

IBM Spectrum Discover
Version 2.0.4

*Concepts, Planning, and Deployment
Guide*



Note

Before using this information and the product it supports, read the information in [“Notices” on page 119](#).

Edition notice

This edition applies to version 2 release 0 modification 4 of the following product, and to all subsequent releases and modifications until otherwise indicated in new editions:

- IBM Spectrum® Discover ordered through Passport Advantage® (product number 5737-I32)
- IBM Spectrum Discover ordered through AAS/eConfig (product number 5641-SG1)

IBM® welcomes your comments; see the topic [“How to send your comments” on page xi](#). When you send information to IBM, you grant IBM a nonexclusive right to use or distribute the information in any way it believes appropriate without incurring any obligation to you.

© **Copyright International Business Machines Corporation 2018, 2020.**

US Government Users Restricted Rights – Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

Contents

Figures.....	v
Tables.....	ix
About this information.....	xi
Prerequisite and related information.....	xi
How to send your comments.....	xi
Summary of changes.....	xiii
Chapter 1. Product overview.....	1
Introduction to IBM Spectrum Discover.....	1
IBM Spectrum Discover architecture.....	2
Role-based access control.....	3
Data source connections.....	4
Cataloging metadata.....	4
Enriching metadata.....	5
Graphical user interface.....	7
Reports for IBM Spectrum Discover.....	9
IBM Spectrum Discover appliance.....	9
Chapter 2. Planning.....	11
Software requirements.....	11
IBM Spectrum Discover deployment models.....	11
CPU and memory requirements for single node trial and single node production IBM Spectrum Discover deployments.....	12
Networking requirements for IBM Spectrum Discover	12
Storage requirements for single node trial and single node production IBM Spectrum Discover deployments.....	13
IBM Spectrum Storage software requirements.....	14
Backup and restore storage requirements for IBM Spectrum Discover.....	16
Single node IBM Spectrum Discover production deployment planning worksheet.....	16
Single node IBM Spectrum Discover trial deployment planning worksheet.....	17
Chapter 3. Deploying and configuring.....	21
Deploy and configure a single node production IBM Spectrum Discover appliance cluster.....	21
Deploying a single node trial or single node production IBM Spectrum Discover virtual appliance.....	21
Configuring storage for a single node trial or single node production of IBM Spectrum Discover virtual appliance.....	27
Configuring CPU and memory allocation for the single node IBM Spectrum Discover virtual appliance.....	36
Configuring networking and perform provisioning of a single node trial or single node production IBM Spectrum Discover virtual appliance.....	40
Known issues with deploying and configuring for single node.....	43
Deploying the IBM Spectrum Discover open virtualization appliance on the Kernel-based Virtual Machine virtualization module.....	43
Configure data source connections.....	44
IBM Spectrum Scale data source connections.....	45

IBM Spectrum Archive data source connections.....	57
IBM Cloud Object Storage data source connection.....	58
IBM Spectrum Discover and S3 object storage data source connections.....	100
Scanning an Elastic Storage Server data source connection.....	105
Creating a Network File System data source connection.....	106
Creating an SMB data source connection.....	107
Creating an IBM Spectrum Protect data source connection.....	108
Editing and using the TimeSinceAccess and Size Range buckets.....	108
Using custom TLS certificate.....	110
Validating code integrity.....	111
Applying the license file.....	111
Chapter 4. Upgrading.....	113
Preparing to run the upgrade tool for IBM Spectrum Discover	113
Running the upgrade tool for IBM Spectrum Discover	114
Upgrading IBM Spectrum Discover to versions 2.0.4 and higher.....	115
Accessibility features for IBM Spectrum Discover.....	117
Accessibility features.....	117
Keyboard navigation.....	117
IBM and accessibility.....	117
Notices.....	119
Trademarks.....	120
Terms and conditions for product documentation.....	120
IBM Online Privacy Statement.....	121
Index.....	123

