

Amazon Bedrock provides API operations to create and manage your Provisioned Throughput resources.

Topics

- [Create Provisioned Throughput \(p. 94\)](#)
- [Run inference using Provisioned Throughput \(p. 95\)](#)
- [Update Provisioned Throughput \(p. 96\)](#)
- [Get Provisioned Throughput \(p. 97\)](#)
- [Delete Provisioned Throughput \(p. 97\)](#)
- [List Provisioned Throughput resources \(p. 98\)](#)

Create Provisioned Throughput

Use the [CreateProvisionedModelThroughput](#) operation to create a Provisioned Throughput for a base or custom model.

When you send a `CreateProvisionedModelThroughput` operation for a base model, Bedrock validates your access permissions for the requested model and checks for capacity availability. If the validations are successful, Bedrock creates the provisioned throughput and returns the ARN of the Provisioned Throughput.

When you send a `CreateProvisionedModelThroughput` operation for a custom model, Bedrock validates your access permissions for the requested custom model, and checks for capacity availability. If the custom model is encrypted using a customer-managed KMS key, your permissions must include `kms:Decrypt` permission for this model.

Note

To see quotas for your options for `commitmentDuration` and `modelUnits`, see [Provisioned Throughput \(p. 90\)](#).

To create a Provisioned Throughput for an Amazon Titan foundation model, use the following IDs instead of the default ones.

Note

Provisioned throughput is currently available for the following models.

Model name	Model ID for provisioned throughput
Titan Text G1 - Express 8K	amazon.titan-text-express-v1:0:8k
Titan Embeddings G1 - Text	amazon.titan-embed-text-v1:2:8k
Anthropic Claude V2 18K	anthropic.claude-v2:0:18k
Anthropic Claude V2 100K	anthropic.claude-v2:0:100k
Anthropic Claude Instant V1 100K	anthropic.claude-instant-v1:2:100K
Stable Diffusion XL 1.0	stability.stable-diffusion-xl-v1:0
Meta Llama 2 Chat 13B	meta.llama2-13b-chat-v1:0:4k

The following code examples demonstrate how to create a Provisioned Throughput using the AWS CLI and the Python SDK.

AWS CLI

Create the Provisioned Throughput using the following command in the AWS CLI:

```
aws bedrock create-provisioned-model-throughput
  --model-units 1
  --commitment-duration SixMonths
  --provisioned-model-name test-model
  --model-id arn:aws:bedrock:us-east-1::foundation-model/anthropic.claude-v2
```

Python (Boto)

The following example demonstrates how to create the Provisioned Throughput using Python:

```
import boto3
bedrock = boto3.client(service_name='bedrock')
bedrock.create_provisioned_model_throughput(
    modelUnits=1,
    commitmentDuration='SixMonths',
    provisionedModelName='test-model',
    modelId='arn:aws:bedrock:us-east-1::foundation-model/anthropic.claude-v2'
)
```

Run inference using Provisioned Throughput

Use the [InvokeModel](#) or [InvokeModelWithResponseStream](#) operation to run inference using Provisioned Throughput. Specify the provisioned model ARN as the `modelId` parameter.

The following code examples show how to run inference for an Anthropic Claude model with a Provisioned Throughput using the AWS CLI and Python. For information about the prompt format and inference parameters, see [Inference parameters for foundation models \(p. 31\)](#).

AWS CLI

The following example demonstrates how to run inference for an Anthropic Claude model with a Provisioned Throughput in the AWS CLI.

```
aws bedrock-runtime invoke-model
  --model-id provisioned-model-arn \
  --body '{"prompt": "\n\nHuman: story of two dogs\n\nAssistant: \""}' invoke-model-output.txt
```

Python (Boto)

The following example demonstrates how to run inference for an Anthropic Claude model with a Provisioned Throughput using Python:

```
import boto3
import json

bedrock = boto3.client(service_name='bedrock-runtime')

body = json.dumps({"prompt": "\n\nHuman: explain black holes to 8th graders\n\nAssistant: \""})
modelId = 'provisioned-model-arn'
accept = 'application/json'
contentType = 'application/json'

response = bedrock.invoke_model(body=body, modelId=modelId, accept=accept,
                                contentType=contentType)
response_body = json.loads(response.get('body').read())

# text
print(response_body.get('results')[0].get('outputText'))

# embedding
print(response_body.get('embedding'))
```

Update Provisioned Throughput

Use the [UpdateProvisionedModelThroughput](#) operation to update the specified Provisioned Throughput.

The following code examples demonstrate how to update a Provisioned Throughput using the AWS CLI and Python.

AWS CLI

Update the Provisioned Throughput using the following command in the AWS CLI:

```
aws bedrock update-provisioned-model-throughput
  --provisioned-model-id provisioned-model-arn | provisioned-model-name
  --model-arn custom-model-arn | foundation-model-arn
```

Python (Boto)

The following example demonstrates how to update the Provisioned Throughput using Python:

```
import boto3
bedrock = boto3.client(service_name='bedrock')

bedrock.update_provisioned_model_throughput(
    provisionedModelName='provisioned-model-arn | provisioned-model-name'
    modelArn='custom-model-arn | foundation-model-arn'
)
```

Get Provisioned Throughput

Use the [GetProvisionedModelThroughput](#) operation to retrieve information about the specified Provisioned Throughput.

The following code examples demonstrates how to retrieve information using the AWS CLI and Python.

AWS CLI

Retrieve information about the Provisioned Throughput using the following command in the AWS CLI:

```
aws bedrock get-provisioned-model-throughput
--provisioned-model-id provisioned-model-arn | provisioned-model-name
```

Python (Boto)

The following example demonstrates how to retrieve information about the Provisioned Throughput using Python:

```
import boto3
bedrock = boto3.client(service_name='bedrock')

bedrock.get_provisioned_model_throughput(
    provisionedModelName='my-provisioned-model-arn | provisioned model name'
)
```

Delete Provisioned Throughput

Use the [DeleteProvisionedModelThroughput](#) operation to delete the specified provisioned throughput.

The following code examples demonstrate how to delete a Provisioned Throughput using the AWS CLI and Python.

AWS CLI

Delete the Provisioned Throughput using the following command in the AWS CLI:

```
aws bedrock delete-provisioned-model-throughput
--provisioned-model-id provisioned-model-arn | provisioned-model-name
```

Python (Boto)

The following example demonstrates how to delete the Provisioned Throughput using Python:

```
import boto3
bedrock = boto3.client(service_name='bedrock')
```



```
bedrock.delete_provisioned_model_throughput(  
    provisionedModelName='my-provisioned-model-arn | provisioned model name'  
)
```

List Provisioned Throughput resources

Use the [ListProvisionedModelThroughputs](#) operation to list the Provisioned Throughput resources that you have defined.

The following code examples demonstrate how to list the Provisioned Throughput resources using the AWS CLI and Python.

AWS CLI

List the Provisioned Throughput resources using the following command in the AWS CLI:

```
aws bedrock list-provisioned-model-throughputs
```

Python (Boto)

The following example demonstrates how to list the Provisioned Throughput resources using Python:

```
import boto3  
bedrock = boto3.client(service_name='bedrock')  
  
bedrock.list_provisioned_model_throughputs()
```

Agents for Amazon Bedrock

Note

Agents for Amazon Bedrock is in limited preview release. To request access, contact your AWS account manager.

Warning

As of October 31, 2023, the payload structure of the Lambda input event and response format for action groups is changing. Refer to [Deprecated Lambda function payload format \(p. 118\)](#) to see the deprecated format.

Agents for Amazon Bedrock enables developers to configure an agent to complete actions based on organization data and user input. Agents orchestrate interactions between foundation models, data sources, software applications, and user conversations, and automatically call APIs to take actions. Developers can easily integrate the agents and accelerate delivery of generative AI applications saving weeks of development effort.

With Agents for Amazon Bedrock, you can automate tasks for your customers. For example, you can create an agent that helps customers process insurance claims or one that helps customers make travel reservations and answer questions related to these tasks. You don't have to worry about provisioning, managing infrastructure, or writing custom code. Agents for Amazon Bedrock manages monitoring, encryption, user permissions, and API invocation.

Agents for Amazon Bedrock can carry out the following tasks:

- Extend foundation models to understand user requests and break down the tasks it needs to perform into smaller steps.
- Collect additional information from a user through natural conversation.
- Take actions to fulfill a customer's request.
- Make API calls to your company systems to carry out actions.
- Augment performance and accuracy by using data sources that you provide to facilitate Retrieval-Augmented Generation (RAG).
- Carry out source attribution.

To take advantage of Agents for Amazon Bedrock, you carry out the following steps:

1. (Optional) Set up a vector database and then create a knowledge base to store your private data in that database. For more information, see [Building a knowledge base \(p. 100\)](#).
2. Create an agent for your use-case, add actions that it can carry out, and attach the knowledge base you created to augment its performance. For more information, see [Building an agent \(p. 112\)](#).
3. Test your agent in the console or through API calls and modify the configurations as necessary. For more information, see [Test your agent \(p. 121\)](#).
4. When you have sufficiently modified your agent and it is ready to be deployed to your application, create an alias to point to a version of your agent. For more information, see [Deploying an agent: versioning and aliases \(p. 124\)](#).
5. Set up your application to make API calls to your agent alias.

Topics

- [Building a knowledge base \(p. 100\)](#)
- [Building an agent \(p. 112\)](#)
- [Test your agent \(p. 121\)](#)

- [Deploying an agent: versioning and aliases \(p. 124\)](#)
- [Using the API \(p. 125\)](#)
- [How Bedrock Agent works with IAM \(p. 126\)](#)

Building a knowledge base

Note

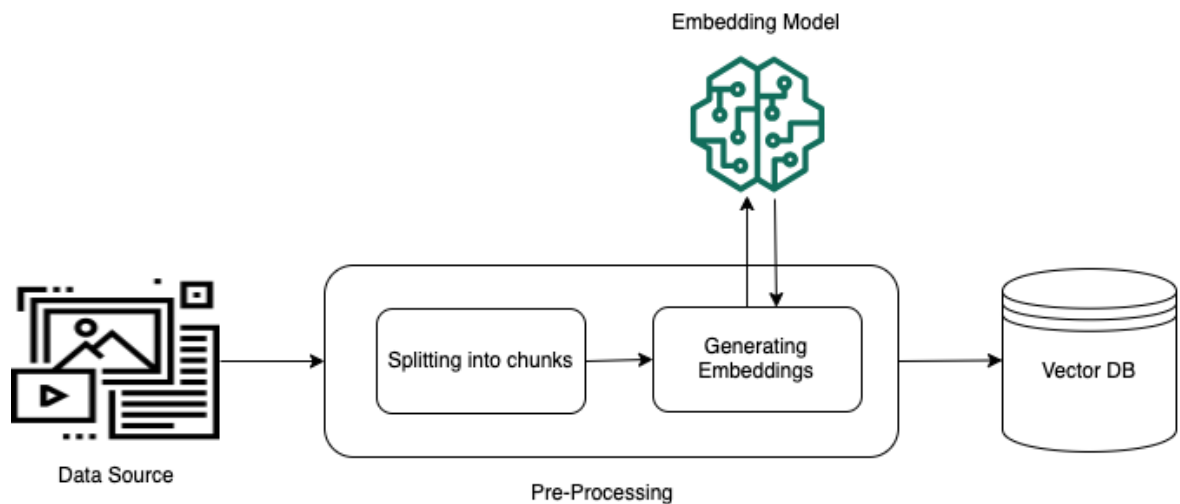
Knowledge base is in limited preview release. To request access, contact your AWS account manager.

With Amazon Bedrock, you can enable a Retrieval-Augmented Generation (RAG) workflow by using knowledge bases to build contextual applications by using the reasoning capabilities of LLMs. RAG is a popular technique that combines the use of private data with Large Language Models (LLMs). The combination of Amazon Bedrock with knowledge bases enables a faster time to market by automating the RAG solution and reducing the build time for your agent. Adding a knowledge base also increases cost-effectiveness by removing the need to continually train your model to be able to leverage your private data.

RAG starts with an initial step to retrieve relevant documents from a data store (most commonly a vector index) based on the user's query. It then employs a language model to generate a response by considering both the retrieved documents and the original query. The following steps in the setup and implementation of RAG are automated for you by the knowledge base service.

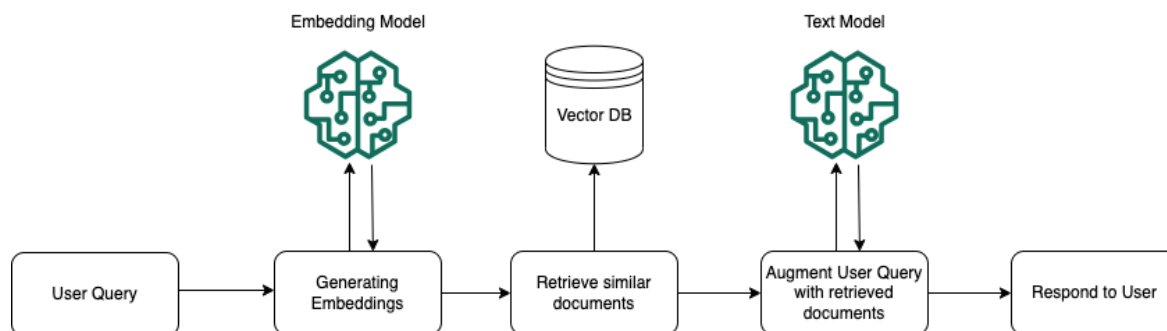
Pre-processing data

To enable effective retrieval from private data, a common practice is to first split the documents into manageable chunks. The following step is to convert the chunks to vectors and then to write them to a vector index while maintaining a mapping to the original document, in order to generate the embeddings. The following image illustrates pre-processing of data for the vector database.



Runtime execution

At runtime, an embedding model is used to convert the user's query to a vector. The vector index is then queried to find documents similar to the user's query by comparing document vectors to the user query vector. In the final step, semantically similar documents retrieved from the vector index are added as context for the original user query. When generating a response for the user, the semantically similar documents are prompted in the text model. The following image illustrates how RAG operates at runtime to augment responses to user queries.



You can build a RAG-based application in Amazon Bedrock by creating a knowledge base and associating it to an agent to augment its generative capabilities with your own data. You use a knowledge base to load your private data into a vector index. A knowledge base reads data from your Amazon S3 bucket, splits it into smaller chunks, generates vector embeddings, and stores the embeddings in a vector index that you provide. You can associate a knowledge base with multiple agents. After you configure an agent with a knowledge base, it can use the information stored in the corresponding vector index to augment its responses to user queries.

Topics

- [Create a service role and configure IAM permissions \(p. 101\)](#)
- [Set up your data for ingestion \(p. 105\)](#)
- [Create a knowledge base \(p. 108\)](#)
- [Manage a knowledge base \(p. 111\)](#)
- [Add a knowledge base to an agent \(p. 111\)](#)

Create a service role and configure IAM permissions

Knowledge bases use service roles to access AWS resources (for more information, see [Creating a role to delegate permissions to an AWS service](#)).

Before you can create a knowledge base, you need to create a service role and attach a trust policy for the role you create.

To create the service role and attach a trust policy

1. Create an IAM role with the prefix `AmazonBedrockExecutionRoleForKnowledgeBase_`. For more information about creating a role, see [Creating a role to delegate permissions to an AWS service](#).
2. Create a trust policy for the role you create. The following shows an example policy you can use. You can restrict the scope of the permission by using one or more global condition context keys. For more information, see [AWS global condition context keys](#). Set the `aws:SourceAccount` value to your account ID. You can use the `ArnEquals` or `ArnLike` condition to restrict the scope to specific knowledge bases.

Note

As a best practice for security purposes, replace the `*` with specific knowledge base IDs after you have created them.

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Principal": {
```

```
        "Service": "bedrock.amazonaws.com"
      },
      "Action": "sts:AssumeRole",
      "Condition": {
        "StringEquals": {
          "aws:SourceAccount": "account-id"
        },
        "ArnLike": {
          "AWS:SourceArn": "arn:aws:bedrock:region:account-id:knowledge-base/*"
        }
      }
    }
  ]
}
```

3. Attach the trust policy to the role.

For a knowledge base, you need to grant Amazon Bedrock the following permissions:

- Access to Amazon Bedrock embedding models
- Access to the Amazon S3 object containing your data sources
- (If you encrypted your Amazon S3 data) Permissions to decrypt and encrypt your customer-managed AWS KMS key for your data sources
- (If you create a vector database in Amazon OpenSearch Service) Access to your OpenSearch Service collection
- (If you create a vector database in Pinecone or Redis Enterprise Cloud) Permissions for AWS Secrets Manager to authenticate your Pinecone or Redis Enterprise Cloud account

Topics

- [Permissions to access Amazon Bedrock models \(p. 102\)](#)
- [Permissions to access your data sources in Amazon S3 \(p. 103\)](#)
- [\(Optional\) Permissions to decrypt your AWS KMS key for your data sources in Amazon S3 \(p. 103\)](#)
- [\(Optional\) Permissions to access your vector database in Amazon OpenSearch Service \(p. 104\)](#)
- [\(Optional\) Permissions to access your vector database in Pinecone or Redis Enterprise Cloud \(p. 104\)](#)
- [\(Optional\) Permissions for AWS to manage a AWS KMS key for transient data storage during data ingestion \(p. 105\)](#)

Permissions to access Amazon Bedrock models

To allow Amazon Bedrock to access Amazon Bedrock models to embed your source data, attach the following policy to your Amazon Bedrock service role. Replace *region* with your region your foundation model is in and *foundation-model-id* with the ID of the foundation model. Currently, only the Titan Embeddings G1 - Text Model *amazon.titan-embed-text-v1* is supported.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "bedrock:ListFoundationModels",
        "bedrock:ListCustomModels"
      ],
      "Resource": "*"
    }
  ]
}
```

```
{
  "Effect": "Allow",
  "Action": [
    "bedrock:InvokeModel"
  ],
  "Resource": [
    "arn:aws:bedrock:region::foundation-model/foundation-model-id"
  ]
}
```

Permissions to access your data sources in Amazon S3

To allow Amazon Bedrock to access your Amazon S3 data, attach the following policy attached to your Amazon Bedrock service role. Replace *bucket/path/to/folder* with the path to the object containing all the data source files for your knowledge base and *account-id* with the account that the Amazon S3 object belongs to.

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Action": [
      "s3:GetObject",
      "s3:ListBucket"
    ],
    "Resource": [
      "arn:aws:s3::bucket/path/to/folder",
      "arn:aws:s3::bucket/path/to/folder/*"
    ],
    "Condition": {
      "StringEquals": {
        "aws:PrincipalAccount": "account-id"
      }
    }
  }]
}
```

(Optional) Permissions to decrypt your AWS KMS key for your data sources in Amazon S3

If you encrypted your data sources in Amazon S3 with a AWS KMS key, attach the following policy to your Amazon Bedrock service role to allow Amazon Bedrock to decrypt your key. Replace *region* and *account-id* with the region and account ID to which the key belongs. Replace *key-id* with the ID of your AWS KMS key.

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Action": [
      "KMS:Decrypt",
    ],
    "Resource": [
      "arn:aws:kms:region:account-id:key/key-id"
    ],
    "Condition": {
      "StringEquals": {
        "kms:ViaService": [
```

```
    "s3. region .amazonaws.com"
  ]
}
}
}
```

(Optional) Permissions to access your vector database in Amazon OpenSearch Service

If you created a vector database in Amazon OpenSearch Service for your knowledge base, attach the following policy to your Amazon Bedrock service role to allow access to the collection. Replace *region* and *account-id* with the region and account ID to which the database belongs. Input the ID of your Amazon OpenSearch Service collection in *collection-id*. You can allow access to multiple collections by adding them to the Resources list.

