



BlueXP workload factory for GenAI documentation

GenAI

NetApp
November 11, 2024

Table of Contents

- BlueXP workload factory for GenAI documentation 1
- Release notes 2
 - What's new 2
- Get started 4
 - Learn about workload factory for GenAI 4
 - Quick start for GenAI 6
 - GenAI requirements 6
 - Identify the data sources to integrate in your knowledge base 9
 - Deploy the GenAI infrastructure 10
 - Use the external example chatbot application 12
 - Components for the NetApp GenAI engine 13
- Use GenAI to build knowledge bases 15
 - Create a knowledge base 15
 - Test a knowledge base 17
 - Activate external authentication for a knowledge base 18
 - Publish a knowledge base and view the unique endpoint 19
 - Create a RAG-based AI application 20
 - What you can do next 20
- Administer GenAI 22
 - Manage the GenAI infrastructure 22
 - Manage knowledge bases 23
 - Manage data sources 27
- Knowledge and support 30
 - Register for support 30
 - Troubleshooting 32
 - Get help 34
- Legal notices 39
 - Copyright 39
 - Trademarks 39
 - Patents 39
 - Privacy policy 39
 - Open source 39

BlueXP workload factory for GenAI documentation

Release notes

What's new

Learn what's new with the Generative AI workloads capability of workload factory.

3 November 2024

Mask Personal Identifiable Information with data guardrails

The Generative AI workload introduces the data guardrails feature, powered by BlueXP classification. The data guardrails feature identifies and masks Personal Identifiable Information (PII) helping you maintain compliance and strengthen security for your sensitive organizational data.

[Create a knowledge base](#)

[Learn about BlueXP classification](#)

29 September 2024

Snapshot and restore support for knowledge base volumes

You can now protect your Generative AI workloads data by taking a point-in-time copy of a knowledge base. This enables you to protect your data against accidental loss or test changes to the settings of the knowledge base. You can restore the previous version of the knowledge base volume at any time.

[Take a snapshot of a knowledge base volume](#)

[Restore a snapshot of a knowledge base volume](#)

Pause scheduled scans

You can now pause scheduled data source scans. By default, Generative AI workloads scans each data source daily to ingest new data into each knowledge base. If you don't want the latest changes to be ingested (during testing or while restoring a snapshot, for example) you can pause the scheduled scans and resume them at any time.

[Manage knowledge bases](#)

Data protection volumes now supported for knowledge bases

When selecting a knowledge base volume, you can now choose a data protection volume that is part of a NetApp SnapMirror replication relationship. This enables you to store knowledge bases on volumes that are already protected by SnapMirror replication.

[Identify the data sources to integrate in your knowledge base](#)

1 September 2024

Additional chunking strategies

Generative AI workloads now supports multi-sentence chunking and overlap-based chunking for data sources.

Dedicated volume for each knowledge base

Generative AI workloads now creates a dedicated Amazon FSx for NetApp ONTAP volume for each new knowledge base, enabling individual snapshot policies for each knowledge base and improved protection against failures and data poisoning.

4 August 2024

Amazon CloudWatch Logs integration

Generative AI workloads is now integrated with Amazon CloudWatch Logs, enabling you to monitor Generative AI workloads log files.

Example chatbot application

The NetApp workload factory GenAI sample application enables you to test authentication and retrieval from your published NetApp workload factory knowledge base by interacting directly with it in a web-based chatbot application.

7 July 2024

Initial release of the workload factory for GenAI

The initial release includes the capability to develop a knowledge base that is customized by embedding your organization's data. The knowledge base can be accessed by a chatbot application for your users. This capability ensures accurate and relevant responses to organization-specific questions, enhancing the satisfaction and productivity for all of your users.

Get started

Learn about workload factory for GenAI

Workload factory for GenAI enables you to integrate Amazon FSx for NetApp ONTAP file systems with GenAI foundation models. This provides you with high performance storage with a rich set of protection, security, and cost optimization features for your AI datasets.

What is workload factory for GenAI?

Workload factory for GenAI enables you to use your enterprise data sources on Amazon FSx for NetApp ONTAP with Generative AI applications. Utilizing the Retrieval Augmentation Generation (RAG) framework, you can quickly connect data sources to Foundation Models available via Amazon Bedrock to develop Generative AI powered applications such as virtual assistants, Q&A chatbots, document summarization, content creation, etc.

Using Generative AI with your organizational data enables you to leverage your own knowledge and expertise, not rely on just the model's intelligence based on public data the models were trained on. Using RAG to customize the models ensures accurate and relevant responses to organization-specific questions, enhancing the productivity and efficiency for the users of your applications using Generative AI.

Developing a GenAI application that is tailored to your organization's data enables you to leverage your own knowledge and expertise. This customization capability ensures accurate and relevant responses to organization-specific questions, enhancing the satisfaction and productivity for all of your users.

For more information about workload factory, refer to the [workload factory overview](#).

Benefits of using GenAI to create generative AI applications

Workload factory for GenAI simplifies the process to deploy infrastructure needed to build Generative AI applications using Retrieval Augmented Generation (RAG). Specifically, GenAI provides the following benefits:

- Without needing a deep knowledge of data infrastructure, foundation and language models, IT administrators and developers can accelerate application development by utilizing the automation provided by GenAI. Data administrators and developers can easily and quickly create enterprise knowledge bases that embed your organization's unstructured data to be used by generative AI applications.
- Enhance security by preserving user permissions in files embedded in the knowledge bases to ensure that data security and privacy is maintained. An application, such as a chatbot, can be developed to provide only the authenticated users with answers based on data the users have access to.
- Keep your enterprise data private and secure within your AWS customer account where your organizational data is never externally exposed.
- Accelerate development of GenAI applications such as a Q&A chatbot using open-source frameworks such as LangChain utilizing the GenAI API to provision and manage knowledge bases, chat with a knowledge base, and store and retrieve chat history.
- Improve data protection and availability posture by deploying the generative AI data infrastructure on FSx for NetApp ONTAP file systems and taking advantage of ONTAP features such as high availability, snapshots for local data protection and recovery, SnapMirror for disaster recovery, and SnapVault for backing up your data infrastructure.
- Reduce overall storage costs for generative AI data infrastructure by taking advantage of ONTAP data efficiency features such as data deduplication, compression and compaction, data tiering, and thin

provisioning.

How GenAI works

GenAI uses your organization's private data to complement the model's intelligence (based on the data it was trained on) to provide customized answers to questions asked by your users in your organization. You first deploy the infrastructure needed for a RAG framework, then build a knowledge base using your organization's data sources and foundation models available via Amazon Bedrock, and then connect an application (such as a Q&A chatbot) to the knowledge base.

How workload factory for GenAI helps to build Generative AI applications

GenAI helps to build generative AI applications using RAG in the following ways:

- Deploys the required infrastructure for Retrieval Augmented Generation (RAG) framework to work with data sources on FSx for ONTAP file systems and Amazon Bedrock. The infrastructure includes the NetApp GenAI Engine instance for managing data, an embedded vector database (LanceDB), and storage on your FSx for ONTAP file system for the vector database.
- Helps connect the data sources to embeddings and language models available via Amazon Bedrock for embedding data sources and retrieving the responses for user queries. The data sources, along with models and their configuration, are presented as FSx for ONTAP knowledge bases.
- Ingests source data into the knowledge base to embed source files on SMB shares and NFS exports on FSx for ONTAP file systems along with storing file permissions for files within SMB shares.
- Automatically builds conversation starter questions based on the content in the knowledge bases.
- Provides a chat simulator for data administrators to test chatting with the knowledge bases.

Tools to use workload factory

You can use BlueXP workload factory with the following tools:

- **Workload factory console:** The Workload factory console provides a visual interface that gives you a holistic view of your applications and projects.
- **BlueXP console:** The BlueXP console provides a hybrid interface experience so that you can use BlueXP workload factory along with other BlueXP services.
- **REST API:** Workload factory REST APIs let you deploy and manage your FSx for ONTAP file systems and other AWS resources.
- **CloudFormation:** AWS CloudFormation code lets you perform the actions you defined in the workload factory console to model, provision, and manage AWS and third-party resources from the CloudFormation stack in your AWS account.
- **Terraform BlueXP workload factory Provider:** Terraform lets you build and manage infrastructure workflows generated in the workload factory console.

Cost

There is no cost for using the GenAI capability of workload factory.

However, you'll need to pay for AWS resources that you deploy in order to support the generative AI infrastructure. For example, you will pay AWS for the Amazon Bedrock AI service, FSx for ONTAP file system

and storage capacity, and the GenAI engine EC2 instance.

Licensing

No special licenses are needed from NetApp to use the AI capabilities of workload factory.

Quick start for GenAI

Get started creating a knowledge base using your organization's data that exists on Amazon FSx for NetApp ONTAP file systems. An application such as a chatbot will access this knowledge base to provide organization-focused responses to end users.

1

Log in to workload factory

You'll need to [set up an account with workload factory](#) and log in using one of the [console experiences](#).

2

Set up your environment to meet GenAI requirements

You'll need AWS credentials to deploy the AWS infrastructure, a deployed and discovered FSx for ONTAP file system, the list of data sources that you want to integrate in your knowledge base, access to the Amazon Bedrock AI service, and more.