Figures

1. IBM Spectrum Discover Architecture.....	2
2. Application SDK architecture.....	3
3. Example of the IBM Spectrum Discover dashboard.....	8
4. Single node deployment.....	11
5. ESXi server menu.....	21
6. Deploy OVF template wizard.....	22
7. Virtual appliance location.....	23
8. Physical server location.....	24
9. Select storage window.....	25
10. Select virtual machine network.....	26
11. Review Settings.....	27
12. Edit Settings.....	29
13. Edit hard disk settings.....	30
14. Hard disk options.....	31
15. Edit Settings.....	32
16. Edit settings dialog box.....	33
17. Hard disk options.....	33
18. Selecting the Edit Settings menu option.....	35
19. Selecting options in the Edit settings dialog box.....	35
20. Selecting hard disk options.....	36
21. Edit settings menu.....	37
22. CPU allocation settings.....	38
23. Memory allocation settings.....	39

24. Reserved memory allocation settings.....	40
25. vSphere client settings menu.....	41
26. Example of the system advanced configuration.....	58
27. Displaying the source names for data source connections.....	59
28. Example of window that shows Data Connections Add data source Connection.....	60
29. Example of the screen for an IBMCOS connection.....	61
30. Data source connections.....	63
31. Selecting a data source connection to scan.....	63
32. Active scans.....	64
33. IBM Cloud Object Storage Replay architecture.....	66
34. Python process count.....	76
35. Settings for three items on the vault configuration page in the net Manager user interface.....	77
36. IBM Cloud Object Storage Scanner progress report.....	85
37. Directory structure from the configuration file.....	89
38. Example of running in debug mode.....	89
39. Example of a log file.....	91
40. Scanner debug.....	92
41. Scanner debug (continued).....	93
42. Scanner debug (continued).....	94
43. Scanner debug (continued).....	95
44. Configurations.....	96
45. Add a storage vault to the configuration.....	98
46. Total number of indexed records.....	100
47. Displaying the source names for Data Source Connections.....	101
48. Example of the Add Data Source Connection GUI window.....	101

49. Selecting S3 Object Storage.....	102
50. Completing the S3 information fields.....	102
51. Data source connections.....	103
52. Selecting a data source connection to scan.....	104
53. Active scans.....	104
54. Example of how to define the settings for a SizeRange bucket.....	109
55. Example of how to modify and define the settings of a bucket that is older than one year old.....	110

Tables

1. IBM Spectrum Discover library information units.....	xi
2. Benefits of IBM Spectrum Discover.....	1
3. Virtual resources for the virtual appliance.....	9
4. Browser requirements for the IBM Spectrum Discover GUI.....	11
5. CPU and memory requirements for single node production.....	12
6. CPU and memory requirements for single node trial.....	12
7. Network parameter example.....	13
8. Storage requirements for single node production.....	13
9. Storage requirements for single node trial.....	14
10. IBM Cloud Object Storage software requirements.....	14
11. IBM Spectrum Scale software requirements.....	15
12. IBM Spectrum Protect software requirements.....	15
13. IBM Spectrum Archive software requirements.....	15
14. Single node IBM Spectrum Discover production deployment planning.....	16
15. Single node IBM Spectrum Discover trial deployment planning.....	17
16. Explanation of the configuration file.....	70
17. Behaviors for Replay for four variables.....	77
18. Description for IBM Cloud Object Storage Scanner progress report.....	85
19. What is reported beneath the report title.....	86
20. List of directories generated by scanner, notifier, and replay.....	87
21. Leaf directory file names.....	89
22. Examples of size ranges and sizes of buckets with user-defined labels.....	109

About this information

IBM Spectrum Discover is a metadata-driven management system for large-scale file and object environments. IBM Spectrum Discover maintains a real-time metadata repository for large-scale enterprise storage environments. Metadata can be searched, enhanced, discovered, and leveraged for data processing by using built-in or custom agents.

IBM Spectrum Discover - Information units

Table 1. IBM Spectrum Discover library information units		
Information unit	Type of information	Intended users
IBM Spectrum Discover: Concepts, Planning, and Deployment Guide	This information unit provides information about the following topics: <ul style="list-style-type: none">• Product Overview• Planning• Deploying and configuring	Users, system administrators, analysts, installers, planners, and programmers of IBM Spectrum Discover.
IBM Spectrum Discover: Administration Guide	This information unit provides information about administration, monitoring, and troubleshooting tasks.	Users, system administrators, analysts, installers, planners, and programmers of IBM Spectrum Discover.
IBM Spectrum Discover: REST API Guide	This information unit provides information about the following topics: <ul style="list-style-type: none">• IBM Spectrum Discover REST APIs• Endpoints for working with a DB2[®] warehouse• Endpoints for working with policy management• Endpoints for working with connection management• Action agent management using APIs• RBAC management using APIs	Users, system administrators, analysts, installers, planners, and programmers of IBM Spectrum Discover.