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Action": [
      "aoss:APIAccessAll"
    ],
    "Resource": [
      "arn:aws:aoss:region:account-id:collection/collection-id"
    ]
  }]
}
```

(Optional) Permissions to access your vector database in Pinecone or Redis Enterprise Cloud

If you created a vector database in Pinecone or Redis Enterprise Cloud for your knowledge base, attach the following policy to your Amazon Bedrock service role to allow AWS Secrets Manager to authenticate your account to access the database. Replace *region* and *account-id* with the region and account ID to which the database belongs. Replace *secret-id* with the ID of your secret.

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Action": [
      "secretsmanager:GetSecretValue"
    ],
    "Resource": [
      "arn:aws:secretsmanager:region:account-id:secret:secret-id"
    ]
  }]
}
```

If your secret is encrypted with a AWS KMS key, attach the following policy to your Amazon Bedrock service role to allow it to decrypt your key. Replace *region* and *account-id* with the region and account ID to which the key belongs. Replace *key-id* with the ID of your AWS KMS key.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
```

```
        "Effect": "Allow",
        "Action": [
            "kms:Decrypt"
        ],
        "Resource": [
            "arn:aws:kms:region:account-id:key/keyId"
        ]
    }
]
```

(Optional) Permissions for AWS to manage a AWS KMS key for transient data storage during data ingestion

To allow the creation of a AWS KMS key for transient data storage in the process of ingesting your data source, attach the following policy to your Amazon Bedrock service role. Replace the *region*, *account-id*, and *key-id* with the appropriate values.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "kms:GenerateDataKey",
        "kms:Decrypt"
      ],
      "Resource": [
        "arn:aws:kms:region:account-id:key/key-id"
      ]
    }
  ]
}
```

Set up your data for ingestion

Before you can create a knowledge base, you need to carry out the following actions:

1. Configure your source data and upload them to an Amazon S3 bucket.
2. Create a vector database and configure fields for Amazon Bedrock to write data to and to access the index.

Configure your source data

By default, knowledge base automatically splits your source data into chunks, such that each chunk contains, at most, 200 tokens. If a document contains less than 200 tokens, then it is not split any further. To split your documents into chunks containing less than 200 tokens, pre-process the documents by splitting them into smaller chunks.

After you pre-process your source data, upload it to an Amazon S3 bucket. To set up an Amazon S3 bucket, see [Getting started with Amazon S3](#). Knowledge base supports the following file formats:

- Plain text (.txt)
- Markdown (.md)
- HyperText Markup Language (.html)
- Microsoft Word document (.doc/.docx)
- Comma-separated values (.csv)

- Microsoft Excel spreadsheet (.xls/.xlsx)
- Portable Document Format (.pdf)

Create a vector database

Create a vector index in one of the following supported options.

- Amazon OpenSearch Service
- Pinecone
- Redis Enterprise Cloud

After processing your data, Amazon Bedrock writes the following information to the index you created.

- Text extracted from your documents.
- The vectors, corresponding to your text, that were generated by the embeddings model.
- The Amazon S3 path of the source file where the text was extracted from.

Topics

- [Create a vector index in Amazon OpenSearch Service \(p. 106\)](#)
- [Create a vector index in Pinecone \(p. 107\)](#)
- [Create a vector index in Redis Enterprise Cloud \(p. 108\)](#)

Create a vector index in Amazon OpenSearch Service

1. Log into Amazon OpenSearch Service and create a collection. Take note of the **Collection ARN**, which you will fill out when you create a knowledge base.
2. Once the collection is created, select it and create a vector index.

For detailed documentation on setting up a vector index in Amazon OpenSearch Service, see [Working with vector search collections](#).

While you set up the vector index, take note of the **Collection ARN**, which you will fill out when you create a knowledge base.

There are additional configurations that you must provide when creating a vector index:

- **Vector index name** – The name of the vector index. Choose any valid name of your choice. Later, when you create your knowledge base, enter the name you choose in the **Vector index name** field.
- **Vector field** – The field where the vector embeddings will be stored. Choose any valid name of your choice. Later, when you create your knowledge base, enter the name you choose in the **Vector field** field.
- **Dimensions** – The number of dimensions in the vector. Choose 1536 if you use the Titan Embeddings Model. Later, when you create your knowledge base, enter this number for the **Dimensions** field.
- **Distance metric** – The metric used to measure the similarity between vectors. We recommend that you experiment with different metrics for your use-case. If you use the Titan Embeddings Model, you can start with **cosine similarity**.
- **Metadata management** – Expand this field and configure the vector index to store additional metadata that a knowledge base can retrieve with vectors. The fields you need to configure are as follows:
 - **Text field** – Amazon Bedrock chunks the raw text in your data and stores the chunks in this field.

- **Mapping field** – The field where the text will be stored. Choose any valid name of your choice. Later, when you create your knowledge base, enter the name you choose for the **Text field name** field.
- **Data type** – Select String.
- **Filterable** – Select False.
- **Bedrock-managed metadata field** – Amazon Bedrock stores metadata related to the data in this field. The metadata includes the following:
 - **Mapping field** – Choose any valid name of your choice. Later, when you create your knowledge base, enter the name you choose for the **Bedrock-managed metadata field name** field.
 - **Data type** – Select String.
 - **Filterable** – Select False.

Security configurations

After you create the knowledge base, you must return to Amazon OpenSearch Service and set up security configurations in your vector database by adjusting the **Network access** and **Data access** settings of your collection. For more information, see [Create a knowledge base \(p. 108\)](#).

Create a vector index in Pinecone

Note

If you use Pinecone, you agree to authorize AWS to access the designated third-party source on your behalf in order to provide the Amazon Bedrock service to you. You're responsible for complying with any third-party terms applicable to use and transfer of data from the third-party service.

For detailed documentation on setting up a vector index in Pinecone, see [Manage indexes](#).

While you set up the vector index, take note of the following information, which you will fill out when you create a knowledge base:

- The endpoint URL for your index management page.
- (Optional) The namespace to be used to write new data to your database. For more information, see [Using namespaces](#).

There are additional configurations that you must provide when creating a Pinecone index:

- **Name** – The name of the vector index. Choose any valid name of your choice. Later, when you create your knowledge base, enter the name you choose in the **Vector index name** field.
- **Dimensions** – The number of dimensions in the vector. Choose 1536 if you use the Titan Embeddings Model. Later, when you create your knowledge base, enter this number in the **Dimensions** field.
- **Distance metric** – The metric used to measure the similarity between vectors. We recommend that you experiment with different metrics for your use-case. If you use the Titan Embeddings Model, you can start with **cosine similarity**.

Configure the Secrets Manager

To access your Pinecone index, you must provide your Pinecone API key to Amazon Bedrock through the AWS Secrets Manager.

To set up a secret for your Pinecone configuration

1. Follow the steps at [Create an AWS Secrets Manager secret](#), setting the key as `apiKey` and the value as the API key to access your Pinecone index.

2. To find your API key, open your [Pinecone console](#) and select **API Keys**.
3. After you create the secret, take note of its ARN. Later, when you create your knowledge base, enter the ARN in the **Credentials secret ARN** field.

Create a vector index in Redis Enterprise Cloud

Note

If you use Redis Enterprise Cloud, you agree to authorize AWS to access the designated third-party source on your behalf in order to provide the Amazon Bedrock service to you. You're responsible for complying with any third-party terms applicable to use and transfer of data from the third-party service.

For detailed documentation on setting up a vector index in Redis Enterprise Cloud, see [Integrating Redis Enterprise Cloud with Amazon Bedrock](#).

While you set up the vector index, take note of the following information, which you will fill out when you create a knowledge base:

- The public endpoint URL for your database.
- The name of the vector index for your database.

Configure the Secrets Manager

To access your Redis Enterprise Cloud cluster, you must provide your Redis Enterprise Cloud security configuration to Amazon Bedrock through the AWS Secrets Manager.

To set up a secret for your Redis Enterprise Cloud configuration

1. Enable TLS to use your database with Amazon Bedrock by following the steps at [Transport Layer Security \(TLS\)](#).
2. Follow the steps at [Create an AWS Secrets Manager secret](#). Set up the following keys with the appropriate values from your Redis Enterprise Cloud configuration in the secret:
 - `username` – The username to access your Redis Enterprise Cloud database. To find your username, look under the **Security** section of your database in the [Redis Console](#).
 - `password` – The password to access your Redis Enterprise Cloud database. To find your password, look under the **Security** section of your database in the [Redis Console](#).
 - `serverCertificate` – The content of the certificate from the Redis Cloud Certificate authority. Download the server certificate from the Redis Admin Console by following the steps at [Download certificates](#).
 - `clientPrivateKey` – The private key of the certificate from the Redis Cloud Certificate authority. Download the server certificate from the Redis Admin Console by following the steps at [Download certificates](#).
 - `clientCertificate` – The public key of the certificate from the Redis Cloud Certificate authority. Download the server certificate from the Redis Admin Console by following the steps at [Download certificates](#).
3. After you create the secret, take note of its ARN. Later, when you create your knowledge base, enter the ARN in the **Credentials secret ARN** field.

Create a knowledge base

To create a knowledge base

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.

2. From the left navigation pane, select **Knowledge base**.
3. In the **Knowledge base** section, select **Create knowledge base**.
4. On the **Provide knowledge base details** page, carry out the following actions:
 - a. In the **Knowledge base details** section, enter a name for the knowledge base and provide an optional description for it.
 - b. In the **IAM permissions** section, choose to let Amazon Bedrock create a new service role or to use an existing service role that allows Amazon Bedrock to access other services on your behalf. For more information, see [Create a service role and configure IAM permissions \(p. 101\)](#).
 - c. If you want to attach any tags to the knowledge base, select **Add new tag** in the **tags** section and add the tags as key-value pairs.
 - d. Select **Next**.
5. On the **Set up data source** page, you provide the information for the data source to add to the knowledge base by carrying out the following actions:
 - a. In the **Data source** section, carry out the following actions:
 - i. Provide a name for the data source and the URI of the Amazon S3 object.
 - ii. If you encrypted your Amazon S3 data, provide the AWS KMS key in the **Customer-managed AWS KMS key for Amazon S3 data** to allow Amazon Bedrock to decrypt it.
 - iii. While converting your data into embeddings, Amazon Bedrock encrypts your transient data with a key that AWS owns and manages, by default. You can select the checkbox labeled **Customize encryption settings (advanced)** under **AWS KMS key for transient data storage**.
 - b. In the **Embeddings model** section, choose an embeddings model to convert the knowledge base from your data into an embedding. Currently, only the Amazon Bedrock Titan embeddings model is available.
 - c. In the **Vector database** section, select the service that contains a vector database that you have already created. Check that your database is already configured with the required fields (for more information, see [Set up your data for ingestion \(p. 105\)](#)). Fill in the fields to allow Amazon Bedrock to map information from the knowledge base to your database, so that it can store, update, and manage embeddings.

Note

If you use a database in Amazon OpenSearch Service, you need to have configured the fields under **Metadata field mapping** beforehand. If you use a database in Pinecone or Redis Enterprise Cloud, you can provide names for these fields here and Amazon Bedrock will dynamically create them in the vector index for you.
 - d. Select **Next**.
6. On the **Review and create** page, check the configuration and details of your knowledge base. Select **Edit** in any section that you need to modify. When you are satisfied, select **Create knowledge base**.
7. The knowledge base creation process begins and the **Status** of the source becomes **In progress**. The time it takes to create the knowledge base depends on the amount of data you provided. When the knowledge base is finished being created, a green success banner appears and the **Status** of the knowledge base changes to **Ready**.

Note

If you chose to store your embeddings in an Amazon OpenSearch Service vector database, remember to set up your security configurations in OpenSearch Service for the knowledge base after it has been created. For more information, see [Security configurations \(p. 107\)](#).

Set up security configurations for your newly created knowledge base. Follow the steps in the tab corresponding to the database that you set up.

OpenSearch Service

To create a data access policy

1. In the OpenSearch Service console, navigate to your collection.
2. Select **Manage data access**.
3. Select **Create access policy** and give the policy a name and an optional description.
4. Choose **JSON** as the policy definition method and paste the following JSON object into the editor, replacing *collection-name* with the name of your collection and *service-role-arn* with the role ARN that you passed when creating your knowledge base.

```
{[
  {
    "Rules": [
      {
        "Resource": [
          "index/collection-name/*"
        ],
        "Permission": [
          "aoss:DescribeIndex",
          "aoss:ReadDocument",
          "aoss:WriteDocument"
        ],
        "ResourceType": "index"
      }
    ],
    "Principal": [
      "service-role-arn"
    ],
    "Description": "Data access policy"
  }
]
```

5. Select **Create** to create the policy.

Pinecone or Redis Enterprise Cloud

To integrate the third-party knowledge base, attach the following policy to your Amazon Bedrock service role, replacing *knowledge-base-arn*.

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Action": [
      "bedrock:AssociateThirdPartyKnowledgeBase"
    ],
    "Resource": [
      "knowledge-base-arn"
    ],
    "Condition": {
      "StringEquals": {
        "bedrock:ThirdPartyKnowledgeBaseCredentialsSecretArn": "secret-arn"
      }
    }
  }]
}
```

Manage a knowledge base

To manage a knowledge base

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.
2. From the left navigation pane, select **Knowledge base**.
3. To view details for a knowledge base, either select the **Name** of the source or choose the radio button next to the source and select **Edit**.
4. On the details page, you can carry out the following actions:
 - To change the details of the knowledge base, select **Edit** in the **Knowledge base overview** section.
 - To update the tags attached to the knowledge base, select **Manage tags** in the **Tags** section.
 - If you update the data source from which the knowledge base was created and need to sync the changes, select **Sync** in the **Data source** section.
 - To view the details of a data source, select a **Data source name**. Within the details, you can choose the radio button next to a sync event in the **Sync history** section and select **View warnings** to see why files in the data ingestion job failed to sync.
 - To manage the embeddings model used for the knowledge base, select **Edit Provisioned Throughput**.
 - Select **Save changes** when you are finished editing.

To delete a knowledge base

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.
2. From the left navigation pane, select **Knowledge base**.
3. To delete a source, either choose the radio button next to the source and select **Delete** or choose the **Name** of the source and then select **Delete** in the top right corner of the details page.
4. Review the warnings for deleting a knowledge base. If you accept these conditions, enter **delete** in the input box and select **Delete** to confirm.
5. Make sure to delete the knowledge base from any agents that it was added to. To do this, carry out the following steps:
 - a. From the left navigation pane, select **Agents**.
 - b. Choose the **Name** of the agent that you want to delete the knowledge base from.
 - c. A red banner appears to warn you to delete the reference to the knowledge base, which no longer exists, from the agent.
 - d. Select the radio button next to the knowledge base that you want to remove. Select **More** and then choose **Delete**.

Add a knowledge base to an agent

To add a knowledge base to an agent

1. Open the Amazon Bedrock console at <https://console.aws.amazon.com/bedrock/>.
2. From the left navigation pane, select **Agents**.
3. Choose the **Name** of the agent that you want to add knowledge bases to.
4. In the **Knowledge base** section, select **Add**.
5. Choose the knowledge base from the dropdown list under **Select knowledge base** and specify the instructions for the agent regarding the knowledge base.

Building an agent

Note

Amazon Bedrock is in limited preview release. To request access, contact your AWS account manager.

To build an agent, you set up the following components:

- The configuration of the agent itself, which defines the purpose of the agent.
- Action groups that define what actions the agent is designed to carry out.
- (Optional) A knowledge base of data sources to augment the generative capabilities of the agent.

Note

If you plan to attach a knowledge base to your agent, first set up your knowledge base by following the steps at [Building a knowledge base \(p. 100\)](#).

You carry out the following steps to set up an agent to interact with your customers:

1. Choose a foundational model that the agent can use for orchestration.
2. Give instructions to the agent, describe the actions it can perform, and provide an API schema for the actions.
3. Add private data to the agent by setting up your data sources in a knowledge base and associating the knowledge base with your agent.
4. Test the agent and iterate on the working draft.
5. After the agent is working as expected, create an alias so that you can integrate the agent with your application.
6. Configure your application to make API calls to the agent alias.
7. Update the working draft of the agent and create new aliases as necessary.

Topics

- [Create a service role and configure IAM permissions \(p. 112\)](#)
- [Create an agent \(p. 115\)](#)
- [Edit your agent \(p. 119\)](#)

Create a service role and configure IAM permissions

Before you can create an agent, you need to create a service role and attach a trust policy for the role you create.

To create the service role and attach a trust policy

1. Create an IAM role with the prefix `AmazonBedrockExecutionRoleForAgents_`. For more information about creating a role, see [Creating a role to delegate permissions to an AWS service](#).
2. Create a trust policy for the role you create. The following shows an example policy you can use. You can restrict the scope of the permission by using one or more global condition context keys. For more information, see [AWS global condition context keys](#). Set the `aws:SourceAccount` value to your account ID. You can use the `ArnEquals` or `ArnLike` condition to restrict the scope to specific agents.

Note

As a best practice for security purposes, replace the `*` with specific agent IDs after you have created them.

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Principal": {
      "Service": "bedrock.amazonaws.com"
    },
    "Action": "sts:AssumeRole",
    "Condition": {
      "StringEquals": {
        "aws:SourceAccount": "account-id"
      },
      "ArnLike": {
        "AWS:SourceArn": "arn:aws:bedrock:region:account-id:agent/*"
      }
    }
  }]
}
```

3. Attach the trust policy to the role.

Grant Amazon Bedrock permissions to access the following resources on your behalf through configuring role-based permissions:

- The Amazon Bedrock base models
- The Amazon S3 objects containing the OpenAPI schemas for the action groups in your agents

You also need to provide permissions for Amazon Bedrock to access the AWS Lambda functions for the action groups in your agents through configuring resource-based policies for your Lambda functions.

If you are going to attach a knowledge base to your agent, you also need to provide permissions for Amazon Bedrock to query the knowledge base.

Topics

- [Permissions to access Amazon Bedrock base models \(p. 113\)](#)
- [Permissions to access your action group API schemas in Amazon S3 \(p. 114\)](#)
- [Permissions to access your action group Lambda functions \(p. 114\)](#)
- [\(Optional\) Permissions to access your knowledge bases \(p. 115\)](#)

Permissions to access Amazon Bedrock base models

To provide permissions for Amazon Bedrock to use foundation models for orchestration, attach the following policy to an IAM role.

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Action": [
      "bedrock:InvokeModel"
    ],
    "Resource": [
      "arn:aws:bedrock:us-east-1::foundation-model/anthropic.claude-v1",
      "arn:aws:bedrock:us-east-1::foundation-model/anthropic.claude-v2",
      "arn:aws:bedrock:us-east-1::foundation-model/anthropic.claude-instant-v1"
    ]
  }]
}
```



```
}  
  }  
}
```

Permissions to access your action group API schemas in Amazon S3

Provide permissions for Amazon Bedrock to access the Amazon S3 URIs of the API schemas for your agent's action groups. Attach the following policy to an IAM role. In the Resource field, you can provide an Amazon S3 object containing the schemas or you can add the URI of each schema to the list.

```
{  
  "Version": "2012-10-17",  
  "Statement": [{  
    "Effect": "Allow",  
    "Action": [  
      "s3:GetObject"  
    ],  
    "Resource": [  
      "arn:aws:s3:::path/to/schema"  
    ],  
    "Condition": {  
      "StringEquals": {  
        "aws:PrincipalAccount": "account-id"  
      }  
    }  
  }  
}]  
}
```

Permissions to access your action group Lambda functions

Provide permissions for Amazon Bedrock to access the Lambda functions for your agent's action groups. Attach the following resource-based policy to a Lambda function. The Condition specifies the ARN of the user account and the ARN of the agent. The Condition field is optional, but is recommended for security purposes. For more information, see [Using resource-based policies for Lambda](#).

```
{  
  "Version": "2012-10-17",  
  "Statement": [{  
    "Sid": "statement-name",  
    "Effect": "Allow",  
    "Principal": {  
      "Service": "bedrock.amazonaws.com"  
    },  
    "Action": [  
      "lambda:InvokeFunction",  
    ],  
    "Condition": {  
      "StringEquals": {  
        "AWS:SourceAccount": "account-id"  
      },  
      "ArnLike": {  
        "AWS:SourceArn": "arn:aws:bedrock:region:account-id:agent/agent-id"  
      }  
    }  
  }  
}]  
}
```

(Optional) Permissions to access your knowledge bases

To provide permissions for Amazon Bedrock to access a knowledge base you have set up, attach the following policy to an IAM role. Replace *account-id* and *knowledge-base-id* with the appropriate values and *statement-name* with the statement name of your choice.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "statement-name",
      "Effect": "Allow",
      "Action": [
        "bedrock:QueryKnowledgeBase"
      ],
      "Resource": [
        "arn:aws:bedrock:*:account-id:knowledge-base/knowledge-base-id"
      ]
    }
  ]
}
```

Create an agent

To create an agent, <https://console.aws.amazon.com/bedrock/> Choose **Agents** in the left navigation pane. Then select **Create** at the top right corner of the **Agents** section.