[Learn more about GenAI requirements.](#)

3

Identify the FSx for ONTAP file system that contains the data sources

The data sources that you'll integrate in your knowledge base can be located on a single FSx for ONTAP file system, or on multiple FSx for ONTAP file systems. These systems can be in different VPCs, but they must be accessible within the same network.

[Learn how to identify data sources.](#)

4

Deploy the knowledge base infrastructure

Launch the infrastructure deployment wizard to deploy the knowledge base infrastructure in your AWS environment. This process deploys an EC2 instance for the NetApp GenAI engine, and a volume on an FSx for ONTAP file system to contain the NetApp AI Engine databases. The volume is used to store the vector database used by the knowledge base.

[Learn how to deploy the knowledge base infrastructure.](#)

What's next

You can now build a knowledge base to provide organization-focused responses to end users.

GenAI requirements

Ensure that workload factory and AWS are set up properly before you build your knowledge base. This includes having your AWS log in credentials, a deployed FSx for

ONTAP file system that contains the data sources you want to integrate in your knowledge base, access to the Amazon Bedrock AI service, and more.

Workload factory login and account

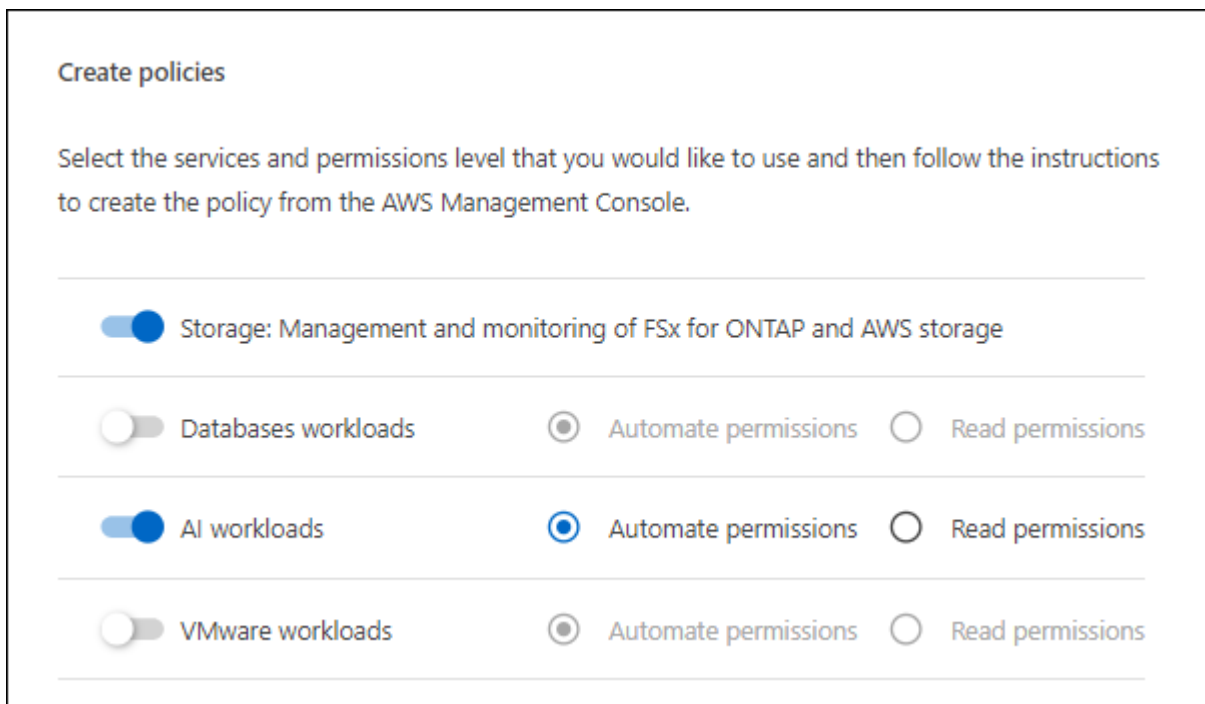
You'll need to [set up an account with workload factory](#) and log in using one of the [console experiences](#).

AWS credentials and permissions

You need to add AWS credentials to workload factory with automate permissions, which means you'll be using workload factory in *Automate* mode for GenAI.

Basic mode and *Read* mode permissions are not supported at this time.

When setting up your credentials, selecting permissions as shown below provides you with full access to manage FSx for ONTAP file systems and to deploy and manage the GenAI EC2 instance and other AWS resources needed for your knowledge base and chatbot.



[Learn how to add AWS credentials to workload factory](#)

Amazon Bedrock

Amazon Bedrock enables you to use foundation models and it provides the capabilities to build generative AI applications.

Before you get started with workload factory for GenAI, you must set up Amazon Bedrock. Your GenAI deployment must be in an AWS region that has Amazon Bedrock enabled.

- [AWS documentation: Set up Amazon Bedrock](#)
- [AWS documentation: Supported regions and models for Knowledge bases for Amazon Bedrock](#)

Embedding model

You must enable the embedding model that you plan to use before creating your knowledge base. The following embedding models are supported:

- Titan Embeddings G1 - Text
- Titan Embedding Text v2
- Titan Multimodal Embeddings G1

[Learn more about Amazon Titan](#)

Chat model

You must enable the foundational chat model that you plan to use before creating your knowledge base. The following Claude chat models are supported:

- Claude 3.5 Sonnet
- Claude 3 Opus
- Claude 3 Haiku
- Claude 3 Sonnet
- Claude 2.1
- Claude 2.0

Since model support varies by AWS region, refer to [this AWS documentation page](#) to verify which models you can use in the regions where you plan to deploy your knowledge base.

Learn more about the available models to help make your selection: [Anthropic's Claude in Amazon Bedrock](#)

FSx for ONTAP file system

You need a minimum of one FSx for ONTAP file system:

- One file system will be used by the NetApp GenAI engine to store the vector database used by the knowledge base.

This FSx for ONTAP file system must use FlexVol volumes. FlexGroup volumes are not supported.

- One or more file systems will contain the data sources that you'll be integrating into your knowledge base.

One FSx for ONTAP file system can be used for both of these purposes, or you can use multiple FSx for ONTAP file systems.

- You'll need to know the AWS region, VPC, and subnet where the AWS FSx for ONTAP file system resides. The file system must be in an AWS region that has Amazon Bedrock enabled.
- You'll need to consider the tag key/value pairs that you want to apply to the AWS resources that are part of this deployment (optional).
- You'll need to know the key pair information that allows you to securely connect to the NetApp AI engine instance.

[Learn how to deploy and manage FSx for ONTAP file systems](#)

Identify the data sources to integrate in your knowledge base

Identify, or create, the documents (data sources) that reside on your FSx for ONTAP file system that you'll integrate in your knowledge base. These data sources enable the knowledge base to provide accurate and personalized answers to user queries based on data that is relevant to your organization.

Maximum number of data sources

The maximum number of supported data sources is 10.

Location of data sources

Data sources can be stored in a single volume, or in a folder within a volume, on an SMB share or NFS export on an Amazon FSx for NetApp ONTAP file system. Data sources can also be stored on Amazon FSx for NetApp ONTAP volumes that are in a NetApp SnapMirror data protection relationship.

You can't select individual documents within a volume or folder, therefore, you should ensure that each volume or folder that contains data sources does not contain extraneous documents that shouldn't be integrated with your knowledge base.

You can add multiple data sources into each knowledge base, but they all need to reside on FSx for ONTAP file systems that are accessible from your AWS account.

The maximum file size for each data source is 50 MB.

Supported protocols

The knowledge base supports data from volumes that use either NFS or SMB/CIFS protocols. When selecting files stored using the SMB protocol, you'll need to enter the Active Directory information so that the knowledge base can access the files on those volumes. This includes the Active Directory Domain, IP address, user name, and password.

When storing your data source on a share (file or directory) accessed over SMB, the data is only accessible by chatbot users or groups who have the permissions to access that share. When this "permission-aware capability" is enabled, the AI system will compare the user email in auth0 to the users allowed to view or use the files on the SMB share. The chatbot will provide answers based on user permissions for the embedded files.

For example, if you have integrated 10 files (data sources) into your knowledge base, and 2 of the files are human resources files that contain restricted information, only chatbot users who are authenticated to access those 2 files will receive responses from the chatbot that include data from those files.

Supported data source file formats

The following data source file formats are currently supported.

File format	Extension
Comma-separated values file	.csv

File format	Extension
Markdown	.md
Microsoft Word	.doc or .docx
Plain text	.txt
Portable Document Format	.pdf

Deploy the GenAI infrastructure

You need to deploy the GenAI infrastructure for RAG framework in your environment before you can build FSx for ONTAP knowledge bases and applications for your organization. The primary infrastructure components are the Amazon Bedrock service, a virtual machine instance for the NetApp GenAI engine, and an FSx for ONTAP file system.

The deployed infrastructure can support multiple knowledge bases and chatbots, so you'll typically only need to perform this task once.

Infrastructure details

Your GenAI deployment must be in an AWS region that has Amazon Bedrock enabled. [View the list of supported regions](#)

The deployment consists of the following components.

Amazon Bedrock service

Amazon Bedrock is a fully-managed service that enables you to use foundation models (FMs) from leading AI companies via a single API. It also provides the capabilities you need to build secure generative AI applications.

[Learn more about Amazon Bedrock](#)

Virtual machine for the NetApp GenAI engine

The NetApp GenAI engine gets deployed during this process. It provides the processing power to ingest the data from your data sources and then write that data in the vector database.