Prerequisite and related information

For updates to this information, see IBM Spectrum Discover in IBM Knowledge Center (<https://www.ibm.com/support/knowledgecenter/SSY8AC>).

How to send your comments

You can add your comments in IBM Knowledge Center. To add comments directly in IBM Knowledge Center, you need to log in with your IBM ID.

You can also send your comments to ibmkc@us.ibm.com.

Summary of changes

The summary of changes compiles a list of changes that are implemented in the IBM Spectrum Discover licensed program and the IBM Spectrum Discover library. Within each topic, these markers (□) surrounding text or illustrations indicate technical changes or additions that are made to the previous edition of the information.

[

Summary of changes for IBM Spectrum Discover version 2.0.4 as updated, December 2020

This release of the IBM Spectrum Discover licensed program and the IBM Spectrum Discover library includes the following improvements. All improvements are available after an upgrade, unless otherwise specified.

Administering

- Search is enhanced with the ability to build visual query for an improved experience. You can also build custom queries to meet your specific search requirements. For more information, see the topic *Searching* in the *IBM Spectrum Discover: Administration Guide*.
- With the IBM Spectrum Discover Import Tags application you can apply a pre-curated set of tags from external CSV file to S3/COS object records. For more information, see the topic *Importing externally curated tags for COS/S3 using the import tags application* in the *IBM Spectrum Discover: Administration Guide*.
- With ScaleAFM application you can copy files/objects between IBM Spectrum Scale and IBM Cloud® Object Storage connections. For more information, see *Copying data using ScaleAFM application* in the *IBM Spectrum Discover: Administration Guide*.
- With integration of IBM Spectrum Discover and Moonwalk, the application automates and manages the movement of data from primary storage locations to lower-cost file systems, object stores or cloud storage services. For more information, see *Data movement with IBM Spectrum Discover and Moonwalk* in the *IBM Spectrum Discover: Administration Guide*.
- The configuration of the Watson™ Knowledge Catalog connector app has been simplified. For more information, see the topic *Configuring Watson Knowledge Catalog connector* in the *IBM Spectrum Discover: Administration Guide*.

GUI changes

A new interface allows you to view policy execution statistics, policy history details and debug log data.

REST API changes

The following REST API documentation is added:

- How to retrieve policy execution history details with: `/policyengine/v1/policyhistory:GET`
- How retrieve policy run log details specific policies with: `/policyengine/v1/policyhistory/<pol_id>/<log_id>: GET`
- How to delete policy execution history and logs with: `/policyengine/v1/policyhistory:DELETE`

The following REST API documentation command requests and responses are updated:

- `api/application/appcatalog/help:DELETE`
- `api/application/appcatalog/help:GET`
- `api/application/appcatalog/help:PATCH`

- api/application/appcatalog/help:POST

]

Chapter 1. Product overview

Introduction to IBM Spectrum Discover

Companies need the ability to use unstructured data to meet their business priorities.

IBM Spectrum Discover is a modern metadata management software that provides data insight for exabyte-scale heterogeneous file, object, backup, and archive storage on premises and in the cloud. The software easily connects to these data sources to rapidly ingest, consolidate, and index metadata for billions of files and objects.

IBM Spectrum Discover provides a rich metadata layer that enables storage administrators, data stewards, and data scientists to efficiently manage, classify, and gain insights from massive amounts of data. It improves storage economics, helps mitigate risk, and accelerates large-scale analytics to create competitive advantage and speed critical research.

Many companies face significant challenges to manage their data. Some difficult challenges that companies face include:

- Pinpointing and activating relevant data for large-scale analytics.
- Lacking the fine-grained visibility needed to map data to business priorities.
- Removing redundant, trivial, and obsolete data.
- Identifying and classifying sensitive data.