Provide agent details

1. In the **Agent name** section, give the agent a name and an optional description.
2. In the **User input** section, select whether you want the agent to request the user for additional information when trying to complete a task. If you select **No**, the agent doesn't request the user for additional details and informs the user that it doesn't have enough information to complete the task.
3. In the **IAM permissions** section, choose an AWS Identity and Access Management (IAM) role that provides Amazon Bedrock permission to access other AWS services. See [Create a service role and configure IAM permissions \(p. 112\)](#) for more information.
4. In the **Idle session timeout** section, choose the duration that Amazon Bedrock keeps a session with a user open. Amazon Bedrock maintains prompt variables for the duration of the session so that your agent can resume an interaction with the same variables.
5. Select **Next** when you are done setting up the agent configuration.

Select model

1. For the agent to perform orchestration, you need to choose a model. Choose a model provider and then choose a model in the dropdown menu to train your agent.
2. In **Instructions for the Agent**, provide details to tell the agent what it should do and how it should interact with users. Refer to the following sample instructions:

You are an office assistant in an insurance agency. You are friendly and polite. You help with managing insurance claims and coordinating pending paperwork.

Add Action groups

1. Provide a name for your action group and describe what the action does in the **Description for action groups**. For example, you can use the following text for an action group called *SetupDNSResources*:

Use this Action whenever you want to setup DNS resources for setting up a website.

2. In **Select Lambda function**, choose a Lambda function that you created in AWS Lambda. The Lambda function provides the business logic that is carried out upon invoking the action. Choose the version of the function to use.

Note

Remember to configure a policy to allow the Amazon Bedrock service principal to access the Lambda function. For more information, see [Create a service role and configure IAM permissions \(p. 112\)](#) for more information.. You can optionally provide the ARN of the agent as the SourceArn when configuring the policy.

The following is the general format of the Lambda input event. Use the input event fields to create your Lambda function. For more information, see [Event-driven invocation](#) in the AWS Lambda documentation.

Warning

As of October 31, 2023, the payload structure of the Lambda input event and response format is changing. Refer to [Deprecated Lambda function payload format \(p. 118\)](#) to see the deprecated format.

```
{
  "messageVersion": "1.0",
  "agent": {
    "name": "string",
    "id": "string",
    "alias": "string",
    "version": "string"
  },
  "inputText": "string",
  "sessionId": "string",
  "actionGroup": "string",
  "apiPath": "string",
  "httpMethod": "string",
  "parameters": [
    {
      "name": "string",
      "type": "string",
      "value": "string"
    },
    ...
  ],
  "requestBody": {
    "content": {
      "<content_type>": {
        "properties": [
          {
            "name": "string",
            "type": "string",
            "value": "string"
          },
          ...
        ]
      }
    }
  },
  ...
}
```

```
"sessionAttributes": {  
  "string": "string",  
}  
}
```

Amazon Bedrock expects a response from your Lambda function in the following format.

```
{  
  "messageVersion": "1.0",  
  "response": {  
    "actionGroup": "string",  
    "apiPath": "string",  
    "httpMethod": "string",  
    "httpStatusCode": number,  
    "responseBody": {  
      "<contentType>": {  
        "body": "string"  
      }  
    }  
  }  
}
```

The following is an example of how the Lambda function can be defined in Python.

```
def lambda_handler(event, context):  
  
    response_body = {  
        'application/json': {  
            'body': "<sample response>"  
        }  
    }  
  
    action_response = {  
        'actionGroup': event['actionGroup'],  
        'apiPath': event['apiPath'],  
        'httpMethod': event['httpMethod'],  
        'httpStatusCode': 200,  
        'responseBody': response_body  
    }  
  
    api_response = {  
        'messageVersion': '1.0',  
        'response': action_response  
    }  
    return mock_api_response
```

3. In **Select API schema**, provide a link to the Amazon S3 URI of the schema with the API description, structure, and parameters for the action group. For examples of the schema, see <https://github.com/OAI/OpenAPI-Specification/tree/main/examples/v3.0>.

Note

Remember to assign the IAM role permissions to access the Amazon S3 URI of the schema. For more information, see [Create a service role and configure IAM permissions \(p. 112\)](#) for more information.

4. Select **Add another Action group** to set up another action group for your agent. When you are done adding action groups, select **Next**.

Add knowledge bases

1. If you have not yet created any knowledge bases, select **Knowledge base** in the left navigation pane and follow the instructions at [Create a knowledge base \(p. 108\)](#) to create one. Otherwise, select a knowledge base from the dropdown menu.
2. Write a prompt in **Knowledge base instructions for the agent** to describe how the agent should use the knowledge base. For example, you can use the following text for a knowledge base called *Domain name system details*:

Use this knowledge base whenever you are creating a DNS record

3. Select **Add another knowledge base** to set up another knowledge base for your agent. When you are done adding knowledge bases, select **Next**.

Review the configuration of your agent and choose **Edit** for any sections you want to change. Select **Create** when you are ready to create the agent. When the process finishes, a green banner appears at the top to inform you that the agent was successfully created.

Deprecated Lambda function payload format

For reference, the deprecated format of the Lambda function input event is as follows. Update your action groups to the updated format (for more information, see [Create an agent \(p. 115\)](#)).

```
{
  "messageVersion": "1.0",
  "agent": {
    "name": "string",
    "id": "string",
    "alias": "string",
    "version": "string"
  },
  "inputText": "string",
  "sessionId": "string",
  "actionGroups": [
    {
      "actionGroup": "string",
      "apiPath": "string",
      "httpMethod": "string",
      "parameters": [
        {
          "name": "string",
          "type": "string",
          "value": "string"
        },
        ...
      ],
      "requestBody": {
        "content": {
          "<content_type>": {
            "properties": [
              {
                "name": "string",
                "type": "string",
                "value": "string"
              },
              ...
            ]
          }
        }
      }
    }
  ]
}
```

```
    },  
    ...  
  ],  
  "sessionAttributes": {  
    "string": "string",  
  }  
}
```

For reference, the deprecated format of the Lambda function response is as follows. Update your action groups to the updated format.

```
{  
  "messageVersion": "1.0",  
  "response": [  
    {  
      "actionGroup": "string",  
      "apiPath": "string",  
      "httpMethod": "string",  
      "httpStatusCode": number,  
      "responseBody": {  
        "<contentType>": {  
          "body": "string"  
        }  
      }  
    }  
  ]  
}
```

Edit your agent

After you create an agent, you can update its configuration as required. The configuration applies to the working draft.

To edit the agent configuration

1. Choose an agent in the **Agents** section.
2. Select **Edit** in the **Agent overview** section.
3. Edit the existing fields as necessary.
4. Select **Save changes**. A green success banner appears at the top.

You might want to try different foundation models for your agent or change the instructions for the agent. These changes apply only to the working draft.

To change the foundation model that your agent uses or the instructions to the agent.

1. Choose an agent in the **Agents** section.
2. Select **View** in the **Working draft** section or choose the working draft.
3. In the **Model details** section, select **Edit**
4. Edit the fields as necessary.
5. Select **Save** to remain in the same window. A green success banner appears at the top.
6. Test the updated agent in the right panel and make changes as necessary.
7. Select **Save and exit** to return to the working draft page.

Manage the action groups of an agent

After creating an agent, you can add more action groups or edit them. Adding and editing take place within the working draft. To carry out these operations, choose an agent from the **Agents** section and then choose the **Working draft** in the **Working Draft** section.

To add an action group

1. Select **Add** in the **Action groups** section.
2. Fill out the action group details
3. Select **Add**. A green success banner appears at the top.

To edit an action group

1. Do one of the following:
 - Choose the radio button next to the action group to edit and select **Edit**.
 - Select an action group to see its details. Choose **Edit** at the top.
2. Edit the existing fields as necessary.
3. Select **Save** to remain in the same window. A green success banner appears if there are no issues and a red banner appears if there are errors in the edit.
4. Test the updated agent in the right panel and make changes as necessary.
5. Select **Save and exit** to return to the working draft page.

Manage the knowledge bases of an agent

After creating an agent, you can add more knowledge bases or edit them. Adding and editing take place within the working draft. To carry out these operations, choose an agent from the **Agents** section and then choose the **Working draft** in the **Working Draft** section.

To add a knowledge base

1. Select **Add** in the **Knowledge bases** section.
2. Choose a knowledge base that you have created and provide instructions for how the agent should interact with it.
3. Select **Add**. A green success banner appears at the top.

To edit a knowledge base

1. Select a knowledge base.
2. In the knowledge base details page, select **Edit**
3. Edit the existing fields as necessary.
4. Select **Save** to remain in the same window. A green success banner appears if there are no issues and a red banner appears if there are errors in the edit.
5. Test the updated agent in the right panel and make changes as necessary.
6. Select **Save and exit** to return to the working draft page.

To delete a knowledge base

1. Select the radio button next to the knowledge base to delete.
2. Select **Delete**.

Test your agent

Once you have created an agent, you can find it in the **Agents** section. When you first create an agent, you will have a *working draft*. The working draft is a version of the agent that you can use to iteratively build the agent. By default, you can interact with the working draft with the AgentTestAlias. You can also select a different alias to test. In the test window, you can opt to show the trace for each response. The trace shows the agent's reasoning process, step-by-step, and is a useful tool for debugging your agent. To learn more about the trace, see [Trace enablement \(p. 121\)](#). To learn how to enable the trace through the API, see [Invoke your agent \(p. 125\)](#).

You have access to a test window to interact with your agent. You can use this function to debug your agent and make the necessary changes for it to work as expected.

To test the agent

1. Choose an agent in the **Agents** section.
2. In any of the pages inside an agent, you can select the left arrow icon at the top right to expand the test window.
3. Use the dropdown menu at the top of the test window to choose an alias and associated version to test.
4. Enter a message to request the agent to perform a task. Use the test window to help debug your agent.
5. When viewing a response, you have the following options.
 - If the agent pulls information from a source, the response contains footnotes. Select a footnote to view the citation for that part of the response.
 - Under the response from the bot, select **Show trace** to view the agent's reasoning process and usage of the attached knowledge bases and action groups, in addition to the **Inference configurations** used. The trace breaks down how the response was formed, step-by-step. Select an arrow next to a step to expand or collapse the trace for that step. For more information, see [Trace enablement \(p. 121\)](#).

You have the option of turning action groups and knowledge bases on and off. Use this feature to debug your agent and to isolate which actions or knowledge bases need to be updated by assessing its behavior with different settings.

To turn an action group or knowledge base on or off

1. Choose an agent in the **Agents** section.
2. Select the working draft in the **Working draft**.
3. In the **Action groups** or **Knowledge base** section, hover over the **State** of the action group whose state you want to change.
4. An edit button appears. Select it and then choose from the dropdown menu whether the action group or knowledge base is **Enabled** or **Disabled**.
5. If an action group is **Disabled**, the agent doesn't try to carry it out. If a knowledge base is **Disabled**, the agent doesn't use it in orchestration. Turn action groups on and off and use the test window to debug your agent.

Trace enablement

Each response from an agent is accompanied by a *trace* that you can view when you test the agent. The trace helps you follow the agent's reasoning process that leads it to the response it gives at that point in the conversation.

In the console test window, the trace is enabled by default. When you use an API, the trace is disabled by default unless you set the `enableTrace` to `true`. For more information about enabling the trace through the API, see [Invoke your agent \(p. 125\)](#).

Use the trace to track the agent's path from the user input to the response it returns. With the trace, you can find information about the inputs to the action groups that the agent invokes and the knowledge bases that it looks up to help it respond to the user, in addition to the outputs that the action groups and knowledge bases return. You can view the reasoning that the agent uses to determine the action that it takes or the query that it makes to a knowledge base. If a step in the trace fails, the trace returns a reason for the failure. Use the detailed information in the trace to debug your agent by identifying steps at which it has trouble or at which it yields unexpected behavior and ways in which you can improve the behavior.

When you show the trace in the test window in the console, a window appears showing a trace for each **Step** in the reasoning process. The **Step** consists of the following four types of traces. Each **Step** contains either a `failureTrace` or any combination of a `rationaleTrace`, `invocationInputTrace`, or `observationTrace`.

- `rationaleTrace` – Contains the reasoning, based on the user input, that the agent uses to justify carrying out an action group or getting information from a knowledge base.

```
{
  "traceId": "string",
  "text": "string"
}
```

The fields are described below.

- `traceId` – The unique identifier of the trace step. Any `rationaleTrace`, `invocationInput`, or `observationTrace` that is in the same step has the same `traceId`.
- `text` – The reasoning or thought process of the agent, based on the user input.
- `invocationInputTrace` – Contains information pertaining to the action group or knowledge base that is being invoked.

```
{
  "traceId": "string",
  "invocationType": "ACTION_GROUP | KNOWLEDGE_BASE",
  "actionGroupInvocationInput": {
    "actionGroupName": "string",
    "apiPath": "string",
    "parameters": [
      {
        "name": "string",
        "type": "string",
        "value": "string",
      },
      ...
    ],
    "requestBody": {
      "content": {
        "string": [
          {
            "name": "string",
            "type": "string",
            "value": "string",
          },
          ...
        ]
      }
    },
    "verb": "string"
  }
}
```

```
    },
    "knowledgeBaseLookupInput": {
      "knowledgeBaseId": "string",
      "text": "string"
    }
  }
}
```

The fields are described below.

- **traceId** – The unique identifier of the trace.
- **invocationType** – Specifies whether the agent is invoking an action group or a knowledge base.
- **actionGroupInvocationInput** – Contains the following metadata about the action group being invoked.
 - **actionGroupName** – The name of the action group that the agent is invoking.
 - **apiPath** – The path to the API to call, based off the action group.
 - **parameters** – The parameters in the Lambda input event.
 - **requestBody** – The parameters in the request body for the Lambda input event.
 - **verb** – The API method being used, based off the action group.
- **knowledgeBaseLookupInput** – Contains the following information about the knowledge base and the search query for the knowledge base.
 - **knowledgeBaseId** – The unique identifier of the knowledge base that the agent is looking up.
 - **text** – The query being made to the knowledge base.
- **observationTrace** – Contains the result or output of an action group or knowledge base, or the response to the user.

```
{
  "traceId": "string",
  "invocationType": "ACTION_GROUP | KNOWLEDGE_BASE | FINISH",
  "actionGroupInvocationOutput": {
    "text": "string"
  },
  "knowledgeBaseLookupOutput": {
    "sourceReferences": {
      "textSourceReferences": [
        {
          "referenceText": "string",
          "sourceLocation": {
            "s3SourceLocation": {
              "s3Uri": "string"
            }
          }
        },
        ...
      ]
    }
  },
  "finalResponse": {
    "text": "string"
  }
}
```

The fields are described below.

- **traceId** – The unique identifier of the trace.
- **invocationType** – Specifies whether the agent is invoking an action group, a knowledge base, or returning a response to the user.
- **actionGroupInvocationOutput** – Contains the JSON-formatted string returned by the API invoked by the action group.

- `knowledgeBaseLookupOutput` – Contains the text in the knowledge base and the S3 location of the data source. Each object in the list of `testSourceReferences` contains the following fields.
 - `referenceText` – The text from the knowledge base that is returned from the knowledge base query.
 - `sourceLocation` – Contains the S3 URI of the data source from which the returned text was found.
- `failureTrace` – Contains one field, a `failureReason`, which describes the error that occurred.

```
{  
  "failureReason": "string"  
}
```

Deploying an agent: versioning and aliases

After you have sufficiently iterated on your working draft and are satisfied with the behavior of your agent, you can set it up for deployment and integration into your application by creating *aliases* of your agent.

To deploy your agent, you create an *alias*. During alias creation, Amazon Bedrock automatically creates a version of your agent. The alias points to this newly created version. You can point the alias to previously created version if necessary. You then configure your application to make API calls to that alias.

The version is like a snapshot that preserves the resource as it exists at the time it was created. You can keep modifying the working draft and create new aliases (and consequently, versions) of your agent as necessary. In Amazon Bedrock, you create a new version of your agent by creating an alias that points to the new version by default. Amazon Bedrock creates versions in numerical order, starting from 1. Because a version acts as a snapshot of your agent at the time you created it, it is immutable.

Aliases let you efficiently switch between different versions of your agent without requiring the application to keep track of the version. For example, you can change an alias to point to a previous version of your agent if there are changes that you need to quickly revert.

The working draft version is DRAFT and the alias that points to it is the `AgentTestAlias`.

To manage versions and aliases of an agent, select **Agents** from the left navigation pane and choose the agent from the **Agents** section.

To create a new alias (and optionally a new version)

1. Select **Create alias** at the top right corner. Alternatively, select the **Deploy** tab and choose **Create** in the **Aliases** section.
2. Enter a unique name for the alias and provide an optional description.
3. Choose one of the following options
 - Create a new version and to associate the alias with it
 - Associate the alias with an existing version. From the dropdown menu, choose the version you want to associate the alias to.
4. Select **Create alias**. A green success banner appears at the top.

To manage versions of the agent, check that you are in the **Build** tab and that the **Versions** section displays underneath.

To view the details of a version

1. Select the version to view from the **Versions** section.

2. You can't modify any part of a version, but you can view details about the model, action groups, or Lambda function by choosing the name of the information you want to view.

To manage aliases for the agent, select the **Deploy** tab and check that the **Aliases** section displays underneath.

To associate an alias to a different version

1. Select the radio button next to the alias you want to edit.
2. Select the **Edit** button.
3. Choose one of the following options.
 - Create a new version and to associate the alias with it
 - Associate the alias with an existing version. From the dropdown menu, choose the version you want to associate the alias to.

Using the API

The following are the service endpoints for the Amazon Bedrock service. To connect programmatically to an AWS service, you use an endpoint. For information about endpoints for other AWS services, see [AWS service endpoints](#) in the AWS General Reference.

The following table provides a list of Region-specific endpoints that Amazon Bedrock supports.

Region Name	Region	Endpoint	Protocol
US East (N. Virginia)	us-east-1	bedrock-agent-runtime.us-east-1.amazonaws.com	HTTPS
US West (Oregon)	us-west-2	bedrock-agent-runtime.us-west-2.amazonaws.com	HTTPS

Currently, the only API available for Amazon Bedrock is InvokeAgent.

Invoke your agent

To interact with your agent, send an InvokeAgent request. Use TSTALIASID as the agentAliasId to invoke the draft version of your agent. You must use the Amazon Bedrock runtime endpoint to call this operation (bedrock-agent-runtime.*region*.amazonaws.com)

The following shows a sample request.

```
POST /agents/agentId/agentAliases/agentAliasId/sessions/sessionId/text HTTP/1.1
Content-type: application/json

{
  "inputText": "string",
  "endSession": boolean,
  "enableTrace": boolean,
  "sessionState": {
    "sessionAttributes": {
      "string" : "string"
    }
  }
}
```

```
}
```

The fields in the request are described below.

- `agentId` – The unique identifier for your agent. You can find the ID in your agent's details page.
- `agentAliasId` – The unique identifier for your agent's alias. You can find the ID in the details page for your agent's alias. To invoke the draft version of your agent, use `TSTALIASID`.
- `sessionId` – The unique identifier for the session. If you reuse this value, you continue an existing session with the agent if the value you set for the idle session timeout hasn't been exceeded.
- `inputText` – The prompt to provide the agent.
- `endSession` – Specifies whether to end the session with the agent or not.
- `enableTrace` – Specifies whether to return the trace or not in the response. By default, the trace is disabled. For more information, see [Trace enablement \(p. 121\)](#).
- `sessionState` – Contains session attributes for the session with the agent.

The response returns the following objects. If there is an error, the response returns an exception.

- `chunk` – The `bytes` field contains the response to the user as a base-64 encoded binary object. The `attribution` object contains a list of citations accompanying the response.
- `trace` – Contains the `agentId`, `agentAliasId`, and `sessionId`, alongside the trace object. For information about the fields in the trace object, see [Trace enablement \(p. 121\)](#). The trace is only returned if `enableTrace` was set as `true` in the request.

How Bedrock Agent works with IAM

Before you use IAM to manage access to Bedrock Agent, learn what IAM features are available to use with Bedrock Agent.

IAM features you can use with Bedrock Agent

IAM feature	Bedrock Agent support
Identity-based policies (p. 127)	Yes
Resource-based policies (p. 127)	No
Policy actions (p. 127)	Yes
Policy resources (p. 128)	Yes
Policy condition keys (p. 128)	No
ACLs (p. 129)	No
ABAC (tags in policies) (p. 129)	No
Temporary credentials (p. 130)	Yes
Principal permissions (p. 130)	Yes
Service roles (p. 130)	Yes
Service-linked roles (p. 130)	No

To get a high-level view of how Bedrock Agent and other AWS services work with most IAM features, see [AWS services that work with IAM](#) in the *IAM User Guide*.