FSx for ONTAP file system

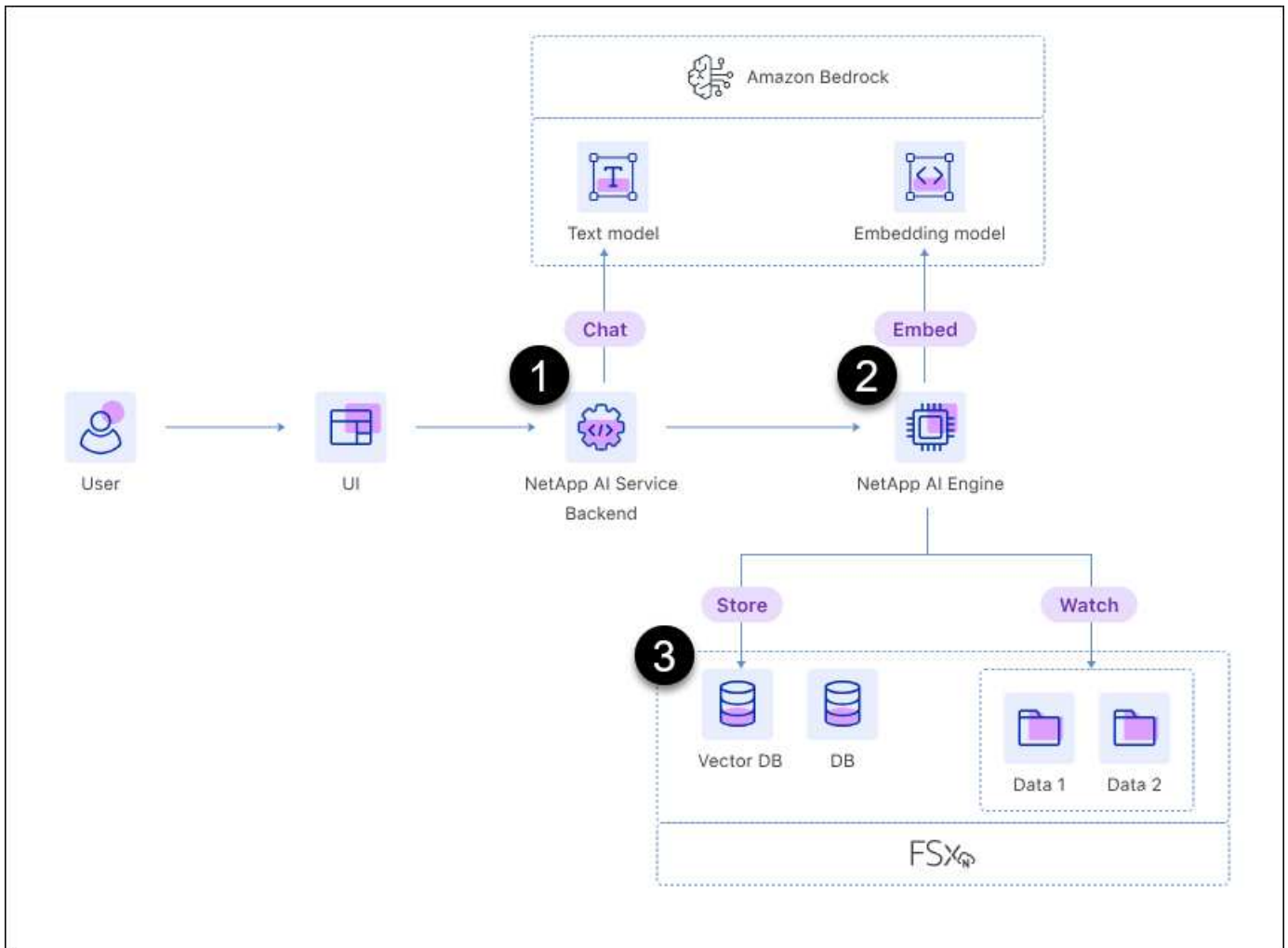
The FSx for ONTAP file system provides the storage for your GenAI system.

A single volume is deployed that will contain the vector database that stores the data that has been generated by the foundational model based on your data sources.

The data sources that you'll integrate into your knowledge base can reside on the same FSx for ONTAP file system, or on a different system.

The NetApp GenAI engine monitors and interacts with both of these volumes.

The following image shows the GenAI infrastructure. The components numbered 1, 2, and 3 are deployed during this procedure. The other elements must be in place before starting the deployment.



Deploy the GenAI infrastructure

You'll need to enter your AWS credentials and select the FSx for ONTAP file system to deploy the Retrieval Augmentation Generation (RAG) infrastructure.

Before you begin

Make sure your environment meets [the requirements](#) before you start this procedure.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. In the AI workloads tile, select **Deploy & manage**.
3. Review the infrastructure diagram and select **Next**.
4. Complete the items in the **AWS settings** section:
 - a. **AWS credentials:** Select or add the AWS credentials that provide permissions to deploy the AWS resources.
 - b. **Location:** Select an AWS region, VPC, and subnet.

The GenAI deployment must be in an AWS region that has Amazon Bedrock enabled. [View the list of supported regions](#)

5. Complete the items in the **Infrastructure settings** section:

- a. **Tags:** Enter any tag key/value pairs that you want to apply to all the AWS resources that are part of this deployment.
 - b. **Key pair:** Select a key pair that allows you to securely connect to the NetApp GenAI engine instance.
6. Select **Deploy** to begin the deployment.

Result

Workload factory starts deploying the chatbot infrastructure. This process can take up to 10 minutes.

During the deployment process, the following items are set up:

- The network is set up along with the private endpoints.
- The IAM role, instance profile, and security group are created.
- The virtual machine instance for the GenAI engine is deployed.
- Amazon Bedrock is configured to send logs to Amazon CloudWatch Logs, using a log group with the prefix `/aws/bedrock/`.
- The GenAI engine is configured to send logs to Amazon CloudWatch Logs, using a log group with the name `/netapp/wlmai/<tenancyAccountId>/randomId`, where `<tenancyAccountId>` is the [BlueXP account ID](#) for the current user.

Use the external example chatbot application

After you configure, activate, and publish a knowledge base, external application developers can configure and run the open source example chatbot application provided by NetApp to interact with your knowledge base and to learn how to use the workload factory API to create their own generative AI applications.

Steps

1. [Create a knowledge base](#).
2. [Activate authentication](#) for the knowledge base that you created.

This enables the knowledge base to authenticate API requests, and makes token validation and ACLs required when using the API endpoints.



External chat applications that integrate with this knowledge base will need to use the same authentication provider (issuer) that you configure in the authentication settings for the knowledge base.

3. [Publish the knowledge base](#) to enable API access for external applications.

After a knowledge base is published, the API endpoints are accessible externally, and you can integrate the knowledge base with an external chat application (such as the example chatbot application).

4. Download the example chatbot application package from [GitHub](#).
5. Install and run the chatbot application by following the instructions in the README file included in the package.
6. Browse to <http://localhost:9091> to log in to the application.

The example chatbot application appears.

Learn more

[workload factory API documentation](#)

Components for the NetApp GenAI engine

When you deploy the GenAI infrastructure, workload factory creates an EC2 instance for the GenAI engine. It also creates an IAM role, security group, and private endpoints for this instance. You might want to understand more details about these components that workload factory creates in your AWS environment.

EC2 instance type

m5.large

IAM role

The GenAI engine instance needs permissions to send chunks of data to the embedding model on Amazon Bedrock and to communicate with the NetApp AI Service Backend. The IAM role includes the following permissions:

```
"ssm:Describe*"
"ssm:Get*"
"ssm:List*"
"ssm:PutInventory"
"ssm:PutComplianceItems"
"ssm:PutConfigurePackageResult"
"ssm:SendCommand"
"ssm:UpdateAssociationStatus"
"ssm:UpdateInstanceAssociationStatus"
"ssm:UpdateInstanceInformation"
"ssmmessages:*"
"ec2messages:GetMessages"
```

Security group

The outbound rules are open to all traffic, while the inbound rules are completely closed.

Private endpoints

If the target VPC doesn't already have them, workload factory creates private endpoints for the GenAI engine EC2 instance so that it can communicate with the following AWS services:

- Amazon Bedrock
- Amazon EC2
- Amazon Elastic Container Registry (ECR)
- Amazon S3
- AWS Systems Manager (SSM)
- Amazon FSx for NetApp ONTAP

- Amazon CloudWatch

Use GenAI to build knowledge bases

Create a knowledge base

After you've deployed the AI infrastructure and identified the data sources that you'll integrate in your knowledge base from your FSx for ONTAP datastores, you are ready to build the knowledge base using workload factory. As part of this step, you'll also define the AI characteristics and create conversation starters.

About this task

Knowledge bases have two data integration modalities - *public mode* and *Enterprise mode*.

Public mode

A knowledge base can be used without integrating data sources from your organization. In this case, an application integrated with the knowledge base will only provide results from publicly available information on the internet. This is known as a *public mode* integration.

Enterprise mode

In most cases you'll want to integrate data sources from your organization into the knowledge base. This is known as an *Enterprise mode* integration because it provides knowledge from your enterprise.

Data sources from your organization may contain Private Identity Information (PII). To safeguard this sensitive information, you can enable *data guardrails* when creating and configuring knowledge bases. Data guardrails, powered by BlueXP classification, identifies and masks PII, making it inaccessible and irretrievable.