Benefits of IBM Spectrum Discover

IBM Spectrum Discover can help you manage your unstructured data by reducing the data storage costs, uncovering hidden data value, and reducing the risk of massive data stores. See [Table 2 on page 1](#).

Table 2. Benefits of IBM Spectrum Discover			
Optimize - Improve storage usage	Analyze - Uncover hidden data value	Govern - Mitigate risk and improve data quality	Data Management
Decreases storage capital expenditure (CaPex) by facilitating data movement to colder, cheaper storage.	Accelerates data identification for large-scale analytics.	Perform data inspection and classification.	Automate tags for custom insight.
Increases storage efficiency by eliminating trivial or redundant data.	Operationalize tasks to reduce the burden of data preparation.	Helps ensure that data is compliant with governance policies by labeling sensitive data.	Create reports for analysis.
Reduces storage operating expenditure (OpEx) by improving storage administrator productivity.	Orchestrates the ML/DL and Platform Symphony® MapReduce process.	Helps reduce risk that is hidden in heterogeneous data sources.	GUI search for real-time results Search content for fast discovery.

IBM Spectrum Discover architecture

IBM Spectrum Discover is an extensible platform that provides exabyte scale data ingest, data visualization, data activation, and business-oriented data mapping.

Exabyte-scale data ingest

- Scan billions of files and objects in a day
- Real-time event notifications
- Automatic indexing

Data Visualization

- Fast queries of billions of records
- Multi-faceted search
- Drilldown Dashboard

Data Activation

- Application software development kit (SDK)
- Extensible architecture
- Solution blueprints

Business-oriented data mapping

- System-level data tagging
- Contextual data tagging
- Policy-driven workflows

The following figure illustrates a high-level view of the IBM Spectrum Discover architecture.

IBM Spectrum Discover Overview

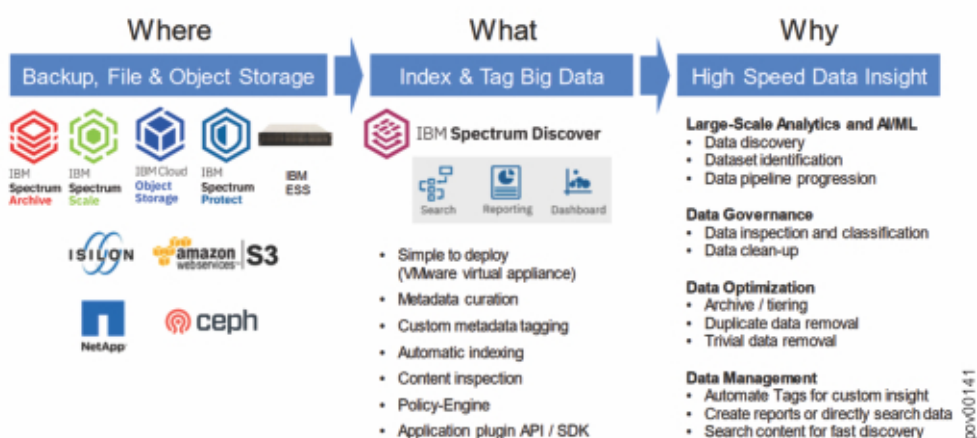


Figure 1. IBM Spectrum Discover Architecture

Data management

IBM Spectrum Discover connects to the data sources shown in the architecture image (*IBM Spectrum Discover architecture*), and automatically harvests and indexes the system metadata where the system metadata refers to certain information. This might include the following information.

- It might include the names of the files and objects.

- It might include the bucket or path where the data resides.
- It might include the size.
- It might include the time the data sources were last modified.

After the data is ingested, analytics are automatically applied to classify and group the data according to the different system metadata attributes. The data can be inspected automatically in IBM Spectrum Discover by using built-in content search capabilities to identify sensitive and personally identifiable information and perform data classification. The content inspection capabilities can also be used by researchers and data scientists to extract content from their data sets. This easy-to-use extraction ability assists with data discovery.

The records that are maintained by IBM Spectrum Discover can also be further enriched with custom metadata tags that map the data to business constructs and further increase the value of the data.