Identity-based policies for Bedrock Agent

Supports identity-based policies	Yes
----------------------------------	-----

Identity-based policies are JSON permissions policy documents that you can attach to an identity, such as an IAM user, group of users, or role. These policies control what actions users and roles can perform, on which resources, and under what conditions. To learn how to create an identity-based policy, see [Creating IAM policies](#) in the *IAM User Guide*.

With IAM identity-based policies, you can specify allowed or denied actions and resources as well as the conditions under which actions are allowed or denied. You can't specify the principal in an identity-based policy because it applies to the user or role to which it is attached. To learn about all of the elements that you can use in a JSON policy, see [IAM JSON policy elements reference](#) in the *IAM User Guide*.

Identity-based policy examples for Bedrock Agent

To view examples of Bedrock Agent identity-based policies, see [Identity-based policy examples for Bedrock Agent \(p. 131\)](#).

Resource-based policies within Bedrock Agent

Supports resource-based policies	No
----------------------------------	----

Resource-based policies are JSON policy documents that you attach to a resource. Examples of resource-based policies are IAM *role trust policies* and Amazon S3 *bucket policies*. In services that support resource-based policies, service administrators can use them to control access to a specific resource. For the resource where the policy is attached, the policy defines what actions a specified principal can perform on that resource and under what conditions. You must [specify a principal](#) in a resource-based policy. Principals can include accounts, users, roles, federated users, or AWS services.

To enable cross-account access, you can specify an entire account or IAM entities in another account as the principal in a resource-based policy. Adding a cross-account principal to a resource-based policy is only half of establishing the trust relationship. When the principal and the resource are in different AWS accounts, an IAM administrator in the trusted account must also grant the principal entity (user or role) permission to access the resource. They grant permission by attaching an identity-based policy to the entity. However, if a resource-based policy grants access to a principal in the same account, no additional identity-based policy is required. For more information, see [How IAM roles differ from resource-based policies](#) in the *IAM User Guide*.

Policy actions for Bedrock Agent

Supports policy actions	Yes
-------------------------	-----

Administrators can use AWS JSON policies to specify who has access to what. That is, which **principal** can perform **actions** on what **resources**, and under what **conditions**.

The Action element of a JSON policy describes the actions that you can use to allow or deny access in a policy. Policy actions usually have the same name as the associated AWS API operation. There are some exceptions, such as *permission-only actions* that don't have a matching API operation. There are also some operations that require multiple actions in a policy. These additional actions are called *dependent actions*.

Include actions in a policy to grant permissions to perform the associated operation.

For a list of Bedrock Agent actions, see [Actions Defined by Bedrock Agent](#) in the *Service Authorization Reference*.

Policy actions in Bedrock Agent use the following prefix before the action:

```
bedrock
```

To specify multiple actions in a single statement, separate them with commas.

```
"Action": [  
    "bedrock:action1",  
    "bedrock:action2"  
]
```

To view examples of Bedrock Agent identity-based policies, see [Identity-based policy examples for Bedrock Agent \(p. 131\)](#).

Policy resources for Bedrock Agent

Supports policy resources	Yes
---------------------------	-----

Administrators can use AWS JSON policies to specify who has access to what. That is, which **principal** can perform **actions** on what **resources**, and under what **conditions**.

The Resource JSON policy element specifies the object or objects to which the action applies. Statements must include either a Resource or a NotResource element. As a best practice, specify a resource using its [Amazon Resource Name \(ARN\)](#). You can do this for actions that support a specific resource type, known as *resource-level permissions*.

For actions that don't support resource-level permissions, such as listing operations, use a wildcard (*) to indicate that the statement applies to all resources.

```
"Resource": "*"
```

To view examples of Bedrock Agent identity-based policies, see [Identity-based policy examples for Bedrock Agent \(p. 131\)](#).

Policy condition keys for Bedrock Agent

Supports service-specific policy condition keys	No
---	----

Administrators can use AWS JSON policies to specify who has access to what. That is, which **principal** can perform **actions** on what **resources**, and under what **conditions**.

The Condition element (or Condition *block*) lets you specify conditions in which a statement is in effect. The Condition element is optional. You can create conditional expressions that use [condition operators](#), such as equals or less than, to match the condition in the policy with values in the request.

If you specify multiple Condition elements in a statement, or multiple keys in a single Condition element, AWS evaluates them using a logical AND operation. If you specify multiple values for a single condition key, AWS evaluates the condition using a logical OR operation. All of the conditions must be met before the statement's permissions are granted.

You can also use placeholder variables when you specify conditions. For example, you can grant an IAM user permission to access a resource only if it is tagged with their IAM user name. For more information, see [IAM policy elements: variables and tags](#) in the *IAM User Guide*.

AWS supports global condition keys and service-specific condition keys. To see all AWS global condition keys, see [AWS global condition context keys](#) in the *IAM User Guide*.

All Bedrock Agent actions support condition keys using agents or aliases as the resource.

To view examples of Bedrock Agent identity-based policies, see [Identity-based policy examples for Bedrock Agent \(p. 131\)](#).

ACLs in Bedrock Agent

Supports ACLs	No
---------------	----

Access control lists (ACLs) control which principals (account members, users, or roles) have permissions to access a resource. ACLs are similar to resource-based policies, although they do not use the JSON policy document format.

ABAC with Bedrock Agent

Supports ABAC (tags in policies)	No
----------------------------------	----

Attribute-based access control (ABAC) is an authorization strategy that defines permissions based on attributes. In AWS, these attributes are called *tags*. You can attach tags to IAM entities (users or roles) and to many AWS resources. Tagging entities and resources is the first step of ABAC. Then you design ABAC policies to allow operations when the principal's tag matches the tag on the resource that they are trying to access.

ABAC is helpful in environments that are growing rapidly and helps with situations where policy management becomes cumbersome.

To control access based on tags, you provide tag information in the [condition element](#) of a policy using the `aws:ResourceTag/key-name`, `aws:RequestTag/key-name`, or `aws:TagKeys` condition keys.

If a service supports all three condition keys for every resource type, then the value is **Yes** for the service. If a service supports all three condition keys for only some resource types, then the value is **Partial**.

For more information about ABAC, see [What is ABAC?](#) in the *IAM User Guide*. To view a tutorial with steps for setting up ABAC, see [Use attribute-based access control \(ABAC\)](#) in the *IAM User Guide*.

Using temporary credentials with Bedrock Agent

Supports temporary credentials	Yes
--------------------------------	-----

Some AWS services don't work when you sign in using temporary credentials. For additional information, including which AWS services work with temporary credentials, see [AWS services that work with IAM](#) in the *IAM User Guide*.

You are using temporary credentials if you sign in to the AWS Management Console using any method except a user name and password. For example, when you access AWS using your company's single sign-on (SSO) link, that process automatically creates temporary credentials. You also automatically create temporary credentials when you sign in to the console as a user and then switch roles. For more information about switching roles, see [Switching to a role \(console\)](#) in the *IAM User Guide*.

You can manually create temporary credentials using the AWS CLI or AWS API. You can then use those temporary credentials to access AWS. AWS recommends that you dynamically generate temporary credentials instead of using long-term access keys. For more information, see [Temporary security credentials in IAM](#).

Cross-service principal permissions for Bedrock Agent

Supports principal permissions	Yes
--------------------------------	-----

When you use an IAM user or role to perform actions in AWS, you are considered a principal. Policies grant permissions to a principal. When you use some services, you might perform an action that then triggers another action in a different service. In this case, you must have permissions to perform both actions. To see whether an action requires additional dependent actions in a policy, see [Actions, Resources, and Condition Keys for Bedrock Agent](#) in the *Service Authorization Reference*.

Service roles for Bedrock Agent

Supports service roles	Yes
------------------------	-----

A service role is an [IAM role](#) that a service assumes to perform actions on your behalf. An IAM administrator can create, modify, and delete a service role from within IAM. For more information, see [Creating a role to delegate permissions to an AWS service](#) in the *IAM User Guide*.

Warning

Changing the permissions for a service role might break Bedrock Agent functionality. Edit service roles only when Bedrock Agent provides guidance to do so.

When you create an agent resource, you must create or choose a role to allow Bedrock Agent to access resources on your behalf. If you have previously created a service role or service-linked role, Bedrock Agent provides you with a list of roles to choose from. It's important to choose a role that allows access to the required actions.

Service-linked roles for Bedrock Agent

Supports service-linked roles	No
-------------------------------	----

A service-linked role is a type of service role that is linked to an AWS service. The service can assume the role to perform an action on your behalf. Service-linked roles appear in your AWS account and are owned by the service. An IAM administrator can view, but not edit the permissions for service-linked roles.

Identity-based policy examples for Bedrock Agent

By default, users and roles don't have permission to create or modify Bedrock Agent resources. They also can't perform tasks by using the AWS Management Console, AWS Command Line Interface (AWS CLI), or AWS API. To grant users permission to perform actions on the resources that they need, an IAM administrator can create IAM policies. The administrator can then add the IAM policies to roles, and users can assume the roles.

To learn how to create an IAM identity-based policy by using these example JSON policy documents, see [Creating IAM policies](#) in the *IAM User Guide*.

For details about actions and resource types defined by Bedrock Agent, including the format of the ARNs for each of the resource types, see [Actions, Resources, and Condition Keys for Bedrock Agent](#) in the *Service Authorization Reference*.

Note

The Bedrock Agent service is available as a limited preview release, so its information is not included in the Service Authorization Reference.

Topics

- [Policy best practices \(p. 131\)](#)
- [Using the Bedrock Agent console \(p. 132\)](#)
- [Allow users to view their own permissions \(p. 132\)](#)
- [Allow users to perform actions on agent and alias resources \(p. 133\)](#)

Policy best practices

Identity-based policies determine whether someone can create, access, or delete Bedrock Agent resources in your account. These actions can incur costs for your AWS account. When you create or edit identity-based policies, follow these guidelines and recommendations:

- **Get started with AWS managed policies and move toward least-privilege permissions** – To get started granting permissions to your users and workloads, use the *AWS managed policies* that grant permissions for many common use cases. They are available in your AWS account. We recommend that you reduce permissions further by defining AWS customer managed policies that are specific to your use cases. For more information, see [AWS managed policies](#) or [AWS managed policies for job functions](#) in the *IAM User Guide*.
- **Apply least-privilege permissions** – When you set permissions with IAM policies, grant only the permissions required to perform a task. You do this by defining the actions that can be taken on specific resources under specific conditions, also known as *least-privilege permissions*. For more information about using IAM to apply permissions, see [Policies and permissions in IAM](#) in the *IAM User Guide*.
- **Use conditions in IAM policies to further restrict access** – You can add a condition to your policies to limit access to actions and resources. For example, you can write a policy condition to specify that all requests must be sent using SSL. You can also use conditions to grant access to service actions if they are used through a specific AWS service, such as AWS CloudFormation. For more information, see [IAM JSON policy elements: Condition](#) in the *IAM User Guide*.
- **Use IAM Access Analyzer to validate your IAM policies to ensure secure and functional permissions** – IAM Access Analyzer validates new and existing policies so that the policies adhere to the IAM policy language (JSON) and IAM best practices. IAM Access Analyzer provides more than 100 policy

checks and actionable recommendations to help you author secure and functional policies. For more information, see [IAM Access Analyzer policy validation](#) in the *IAM User Guide*.

- **Require multi-factor authentication (MFA)** – If you have a scenario that requires IAM users or a root user in your AWS account, turn on MFA for additional security. To require MFA when API operations are called, add MFA conditions to your policies. For more information, see [Configuring MFA-protected API access](#) in the *IAM User Guide*.

For more information about best practices in IAM, see [Security best practices in IAM](#) in the *IAM User Guide*.

Using the Bedrock Agent console

To access the Bedrock Agent console, you must have a minimum set of permissions. These permissions must allow you to list and view details about the Bedrock Agent resources in your AWS account. If you create an identity-based policy that is more restrictive than the minimum required permissions, the console won't function as intended for entities (users or roles) with that policy.

You don't need to allow minimum console permissions for users that are making calls only to the AWS CLI or the AWS API. Instead, allow access to only the actions that match the API operation that they're trying to perform.

To provide access to the Bedrock Agent console, attach the following policy to the roles or entities that need access.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "BedrockConsole",
      "Effect": "Allow",
      "Action": [
        "bedrock:CreateAgent",
        "bedrock:UpdateAgent",
        "bedrock:GetAgent",
        "bedrock:ListAgents",
        "bedrock:CreateActionGroup",
        "bedrock:UpdateActionGroup",
        "bedrock:GetActionGroup",
        "bedrock:ListActionGroups",
        "bedrock:CreateAgentDraftSnapshot",
        "bedrock:GetAgentVersion",
        "bedrock:ListAgentVersions",
        "bedrock:CreateAgentAlias",
        "bedrock:UpdateAgentAlias",
        "bedrock:GetAgentAlias",
        "bedrock:ListAgentAliases",
        "bedrock:InvokeAgent"
      ],
      "Resource": "*"
    }
  ]
}
```

For more information, see [Adding permissions to a user](#) in the *IAM User Guide*.

Allow users to view their own permissions

This example shows how you might create a policy that allows IAM users to view the inline and managed policies that are attached to their user identity. This policy includes permissions to complete this action on the console or programmatically using the AWS CLI or AWS API.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "ViewOwnUserInfo",
      "Effect": "Allow",
      "Action": [
        "iam:GetUserPolicy",
        "iam:ListGroupForUser",
        "iam:ListAttachedUserPolicies",
        "iam:ListUserPolicies",
        "iam:GetUser"
      ],
      "Resource": ["arn:aws:iam:*:user/${aws:username}"]
    },
    {
      "Sid": "NavigateInConsole",
      "Effect": "Allow",
      "Action": [
        "iam:GetGroupPolicy",
        "iam:GetPolicyVersion",
        "iam:GetPolicy",
        "iam:ListAttachedGroupPolicies",
        "iam:ListGroupPolicies",
        "iam:ListPolicyVersions",
        "iam:ListPolicies",
        "iam:ListUsers"
      ],
      "Resource": "*"
    }
  ]
}
```

Allow users to perform actions on agent and alias resources

You can provision identities with permissions to perform actions on agent and alias resources. The ARNs of these resources are formatted as follows.

- Agents `arn:aws:bedrock:region:account-id:agent/AGENTID`
- Aliases `arn:aws:bedrock:region:account-id:agent-alias/AGENTID/ALIASID`

A role can call API operations on specific resources. For example, the `InvokeAgent` operation can only be used on alias resources and the `UpdateAgent` operation can only be used on agent resources. If you specify an operation in a policy that can't be used on the resource specified in the policy, Bedrock Agent returns an error. For a list of operations and the resources that they can be used with, see the following table. `CreateAgent` and `ListAgents` are not performed on a specific resource.

Operation	Resource
CreateAgent	N/A
UpdateAgent	Agent
GetAgent	Agent
ListAgents	N/A
CreateActionGroup	Agent
UpdateActionGroup	Agent

Operation	Resource
GetActionGroup	Agent
ListActionGroups	Agent
CreateAgentDraftSnapshot	Agent
GetAgentVersion	Agent
ListAgentVersions	Agent
CreateAgentAlias	Agent
UpdateAgentAlias	Alias
GetAgentAlias	Alias
ListAgentAliases	Agent
InvokeAgent	Alias