[Learn about BlueXP classification.](#)



Data guardrails can be enabled or disabled at any time. If you switch data guardrails enablement, Workload Factory scans the entire knowledge base from scratch, which incurs a cost.

Create and configure the knowledge base

The knowledge base defines characteristics such as the Bedrock AI models and embedding format that you want to use to create your knowledge base.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. In the AI workloads tile, select **Deploy & manage**.
3. From the Knowledge bases tab, select **Add knowledge base**.
4. On the Define knowledge base page, configure the knowledge base settings:
 - a. **Name:** Enter the name you want to use for the knowledge base.
 - b. **Description:** Enter a detailed description for the knowledge base.
 - c. **Embedding model:** The embedding model defines how your data will be converted into vector embeddings for the knowledge base. Workload factory supports the following models:
 - Titan Embeddings G1 - Text

- Titan Embedding Text v2
- Titan Multimodal Embeddings G1

Note that you must have already enabled the embedding model from Amazon Bedrock.

[Learn more about Amazon Titan](#)

- d. **Chat model:** Choose from Claude chat models that are integrated in Amazon Bedrock. Note that you must have already enabled the chat model from Amazon Bedrock.

Learn more about the available models to help make your selection: [Anthropic's Claude in Amazon Bedrock](#)

- e. **Data guardrails:** Choose whether you want to enable or disable data guardrails. [Learn about data guardrails, powered by BlueXP classification.](#)

The following prerequisites must be complete to enable data guardrails.

- A service account is required to communicate with BlueXP classification. You must have the *Organization admin* role on your BlueXP tenancy account for service account creation. A member who has the Organization admin role can complete all actions in BlueXP. [Learn how to add a role to a member in BlueXP](#)
- The AI engine must have access to the [BlueXP console](#)
- You'll need to do the following as described in [BlueXP classification documentation](#):
 1. Create a BlueXP Connector
 2. Ensure that your environment can meet the prerequisites
 3. Deploy BlueXP classification

- f. **Conversation starters:** Choose whether you want to provide up to four conversation starter prompts that are displayed to users who interact with a chatbot that uses this knowledge base. We recommend that you enable this setting.

If you activate conversation starters, "Automatic mode" is selected by default. "Manual mode" can be enabled only after you've added data sources to your knowledge base. [Learn how to modify knowledge base settings.](#)

- g. **FSx for ONTAP file system:** When you define a new knowledge base, Workload factory creates a new Amazon FSx for NetApp ONTAP volume to store it. Choose an existing file system name and SVM (also called a storage VM) where the new volume will be created.

- h. **Snapshot policy:** Choose a snapshot policy from the list of existing policies defined in the workload factory storage inventory. Recurring snapshots of the knowledge base will automatically be created at a frequency based on the snapshot policy you select.

If the snapshot policy you need doesn't exist, you can [create a snapshot policy](#) on the storage VM that contains the volume.

5. Select **Create knowledge base** to add the knowledge base to GenAI.

A progress indicator appears while the knowledge base is created.

After the knowledge base is created, you have the option to add a data source to your new knowledge base or to end the process without adding a data source. We recommend that you select **Add data source** and add one or more data sources now.

Add data sources to the knowledge base

You can add one or more data sources to populate the knowledge base with your organization's data.

About this task

The maximum number of supported data sources is 10.

Steps

1. After you select **Add data source**, the **Select a file system** page displays.
2. **Select a file system**: Select the FSx for ONTAP file system where your data source files reside and select **Next**.
3. **Select a volume**: Select the volume on which your data source files reside and select **Next**.

When selecting files stored using the SMB protocol, you'll need to enter the Active Directory information, which includes the domain, IP address, user name, and password.

4. **Select a data source**: Select the data source location based on where you have saved the files. This can be an entire volume, or just a specific folder or sub-folder in the volume, and select **Next**.
5. **Define AI parameters**: In the **Chunking strategy** section, define the how the GenAI engine splits data source content into chunks when the data source is integrated with a knowledge base. You can choose one of the following strategies:
 - **Multi-sentence chunking**: Organizes information from your data source into sentence-defined chunks. You can choose how many sentences make up each chunk (up to 100).
 - **Overlap-based chunking**: Organizes information from your data source into character-defined chunks that can overlap neighboring chunks. You can choose the size of each chunk in characters, and how much each chunk overlaps with adjacent chunks. You can configure a chunk size of between 50 and 3000 characters, and an overlap percentage of between 1 and 99%.



Choosing a high overlap percentage can greatly increase storage requirements with only slight improvements in retrieval accuracy.

6. In the **Permission aware** section, which is available only when the data source you selected is on a volume that uses the SMB protocol, you can enable or disable the selection:
 - **Enabled**: Users of the chatbot who access this knowledge base will only get responses to queries from data sources to which they have access.
 - **Disabled**: Users of the chatbot will receive responses using content from all integrated data sources.
7. Select **Add** to add this data source to your knowledge base.

Result

The data source starts to be embedded into your knowledge base. The status changes from "Embedding" to "Embedded" when the data source is completely embedded.

After you add a single data source to the knowledge base, you can test it locally in the chatbot simulator window and make any required changes before you make the chatbot available to your users. You can also follow the same steps to add additional data sources to the knowledge base.

Test a knowledge base

After you create the knowledge base, you'll be able to test it locally using the chatbot

simulator and make any required changes before you make the knowledge base available to your users through a chatbot application.

About this task

You test your knowledge base to make sure it performs as you expect, and you can customize the conversation starters that you want to be available by default for chatbot users of this knowledge base. The chatbot simulator runs against all the data sources that have been embedded in the knowledge base.

You can test a knowledge base by chatting with your embedded data sources in the chatbot simulator. Note that none of the interaction or insights are captured in the GenAI vector database when testing the knowledge base locally.

You'll perform most of your testing within workload factory before you deploy the knowledge base in an application for your users. If you need to make changes to your data source or the chatbot operation, you'll want to do it now before you publish your knowledge base.


Some of the tasks you'll want to perform to test your chatbot are:

- Enter a large number of questions that are relevant to your organization to make sure the answers are as expected.
- Customize the conversation starters that you want to be available by default for your users in the chatbot application.
- Make sure that the attributed content that is provided at the bottom of the chatbot answers contain the correct references.

Steps

1. From the Knowledge bases inventory page, select the knowledge base that you want to test.


The chatbot simulator appears in the right pane. If defined, existing conversation starters are displayed as well.

2. In the chatbot entry field, enter a prompt or question and select  to see how your chatbot responds with your organizational knowledge.
3. If you need to update any of your data sources so that your knowledge base provides more focused answers, make those changes now and then retest the knowledge base.

Activate external authentication for a knowledge base

Activate authentication for a knowledge base so that token validation and ACLs are required when using the API endpoints to integrate a knowledge base with a chatbot application. When you activate authentication, you configure settings for a JSON Web Token that will be used for API requests to a knowledge base from chatbot clients.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, select the knowledge base for which you want to activate authentication.
3. Select  and select **Manage knowledge base**.
4. Select the **Actions** menu and select **Manage authentication settings**.

5. Set up authentication:

a. Select **Activate authentication settings**.

b. Provide the required information. Examples are provided, but you should obtain the values for these fields from your authentication provider:

- **Algorithms:** The signing algorithm that your authentication provider uses.
- **Audience (Optional):** A string containing the intended recipient of the token (sometimes a URL).
- **Issuer:** A string that identifies the provider that issued the token.

For example, Amazon Cognito uses issuer strings with the following format:

```
https://cognito-idp-<region>.amazonaws.com/<UserPoolID>
```

Where <region> is the AWS region containing the user pool, and <UserPoolID> is your user pool ID. You can retrieve your user pool ID using the following command:

```
aws cognito-idp list-user-pools --max-results=60 --output=table
```

- **JWKS URI:** The URI string that provides public keys needed to verify signatures of this token.

For example, Amazon Cognito uses JWKS URI strings with the following format:

```
https://cognito-idp.<region>.amazonaws.com/<userPoolId>/well-known/jwks.json
```

Where <region> is the AWS region containing the user pool, and <UserPoolID> is your user pool ID. You can retrieve your user pool ID using the following command:

```
aws cognito-idp list-user-pools --max-results=60 --output=table
```

6. Select **Save**.

Result

Authentication for the knowledge base is now active, and you can use API endpoints to interact with the knowledge base and integrate the knowledge base with a chatbot application.

Publish a knowledge base and view the unique endpoint

After you've built and tested your knowledge base locally, you can publish the knowledge base so that it can be integrated with a chatbot application that will enable your users to query the knowledge base.

About this task

Publishing the knowledge base enables you to use it in chat applications. The publish action triggers the

workload factory API to generate and publish unique endpoints. After publishing, the knowledge base becomes accessible for chat applications, and the API endpoints are ready for integration.

Each knowledge base that you publish has unique endpoints.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, select the knowledge base that you want to publish.
3. Select **...** and select **Manage knowledge base**.

This page displays the published status, embedding status of the data sources, embedding mode, and the list of all embedded data sources.

4. Select the **Actions** menu and select **Publish**.

Workload factory publishes the knowledge base. On the details page for the knowledge base, the status changes from **Unpublished** to **Published**.

You can now get details about the unique endpoint for the knowledge base.

5. Next to the published status, select **View**.

Details about how to access the knowledge base using the workload factory API is displayed.