You can use IBM Spectrum Discover catalog to gain insight about your data and to find your data easily.

The IBM Spectrum Discover architecture also supports a community-supported catalog of open source applications that enhance and customize the capabilities of IBM Spectrum Discover with third-party extensions. Users can find and install available applications and can develop and share new applications that use an SDK that contains sample code and a fully published API. For more information, see the topic *Creating your own applications to use in the IBM Spectrum Discover application catalog* in the *IBM Spectrum Discover: Administration Guide*.

Figure 2 on page 3 shows an example of the Application SDK architecture.

Extensible Foundation for Data Insight

- Action Agent SDK extends capabilities via well defined API
- Customize actions taken based on Discover metadata
 - ❖ Content indexing
 - ❖ Data movement (tiering)
 - ❖ Classification
 - ❖ Sensitive data identification
 - ❖ RDT Detection/Disposal
 - ❖ Etc...
- Integrate with upstream information management applications

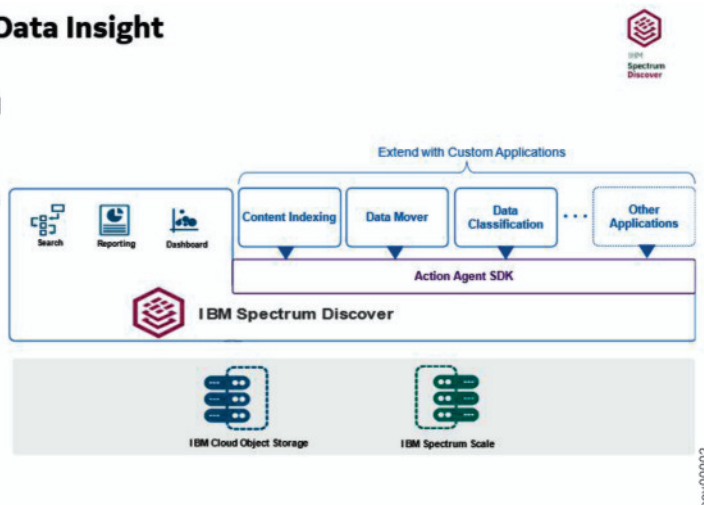


Figure 2. Application SDK architecture

Role-based access control

IBM Spectrum Discover provides access to resources based on roles. You can restrict access to information based on roles.

The role that is assigned to a user or group determines the privileges for that user or group. Users and groups can be associated with collections, which use policies that determine the metadata that is available to view.

User and group access can be authenticated by IBM Spectrum Discover, an LDAP server, or the IBM Cloud Object Storage System. The administrator can manage the user access functions.

Roles

Roles determine how users and groups access records or the IBM Spectrum Discover environment.

Remember: If a user or group is assigned to multiple roles, the least restrictive role is applicable.

For example, if you are assigned the role of **Data User**, and you are also assigned the role of a **Data Admin**, you have the privileges of a **Data Admin**.

Admin

An Admin can create users, groups, collections, manage LDAP, and IBM Cloud Object Storage connections for user access management.

Data Admin

Users with the **Data Admin** role can access all metadata that is collected by IBM Spectrum Discover and is not restricted by collections.

Collection Admin

The **Collection Admin** role is as a bridge between the **Data Admin** role and the **Data User** role:

- Users with the **Collection Admin** role can list any type of tag and create or modify Characteristic tags. Users with the **Collection Admin** role cannot create, modify, or delete Open and Restricted tags. These permissions are the same permissions as the **Data User** role.

Note: The built-in Collection tag is a special tag that can be set only by users with the **Data Admin** role. All other tags can be set by any user with the **Data User** or **Data Admin** or **Collection Admin** role.

- Users with the **Collection Admin** role can
 - Create, update, and delete the policies for the collections they administer.
 - View, update, and delete policies of data users for the collections they administer. They cannot delete a policy if it has a collection that they do not administer.
 - Add users to collections that they administer. These data users can access a particular collection, which means that they can access the records that are marked with that collection value.

Data User

Users with the **Data User** role can access metadata that is collected by IBM Spectrum Discover. Metadata access might be restricted by policies in the collections that are assigned to users in this role. A user with the **Data User** role can also define tags and policies based on the collections to which the role is assigned.