The following is a sample policy that you can attach to an IAM role to allow it to call Bedrock Agent API operations to get information about an agent, update an agent alias, and interact with an agent. Replace the *sid* with a policy identifier of your choice, the *account-id* with the account ID to which the agent belongs, the *AGENTID* with the ID of the agent, and the *ALIASID* with the alias of the agent. You can find agent and alias IDs with the `ListAgents` API operation or in the agent and alias details pages in the console.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "sid",
      "Effect": "Allow",
      "Action": "bedrock:GetAgent",
      "Resource": "arn:aws:bedrock:region:account-id:agent/AGENTID"
    },
    {
      "Sid": "sid",
      "Effect": "Allow",
      "Action": [
        "bedrock:UpdateAgentAlias",
        "bedrock:InvokeAgent"
      ],
      "Resource": [
        "arn:aws:bedrock:region:account-id:agent-alias/AGENTID/ALIASID"
      ]
    }
  ]
}
```

Tag resources

To help you manage your Amazon Bedrock resources, you can assign metadata to each resource as tags. A tag is a label that you assign to an AWS resource. Each tag consists of a key and a value.

The Amazon Bedrock resources that you can tag are:

- Custom models
- Model customization jobs
- Provisioned models

Tags enable you to categorize your AWS resources in different ways, for example, by purpose, owner, or application. Tags help you to do the following:

- Identify and organize your AWS resources. Many AWS resources support tagging, so you can assign the same tag to resources in different services to indicate that the resources are the same.
- Allocate costs. You activate tags on the AWS Billing and Cost Management dashboard. AWS uses the tags to categorize your costs and deliver a monthly cost allocation report to you. For more information, see [Use cost allocation tags](#) in the *AWS Billing and Cost Management User Guide*.
- Control access to your resources. You can use tags with Amazon Bedrock to create policies to control access to Amazon Bedrock resources. These policies can be attached to an IAM role or user to enable tag-based access control.

Topics

- [Use the console \(p. 135\)](#)
- [Use APIs \(p. 136\)](#)
- [Tag restrictions \(p. 136\)](#)

Use the console

Tagging is supported for custom models and fine-tuning jobs. You can add and modify tags at any time in the fine-tuning process.

To add tags when submitting a new fine-tuning job:

1. Select **Fine-tune** from the left sidebar.
2. In the **Models** section, choose **Fine-tune model**.
3. In the **Model configuration** or **Job configuration** section, expand the **Tags** section.
4. Select **Add new tag** and fill in the **Key** and **Value** for the tag. You can add up to 50 tags in total.
5. Select **Remove tag** next to a key-value pair to remove that tag.

To add tags to an existing custom model or fine-tuning job:

1. Select **Fine-tune** from the left sidebar.
2. In the **Models** or **Training jobs** section, choose the model or job that you want to add tags to.
3. In the **Tags** section at the bottom, select **Manage Tags**.
4. Select **Add new tag** and fill in the **Key** and **Value** for the tag. You can add up to 50 tags in total.

5. Select **Remove tag** next to a key-value pair to remove that tag.

Use APIs

To carry out tagging operations, you need the Amazon Resource Name (ARN) of the resource on which you want to carry out a tagging operation. Use the [TagResource](#) and [UntagResource](#) operations to tag and untag resources. To list the tags for a resource, use the [ListTagsForResource](#) operation. For tagging examples, see [Tag resources \(p. 29\)](#).

Tag restrictions

The following basic restrictions apply to tags on Amazon Bedrock resources:

- Maximum number of tags per resource – 50.
- Maximum number of keys – 50.
- Maximum key length – 128 characters.
- Maximum value length – 256 characters.
- Valid characters for key and value – a-z, A-Z, 0-9, space, and the following characters: `_:/=+-` and `@`
- Keys and values are case-sensitive.
- Don't use `aws :` as a prefix for keys. This is reserved for AWS use.

Security in Amazon Bedrock

Cloud security at AWS is the highest priority. As an AWS customer, you benefit from data centers and network architectures that are built to meet the requirements of the most security-sensitive organizations.

Security is a shared responsibility between AWS and you. The [shared responsibility model](#) describes this as security *of* the cloud and security *in* the cloud:

- **Security of the cloud** – AWS is responsible for protecting the infrastructure that runs AWS services in the AWS Cloud. AWS also provides you with services that you can use securely. Third-party auditors regularly test and verify the effectiveness of our security as part of the [AWS Compliance Programs](#). To learn about the compliance programs that apply to Amazon Bedrock, see [AWS Services in Scope by Compliance Program](#).
- **Security in the cloud** – Your responsibility is determined by the AWS service that you use. You are also responsible for other factors including the sensitivity of your data, your company's requirements, and applicable laws and regulations.

This documentation helps you understand how to apply the shared responsibility model when using Amazon Bedrock. The following topics show you how to configure Amazon Bedrock to meet your security and compliance objectives. You also learn how to use other AWS services that help you to monitor and secure your Amazon Bedrock resources.

Topics

- [Data protection in Amazon Bedrock \(p. 137\)](#)
- [Identity and access management for Amazon Bedrock \(p. 146\)](#)
- [Compliance validation for Amazon Bedrock \(p. 161\)](#)
- [Incident response in Amazon Bedrock \(p. 162\)](#)
- [Resilience in Amazon Bedrock \(p. 162\)](#)
- [Infrastructure security in Amazon Bedrock \(p. 162\)](#)
- [Cross-service confused deputy prevention \(p. 163\)](#)
- [Configuration and vulnerability analysis in Amazon Bedrock \(p. 164\)](#)

Data protection in Amazon Bedrock

The AWS [shared responsibility model](#) applies to data protection in Amazon Bedrock. As described in this model, AWS is responsible for protecting the global infrastructure that runs all of the AWS Cloud. You are responsible for maintaining control over your content that is hosted on this infrastructure. This content includes the security configuration and management tasks for the AWS services that you use. For more information about data privacy, see the [Data Privacy FAQ](#). For information about data protection in Europe, see the [AWS Shared Responsibility Model and GDPR](#) blog post on the *AWS Security Blog*.

For data protection purposes, we recommend that you protect AWS account credentials and set up individual users with AWS IAM Identity Center or AWS Identity and Access Management (IAM). That way, each user is given only the permissions necessary to fulfill their job duties. We also recommend that you secure your data in the following ways:

- Use multi-factor authentication (MFA) with each account.
- Use SSL/TLS to communicate with AWS resources. We require TLS 1.2 and recommend TLS 1.3.
- Set up API and user activity logging with AWS CloudTrail.
- Use AWS encryption solutions, along with all default security controls within AWS services.
- Use advanced managed security services such as Amazon Macie, which assists in discovering and securing sensitive data that is stored in Amazon S3.
- If you require FIPS 140-2 validated cryptographic modules when accessing AWS through a command line interface or an API, use a FIPS endpoint. For more information about the available FIPS endpoints, see [Federal Information Processing Standard \(FIPS\) 140-2](#).

We strongly recommend that you never put confidential or sensitive information, such as your customers' email addresses, into tags or free-form text fields such as a **Name** field. This includes when you work with Amazon Bedrock or other AWS services using the console, API, AWS CLI, or AWS SDKs. Any data that you enter into tags or free-form text fields used for names may be used for billing or diagnostic logs. If you provide a URL to an external server, we strongly recommend that you do not include credentials information in the URL to validate your request to that server.

Amazon Bedrock doesn't use your prompts and continuations to train any AWS models or distribute them to third parties. Your training data isn't used to train the base Amazon Titan models or distributed to third parties. Other usage data, such as usage timestamps, logged account IDs, and other information logged by the service, is also not used to train the models.

Amazon Bedrock uses the fine tuning data you provide only for fine tuning an Amazon Titan model. Amazon Bedrock doesn't use fine tuning data for any other purpose, such as training base foundation models.

Each model provider has an escrow account that they upload their models to. The Amazon Bedrock inference account has permissions to call these models, but the escrow accounts themselves don't have outbound permissions to Amazon Bedrock accounts. Additionally, model providers don't have access to Amazon Bedrock logs or access to customer prompts and continuations.

Amazon Bedrock doesn't store or log your data in its service logs.

Bedrock uses your training data with the [CreateModelCustomizationJob](#) action, or with the [console \(p. 75\)](#), to create a custom model which is a fine tuned version of an Amazon Bedrock foundational model. Your custom models are managed and stored by AWS. By default, custom models are encrypted with AWS Key Management Service keys that AWS owns, but you can use your own AWS KMS keys to encrypt your custom models. You encrypt a custom model when you submit a fine tuning job with the console or programmatically with the `CreateModelCustomizationJob` action.

None of the training or validation data you provide for fine tuning is stored in Amazon Bedrock accounts, once the fine tuning job completes. During training, your data exists in AWS Service Management Connector instance memory, but is encrypted on these machines using an XTS-AES-256 cipher that is implemented on a hardware module, on the instance itself.

We don't recommend using confidential data to train a custom model as the model might generate inference responses based on that confidential data. If you use confidential data to train a custom model, the only way to prevent responses based on that data is to delete the custom model, remove the confidential data from your training dataset, and retrain the custom model.

Custom model metadata (name and Amazon Resource Name) and a provisioned model's metadata is stored in an Amazon DynamoDB table that is encrypted with a key that the Amazon Bedrock service owns.

Topics

- [Data encryption \(p. 139\)](#)

- [Protect jobs using a VPC \(p. 141\)](#)

Data encryption

Amazon Bedrock uses encryption to protect data at rest and data in transit.

Topics

- [Encryption in transit \(p. 139\)](#)
- [Encryption at rest \(p. 139\)](#)
- [Key management \(p. 139\)](#)
- [Encryption of model customization jobs \(p. 139\)](#)

Encryption in transit

Within AWS, all inter-network data in transit supports TLS 1.2 encryption.

Requests to the Amazon Bedrock API and console are made over a secure (SSL) connection. You pass AWS Identity and Access Management (IAM) roles to Amazon Bedrock to provide permissions to access resources on your behalf for training and deployment.

Encryption at rest

Amazon Bedrock provides [Encryption of model customization jobs \(p. 139\)](#) at rest.

Key management

You can use AWS Key Management Service (KMS) to manage the keys that you use with model customization jobs. For more information, see [AWS Key Management Service concepts](#).

Encryption of model customization jobs

Amazon Bedrock encrypts the model artifacts from your model customization jobs. By default, Amazon Bedrock encrypts this data using an AWS managed key. Optionally, you can encrypt the model artifacts using a customer managed key.

For more information about AWS KMS keys, see [Customer managed keys](#) in the *AWS Key Management Service Developer Guide*.

Topics

- [Model customization inputs \(p. 139\)](#)
- [Model customization outputs \(p. 140\)](#)
- [Create a customer managed key \(p. 140\)](#)
- [Use a customer managed key to run customization jobs \(p. 140\)](#)
- [Use a customer managed key during inference \(p. 141\)](#)

Model customization inputs

When you use Amazon Bedrock to run a model customization job, you store the input documents (training/validation data) in your Amazon S3 bucket. To encrypt these documents at rest, you can use the Amazon S3 SSE-S3 server-side encryption option. With this option, objects are encrypted with service keys managed by the Amazon S3 service.

For more information, see [Protecting data using server-side encryption with Amazon S3-managed encryption keys \(SSE-S3\)](#) in the *Amazon Simple Storage Service User Guide*.

Model customization outputs

When Amazon Bedrock completes a model customization job, it stores the job metrics in the Amazon S3 location that you specified when you created the job. To encrypt the metrics, you can use the Amazon S3 SSE-S3 server-side encryption option described in the previous section.

Amazon Bedrock stores the custom model artifacts in an Amazon S3 bucket controlled by AWS. By default, Amazon Bedrock encrypts this data using an *AWS managed key*. This type of KMS key is created by AWS, so you don't manage this KMS key yourself. AWS manages the key and uses it on your behalf.

Optionally, you can choose to encrypt the custom model artifacts with a *customer managed key*. This is a KMS key that you create, own, and manage in your AWS account.

Before you can use your own KMS key, configure the policies and permissions as described in the following sections.

Create a customer managed key

Any user with `CreateKey` permissions can create customer managed keys using either the AWS Key Management Service (AWS KMS) console or the [CreateKey](#) API operation. Make sure to create a symmetric encryption key.

Create a key policy and add the following policy statements to grant permissions to custom model builders and users.

```
{
  "Version": "2012-10-17",
  "Id": "KMS Key Policy",
  "Statement": [
    {
      "Sid": "Permissions for custom model users",
      "Effect": "Allow",
      "Principal": {
        "AWS": "arn:aws:iam::111122223333:user/CustomModelCaller"
      },
      "Action": "kms:Decrypt",
      "Resource": "*"
    },
    {
      "Sid": "Permissions for custom model builders",
      "Effect": "Allow",
      "Principal": {
        "AWS": "arn:aws:iam::111122223333:user/CustomModelBuilder"
      },
      "Action": [
        "kms:Decrypt",
        "kms:GenerateDataKey",
        "kms:DescribeKey",
        "kms:CreateGrant"
      ],
      "Resource": "*"
    }
  ]
}
```

Use a customer managed key to run customization jobs

For users to create a model customization job, the user or role needs the following permissions to use the customer managed key.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "Bedrock custom model builder policy",
      "Effect": "Allow",
      "Action": [
        "kms:Decrypt",
        "kms:DescribeKey",
        "kms:GenerateDataKey",
        "kms:CreateGrant"
      ],
      "Resource": "arn:aws:kms:region:111122223333:key/key id"
    }
  ]
}
```

Use a customer managed key during inference

For users to run inference on a custom model, the user or role needs the following permissions to use the customer managed key.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "Bedrock Customer Invocation Policy",
      "Effect": "Allow",
      "Action": [
        "kms:Decrypt"
      ],
      "Resource": "arn:aws:kms:region:111122223333:key/key id"
    }
  ]
}
```

When you invoke a model to run inference (for a custom model encrypted with a customer managed key), if you don't have `kms:Decrypt` permissions for that key, the request fails with the following error message:

You don't have sufficient access to the model's KMS key. Ensure that `kms:Decrypt` permissions are correctly configured

Protect jobs using a VPC

When you run a model customization job, the job accesses your Amazon S3 bucket to download the input data and to upload job metrics.

To control access to your data, we recommend that you create a virtual private cloud (VPC). Configure it so that your training data is not accessible over the internet. For information about creating and configuring a VPC, see [Getting Started With Amazon VPC](#) in the *Amazon VPC User Guide*.

Using a VPC protects your data and lets you monitor all network traffic in and out of the AWS job containers by using VPC Flow Logs. For more information, see [VPC Flow Logs](#) in the *Amazon VPC User Guide*.

If you configure your VPC with no internet access, you must create a VPC endpoint that allows the customization job to access these S3 buckets.

When you configure your customization job to use VPC, the job creates an elastic network interface (ENI) that uses your VPC endpoint to access your S3 buckets. For information about ENIs, see [Elastic Network Interfaces](#) in the *Amazon VPC User Guide*.

Topics

- [Configure a VPC for Amazon Bedrock \(p. 142\)](#)
- [Configure your model customization job to use VPC \(p. 145\)](#)

Configure a VPC for Amazon Bedrock

When you configure the VPC for an Amazon Bedrock model customization job, use the following guidelines. For information about setting up a VPC, see [Working with VPCs and Subnets](#) in the *Amazon VPC User Guide*.

Topics

- [Create an Amazon S3 VPC Endpoint \(p. 142\)](#)
- [Use a custom endpoint policy to restrict access to Amazon S3 \(p. 142\)](#)
- [VPC permissions for customization job role \(p. 143\)](#)
- [Configure route tables \(p. 145\)](#)

Create an Amazon S3 VPC Endpoint

If you configure your VPC with no internet access, you need create a VPC endpoint. The endpoint allows access to the S3 buckets that contain your training data and the training loss metrics data that's stored by the job.

By creating a VPC endpoint, you allow your model customization jobs to access the buckets where you store your data and model artifacts.

We recommend that you also create a custom policy that allows only requests from your private VPC to access your S3 buckets. For more information, see [Use interface VPC endpoints \(AWS PrivateLink\) \(p. 174\)](#).

To create an Amazon S3 VPC endpoint

1. Open the Amazon VPC console at <https://console.aws.amazon.com/vpc/>.
2. In the navigation pane, choose **Endpoints**, then choose **Create Endpoint**.
3. For **Service Name**, search for `com.amazonaws.region.s3`. Replace *region* with the name of the Region where your VPC resides.
4. Choose the **Gateway** type.
5. For **VPC**, choose the VPC you want to use for this endpoint.
6. For **Configure route tables**, select the route tables to be used by the endpoint. Amazon VPC automatically adds a route to each selected route table that points any Amazon S3 traffic to the new endpoint.
7. For **Policy**, choose **Full Access** to allow full access to Amazon S3 by any user or service within the VPC. Choose **Custom** to restrict access further. For information, see [Use a custom endpoint policy to restrict access to Amazon S3 \(p. 142\)](#).

Use a custom endpoint policy to restrict access to Amazon S3

The default endpoint policy allows full access to Amazon S3 for any user or service in your VPC. To further restrict access to Amazon S3, create a custom endpoint policy. For more information, see [Using](#)

[Endpoint Policies for Amazon S3](#). You can also use a bucket policy to restrict access to your S3 buckets to only traffic that comes from your VPC.

The following policy allows access to S3 buckets. Edit this policy to allow access to only the resources that your job needs.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "Allow access to output location",
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:PutObject",
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3::my_output_bucket/myfolder"
      ]
    },
    {
      "Sid": "Allow access to input data",
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:PutObject",
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3::my_training_data_bucket/myfolder",
        "arn:aws:s3::my_validation_data_bucket/myfolder"
      ]
    }
  ]
}
```

For information, see [Using Amazon S3 Bucket Policies](#).

VPC permissions for customization job role

When you use a VPC with your model customization job, the data access role that you provide in the `CreateModelCustomizationJob` request must include the following permissions. Edit this policy to allow access to only the VPC resources that your job needs.

```
{
  "Effect": "Allow",
  "Action": [
    "ec2:DescribeNetworkInterfaces",
    "ec2:DescribeVpcs",
    "ec2:DescribeDhcpOptions",
    "ec2:DescribeSubnets",
    "ec2:DescribeSecurityGroups"
  ],
  "Resource": "*"
}, {
  "Effect": "Allow",
  "Action": [
    "ec2:CreateNetworkInterface"
  ],
  "Resource": [
    "arn:aws:ec2:{region}:{AccountId}:network-interface/*"
  ]
}
```

```

    ],
    "Condition": {
      "StringEquals": {
        "aws:RequestTag/BedrockManaged": ["true"]
      },
      "ArnEquals": {
        "aws:RequestTag/BedrockModelCustomizationJobArn": ["arn:aws:bedrock:
{region}:{AccountId}:model-customization-job/*"]
      }
    }
  }, {
    "Effect": "Allow",
    "Action": [
      "ec2:CreateNetworkInterface",
    ],
    "Resource": [
      "arn:aws:ec2:{region}:{AccountId}:subnet/{subnet-Id}",
      "arn:aws:ec2:{region}:{AccountId}:subnet/{subnet-Id2}",
      "arn:aws:ec2:{region}:{AccountId}:security-group/{sg-Id}",
      "arn:aws:ec2:{region}:{AccountId}:security-group/{sg-Id2}"
    ],
  }, {
    "Effect": "Allow",
    "Action": [
      "ec2:CreateNetworkInterfacePermission",
      "ec2:DeleteNetworkInterface",
      "ec2:DeleteNetworkInterfacePermission",
    ],
    "Resource": "*",
    "Condition": {
      "ArnEquals": {
        "ec2:Subnet": [
          "arn:aws:ec2:{region}:{AccountId}:subnet/{subnet-Id}",
          "arn:aws:ec2:{region}:{AccountId}:subnet/{subnet-Id2}"
        ],
        "ec2:ResourceTag/BedrockModelCustomizationJobArn": ["arn:aws:bedrock:
{region}:{AccountId}:model-customization-job/*"]
      },
      "StringEquals": {
        "ec2:ResourceTag/BedrockManaged": "true",
      }
    }
  }, {
    "Effect": "Allow",
    "Action": [
      "ec2:CreateTags"
    ],
    "Resource": "arn:aws:ec2:{region}:{AccountId}:network-interface/*",
    "Condition": {
      "StringEquals": {
        "ec2:CreateAction": [
          "CreateNetworkInterface"
        ]
      },
      "ForAllValues:StringEquals": {
        "aws:TagKeys": [
          "BedrockManaged",
          "BedrockModelCustomizationJobArn"
        ]
      }
    }
  }
}

```

For more information about configuring the data access role for model customization jobs, see [Set up an IAM role for model customization \(p. 79\)](#).

Configure route tables

Use default DNS settings for your endpoint route table, so that standard Amazon S3 URLs (for example, `http://s3-aws-region.amazonaws.com/MyBucket`) resolve. If you don't use default DNS settings, ensure that the URLs for the locations of the data in your training jobs resolve. Do this by configuring the endpoint route tables. For information about VPC endpoint route tables, see [Routing for Gateway Endpoints](#) in the *Amazon VPC User Guide*.

Configure your model customization job to use VPC

After you configure the VPC and the required roles and permissions as described in the previous sections, you can create a model customization job that uses this VPC.

When you specify the VPC subnets and security groups for a job, Amazon Bedrock creates *elastic network interfaces* (ENIs) that are associated with your security groups in one of the subnets. ENIs allow the Amazon Bedrock job to connect to resources in your VPC. For information about ENIs, see [Elastic Network Interfaces](#) in the *Amazon VPC User Guide*. Amazon Bedrock tags ENIs that it creates with `BedrockManaged` and `BedrockModelCustomizationJobArn` tags.

We recommend that you provide at least one subnet in each Availability Zone.

You can use security groups to establish rules for controlling Amazon Bedrock access to your VPC resources.

Use the API

For the Amazon Bedrock API, you specify VPC subnets and security groups in the `VpcConfig` request parameter. The following is an example of the `VpcConfig` parameter that you include in your API call:

```
"VpcConfig": {
  "SecurityGroupIds": [
    "sg-0123456789abcdef0"
  ],
  "Subnets": [
    "subnet-0123456789abcdef0",
    "subnet-0123456789abcdef1",
    "subnet-0123456789abcdef2"
  ]
}
```

For more information about configuring jobs using the API, see [Submit a job \(API\) \(p. 81\)](#).

Use the console

For the Amazon Bedrock console, you specify VPC subnets and security groups in the optional **VPC settings** section when you create the model customization job. For more information about configuring jobs using the console, see [Submit a job \(console\) \(p. 75\)](#).

Note

For a job that includes VPC configuration, the console cannot create a new service role for the job. Create the service role using the following example.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "bedrock.amazonaws.com"
      }
    }
  ]
}
```