6. From the **View published info** dialog box, copy the API endpoints that you can use to integrate the knowledge base with an application.

To learn more about the API endpoints, go to the [API documentation](#) and select **AI > External**.

Before you can use these endpoints, you need to obtain a user token from your authentication provider.

Result

You now have a published knowledge base and the unique endpoint that you can use to integrate the knowledge base with a chatbot application.

Create a RAG-based AI application

After you build your knowledge base and test your chatbot, you are ready to set up the application that will enable your users to query the chatbot.

[Learn how to create a RAG-based AI application on FSx for ONTAP](#)

What you can do next

Now that you've created a knowledge base using your enterprise data and deployed it for your users, you can manage the knowledge base, data sources, and the RAG infrastructure, including FSx for ONTAP file systems.

Some of the tasks you can perform to manage your knowledge base components are:

- Update the content of your data sources, or add new data sources, and sync those changes with your

knowledge base and chatbot.

- Manage your data source settings, including the chunking strategy and permission awareness (for SMB file access).
- Manage your knowledge base settings, including the chat model and conversation starters.
- Unpublish a knowledge base or republish it after making changes.
- Back up and protect the important data on your FSx for ONTAP file system to make sure your knowledge base data and other infrastructure components are always available.

For information about managing your FSx for ONTAP file system, go to the [workload factory for Amazon FSx for NetApp ONTAP documentation](#) to view the backup and protection capabilities you can use.

Administer GenAI

Manage the GenAI infrastructure

You can view details about your deployed GenAI RAG infrastructure or remove the chatbot infrastructure if you no longer need it.

View information about the infrastructure

You can view information about the chatbot infrastructure.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the workload factory navigation menu, select **AI**.
3. Select the **Infrastructure** tab.
4. View information about the infrastructure, which includes details about the following components:
 - AWS settings
 - Infrastructure settings
 - The AI engine
 - The vector database

Remove the infrastructure

If you no longer need the chatbot infrastructure that you deployed for one or more chatbots, you can remove it from workload factory.



All chatbots that have been deployed on this infrastructure will be disabled and all chat history will be deleted.

This operation only removes the links to the AI infrastructure from workload factory—it does not remove all the components from AWS. You'll need to manually delete the following infrastructure components from AWS:

- The VM instance
- Private endpoints
- The volume on the FSx for ONTAP file system that contains the AI databases
- The IAM role
- The policy
- The security group

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the workload factory navigation menu, select **AI**.
3. Select the **Infrastructure** tab.
4. Select **...** and select **Remove chatbot infrastructure**.

5. Confirm that you want to delete the infrastructure and select **Remove**.

Result

The chatbot infrastructure components are removed from workload factory.

Manage knowledge bases

After you create a knowledge base, you can view the knowledge base details, modify the knowledge base, integrate additional data sources, or delete the knowledge base.

View information about a knowledge base

You can view information about the settings for a knowledge base and the data source that are integrated.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the workload factory navigation menu, select **AI**.
3. Select the knowledge base that you want to view.

If defined, the conversation starters that are currently being used display in the right pane.

4. To view knowledge base details, select **...** and select **Manage knowledge base**.

This page displays the published status, embedding status of the data sources, embedding mode, the list of all embedded data sources, and more.

The **Actions** menu enables you to manage the knowledge base if you want to make any changes.

Edit a knowledge base

You can update a knowledge base by changing some settings, or you can add or remove data sources.

Each time you add, modify, or remove data sources from the knowledge base, you must sync the data source so that it is re-indexed to the knowledge base. Syncing is incremental, so Amazon Bedrock only processes the objects in your FSx for ONTAP volume that have been added, modified, or deleted since the last sync.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, select the knowledge base that you want to update.
3. Select **...** and select **Manage knowledge base**.

This page displays the published status, embedding status of the data sources, embedding mode, the list of all embedded data sources, and more.

4. Select the **Actions** menu and select **Edit knowledge base**.
5. In the Edit knowledge base page, you can change the knowledge base name, description, embedding model, chat model, data guardrails enablement, choose whether conversation starters are created automatically or manually, and the snapshot policy used for the volume that contains the knowledge base.

If you use Manual mode for conversation starters, you can change conversation starters here as well.



Every knowledge base scan, which includes embedding, costs. If data guardrails is enabled after a knowledge base has been created, then the knowledge base gets scanned again and incurs costs.

6. Select **Save** after you have made your changes.

Protect a knowledge base with snapshots

You can protect your knowledge base data by taking and restoring snapshots of your knowledge base volumes. You can restore from a snapshot to revert to the previous version of the knowledge base at any time.

Snapshots can be faster and more storage-efficient than backups, and enable you to protect each knowledge base using a different protection policy. Some of the scenarios where snapshots can be useful are:

- Accidental data loss or corruption
- Recovering from incorrect data being ingested into the knowledge base
- Testing different data sources or chunking strategies, and quickly reverting when the testing is complete

Take a snapshot of a knowledge base volume

You can save the state of a knowledge base by taking a manual snapshot of the knowledge base volume.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, select the knowledge base that you want to protect.
3. Select **...** and select **Manage knowledge base**.

This page displays the published status, embedding status of the data sources, embedding mode, the list of all embedded data sources, and more.

4. Select the **Actions** menu and select **Snapshot > Create new snapshot**.

A snapshot of the knowledge base is created.

Restore a snapshot of a knowledge base volume

You can restore a manual or scheduled snapshot of a knowledge base volume at any time.



You cannot restore a snapshot using the Generative AI workloads UI if the database stored on the volume is corrupt or has been deleted. As a workaround, you can restore the snapshot using the [ONTAP CLI](#) on the ONTAP cluster where the volume is hosted.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, select the knowledge base that you want to restore.
3. Select **...** and select **Manage knowledge base**.

This page displays the published status, embedding status of the data sources, embedding mode, the list of all embedded data sources, and more.

4. Select the **Actions** menu and select **Snapshot > Restore snapshot**.

The snapshot selection dialog appears, where you can see a list of the snapshots that have been created for this knowledge base.

5. (Optional) Deselect the **Pause running and scheduled scans after restoring the snapshot** option if you want scheduled and currently running data source scans to continue after the snapshot is restored.

This option is enabled by default to ensure that a scan doesn't happen while the knowledge base is in a partially restored state, or that a scan doesn't update a freshly restored knowledge base with older data.

6. Select the snapshot you want to restore from the list.
7. Select **Restore**.

Add additional data sources to a knowledge base

You can embed additional data sources in your knowledge base to populate it with additional organization data.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, select the knowledge base where you want to add the data source.
3. Select **...** and select **Add data source**.
4. **Select a file system:** Select the FSx for ONTAP file system where your data source files reside and select **Next**.
5. **Select a volume:** Select the volume on which your data source files reside and select **Next**.

When selecting files stored using the SMB protocol, you'll need to enter the Active Directory information, which includes the domain, IP address, user name, and password.

6. **Select a data source:** Select the data source location based on where you have saved the files. This can be an entire volume, or just a specific folder or sub-folder in the volume, and select **Next**.
7. **Define AI parameters:** In the **Chunking strategy** section, define the how the GenAI engine splits data source content into chunks when the data source is integrated with a knowledge base. You can choose one of the following strategies:
 - **Multi-sentence chunking:** Organizes information from your data source into sentence-defined chunks. You can choose how many sentences make up each chunk (up to 100).
 - **Overlap-based chunking:** Organizes information from your data source into character-defined chunks that can overlap neighboring chunks. You can choose the size of each chunk in characters, and how much each chunk overlaps with adjacent chunks. You can configure a chunk size of between 50 and 3000 characters, and an overlap percentage of between 1 and 99%.



Choosing a high overlap percentage can greatly increase storage requirements with only slight improvements in retrieval accuracy.

8. In the **Permission aware** section, which is available only when the data source you selected is on a volume that uses the SMB protocol, you can enable or disable the selection:
 - **Enabled:** Users of the chatbot who access this knowledge base will only get responses to queries from data sources to which they have access.

- **Disabled:** Users of the chatbot will receive responses using content from all integrated data sources.

9. Select **Add** to add this data source to your knowledge base.

Result

The data source is integrated into your knowledge base.

Synchronize your data sources with a knowledge base

Data sources are synchronized with the associated knowledge base automatically once a day so that any data source changes are reflected in the chatbot. If you make changes to any of your data sources and you'd like to synchronize the data immediately, you can perform an on-demand synchronization.

Syncing is incremental, so Amazon Bedrock only processes the objects in your data sources that have been added, modified, or deleted since the last sync.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, select the knowledge base that you want to synchronize.
3. Select **...** and select **Manage knowledge base**.
4. Select the **Actions** menu and select **Scan now**.

You'll see a message that your data sources are being scanned, and a final message when the scan is complete.

Result

The knowledge base is synchronized with the attached data sources and any active chatbot will start using the newest information from your data sources.

Evaluate chat models before creating a knowledge base

You can evaluate the available foundational chat models before creating a knowledge base so you can see which model works best for your implementation. Since model support varies by AWS region, refer to [this AWS documentation page](#) to verify which models you can use in the regions where you plan to deploy your knowledge base.



This functionality is available only when no knowledge bases have been created — when no knowledge bases exist in the Knowledge bases inventory page.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, you'll see the option to select the chat model on the right side of the page for the Chatbot.
3. Select the chat model from the list and enter a set of questions in the prompt area to see how the chatbot responds.
4. Try multiple models to see which model is best for your implementation.