```
    },
    "Action": "sts:AssumeRole",
    "Condition": {
      "StringEquals": {
        "aws:SourceAccount": "111122223333"
      },
      "ArnEquals": {
        "aws:SourceArn": "arn:aws:bedrock:us-east-1:111122223333:model-
customization-job/*"
      }
    }
  }
}
```

Identity and access management for Amazon Bedrock

AWS Identity and Access Management (IAM) is an AWS service that helps an administrator securely control access to AWS resources. IAM administrators control who can be *authenticated* (signed in) and *authorized* (have permissions) to use Amazon Bedrock resources. IAM is an AWS service that you can use with no additional charge.

Topics

- [Audience \(p. 146\)](#)
- [Authenticating with identities \(p. 147\)](#)
- [Managing access using policies \(p. 149\)](#)
- [How Amazon Bedrock works with IAM \(p. 150\)](#)
- [Identity-based policy examples for Amazon Bedrock \(p. 155\)](#)
- [Service role \(p. 160\)](#)
- [Troubleshooting Amazon Bedrock identity and access \(p. 160\)](#)

Audience

How you use AWS Identity and Access Management (IAM) differs, depending on the work that you do in Amazon Bedrock.

Service user – If you use the Amazon Bedrock service to do your job, then your administrator provides you with the credentials and permissions that you need. As you use more Amazon Bedrock features to do your work, you might need additional permissions. Understanding how access is managed can help you request the right permissions from your administrator. If you cannot access a feature in Amazon Bedrock, see [Troubleshooting Amazon Bedrock identity and access \(p. 160\)](#).

Service administrator – If you're in charge of Amazon Bedrock resources at your company, you probably have full access to Amazon Bedrock. It's your job to determine which Amazon Bedrock features and resources your service users should access. You must then submit requests to your IAM administrator to change the permissions of your service users. Review the information on this page to understand the basic concepts of IAM. To learn more about how your company can use IAM with Amazon Bedrock, see [How Amazon Bedrock works with IAM \(p. 150\)](#).

IAM administrator – If you're an IAM administrator, you might want to learn details about how you can write policies to manage access to Amazon Bedrock. To view example Amazon Bedrock identity-based policies that you can use in IAM, see [Identity-based policy examples for Amazon Bedrock \(p. 155\)](#).

Authenticating with identities

Authentication is how you sign in to AWS using your identity credentials. You must be *authenticated* (signed in to AWS) as the AWS account root user, as an IAM user, or by assuming an IAM role.

You can sign in to AWS as a federated identity by using credentials provided through an identity source. AWS IAM Identity Center (IAM Identity Center) users, your company's single sign-on authentication, and your Google or Facebook credentials are examples of federated identities. When you sign in as a federated identity, your administrator previously set up identity federation using IAM roles. When you access AWS by using federation, you are indirectly assuming a role.

Depending on the type of user you are, you can sign in to the AWS Management Console or the AWS access portal. For more information about signing in to AWS, see [How to sign in to your AWS account](#) in the *AWS Sign-In User Guide*.

If you access AWS programmatically, AWS provides a software development kit (SDK) and a command line interface (CLI) to cryptographically sign your requests by using your credentials. If you don't use AWS tools, you must sign requests yourself. For more information about using the recommended method to sign requests yourself, see [Signing AWS API requests](#) in the *IAM User Guide*.

Regardless of the authentication method that you use, you might be required to provide additional security information. For example, AWS recommends that you use multi-factor authentication (MFA) to increase the security of your account. To learn more, see [Multi-factor authentication](#) in the *AWS IAM Identity Center User Guide* and [Using multi-factor authentication \(MFA\) in AWS](#) in the *IAM User Guide*.

AWS account root user

When you create an AWS account, you begin with one sign-in identity that has complete access to all AWS services and resources in the account. This identity is called the AWS account *root user* and is accessed by signing in with the email address and password that you used to create the account. We strongly recommend that you don't use the root user for your everyday tasks. Safeguard your root user credentials and use them to perform the tasks that only the root user can perform. For the complete list of tasks that require you to sign in as the root user, see [Tasks that require root user credentials](#) in the *IAM User Guide*.

Federated identity

As a best practice, require human users, including users that require administrator access, to use federation with an identity provider to access AWS services by using temporary credentials.

A *federated identity* is a user from your enterprise user directory, a web identity provider, the AWS Directory Service, the Identity Center directory, or any user that accesses AWS services by using credentials provided through an identity source. When federated identities access AWS accounts, they assume roles, and the roles provide temporary credentials.

For centralized access management, we recommend that you use AWS IAM Identity Center. You can create users and groups in IAM Identity Center, or you can connect and synchronize to a set of users and groups in your own identity source for use across all your AWS accounts and applications. For information about IAM Identity Center, see [What is IAM Identity Center?](#) in the *AWS IAM Identity Center User Guide*.

IAM users and groups

An *IAM user* is an identity within your AWS account that has specific permissions for a single person or application. Where possible, we recommend relying on temporary credentials instead of creating IAM

users who have long-term credentials such as passwords and access keys. However, if you have specific use cases that require long-term credentials with IAM users, we recommend that you rotate access keys. For more information, see [Rotate access keys regularly for use cases that require long-term credentials](#) in the *IAM User Guide*.

An [IAM group](#) is an identity that specifies a collection of IAM users. You can't sign in as a group. You can use groups to specify permissions for multiple users at a time. Groups make permissions easier to manage for large sets of users. For example, you could have a group named *IAMAdmins* and give that group permissions to administer IAM resources.

Users are different from roles. A user is uniquely associated with one person or application, but a role is intended to be assumable by anyone who needs it. Users have permanent long-term credentials, but roles provide temporary credentials. To learn more, see [When to create an IAM user \(instead of a role\)](#) in the *IAM User Guide*.

IAM roles

An [IAM role](#) is an identity within your AWS account that has specific permissions. It is similar to an IAM user, but is not associated with a specific person. You can temporarily assume an IAM role in the AWS Management Console by [switching roles](#). You can assume a role by calling an AWS CLI or AWS API operation or by using a custom URL. For more information about methods for using roles, see [Using IAM roles](#) in the *IAM User Guide*.

IAM roles with temporary credentials are useful in the following situations:

- **Federated user access** – To assign permissions to a federated identity, you create a role and define permissions for the role. When a federated identity authenticates, the identity is associated with the role and is granted the permissions that are defined by the role. For information about roles for federation, see [Creating a role for a third-party Identity Provider](#) in the *IAM User Guide*. If you use IAM Identity Center, you configure a permission set. To control what your identities can access after they authenticate, IAM Identity Center correlates the permission set to a role in IAM. For information about permissions sets, see [Permission sets](#) in the *AWS IAM Identity Center User Guide*.
- **Temporary IAM user permissions** – An IAM user or role can assume an IAM role to temporarily take on different permissions for a specific task.
- **Cross-account access** – You can use an IAM role to allow someone (a trusted principal) in a different account to access resources in your account. Roles are the primary way to grant cross-account access. However, with some AWS services, you can attach a policy directly to a resource (instead of using a role as a proxy). To learn the difference between roles and resource-based policies for cross-account access, see [How IAM roles differ from resource-based policies](#) in the *IAM User Guide*.
- **Cross-service access** – Some AWS services use features in other AWS services. For example, when you make a call in a service, it's common for that service to run applications in Amazon EC2 or store objects in Amazon S3. A service might do this using the calling principal's permissions, using a service role, or using a service-linked role.
- **Principal permissions** – When you use an IAM user or role to perform actions in AWS, you are considered a principal. Policies grant permissions to a principal. When you use some services, you might perform an action that then triggers another action in a different service. In this case, you must have permissions to perform both actions. To see whether an action requires additional dependent actions in a policy, see [Actions, Resources, and Condition Keys for Amazon Bedrock](#) in the *Service Authorization Reference*.
- **Service role** – A service role is an [IAM role](#) that a service assumes to perform actions on your behalf. An IAM administrator can create, modify, and delete a service role from within IAM. For more information, see [Creating a role to delegate permissions to an AWS service](#) in the *IAM User Guide*.
- **Service-linked role** – A service-linked role is a type of service role that is linked to an AWS service. The service can assume the role to perform an action on your behalf. Service-linked roles appear in your AWS account and are owned by the service. An IAM administrator can view, but not edit the permissions for service-linked roles.

- **Applications running on Amazon EC2** – You can use an IAM role to manage temporary credentials for applications that are running on an EC2 instance and making AWS CLI or AWS API requests. This is preferable to storing access keys within the EC2 instance. To assign an AWS role to an EC2 instance and make it available to all of its applications, you create an instance profile that is attached to the instance. An instance profile contains the role and enables programs that are running on the EC2 instance to get temporary credentials. For more information, see [Using an IAM role to grant permissions to applications running on Amazon EC2 instances](#) in the *IAM User Guide*.

To learn whether to use IAM roles or IAM users, see [When to create an IAM role \(instead of a user\)](#) in the *IAM User Guide*.

Managing access using policies

You control access in AWS by creating policies and attaching them to AWS identities or resources. A policy is an object in AWS that, when associated with an identity or resource, defines their permissions. AWS evaluates these policies when a principal (user, root user, or role session) makes a request. Permissions in the policies determine whether the request is allowed or denied. Most policies are stored in AWS as JSON documents. For more information about the structure and contents of JSON policy documents, see [Overview of JSON policies](#) in the *IAM User Guide*.

Administrators can use AWS JSON policies to specify who has access to what. That is, which **principal** can perform **actions** on what **resources**, and under what **conditions**.

By default, users and roles have no permissions. To grant users permission to perform actions on the resources that they need, an IAM administrator can create IAM policies. The administrator can then add the IAM policies to roles, and users can assume the roles.

IAM policies define permissions for an action regardless of the method that you use to perform the operation. For example, suppose that you have a policy that allows the `iam:GetRole` action. A user with that policy can get role information from the AWS Management Console, the AWS CLI, or the AWS API.

Identity-based policies

Identity-based policies are JSON permissions policy documents that you can attach to an identity, such as an IAM user, group of users, or role. These policies control what actions users and roles can perform, on which resources, and under what conditions. To learn how to create an identity-based policy, see [Creating IAM policies](#) in the *IAM User Guide*.

Identity-based policies can be further categorized as *inline policies* or *managed policies*. Inline policies are embedded directly into a single user, group, or role. Managed policies are standalone policies that you can attach to multiple users, groups, and roles in your AWS account. Managed policies include AWS managed policies and customer managed policies. To learn how to choose between a managed policy or an inline policy, see [Choosing between managed policies and inline policies](#) in the *IAM User Guide*.

Resource-based policies

Resource-based policies are JSON policy documents that you attach to a resource. Examples of resource-based policies are IAM *role trust policies* and Amazon S3 *bucket policies*. In services that support resource-based policies, service administrators can use them to control access to a specific resource. For the resource where the policy is attached, the policy defines what actions a specified principal can perform on that resource and under what conditions. You must [specify a principal](#) in a resource-based policy. Principals can include accounts, users, roles, federated users, or AWS services.

Resource-based policies are inline policies that are located in that service. You can't use AWS managed policies from IAM in a resource-based policy.

Access control lists (ACLs)

Access control lists (ACLs) control which principals (account members, users, or roles) have permissions to access a resource. ACLs are similar to resource-based policies, although they do not use the JSON policy document format.

Amazon S3, AWS WAF, and Amazon VPC are examples of services that support ACLs. To learn more about ACLs, see [Access control list \(ACL\) overview](#) in the *Amazon Simple Storage Service Developer Guide*.

Other policy types

AWS supports additional, less-common policy types. These policy types can set the maximum permissions granted to you by the more common policy types.

- **Permissions boundaries** – A permissions boundary is an advanced feature in which you set the maximum permissions that an identity-based policy can grant to an IAM entity (IAM user or role). You can set a permissions boundary for an entity. The resulting permissions are the intersection of an entity's identity-based policies and its permissions boundaries. Resource-based policies that specify the user or role in the `Principal` field are not limited by the permissions boundary. An explicit deny in any of these policies overrides the allow. For more information about permissions boundaries, see [Permissions boundaries for IAM entities](#) in the *IAM User Guide*.
- **Service control policies (SCPs)** – SCPs are JSON policies that specify the maximum permissions for an organization or organizational unit (OU) in AWS Organizations. AWS Organizations is a service for grouping and centrally managing multiple AWS accounts that your business owns. If you enable all features in an organization, then you can apply service control policies (SCPs) to any or all of your accounts. The SCP limits permissions for entities in member accounts, including each AWS account root user. For more information about Organizations and SCPs, see [How SCPs work](#) in the *AWS Organizations User Guide*.
- **Session policies** – Session policies are advanced policies that you pass as a parameter when you programmatically create a temporary session for a role or federated user. The resulting session's permissions are the intersection of the user or role's identity-based policies and the session policies. Permissions can also come from a resource-based policy. An explicit deny in any of these policies overrides the allow. For more information, see [Session policies](#) in the *IAM User Guide*.

Multiple policy types

When multiple types of policies apply to a request, the resulting permissions are more complicated to understand. To learn how AWS determines whether to allow a request when multiple policy types are involved, see [Policy evaluation logic](#) in the *IAM User Guide*.

How Amazon Bedrock works with IAM

Before you use IAM to manage access to Amazon Bedrock, learn what IAM features are available to use with Amazon Bedrock.

IAM features you can use with Amazon Bedrock

IAM feature	Amazon Bedrock support
Identity-based policies (p. 151)	Yes
Resource-based policies (p. 151)	No

IAM feature	Amazon Bedrock support
Policy actions (p. 152)	Yes
Policy resources (p. 152)	Yes
Policy condition keys (p. 153)	Yes
ACLs (p. 153)	No
ABAC (tags in policies) (p. 153)	Partial
Temporary credentials (p. 154)	Yes
Principal permissions (p. 154)	Yes
Service roles (p. 154)	Yes
Service-linked roles (p. 155)	No

To get a high-level view of how Amazon Bedrock and other AWS services work with most IAM features, see [AWS services that work with IAM](#) in the *IAM User Guide*.

Identity-based policies for Amazon Bedrock

Supports identity-based policies	Yes
----------------------------------	-----

Identity-based policies are JSON permissions policy documents that you can attach to an identity, such as an IAM user, group of users, or role. These policies control what actions users and roles can perform, on which resources, and under what conditions. To learn how to create an identity-based policy, see [Creating IAM policies](#) in the *IAM User Guide*.

With IAM identity-based policies, you can specify allowed or denied actions and resources as well as the conditions under which actions are allowed or denied. You can't specify the principal in an identity-based policy because it applies to the user or role to which it is attached. To learn about all of the elements that you can use in a JSON policy, see [IAM JSON policy elements reference](#) in the *IAM User Guide*.

Identity-based policy examples for Amazon Bedrock

To view examples of Amazon Bedrock identity-based policies, see [Identity-based policy examples for Amazon Bedrock \(p. 155\)](#).

Resource-based policies within Amazon Bedrock

Supports resource-based policies	No
----------------------------------	----

Resource-based policies are JSON policy documents that you attach to a resource. Examples of resource-based policies are IAM *role trust policies* and Amazon S3 *bucket policies*. In services that support resource-based policies, service administrators can use them to control access to a specific resource. For the resource where the policy is attached, the policy defines what actions a specified principal can perform on that resource and under what conditions. You must [specify a principal](#) in a resource-based policy. Principals can include accounts, users, roles, federated users, or AWS services.

To enable cross-account access, you can specify an entire account or IAM entities in another account as the principal in a resource-based policy. Adding a cross-account principal to a resource-based policy is only half of establishing the trust relationship. When the principal and the resource are in different AWS accounts, an IAM administrator in the trusted account must also grant the principal entity (user or role) permission to access the resource. They grant permission by attaching an identity-based policy to the entity. However, if a resource-based policy grants access to a principal in the same account, no additional identity-based policy is required. For more information, see [How IAM roles differ from resource-based policies](#) in the *IAM User Guide*.

Policy actions for Amazon Bedrock

Supports policy actions	Yes
-------------------------	-----

Administrators can use AWS JSON policies to specify who has access to what. That is, which **principal** can perform **actions** on what **resources**, and under what **conditions**.

The Action element of a JSON policy describes the actions that you can use to allow or deny access in a policy. Policy actions usually have the same name as the associated AWS API operation. There are some exceptions, such as *permission-only actions* that don't have a matching API operation. There are also some operations that require multiple actions in a policy. These additional actions are called *dependent actions*.

Include actions in a policy to grant permissions to perform the associated operation.

To see a list of Amazon Bedrock actions, see [Actions Defined by Amazon Bedrock](#) in the *Service Authorization Reference*.

Policy actions in Amazon Bedrock use the following prefix before the action:

```
bedrock
```

To specify multiple actions in a single statement, separate them with commas.

```
"Action": [  
    "bedrock:action1",  
    "bedrock:action2"  
]
```

To view examples of Amazon Bedrock identity-based policies, see [Identity-based policy examples for Amazon Bedrock \(p. 155\)](#).

Policy resources for Amazon Bedrock

Supports policy resources	Yes
---------------------------	-----

Administrators can use AWS JSON policies to specify who has access to what. That is, which **principal** can perform **actions** on what **resources**, and under what **conditions**.

The Resource JSON policy element specifies the object or objects to which the action applies. Statements must include either a Resource or a NotResource element. As a best practice, specify

a resource using its [Amazon Resource Name \(ARN\)](#). You can do this for actions that support a specific resource type, known as *resource-level permissions*.

For actions that don't support resource-level permissions, such as listing operations, use a wildcard (*) to indicate that the statement applies to all resources.

```
"Resource": "*"

```

To view examples of Amazon Bedrock identity-based policies, see [Identity-based policy examples for Amazon Bedrock \(p. 155\)](#).

Policy condition keys for Amazon Bedrock

Supports service-specific policy condition keys	Yes
---	-----

Administrators can use AWS JSON policies to specify who has access to what. That is, which **principal** can perform **actions** on what **resources**, and under what **conditions**.

The Condition element (or Condition *block*) lets you specify conditions in which a statement is in effect. The Condition element is optional. You can create conditional expressions that use [condition operators](#), such as equals or less than, to match the condition in the policy with values in the request.

If you specify multiple Condition elements in a statement, or multiple keys in a single Condition element, AWS evaluates them using a logical AND operation. If you specify multiple values for a single condition key, AWS evaluates the condition using a logical OR operation. All of the conditions must be met before the statement's permissions are granted.

You can also use placeholder variables when you specify conditions. For example, you can grant an IAM user permission to access a resource only if it is tagged with their IAM user name. For more information, see [IAM policy elements: variables and tags](#) in the *IAM User Guide*.

AWS supports global condition keys and service-specific condition keys. To see all AWS global condition keys, see [AWS global condition context keys](#) in the *IAM User Guide*.

All Amazon Bedrock actions support condition keys using Amazon Bedrock models as the resource.

To view examples of Amazon Bedrock identity-based policies, see [Identity-based policy examples for Amazon Bedrock \(p. 155\)](#).

ACLs in Amazon Bedrock

Supports ACLs	No
---------------	----

Access control lists (ACLs) control which principals (account members, users, or roles) have permissions to access a resource. ACLs are similar to resource-based policies, although they do not use the JSON policy document format.

ABAC with Amazon Bedrock

Supports ABAC (tags in policies)	Partial
----------------------------------	---------

Attribute-based access control (ABAC) is an authorization strategy that defines permissions based on attributes. In AWS, these attributes are called *tags*. You can attach tags to IAM entities (users or roles) and to many AWS resources. Tagging entities and resources is the first step of ABAC. Then you design ABAC policies to allow operations when the principal's tag matches the tag on the resource that they are trying to access.

ABAC is helpful in environments that are growing rapidly and helps with situations where policy management becomes cumbersome.

To control access based on tags, you provide tag information in the [condition element](#) of a policy using the `aws:ResourceTag/key-name`, `aws:RequestTag/key-name`, or `aws:TagKeys` condition keys.

If a service supports all three condition keys for every resource type, then the value is **Yes** for the service. If a service supports all three condition keys for only some resource types, then the value is **Partial**.

For more information about ABAC, see [What is ABAC?](#) in the *IAM User Guide*. To view a tutorial with steps for setting up ABAC, see [Use attribute-based access control \(ABAC\)](#) in the *IAM User Guide*.

Using temporary credentials with Amazon Bedrock

Supports temporary credentials	Yes
--------------------------------	-----

Some AWS services don't work when you sign in using temporary credentials. For additional information, including which AWS services work with temporary credentials, see [AWS services that work with IAM](#) in the *IAM User Guide*.

You are using temporary credentials if you sign in to the AWS Management Console using any method except a user name and password. For example, when you access AWS using your company's single sign-on (SSO) link, that process automatically creates temporary credentials. You also automatically create temporary credentials when you sign in to the console as a user and then switch roles. For more information about switching roles, see [Switching to a role \(console\)](#) in the *IAM User Guide*.

You can manually create temporary credentials using the AWS CLI or AWS API. You can then use those temporary credentials to access AWS. AWS recommends that you dynamically generate temporary credentials instead of using long-term access keys. For more information, see [Temporary security credentials in IAM](#).

Cross-service principal permissions for Amazon Bedrock

Supports principal permissions	Yes
--------------------------------	-----

When you use an IAM user or role to perform actions in AWS, you are considered a principal. Policies grant permissions to a principal. When you use some services, you might perform an action that then triggers another action in a different service. In this case, you must have permissions to perform both actions. To see whether an action requires additional dependent actions in a policy, see [Actions, Resources, and Condition Keys for Amazon Bedrock](#) in the *Service Authorization Reference*.

Service roles for Amazon Bedrock

Supports service roles	Yes
------------------------	-----

A service role is an [IAM role](#) that a service assumes to perform actions on your behalf. An IAM administrator can create, modify, and delete a service role from within IAM. For more information, see [Creating a role to delegate permissions to an AWS service](#) in the *IAM User Guide*.

Warning

Changing the permissions for a service role might break Amazon Bedrock functionality. Edit service roles only when Amazon Bedrock provides guidance to do so.

Service-linked roles for Amazon Bedrock

Supports service-linked roles	No
-------------------------------	----

A service-linked role is a type of service role that is linked to an AWS service. The service can assume the role to perform an action on your behalf. Service-linked roles appear in your AWS account and are owned by the service. An IAM administrator can view, but not edit the permissions for service-linked roles.

Identity-based policy examples for Amazon Bedrock

By default, users and roles don't have permission to create or modify Amazon Bedrock resources. They also can't perform tasks by using the AWS Management Console, AWS Command Line Interface (AWS CLI), or AWS API. To grant users permission to perform actions on the resources that they need, an IAM administrator can create IAM policies. The administrator can then add the IAM policies to roles, and users can assume the roles.

To learn how to create an IAM identity-based policy by using these example JSON policy documents, see [Creating IAM policies](#) in the *IAM User Guide*.

For details about actions and resource types defined by Amazon Bedrock, including the format of the ARNs for each of the resource types, see [Actions, Resources, and Condition Keys for Amazon Bedrock](#) in the *Service Authorization Reference*.

Note

The Amazon Bedrock service is available as a limited preview release, so its information is not included in the Service Authorization Reference.

Topics

- [Policy best practices \(p. 155\)](#)
- [Use the Amazon Bedrock console \(p. 156\)](#)
- [Allow users to view their own permissions \(p. 157\)](#)
- [Allow access to third-party model subscriptions \(p. 158\)](#)
- [Deny access for inference on specific models \(p. 158\)](#)
- [Grant custom jobs access to your training data \(p. 159\)](#)
- [Permissions for using KMS keys with model customization \(p. 159\)](#)

Policy best practices

Identity-based policies determine whether someone can create, access, or delete Amazon Bedrock resources in your account. These actions can incur costs for your AWS account. When you create or edit identity-based policies, follow these guidelines and recommendations:

- **Get started with AWS managed policies and move toward least-privilege permissions** – To get started granting permissions to your users and workloads, use the *AWS managed policies* that grant

permissions for many common use cases. They are available in your AWS account. We recommend that you reduce permissions further by defining AWS customer managed policies that are specific to your use cases. For more information, see [AWS managed policies](#) or [AWS managed policies for job functions](#) in the *IAM User Guide*.

- **Apply least-privilege permissions** – When you set permissions with IAM policies, grant only the permissions required to perform a task. You do this by defining the actions that can be taken on specific resources under specific conditions, also known as *least-privilege permissions*. For more information about using IAM to apply permissions, see [Policies and permissions in IAM](#) in the *IAM User Guide*.
- **Use conditions in IAM policies to further restrict access** – You can add a condition to your policies to limit access to actions and resources. For example, you can write a policy condition to specify that all requests must be sent using SSL. You can also use conditions to grant access to service actions if they are used through a specific AWS service, such as AWS CloudFormation. For more information, see [IAM JSON policy elements: Condition](#) in the *IAM User Guide*.
- **Use IAM Access Analyzer to validate your IAM policies to ensure secure and functional permissions** – IAM Access Analyzer validates new and existing policies so that the policies adhere to the IAM policy language (JSON) and IAM best practices. IAM Access Analyzer provides more than 100 policy checks and actionable recommendations to help you author secure and functional policies. For more information, see [IAM Access Analyzer policy validation](#) in the *IAM User Guide*.
- **Require multi-factor authentication (MFA)** – If you have a scenario that requires IAM users or a root user in your AWS account, turn on MFA for additional security. To require MFA when API operations are called, add MFA conditions to your policies. For more information, see [Configuring MFA-protected API access](#) in the *IAM User Guide*.

For more information about best practices in IAM, see [Security best practices in IAM](#) in the *IAM User Guide*.

Use the Amazon Bedrock console

To access the Amazon Bedrock console, you must have a minimum set of permissions. These permissions must allow you to list and view details about the Amazon Bedrock resources in your AWS account. If you create an identity-based policy that is more restrictive than the minimum required permissions, the console won't function as intended for entities (users or roles) with that policy.

You don't need to allow minimum console permissions for users that are making calls only to the AWS CLI or the AWS API. Instead, allow access to only the actions that match the API operation that they're trying to perform.

To provide access to the Amazon Bedrock console, attach the following policy to the roles or entities that need access. A few of the actions in the policy are for use by the Amazon Bedrock console only, so IAM might raise Invalid Action errors. You can ignore these errors and choose **Next** to complete the task.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "BedrockConsole",
      "Effect": "Allow",
      "Action": [
        "bedrock:ListFoundationModels",
        "bedrock:GetFoundationModel",
        "bedrock:InvokeModel",
        "bedrock:InvokeModelWithResponseStream",
        "bedrock:CreateModelCustomizationJob",
        "bedrock:GetModelCustomizationJob",
        "bedrock:GetFoundationModelAvailability",
        "bedrock:ListModelCustomizationJobs",

```

```
        "bedrock:StopModelCustomizationJob",
        "bedrock:GetCustomModel",
        "bedrock:ListCustomModels",
        "bedrock>DeleteCustomModel",
        "bedrock>CreateProvisionedModelThroughput",
        "bedrock:UpdateProvisionedModelThroughput",
        "bedrock:GetProvisionedModelThroughput",
        "bedrock>DeleteProvisionedModelThroughput",
        "bedrock:ListProvisionedModelThroughputs",
        "bedrock:ListTagsForResource",
        "bedrock:UntagResource",
        "bedrock:TagResource",
        "bedrock:PutFoundationModelEntitlement",
        "bedrock:GetModelInvocationLoggingConfiguration",
        "bedrock:PutModelInvocationLoggingConfiguration",
        "bedrock>CreateFoundationModelAgreement",
        "bedrock>DeleteFoundationModelAgreement",
        "bedrock:ListFoundationModelAgreementOffers",
        "bedrock:GetUseCaseForModelAccess",
        "bedrock:PutUseCaseForModelAccess"
    ],
    "Resource": "*"
}
]
```

For more information, see [Adding permissions to a user](#) in the *IAM User Guide*.

Allow users to view their own permissions

This example shows how you might create a policy that allows IAM users to view the inline and managed policies that are attached to their user identity. This policy includes permissions to complete this action on the console or programmatically using the AWS CLI or AWS API.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "ViewOwnUserInfo",
      "Effect": "Allow",
      "Action": [
        "iam:GetUserPolicy",
        "iam:ListGroupsForUser",
        "iam:ListAttachedUserPolicies",
        "iam:ListUserPolicies",
        "iam:GetUser"
      ],
      "Resource": ["arn:aws:iam::*:user/${aws:username}"]
    },
    {
      "Sid": "NavigateInConsole",
      "Effect": "Allow",
      "Action": [
        "iam:GetGroupPolicy",
        "iam:GetPolicyVersion",
        "iam:GetPolicy",
        "iam:ListAttachedGroupPolicies",
        "iam:ListGroupPolicies",
        "iam:ListPolicyVersions",
        "iam:ListPolicies",
        "iam:ListUsers"
      ],
      "Resource": "*"
    }
  ]
}
```

```
}  
  ]  
}
```

Allow access to third-party model subscriptions

To access the Amazon Bedrock models for the first time, you use the Amazon Bedrock console to subscribe to third-party models. Your IAM user or role that the console user assumes requires permission to access the subscription API operations.

The following example shows an identity-based policy to allow access to the subscription API operations. The example includes a condition key that limits the scope of the policy to the Amazon Bedrock products in the Marketplace.

```
{  
  "Version": "2012-10-17",  
  "Statement": [  
    {  
      "Effect": "Allow",  
      "Action": [  
        "aws-marketplace:Subscribe"  
      ],  
      "Resource": "*",  
      "Condition": {  
        "ForAnyValue:StringEquals": {  
          "aws-marketplace:ProductId": [  
            "c468b48a-84df-43a4-8c46-8870630108a7",  
            "99d90be8-b43e-49b7-91e4-752f3866c8c7",  
            "b0eb9475-3a2c-43d1-94d3-56756fd43737",  
            "1d288c71-65f9-489a-a3e2-9c7f4f6e6a85",  
            "cc0bdd50-279a-40d8-829c-4009b77a1fcc",  
            "d0123e8d-50d6-4dba-8a26-3fed4899f388",  
            "a61c46fe-1747-41aa-9af0-2e0ae8a9ce05"  
          ]  
        }  
      }  
    },  
    {  
      "Effect": "Allow",  
      "Action": [  
        "aws-marketplace:Unsubscribe",  
        "aws-marketplace:ViewSubscriptions"  
      ],  
      "Resource": "*"   
    }  
  ]  
}
```

Deny access for inference on specific models

The following example shows an identity-based policy that denies access to running inference on a specific model.

```
{  
  "Version": "2012-10-17",  
  "Statement": {  
    "Sid": "DenyInference",  
    "Effect": "Deny",  
    "Action": [  
      "bedrock:InvokeModel",  
      "bedrock:InvokeModelWithResponseStream"  
    ]  
  }  
}
```

```
    ]
    "Resource": "arn:aws:bedrock:*::foundation-model/model-id-of-model-to-deny"
  }
}
```

Grant custom jobs access to your training data

The following example grants access to the Amazon S3 locations of the input and output data.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:PutObject",
        "s3:ListBucket",
        "s3:ListObjects"
      ],
      "Resource": [
        "arn:aws:s3::my_training_data_bucket/myfolder",
        "arn:aws:s3::my_validation_bucket/myfolder",
        "arn:aws:s3::my_output_bucket/myfolder"
      ]
    }
  ]
}
```

The following example creates a trust policy to allow model-customization jobs to assume a role, if the job was created in the same account.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "bedrock.amazonaws.com"
      },
      "Action": "sts:AssumeRole",
      "Condition": {
        "StringEquals": {
          "aws:SourceAccount": "111122223333"
        },
        "ArnEquals": {
          "aws:SourceArn": "arn:aws:bedrock:us-east-1:111122223333:model-
customization-job/*"
        }
      }
    }
  ]
}
```

Permissions for using KMS keys with model customization

Provide the appropriate AWS KMS permissions for users and roles that create or use custom models. For details, see [Use a customer managed key to run customization jobs \(p. 140\)](#).

Service role

You can set up a service role that lets Amazon Bedrock run on your behalf.

1. Create an AmazonBedrock-ExecutionRole (RolePolicy) with the permissions needed for training (S3, KMS) and a Trusted Entity policy which allows Amazon Bedrock to assume this role (Reference role). Note the Amazon Resource Name (ARN) of the role policy.
2. In the role policy, to set the context keys for `aws:SourceArn` and `aws:SourceAccount` to the Amazon Bedrock job Amazon Resource Name (ARN) and your AWS Account ID when Amazon Bedrock tries to assume these credentials.
3. Create an IAM entity with permissions to access Amazon Bedrock and pass role permissions to the RoleArn you noted in step 1.

You can then run a fine tuning job with the console, or with the `CreateFineTuningJob` operation, passing the `>RoleArn` and the Amazon S3 URI of your training data.

Amazon Bedrock performs `PassRole` checks in the `CreateFineTuningJob` operation to prevent confused deputy scenarios. These `PassRole` checks validate that you are passing in a role that has permissions to pass that role into an AWS service to act on their behalf. Amazon Bedrock then uses the provided Role ARN to generate credentials that it uses to download the training and validation data stored in your account, as well as upload training metrics to your account.

Troubleshooting Amazon Bedrock identity and access

Use the following information to help you diagnose and fix common issues that you might encounter when working with Amazon Bedrock and IAM.

Topics

- [I am not authorized to perform an action in Amazon Bedrock \(p. 160\)](#)
- [I am not authorized to perform iam:PassRole \(p. 161\)](#)
- [I want to allow people outside of my AWS account to access my Amazon Bedrock resources \(p. 161\)](#)

I am not authorized to perform an action in Amazon Bedrock

If you receive an error that you're not authorized to perform an action, your policies must be updated to allow you to perform the action.

The following example error occurs when the `mateojackson` IAM user tries to use the console to view details about a fictional `my-example-widget` resource but doesn't have the fictional `bedrock:GetWidget` permissions.

```
User: arn:aws:iam::123456789012:user/mateojackson is not authorized to perform:
bedrock:GetWidget on resource: my-example-widget
```

In this case, the policy for the `mateojackson` user must be updated to allow access to the `my-example-widget` resource by using the `bedrock:GetWidget` action.

If you need help, contact your AWS administrator. Your administrator is the person who provided you with your sign-in credentials.

I am not authorized to perform iam:PassRole

If you receive an error that you're not authorized to perform the `iam:PassRole` action, your policies must be updated to allow you to pass a role to Amazon Bedrock.

Some AWS services allow you to pass an existing role to that service instead of creating a new service role or service-linked role. To do this, you must have permissions to pass the role to the service.

The following example error occurs when an IAM user named `marymajor` tries to use the console to perform an action in Amazon Bedrock. However, the action requires the service to have permissions that are granted by a service role. Mary does not have permissions to pass the role to the service.

```
User: arn:aws:iam::123456789012:user/marymajor is not authorized to perform: iam:PassRole
```

In this case, Mary's policies must be updated to allow her to perform the `iam:PassRole` action.

If you need help, contact your AWS administrator. Your administrator is the person who provided you with your sign-in credentials.

I want to allow people outside of my AWS account to access my Amazon Bedrock resources

You can create a role that users in other accounts or people outside of your organization can use to access your resources. You can specify who is trusted to assume the role. For services that support resource-based policies or access control lists (ACLs), you can use those policies to grant people access to your resources.

To learn more, consult the following:

- To learn whether Amazon Bedrock supports these features, see [How Amazon Bedrock works with IAM \(p. 150\)](#).
- To learn how to provide access to your resources across AWS accounts that you own, see [Providing access to an IAM user in another AWS account that you own](#) in the *IAM User Guide*.
- To learn how to provide access to your resources to third-party AWS accounts, see [Providing access to AWS accounts owned by third parties](#) in the *IAM User Guide*.
- To learn how to provide access through identity federation, see [Providing access to externally authenticated users \(identity federation\)](#) in the *IAM User Guide*.
- To learn the difference between using roles and resource-based policies for cross-account access, see [How IAM roles differ from resource-based policies](#) in the *IAM User Guide*.

Compliance validation for Amazon Bedrock

To learn whether an AWS service is within the scope of specific compliance programs, see [AWS services in Scope by Compliance Program](#) and choose the compliance program that you are interested in. For general information, see [AWS Compliance Programs](#).

You can download third-party audit reports using AWS Artifact. For more information, see [Downloading Reports in AWS Artifact](#).

Your compliance responsibility when using AWS services is determined by the sensitivity of your data, your company's compliance objectives, and applicable laws and regulations. AWS provides the following resources to help with compliance:

- [Security and Compliance Quick Start Guides](#) – These deployment guides discuss architectural considerations and provide steps for deploying baseline environments on AWS that are security and compliance focused.
- [Architecting for HIPAA Security and Compliance on Amazon Web Services](#) – This whitepaper describes how companies can use AWS to create HIPAA-eligible applications.

Note

Not all AWS services are HIPAA eligible. For more information, see the [HIPAA Eligible Services Reference](#).

- [AWS Compliance Resources](#) – This collection of workbooks and guides might apply to your industry and location.
- [Evaluating Resources with Rules](#) in the *AWS Config Developer Guide* – The AWS Config service assesses how well your resource configurations comply with internal practices, industry guidelines, and regulations.
- [AWS Security Hub](#) – This AWS service provides a comprehensive view of your security state within AWS. Security Hub uses security controls to evaluate your AWS resources and to check your compliance against security industry standards and best practices. For a list of supported services and controls, see [Security Hub controls reference](#).
- [AWS Audit Manager](#) – This AWS service helps you continuously audit your AWS usage to simplify how you manage risk and compliance with regulations and industry standards.

Incident response in Amazon Bedrock

Incident response for Amazon Bedrock is an AWS responsibility. AWS has a formal, documented policy and program that governs incident response.

AWS operational issues with broad impact are posted on the [AWS Service Health Dashboard](#).

Operational issues are also posted to individual accounts via the AWS Health Dashboard. For information on how to use the AWS Health Dashboard, see the [AWS Health User Guide](#).

Resilience in Amazon Bedrock

The AWS global infrastructure is built around AWS Regions and Availability Zones. AWS Regions provide multiple physically separated and isolated Availability Zones, which are connected with low-latency, high-throughput, and highly redundant networking. With Availability Zones, you can design and operate applications and databases that automatically fail over between zones without interruption. Availability Zones are more highly available, fault tolerant, and scalable than traditional single or multiple data center infrastructures.

For more information about AWS Regions and Availability Zones, see [AWS Global Infrastructure](#).

Infrastructure security in Amazon Bedrock

As a managed service, Amazon Bedrock is protected by the AWS global network security. For information about AWS security services and how AWS protects infrastructure, see [AWS Cloud Security](#). To design your AWS environment using the best practices for infrastructure security, see [Infrastructure Protection](#) in *Security Pillar AWS Well-Architected Framework*.

You use AWS published API calls to access Amazon Bedrock through the network. Clients must support the following:

- Transport Layer Security (TLS). We require TLS 1.2 and recommend TLS 1.3.
- Cipher suites with perfect forward secrecy (PFS) such as DHE (Ephemeral Diffie-Hellman) or ECDHE (Elliptic Curve Ephemeral Diffie-Hellman). Most modern systems such as Java 7 and later support these modes.

Additionally, requests must be signed by using an access key ID and a secret access key that is associated with an IAM principal. Or you can use the [AWS Security Token Service](#) (AWS STS) to generate temporary security credentials to sign requests.

Cross-service confused deputy prevention

The confused deputy problem is a security issue where an entity that doesn't have permission to perform an action can coerce a more-privileged entity to perform the action. In AWS, cross-service impersonation can result in the confused deputy problem. Cross-service impersonation can occur when one service (the *calling service*) calls another service (the *called service*). The calling service can be manipulated to use its permissions to act on another customer's resources in a way it should not otherwise have permission to access. To prevent this, AWS provides tools that help you protect your data for all services with service principals that have been given access to resources in your account.

We recommend using the [aws:SourceArn](#) and [aws:SourceAccount](#) global condition context keys in resource policies to limit the permissions that Amazon Bedrock gives another service to the resource. Use `aws:SourceArn` if you want only one resource to be associated with the cross-service access. Use `aws:SourceAccount` if you want to allow any resource in that account to be associated with the cross-service use.

The most effective way to protect against the confused deputy problem is to use the `aws:SourceArn` global condition context key with the full ARN of the resource. If you don't know the full ARN of the resource or if you are specifying multiple resources, use the `aws:SourceArn` global context condition key with wildcard characters (*) for the unknown portions of the ARN. For example, `arn:aws:bedrock:*:123456789012:*`.

If the `aws:SourceArn` value does not contain the account ID, such as an Amazon S3 bucket ARN, you must use both global condition context keys to limit permissions.

The value of `aws:SourceArn` must be `ResourceDescription`.

The following example shows how you can use the `aws:SourceArn` and `aws:SourceAccount` global condition context keys in Bedrock to prevent the confused deputy problem.

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "bedrock.amazonaws.com"
      },
      "Action": "sts:AssumeRole",
      "Condition": {
        "StringEquals": {
          "aws:SourceAccount": "111122223333"
        },
        "ArnEquals": {
          "aws:SourceArn": "arn:aws:bedrock:us-east-1:111122223333:model-
customization-job/*"
        }
      }
    }
  ]
}
```

```
}  
  ]  
}
```

Configuration and vulnerability analysis in Amazon Bedrock

Configuration and IT controls are a shared responsibility between AWS and you, our customer. For more information, see the AWS [shared responsibility model](#).

Monitor

You can monitor Amazon Bedrock with Amazon CloudWatch and with Amazon EventBridge.

Topics

- [Monitor Amazon Bedrock with Amazon CloudWatch \(p. 165\)](#)
- [Monitor Amazon Bedrock events in Amazon EventBridge \(p. 167\)](#)
- [Log Amazon Bedrock API calls using AWS CloudTrail \(p. 170\)](#)

Monitor Amazon Bedrock with Amazon CloudWatch

You can monitor Amazon Bedrock using Amazon CloudWatch, which collects raw data and processes it into readable, near real-time metrics. You can graph the metrics using the CloudWatch console. You can also set alarms that watch for certain thresholds, and send notifications or take actions when values exceed those thresholds.

For more information, see [What is Amazon CloudWatch](#) in the *Amazon CloudWatch User Guide*.

Topics

- [Runtime metrics \(p. 165\)](#)
- [Logging CloudWatch metrics \(p. 166\)](#)
- [Use CloudWatch metrics for Amazon Bedrock \(p. 166\)](#)
- [View Amazon Bedrock metrics \(p. 166\)](#)

Runtime metrics

The following table describes runtime metrics provided by Amazon Bedrock.

Metric name	Unit	Description
Invocations	SampleCount	Number of requests to the InvokeModel or InvokeModelWithResponseStream API operations.
InvocationLatency	Milliseconds	Latency of the invocations.
InvocationClientErrors	SampleCount	Number of invocations that result in client-side errors.
InvocationServerErrors	SampleCount	Number of invocations that result in AWS server-side errors.

Metric name	Unit	Description
InvocationThrottles	SampleCount	Number of invocations that the system throttled.
InputTokenCount	SampleCount	Number of tokens of text input.
OutputTokenCount	SampleCount	Number of tokens of text output.
ContentFilteredCount	SampleCount	Number of times the text output content was filtered.
OutputImageCount	SampleCount	Number of output images.

Logging CloudWatch metrics

For each delivery success or failure attempt, the following Amazon CloudWatch metrics are emitted under the namespace AWS/Bedrock, and Across all model IDs dimension:

- ModelInvocationLogsCloudWatchDeliverySuccess
- ModelInvocationLogsCloudWatchDeliveryFailure
- ModelInvocationLogsS3DeliverySuccess
- ModelInvocationLogsS3DeliveryFailure
- ModelInvocationLargeDataS3DeliverySuccess
- ModelInvocationLargeDataS3DeliveryFailure

If logs fail to deliver due to permission misconfiguration or transient failures, the delivery is retried periodically for up to 24 hours.

Use CloudWatch metrics for Amazon Bedrock

To retrieve metrics for your Amazon Bedrock operations, you specify the following information:

- The metric dimension. A *dimension* is a set of name-value pairs that you use to identify a metric. Amazon Bedrock supports the following dimensions:
 - ModelId – all metrics
 - ModelId + ImageSize + BucketedStepSize – OutputImageCount
- The metric name, such as InvocationClientErrors.