Result

Use that chat model when you create your knowledge base.

Unpublish your knowledge base

After you've published your knowledge base so that it can be integrated with a chatbot application, you can unpublish it if you want to disable the chatbot application from accessing the knowledge base.

Unpublishing the knowledge base stops any chat applications from working. The unique API endpoint at which the knowledge base was accessible is disabled.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, select the knowledge base that you want to unpublish.
3. Select **...** and select **Manage knowledge base**.

This page displays the published status, embedding status of the data sources, embedding mode, and the list of all embedded data sources.

4. Select the **Actions** menu and select **Unpublish**.

Result

The knowledge base is disabled and is no longer accessible by a chatbot application.

Delete a knowledge base

If you no longer need a knowledge base, you can delete it. When you delete a knowledge base, it is removed from workload factory and the volume that contains the knowledge base is deleted. Any applications or chatbots that are using the knowledge base will stop working. Deleting a knowledge base is not reversible.

When you delete a knowledge base, you should also disassociate the knowledge base from any agents it is associated with to fully delete all resources associated with the knowledge base.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the Knowledge bases inventory page, select the knowledge base that you want to delete.
3. Select **...** and select **Manage knowledge base**.
4. Select the **Actions** menu and select **Delete knowledge base**.
5. In the Delete knowledge base dialog, confirm that you want to delete it and select **Delete**.

Result

The knowledge base is removed from workload factory and its associated volume is deleted.

Manage data sources

After you create a knowledge base using data sources on your FSx for ONTAP file system, you can view the data source details, update or change the data source contents, edit data source settings, or delete the data source.

View information about a data source

You can view information about the contents of a data source and you can view its embedding status with the

knowledge base. Since a data source is part of a knowledge base, you'll need to open the knowledge base first before you can view the data source details.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the workload factory navigation menu, select **AI**.
3. Select the knowledge base where the data source resides, and then select **...** and select **Manage knowledge base**.

The bottom part of the page lists the data sources that are part of this knowledge base.

4. Expand each row by selecting the **▼** to view detailed information about each data source, such as the FSx for ONTAP file system, the volume, and the path where the data source resides.

It also lists the embedding information and whether that data source is currently embedded in the knowledge base.

Edit data source settings

You can edit information about a data source that you've integrated in a knowledge base. Most of the information is fixed after you've added a data source to a knowledge base, but you can make changes to the chunking definition and configuration, and to the permission awareness (when using SMB protocol on the data source volume).

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the knowledge bases inventory page, select the knowledge base where the data source resides, and then select **...** and select **Manage knowledge base**.

The bottom part of the page lists the data sources that are part of this knowledge base.

3. In the row of the data source that you want to edit, select **...** and select **Edit data source**.
4. In the Edit data source page, select **▼** to expand the row for chunk definition.
5. Update the settings for the chunking strategy and configuration, and permission awareness (for SMB volumes) and select **Save**.

Result

The data source settings are updated and the AI system synchronizes the data source so that it is re-indexed to the knowledge base.

Update the contents of an existing data source

You can change the contents of a data source at any time to add or update your organizational data. If this data source is being actively used in a knowledge base, you must sync the data source so that it is re-indexed to the knowledge base. Syncing is incremental, so Amazon Bedrock only processes the objects in your FSx for ONTAP volume that have been added, modified, or deleted since the last sync.

Data sources are synchronized with the knowledge base automatically once a day so that any data source changes are reflected in the chatbot. If you make changes to a data source and you'd like to synchronize the data immediately, you can [perform an on-demand synchronization](#).

Delete a data source

If you no longer need a data source to be part of your knowledge base, you can delete it.

Steps

1. Log in to workload factory using one of the [console experiences](#).
2. From the knowledge base inventory page, select the knowledge base where the data source resides, and then select **...** and select **Manage knowledge base**.

The bottom part of the page lists the data sources that are part of this knowledge base.

3. In the row of the data source that you want to delete, select **...** and select **Delete data source**.
4. In the Delete data source dialog, confirm that you want to delete it and select **Confirm**.

Result

The data source is removed from the knowledge base and the AI system removes the indexed information about this data source from the knowledge base. Any information from that data source will no longer be available to chatbots that are using the knowledge base.

Knowledge and support

Register for support

Support registration is required to receive technical support specific to Workload Factory and its storage solutions and services. You must register for support from the BlueXP console, which is a separate web-based console from Workload Factory.

Registering for support does not enable NetApp support for a cloud provider file service. For technical support related to a cloud provider file service, its infrastructure, or any solution using the service, refer to "Getting help" in the Workload Factory documentation for that product.

[Amazon FSx for ONTAP](#)

Support registration overview

Registering your account ID support subscription (your 20 digit 960xxxxxxx serial number located on the Support Resources page in BlueXP) serves as your single support subscription ID. Each BlueXP account-level support subscription must be registered.

Registering enables capabilities like opening support tickets and automatic case generation. Registration is completed by adding NetApp Support Site (NSS) accounts to BlueXP as described below.

Register your account for NetApp support

To register for support and activate support entitlement, one user in your account must associate a NetApp Support Site account with their BlueXP login. How you register for NetApp support depends on whether you already have a NetApp Support Site (NSS) account.

Existing customer with an NSS account

If you're a NetApp customer with an NSS account, you simply need to register for support through BlueXP.

Steps

1. In the upper right of the Workload Factory console, select **Help > Support**.

Selecting this option opens the BlueXP console a new browser tab and loads the Support dashboard.

2. In the upper right of the BlueXP console, select the Settings icon, and select **Credentials**.
3. Select **User Credentials**.
4. Select **Add NSS credentials** and follow the NetApp Support Site (NSS) Authentication prompt.
5. To confirm that the registration process was successful, select the Help icon, and select **Support**.

The **Resources** page should show that your account is registered for support.



Note that other BlueXP users will not see this same support registration status if they have not associated a NetApp Support Site account with their BlueXP login. However, that doesn't mean that your BlueXP account is not registered for support. As long as one user in the account has followed these steps, then your account has been registered.

Existing customer but no NSS account

If you're an existing NetApp customer with existing licenses and serial numbers but *no* NSS account, you need to create an NSS account and associate it with your BlueXP login.

Steps

1. Create a NetApp Support Site account by completing the [NetApp Support Site User Registration form](#)
 - a. Be sure to select the appropriate User Level, which is typically **NetApp Customer/End User**.
 - b. Be sure to copy the BlueXP account serial number (960xxxx) used above for the serial number field. This will speed up the account processing.
2. Associate your new NSS account with your BlueXP login by completing the steps under [Existing customer with an NSS account](#).

Brand new to NetApp

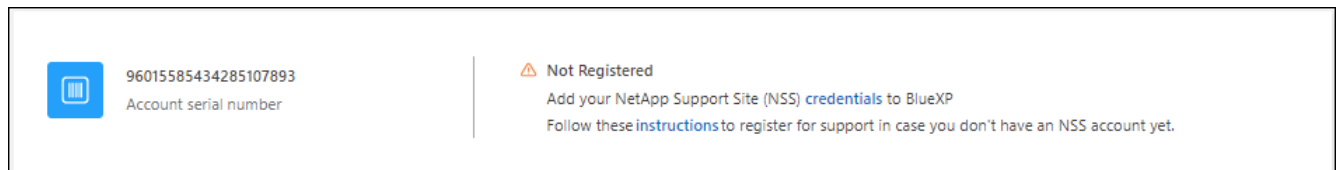
If you are brand new to NetApp and you don't have an NSS account, follow each step below.

Steps

1. In the upper right of the Workload Factory console, select **Help > Support**.

Selecting this option opens the BlueXP console a new browser tab and loads the Support dashboard.

2. Locate your account ID serial number from the Support Resources page.



3. Navigate to [NetApp's support registration site](#) and select **I am not a registered NetApp Customer**.
4. Fill out the mandatory fields (those with red asterisks).
5. In the **Product Line** field, select **Cloud Manager** and then select your applicable billing provider.
6. Copy your account serial number from step 2 above, complete the security check, and then confirm that you read NetApp's Global Data Privacy Policy.

An email is immediately sent to the mailbox provided to finalize this secure transaction. Be sure to check your spam folders if the validation email doesn't arrive in few minutes.

7. Confirm the action from within the email.

Confirming submits your request to NetApp and recommends that you create a NetApp Support Site account.

8. Create a NetApp Support Site account by completing the [NetApp Support Site User Registration form](#)
 - a. Be sure to select the appropriate User Level, which is typically **NetApp Customer/End User**.

- b. Be sure to copy the account serial number (960xxxx) used above for the serial number field. This will speed up the account processing.

After you finish

NetApp should reach out to you during this process. This is a one-time onboarding exercise for new users.

Once you have your NetApp Support Site account, associate the account with your BlueXP login by completing the steps under [Existing customer with an NSS account](#).

Troubleshooting

Learn how to work around some common problems you might encounter.

Common issues and solutions

If you have one of these issues, you can use the steps in the Workaround column to try to resolve it.

Area	Issue	Cause	Workaround
Deployment	Deployment fails because the volume already exists.	Workload factory for GenAI needs to create a new volume during the deployment process, but a volume already exists using the name you have specified.	Specify a unique name to use for the new volume, and try deploying again.
Deployment	The deployment fails because workload factory for GenAI is unable to mount the volume.	One or more of the inbound ports required for FSx for NetApp ONTAP are closed or filtered.	Open the ports listed in Security group rules for FSx for ONTAP .

Area	Issue	Cause	Workaround
Maintenance	The backend Docker instance used by workload factory for GenAI failed to start.	The volume was deleted and the EC2 instance was restarted.	<p>Use the following recovery steps:</p> <p>Steps</p> <ol style="list-style-type: none"> 1. Create a new volume on FSx for NetApp ONTAP. For example, the volume name can be <code>netapp_ai</code> and the volume path can be <code>/netapp_ai</code>. 2. SSH to the Amazon EC2 instance. 3. List the volumes: <pre>docker volume list</pre> <ol style="list-style-type: none"> 4. Remove the old volume: <pre>docker volume rm ec2-user_persistent_folder</pre> <ol style="list-style-type: none"> 5. Open the <code>docker-compose.yml</code> file using a text editor. 6. In the <code>volumes</code> section, change the device path to the new volume path. For example: <pre>volumes: persistent_folder: driver_opts: type: 'nfs' o: "addr=svm-0df66b96a890d8a72.\fs-0d673008aaca12bc3.\fsx.us-east-1.amazonaws.com,nolock,soft,rw" device: ':/netapp_ai' # Path to new volume</pre>
Maintenance	The backend Docker instance used by workload factory for GenAI failed to start.	The root volume was deleted.	Create a volume with a name and path, and then restart the backend Docker instance from Amazon EC2.

Get help

NetApp provides support for Workload Factory and its cloud services in a variety of ways. Extensive free self-support options are available 24x7, such as knowledgebase (KB) articles and a community forum. Your support registration includes remote technical support via web ticketing.

Get support for FSx for ONTAP

For technical support related to FSx for ONTAP, its infrastructure, or any solution using the service, refer to "Getting help" in the Workload Factory documentation for that product.

[Amazon FSx for ONTAP](#)

To receive technical support specific to Workload Factory and its storage solutions and services, use the support options described below.

Use self-support options

These options are available for free, 24 hours a day, 7 days a week:

- Documentation

The Workload Factory documentation that you're currently viewing.

- [Knowledge base](#)

Search through the Workload Factory knowledge base to find helpful articles to troubleshoot issues.

- [Communities](#)

Join the Workload Factory community to follow ongoing discussions or create new ones.

Create a case with NetApp support

In addition to the self-support options above, you can work with a NetApp Support specialist to resolve any issues after you activate support.

Before you get started

To use the **Create a Case** capability, you must first register for support. Associate your NetApp Support Site credentials with your Workload Factory login. [Learn how to register for support.](#)

Steps

1. In the upper right of the Workload Factory console, select **Help > Support**.

Selecting this option opens the BlueXP console a new browser tab and loads the Support dashboard.

2. On the **Resources** page, choose one of the available options under Technical Support:

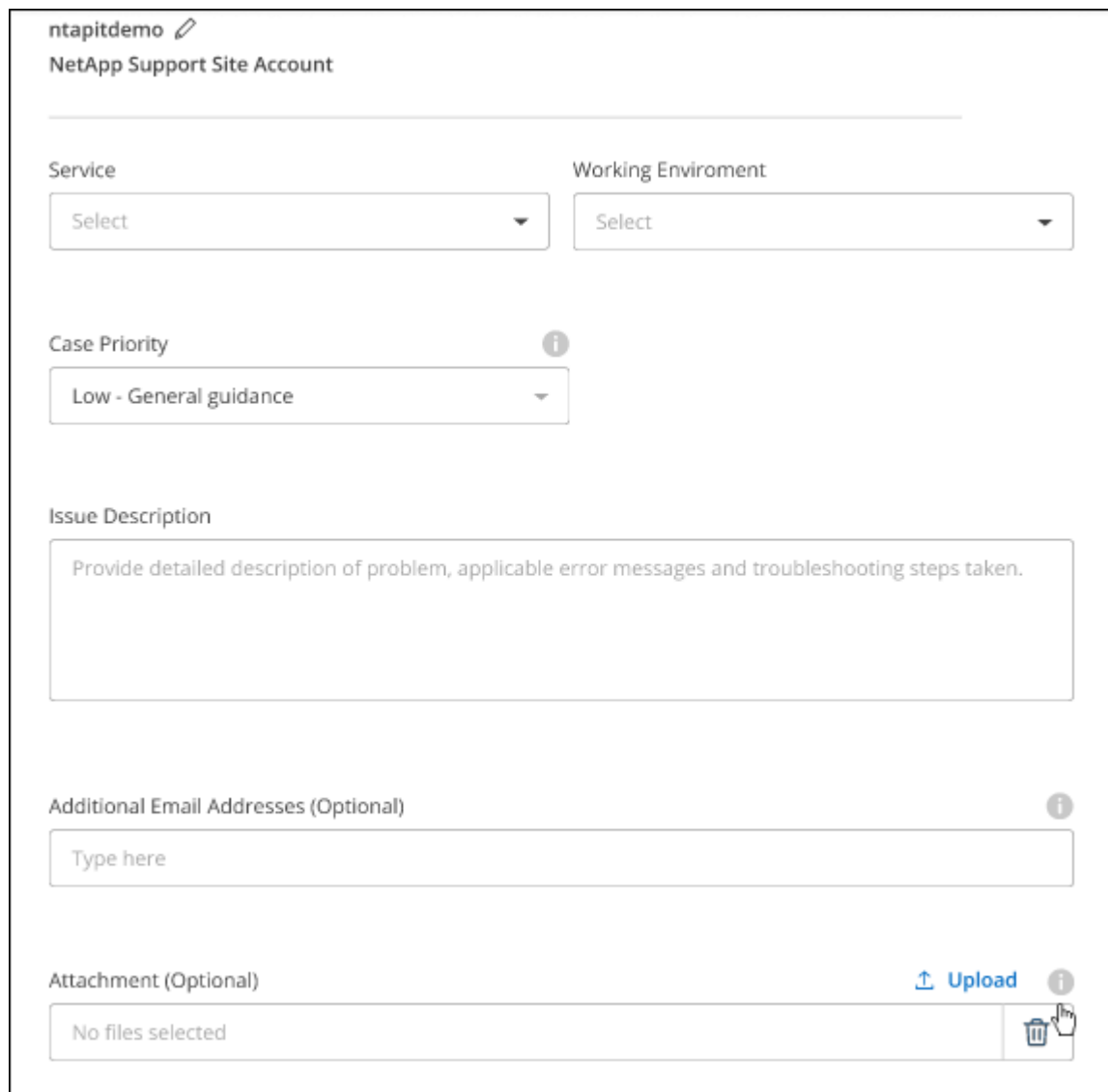
- a. Select **Call Us** if you'd like to speak with someone on the phone. You'll be directed to a page on netapp.com that lists the phone numbers that you can call.
- b. Select **Create a Case** to open a ticket with a NetApp Support specialist:

- **Service:** Select **Workload Factory**.
- **Case Priority:** Choose the priority for the case, which can be Low, Medium, High, or Critical.

To learn more details about these priorities, hover your mouse over the information icon next to the field name.

- **Issue Description:** Provide a detailed description of your problem, including any applicable error messages or troubleshooting steps that you performed.
- **Additional Email Addresses:** Enter additional email addresses if you'd like to make someone else aware of this issue.
- **Attachment (Optional):** Upload up to five attachments, one at a time.

Attachments are limited to 25 MB per file. The following file extensions are supported: txt, log, pdf, jpg/jpeg, rtf, doc/docx, xls/xlsx, and csv.



The screenshot shows a web form titled "NetApp Support Site Account" for a user named "ntapitdemo". The form contains several fields:

- Service:** A dropdown menu with "Select" as the current option.
- Working Environment:** A dropdown menu with "Select" as the current option.
- Case Priority:** A dropdown menu with "Low - General guidance" as the current option. An information icon (i) is located to the right of the label.
- Issue Description:** A large text area with the placeholder text "Provide detailed description of problem, applicable error messages and troubleshooting steps taken."
- Additional Email Addresses (Optional):** A text input field with the placeholder text "Type here". An information icon (i) is located to the right of the label.
- Attachment (Optional):** A file upload area showing "No files selected". It includes an "Upload" button with an upward arrow icon and an information icon (i). A trash can icon is also visible next to the file selection area.

After you finish

A pop-up will appear with your support case number. A NetApp Support specialist will review your case and get back to you soon.

For a history of your support cases, you can select **Settings > Timeline** and look for actions named "create support case." A button to the far right lets you expand the action to see details.

It's possible that you might encounter the following error message when trying to create a case:

"You are not authorized to Create a Case against the selected service"

This error could mean that the NSS account and the company of record it's associated with is not the same company of record for the BlueXP account serial number (ie. 960xxxx) or the working environment serial number. You can seek assistance using one of the following options:

- Use the in-product chat
- Submit a non-technical case at <https://mysupport.netapp.com/site/help>

Manage your support cases (Preview)

You can view and manage active and resolved support cases directly from BlueXP. You can manage the cases associated with your NSS account and with your company.

Case management is available as a Preview. We plan to refine this experience and add enhancements in upcoming releases. Please send us feedback by using the in-product chat.

Note the following:

- The case management dashboard at the top of the page offers two views:
 - The view on the left shows the total cases opened in the past 3 months by the user NSS account you provided.
 - The view on the right shows the total cases opened in the past 3 months at your company level based on your user NSS account.