You can get metrics for Amazon Bedrock with the AWS Management Console, the AWS CLI, or the CloudWatch API. You can use the CloudWatch API through one of the AWS Software Development Kits (SDKs) or the CloudWatch API tools.

You must have the appropriate CloudWatch permissions to monitor Amazon Bedrock with CloudWatch. For more information, see [Authentication and Access Control for Amazon CloudWatch](#) in the *Amazon CloudWatch User Guide*.

View Amazon Bedrock metrics

View Amazon Bedrock metrics in the CloudWatch console.

To view metrics (CloudWatch console)

1. Sign in to the AWS Management Console and open the CloudWatch console at <https://console.aws.amazon.com/cloudwatch/>.
2. Choose **Metrics**, choose **All Metrics**, and then search for **ModelId**.

Monitor Amazon Bedrock events in Amazon EventBridge

You can use Amazon EventBridge to monitor status change events in Amazon Bedrock. With Amazon EventBridge, you can configure Amazon SageMaker to respond automatically to a model customization job status change in Amazon Bedrock. Events from Amazon Bedrock are delivered to Amazon EventBridge in near real time. You can write simple rules to automate actions when an event matches a rule. If you use Amazon EventBridge with Amazon Bedrock, you can:

- Publish notifications whenever there is a state change event in the model customization you have triggered, whether you add new asynchronous workflows in the future. The event published should give you enough information to react to events in downstream workflows.
- Deliver job status updates without invoking the `GetModelCustomizationJob` API, which can mean handling API rate limit issues, API updates, and reduction in additional compute resources.

There is no cost to receive AWS events from Amazon EventBridge. For more information about, Amazon EventBridge, see [Amazon EventBridge](#)

Note

- Amazon Bedrock emits events on a best-effort basis. Events are delivered to Amazon EventBridge in near real time. With Amazon EventBridge, you can create rules that trigger programmatic actions in response to an event. For example, you can configure a rule that invokes an SNS topic to send an email notification or invokes a function to take some action. For more information, see the *Amazon EventBridge User Guide*.
- AWS Bedrock creates a new event every time there is a state change in a model customization job that you trigger and make best-effort delivery of such event.

Topics

- [How it works \(p. 167\)](#)
- [EventBridge schema \(p. 168\)](#)
- [Rules and targets \(p. 169\)](#)
- [Create a rule to handle AWS Bedrock events \(p. 169\)](#)

How it works

To receive events from Amazon Bedrock, you need to create rules and targets to match, receive, and handle state change data through Amazon EventBridge. Amazon EventBridge is a serverless event bus that ingests change state events from AWS services, SaaS partners, and customer applications. It processes events based on rules or patterns that you create, and routes these events to one or more “targets” that you choose, such as AWS Lambda, Amazon Simple Queue Service, and Amazon Simple Notification Service.

Amazon Bedrock publishes your events via Amazon EventBridge whenever there is a change in the state of a model customization job. In each case, a new event is created and sent to Amazon EventBridge,

which then sends the event to your default event-bus. The event shows which customization job's state has changed, and the current state of the job. When Amazon EventBridge receives an event that matches a rule that you created, Amazon EventBridge routes it to the target that you specified. When you create a rule, you can configure these targets as well as downstream workflows based on the contents of the event.

EventBridge schema

The following event fields in the EventBridge event schema are specific to Amazon Bedrock.

- `jobArn` — The ARN of the model customization job.
- `outputModelArn` — The ARN of the output model. Published when the training job has completed.
- `jobStatus` — The current status of the job.
- `FailureMessage` — A failure message. Published when the training job has failed.

Event example

The following is example event JSON for a failed model customization job.

```
{
  "version": "0",
  "id": "UUID",
  "detail-type": "Model Customization Job State Change",
  "source": "aws.bedrock",
  "account": "123412341234",
  "time": "2023-08-11T12:34:56Z",
  "region": "us-east-1",
  "resources": [ "arn:aws:bedrock:us-east-1:123412341234:model-customization-job/
abcdefghwxyz" ],
  "detail": {
    "version": "0.0",
    "jobName": "abcd-wxyz",
    "jobArn": "arn:aws:bedrock:us-east-1:123412341234:model-customization-job/
abcdefghwxyz",
    "outputModelName": "dummy-output-model-name",
    "outputModelArn": "arn:aws:bedrock:us-east-1:123412341234:dummy-output-model-name",
    "roleArn": "arn:aws:iam::123412341234:role/JobExecutionRole",
    "jobStatus": "Failed",
    "failureMessage": "Failure Message here.",
    "creationTime": "2023-08-11T10:11:12Z",
    "lastModifiedTime": "2023-08-11T12:34:56Z",
    "endTime": "2023-08-11T12:34:56Z",
    "baseModelArn": "arn:aws:bedrock:us-east-1:123412341234:base-model-name",
    "hyperParameters": {
      "batchSize": "batchSizeNumberUsed",
      "epochCount": "epochCountNumberUsed",
      "learningRate": "learningRateUsed",
      "learningRateWarmupSteps": "learningRateWarmupStepsUsed"
    },
    "trainingDataConfig": {
      "s3Uri": "s3://bucket/key",
    },
    "validationDataConfig": {
      "s3Uri": "s3://bucket/key",
    },
    "outputDataConfig": {
      "s3Uri": "s3://bucket/key",
    }
  }
}
```

```
}
```

Rules and targets

When an incoming event matches a rule that you created, the event is routed to the target that you specified for that rule, and the target processes these events. Targets support JSON format and can include AWS services such as Amazon EC2 instances, Lambda functions, Kinesis streams, Amazon ECS tasks, Step Functions, Amazon SNS topics, and Amazon SQS. To receive and process events correctly, you need to create rules and targets for matching, receiving, and correctly handling event data. You can create these rules and targets either through the Amazon EventBridge console, or through the AWS CLI.

Example rule

This rule matches an event pattern emitted by: `source ["aws.bedrock"]`. The rule captures all events sent by Amazon EventBridge that have source "aws.bedrock" to your default event bus.

```
{  
  "source": ["aws.bedrock"]  
}
```

Target

When creating a rule in Amazon EventBridge, you need to specify a target where EventBridge sends the event that matches your rule pattern. These targets can be a SageMaker pipeline, a Lambda function, an SNS topic, an SQS queue or any of the other targets that EventBridge currently supports. You can refer to the *Amazon EventBridge* documentation to learn how to set targets for events. For a procedure that shows how to use Amazon Simple Notification Service as a target, see [Create a rule to handle AWS Bedrock events \(p. 169\)](#).

Create a rule to handle AWS Bedrock events

Complete the following procedures in order to receive email notifications about your AWS Bedrock events.

Create an Amazon Simple Notification Service topic

1. Open the Amazon SNS console at <https://console.aws.amazon.com/sns/v3/home>.
2. In the navigation pane, choose **Topics**.
3. Choose **Create topic**.
4. For **Type**, choose **Standard**.
5. For **Name**, enter a name for your topic.
6. Choose **Create topic**.
7. Choose **Create subscription**.
8. For **Protocol**, choose **Email**.
9. For **Endpoint**, enter the email address that receives the notifications.
10. Choose **Create subscription**.
11. You'll receive an email message with the following subject line: **AWS Notification - Subscription Confirmation**. Follow the directions to confirm your subscription.

Use the following procedure to create a rule to handle your AWS Bedrock events.

To create a rule to handle AWS Bedrock events

1. Open the Amazon EventBridge console at <https://console.aws.amazon.com/events/>.
2. Choose **Create rule**.
3. For **Name**, enter a name for your rule.
4. For **Rule type**, choose **Rule with an event pattern**.
5. Choose **Next**.
6. For Event pattern, do the following:
 - a. For **Event source**, choose **AWS services**.
 - b. For **AWS service**, choose **Amazon Bedrock**.
 - c. For **Event type**, choose **Model Customization Job State Change**.
 - d. By default, we send notifications for every event. If you prefer, you can create an event pattern that filters events for a specific job state.
 - e. Choose **Next**.
7. Specify a target as follows:
 - a. For **Target types**, choose **AWS service**.
 - b. For **Select a target**, choose **SNS topic**.
 - c. For **Topic**, choose the SNS topic that you created for notifications.
 - d. Choose **Next**.
8. (Optional) Add tags to your rule.
9. Choose **Next**.
10. Choose **Create rule**.

Log Amazon Bedrock API calls using AWS CloudTrail

Amazon Bedrock is integrated with AWS CloudTrail, a service that provides a record of actions taken by a user, role, or an AWS service in Amazon Bedrock. CloudTrail captures all API calls for Amazon Bedrock as events. The calls captured include calls from the Amazon Bedrock console and code calls to the Amazon Bedrock API operations. If you create a trail, you can enable continuous delivery of CloudTrail events to an Amazon S3 bucket, including events for Amazon Bedrock. If you don't configure a trail, you can still view the most recent events in the CloudTrail console in **Event history**. Using the information collected by CloudTrail, you can determine the request that was made to Amazon Bedrock, the IP address from which the request was made, who made the request, when it was made, and additional details.

To learn more about CloudTrail, see the [AWS CloudTrail User Guide](#).

Amazon Bedrock information in CloudTrail

CloudTrail is enabled on your AWS account when you create the account. When activity occurs in Amazon Bedrock, that activity is recorded in a CloudTrail event along with other AWS service events in **Event history**. You can view, search, and download recent events in your AWS account. For more information, see [Viewing events with CloudTrail Event history](#).

For an ongoing record of events in your AWS account, including events for Amazon Bedrock, create a trail. A *trail* enables CloudTrail to deliver log files to an Amazon S3 bucket. By default, when you create a trail in the console, the trail applies to all AWS Regions. The trail logs events from all Regions in the AWS partition and delivers the log files to the Amazon S3 bucket that you specify. Additionally, you can

configure other AWS services to further analyze and act upon the event data collected in CloudTrail logs. For more information, see the following:

- [Overview for creating a trail](#)
- [CloudTrail supported services and integrations](#)
- [Configuring Amazon SNS notifications for CloudTrail](#)
- [Receiving CloudTrail log files from multiple regions](#) and [Receiving CloudTrail log files from multiple accounts](#)

All Amazon Bedrock actions are logged by CloudTrail and are documented in the [Amazon Bedrock API Reference](#). For example, calls to the `InvokeModel`, `ListFoundationModels` and `StopModelCustomizationJob` actions generate entries in the CloudTrail log files.

Every event or log entry contains information about who generated the request. The identity information helps you determine the following:

- Whether the request was made with root or AWS Identity and Access Management (IAM) user credentials.
- Whether the request was made with temporary security credentials for a role or federated user.
- Whether the request was made by another AWS service.

For more information, see the [CloudTrail userIdentity element](#).

Understanding Amazon Bedrock log file entries

A trail is a configuration that enables delivery of events as log files to an Amazon S3 bucket that you specify. CloudTrail log files contain one or more log entries. An event represents a single request from any source and includes information about the requested action, the date and time of the action, request parameters, and so on. CloudTrail log files aren't an ordered stack trace of the public API calls, so they don't appear in any specific order.

The following example shows a CloudTrail log entry that demonstrates the `InvokeModel` action.

```
{
  "eventVersion": "1.08",
  "userIdentity": {
    "type": "IAMUser",
    "principalId": "AROAI CFHPEXAMPLE",
    "arn": "arn:aws:iam::111122223333:user/userxyz",
    "accountId": "111122223333",
    "accessKeyId": "AKIAIOSFODNN7EXAMPLE",
    "userName": "userxyz"
  },
  "eventTime": "2023-10-11T21:58:59Z",
  "eventSource": "bedrock.amazonaws.com",
  "eventName": "InvokeModel",
  "awsRegion": "us-west-2",
  "sourceIPAddress": "192.0.2.0",
  "userAgent": "Boto3/1.28.62 md/Botocore#1.31.62 ua/2.0 os/macos#22.6.0 md/arch#arm64
lang/python#3.9.6 md/pyimpl#CPython cfg/retry-mode#legacy Botocore/1.31.62",
  "requestParameters": {
    "modelId": "stability.stable-diffusion-xl-v0"
  },
  "responseElements": null,
  "requestID": "a1b2c3d4-5678-90ab-cdef-EXAMPLE22222",
  "eventID": "a1b2c3d4-5678-90ab-cdef-EXAMPLE11111 ",
  "readOnly": false,
  "eventType": "AwsApiCall",
}
```

```
"managementEvent": true,  
"recipientAccountId": "111122223333",  
"eventCategory": "Management",  
"tlsDetails": {  
  "tlsVersion": "TLSv1.2",  
  "cipherSuite": "cipher suite",  
  "clientProvidedHostHeader": "bedrock-runtime.us-west-2.amazonaws.com"  
}  
}
```

Amazon Bedrock abuse detection

AWS is committed to the responsible use of AI. To help prevent potential misuse, Amazon Bedrock implements automated abuse detection mechanisms to identify and mitigate potential violations of AWS's [Acceptable Use Policy](#) (AUP) and [Responsible AI Policy](#) or a third-party model provider's AUP.

Our abuse detection mechanisms are fully automated, so there is no human review of, or access to, user inputs or model outputs.

Automated abuse detection includes several components:

- **Categorize content** — We use classifiers to detect harmful content (such as content that incites violence) in user inputs and model outputs. A classifier is an algorithm that processes model inputs and outputs, and assigns type of harm and level of confidence. We may run these classifiers on both Amazon Titan and third-party model usage. The classification process is automated and does not involve human review of user inputs or model outputs.
- **Identify patterns** — We use classifier metrics to identify potential violations and recurring behavior. We may compile and share aggregated and anonymized classifier metrics with third-party model providers. Amazon Bedrock does not store user input or model output and does not share these with third-party model providers.
- **Contact customers** — We may request information about customers' use of Bedrock and compliance with our policies. In the event that a customer continues to use the service in a manner that may violate AWS's policies or a third-party model provider's AUP, AWS may suspend access to Bedrock, taking into account severity, recurrence of the activity, customer cooperation, and whether customers have mechanisms in place to prevent misuse of the service.

Contact AWS Support if you have additional questions. For more information, see the [Amazon Bedrock FAQs](#).

Use interface VPC endpoints (AWS PrivateLink)

You can use AWS PrivateLink to create a private connection between your VPC and Amazon Bedrock. You can access Amazon Bedrock as if it were in your VPC, without the use of an internet gateway, NAT device, VPN connection, or AWS Direct Connect connection. Instances in your VPC don't need public IP addresses to access Amazon Bedrock.

You establish this private connection by creating an *interface endpoint*, powered by AWS PrivateLink. We create an endpoint network interface in each subnet that you enable for the interface endpoint. These are requester-managed network interfaces that serve as the entry point for traffic destined for Amazon Bedrock.

For more information, see [Access AWS services through AWS PrivateLink](#) in the *AWS PrivateLink Guide*.

Considerations for Amazon Bedrock VPC endpoints

Before you set up an interface endpoint for Amazon Bedrock, review [Considerations](#) in the *AWS PrivateLink Guide*.

Amazon Bedrock supports making calls to [Runtime API actions](#) through the VPC endpoint.

Amazon Bedrock endpoints are not available in all Availability Zones in a Region. When you create the endpoint, use the following command to list the Availability Zones.

```
aws ec2 describe-vpc-endpoint-services \
  --service-names com.amazonaws.region.bedrock-runtime --region region
```

Important

The service name `com.amazonaws.region.bedrock` is deprecated.

Create an interface endpoint for Amazon Bedrock

You can create an interface endpoint for Amazon Bedrock using either the Amazon VPC console or the AWS Command Line Interface (AWS CLI). For more information, see [Create an interface endpoint](#) in the *AWS PrivateLink Guide*.

Create an interface endpoint for Amazon Bedrock using the following service name:

```
com.amazonaws.region.bedrock-runtime
```

If you enable private DNS for the interface endpoint, you can make API requests to Amazon Bedrock using its default Regional DNS name. For example, `bedrock-runtime.us-east-1.amazonaws.com`.

Important

The private DNS `bedrock.region.amazonaws.com` is deprecated.

Create an endpoint policy for your interface endpoint

An endpoint policy is an IAM resource that you can attach to an interface endpoint. The default endpoint policy allows full access to Amazon Bedrock through the interface endpoint. To control the access allowed to Amazon Bedrock from your VPC, attach a custom endpoint policy to the interface endpoint.

An endpoint policy specifies the following information:

- The principals that can perform actions (AWS accounts, IAM users, and IAM roles).
- The actions that can be performed.
- The resources on which the actions can be performed.

For more information, see [Control access to services using endpoint policies](#) in the *AWS PrivateLink Guide*.

Example: VPC endpoint policy for Amazon Bedrock actions

The following is an example of a custom endpoint policy. When you attach this policy to your interface endpoint, it grants access to the listed Amazon Bedrock actions for all principals on all resources.

```
{
  "Statement": [
    {
      "Principal": "*",
      "Effect": "Allow",
      "Action": [
        "bedrock:InvokeModel",
        "bedrock:InvokeModelWithResponseStream"
      ],
      "Resource": "*"
    }
  ]
}
```

Quotas for Amazon Bedrock

Your AWS account has default quotas, formerly referred to as limits, for each AWS service. Unless otherwise noted, each quota is Region-specific. You can request increases for some quotas, and other quotas cannot be increased.

Some service quotas can be adjusted or increased. Refer to the **Adjustable** column in the following tables to see whether a quota can be adjusted. To request a quota increase, use the [limit increase form](#).

Your AWS account has the following quotas related to Amazon Bedrock.

Runtime quotas

The following quotas are enforced when you use a model for inference.

Model	Requests processed per minute	Tokens processed per minute	Adjustable
Amazon Titan Express	400	300,000	No
Amazon Titan Text Embeddings	2,000	300,000	No
Anthropic Claude Instant	400	300,000	No
Anthropic Claude V2	100	200,000	No
AI21 Labs Jurassic-2 Mid	400	300,000	No
AI21 Labs Jurassic-2 Ultra	100	300,000	No
Cohere Command	400	300,000	No
Cohere Command Light	400	300,000	No
Cohere Embed English	400	300,000	No
Cohere Embed Multilingual	400	300,000	No
Meta Llama2 Chat 13B	800	300,000	No
Stability.ai Diffusion XL 1.0	60	N/A	No

Model customization quotas

The following quotas apply to model customization.

Model quotas

The following quotas are enforced for Titan Express.

Description	Maximum value	Adjustable
Sum of input and output tokens when batch size is 1	4096	No
Sum of input and output tokens when batch size is between 2 and 4	2048	No
Character quota	Token quota x 6	No

Fine-tuning quotas

The following quotas are enforced for fine-tuning.

Description	Default	Adjustable
Number of training records in a dataset	10,000	Yes
Number of validation records in a dataset	1,000	Yes
Training dataset file size	1 GB	Yes
Validation dataset file size	100 MB	Yes

Training quotas

The following quotas are enforced for model training.

Description	Default	Adjustable
Number of submitted model customization jobs at one time per account	2	Yes
Number of customized models	100	Yes
Number of tags attached to a resource	See AWS Resource Groups and Tagging endpoints and quotas	Yes

Provisioned Throughput quotas

The following quotas are enforced for Provisioned Throughput.

Description	Default	Adjustable
Number of model units for a Provisioned Throughput for a foundation model	0	Yes
Number of model units for a Provisioned Throughput for a custom model	2	Yes

Document history for the Amazon Bedrock User Guide

- **Latest documentation update:** November 13, 2023

The following table describes important changes in each release of Amazon Bedrock. For notification about updates to this documentation, you can subscribe to an RSS feed.

Change	Description	Date
New model releases (p. 179)	Updated documentation to cover new Meta and Cohere models. For more information, see Amazon Bedrock .	November 13, 2023
Documentation localization (p. 179)	Amazon Bedrock documentation is now available in Japanese and German .	October 20, 2023
Region expansion	Amazon Bedrock is now available in Europe (Frankfurt) (eu-central-1). For information on endpoints, see Amazon Bedrock endpoints and quotas .	October 19, 2023
Region expansion	Amazon Bedrock is now available in Asia Pacific (Tokyo) (ap-northeast-1). For information on endpoints, see Amazon Bedrock endpoints and quotas .	October 3, 2023
Gated general release	Gated general release of the Amazon Bedrock service. For more information, see Amazon Bedrock .	September 28, 2023

AWS Glossary

For the latest AWS terminology, see the [AWS glossary](#) in the *AWS Glossary Reference*.