The results in the table reflect the cases related to the view that you selected.

- You can add or remove columns of interest and you can filter the contents of columns like Priority and Status. Other columns provide just sorting capabilities.

View the steps below for more details.

- At a per-case level, we offer the ability to update case notes or close a case that is not already in Closed or Pending Closed status.

Steps

1. In the upper right of the Workload Factory console, select **Help > Support**.

Selecting this option opens the BlueXP console a new browser tab and loads the Support dashboard.

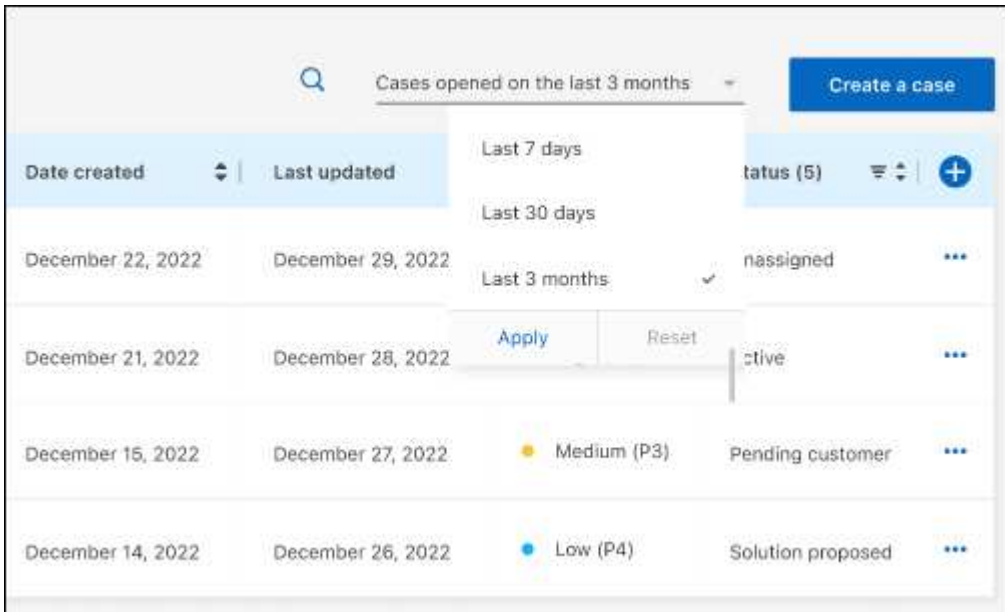
2. Select **Case Management** and if you're prompted, add your NSS account to BlueXP.

The **Case management** page shows open cases related to the NSS account that is associated with your BlueXP user account. This is the same NSS account that appears at the top of the **NSS management** page.

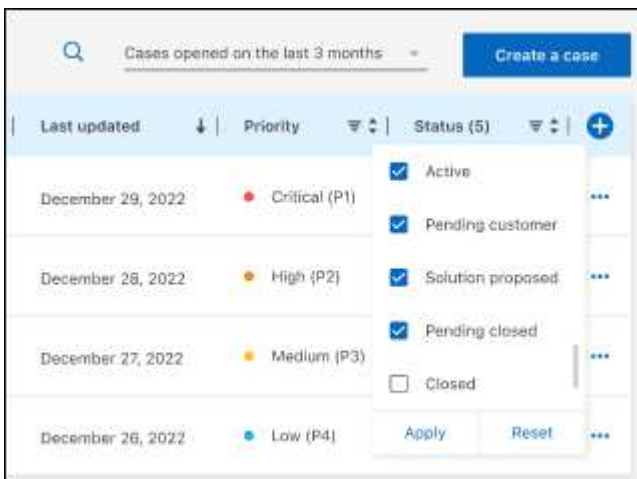
3. Optionally modify the information that displays in the table:


- Under **Organization's cases**, select **View** to view all cases associated with your company.

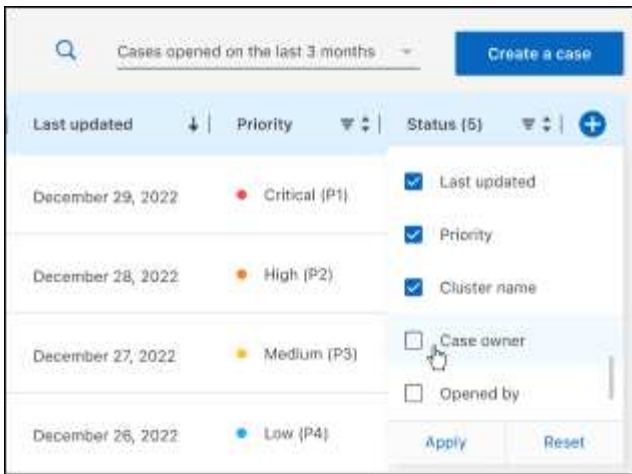
- Modify the date range by choosing an exact date range or by choosing a different time frame.



- Filter the contents of the columns.



- Change the columns that appear in the table by selecting  and then choosing the columns that you'd like to display.

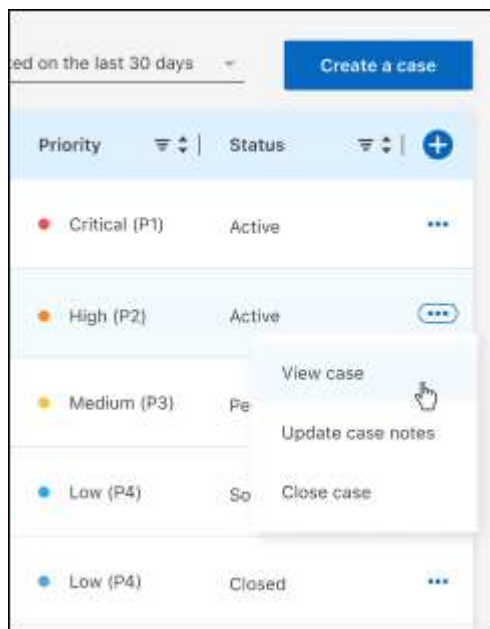


4. Manage an existing case by selecting **...** and selecting one of the available options:

- **View case:** View full details about a specific case.
- **Update case notes:** Provide additional details about your problem or select **Upload files** to attach up to a maximum of five files.

Attachments are limited to 25 MB per file. The following file extensions are supported: txt, log, pdf, jpg/jpeg, rtf, doc/docx, xls/xlsx, and csv.

- **Close case:** Provide details about why you're closing the case and select **Close case**.



Legal notices

Legal notices provide access to copyright statements, trademarks, patents, and more.

Copyright

<https://www.netapp.com/company/legal/copyright/>

Trademarks

NETAPP, the NETAPP logo, and the marks listed on the NetApp Trademarks page are trademarks of NetApp, Inc. Other company and product names may be trademarks of their respective owners.

<https://www.netapp.com/company/legal/trademarks/>

Patents

A current list of NetApp owned patents can be found at:

<https://www.netapp.com/pdf.html?item=/media/11887-patentspage.pdf>

Privacy policy

<https://www.netapp.com/company/legal/privacy-policy/>

Open source

Notice files provide information about third-party copyright and licenses used in NetApp software.

- [Workload Factory](#)
- [Workload Factory for Databases](#)
- [Workload Factory for GenAI](#)
- [Workload Factory for VMware](#)

Copyright information

Copyright © 2024 NetApp, Inc. All Rights Reserved. Printed in the U.S. No part of this document covered by copyright may be reproduced in any form or by any means—graphic, electronic, or mechanical, including photocopying, recording, taping, or storage in an electronic retrieval system—without prior written permission of the copyright owner.

Software derived from copyrighted NetApp material is subject to the following license and disclaimer:

THIS SOFTWARE IS PROVIDED BY NETAPP "AS IS" AND WITHOUT ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE, WHICH ARE HEREBY DISCLAIMED. IN NO EVENT SHALL NETAPP BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

NetApp reserves the right to change any products described herein at any time, and without notice. NetApp assumes no responsibility or liability arising from the use of products described herein, except as expressly agreed to in writing by NetApp. The use or purchase of this product does not convey a license under any patent rights, trademark rights, or any other intellectual property rights of NetApp.

The product described in this manual may be protected by one or more U.S. patents, foreign patents, or pending applications.

LIMITED RIGHTS LEGEND: Use, duplication, or disclosure by the government is subject to restrictions as set forth in subparagraph (b)(3) of the Rights in Technical Data -Noncommercial Items at DFARS 252.227-7013 (FEB 2014) and FAR 52.227-19 (DEC 2007).

Data contained herein pertains to a commercial product and/or commercial service (as defined in FAR 2.101) and is proprietary to NetApp, Inc. All NetApp technical data and computer software provided under this Agreement is commercial in nature and developed solely at private expense. The U.S. Government has a non-exclusive, non-transferrable, nonsublicensable, worldwide, limited irrevocable license to use the Data only in connection with and in support of the U.S. Government contract under which the Data was delivered. Except as provided herein, the Data may not be used, disclosed, reproduced, modified, performed, or displayed without the prior written approval of NetApp, Inc. United States Government license rights for the Department of Defense are limited to those rights identified in DFARS clause 252.227-7015(b) (FEB 2014).

Trademark information

NETAPP, the NETAPP logo, and the marks listed at <http://www.netapp.com/TM> are trademarks of NetApp, Inc. Other company and product names may be trademarks of their respective owners.