Service User

The **Service User** role is assigned to accounts for IBM service and support personnel.

Data source connections

A data source connection specifies the parameters for cataloging of metadata from a source system to IBM Spectrum Discover.

Without the proper connection information, ingesting metadata from a connected system fails. You can use the **Data Source Connections** page to view connection information for the data sources that are connected to your environment.

For more information, see [“Configure data source connections” on page 44](#).

Cataloging metadata

System metadata is created and updated by the host system, and not the application software. IBM Spectrum Discover allows the addition of tags that can capture non-system metadata-specific attributes, which are stored in the IBM Spectrum Discover catalog.

Scans are jobs that are scheduled or on demand, and occur at a data source level. For example, a file system or object vault. A set of metadata records is generated with each record that captures the state of an individual file or object within the data source at the time of the scan.

IBM Spectrum Discover supports scanning the following data sources:

- IBM Spectrum Scale

- IBM Elastic Storage® Server
- IBM Cloud Object Storage
- IBM Spectrum Protect
- IBM Spectrum Archive
- Red Hat® Ceph® Storage
- NetApp Storage Solutions
- Dell EMC Isilon Scale Out Network Attached Storage
- Amazon Simple Storage Service (Amazon S3)

Live event notifications are triggered by user actions on the source data. Examples are reading, writing, moving, deleting data, changing permissions, or ownership. The events generate a metadata record in real time that is stored in IBM Spectrum Discover. IBM Spectrum Discover supports live event notifications for the following data sources:

- IBM Cloud Object Storage
- IBM Elastic Storage Server
- IBM Spectrum Scale
- Red Hat Ceph Storage

With IBM Spectrum Scale you can enable live events to start a watch folder on the specified file system. The IBM Spectrum Scale watch folder works with IBM Spectrum Discover to capture file system event notifications and deliver them to IBM Spectrum Discover by using Kafka.

Important: If the connection from IBM Spectrum Scale to IBM Spectrum Discover is interrupted, the watch suspends. Additionally, events are no longer captured in IBM Spectrum Discover, which requires a file system rescan to capture the lost updates.

Enriching metadata

IBM Spectrum Discover can enrich the metadata from supported platforms with additional information by using policies, content inspection, custom tags, and custom applications.

Policies

Policies are used to add additional information about the source data that is indexed in IBM Spectrum Discover. A policy determines the set of files to add tag values to, or to send to the built-in content inspection capabilities of IBM Spectrum Discover, or to a custom application through filtering criteria. The policies give you the ability to run actions one time or on a set schedule. Policies work in batches and can be paused, resumed, stopped, or restarted. You can control the load on the IBM Spectrum Discover system or source storage system for content inspection policies and policies that start custom applications.

Applications

A deep inspect application extracts information from source data records and returns it to IBM Spectrum Discover to be indexed. For example, by using a custom application, you might create a DEEP-INSPECT policy to extract key characteristics from files of a certain type. The characteristics are applied to the metadata records for the files in IBM Spectrum Discover as custom tags and made searchable. You can search for data by name, size, and content.

You can use custom applications to extend the capabilities that are performed by IBM Spectrum Discover.

Policy engine

Policies offer a method whereby you can schedule one-time or repetitive actions on a filtered set of records.

The policy management API service is a RESTful web service that is designed to create, list, update, and delete policies. You can use a policy to initiate action on a select set of indexed documents or data. You can do a task immediately or on a set schedule.

Several types of policies that are supported by IBM Spectrum Discover enrich the metadata records. You can create policies with information to determine which set of documents to run, the action to take, and when to run policies periodically.

A policy includes

Policy ID

Name of the policy.

Filter

Selects a set of documents to work.

Action

ID, parameters, and schedule.

The following list is a description of the policies.

AUTOTAG

A policy that tags a set of records based on filter criteria with a pre-defined set of tags.

CONTENT SEARCH

A policy that uses the built-in content inspection capabilities of IBM Spectrum Discover to extract content from source data and index it automatically into the IBM Spectrum Discover catalog.

DEEP-INSPECT

A policy that passes lists of files based on filter criteria to the analytics application that opens the source data file and extracts metadata information from it. The policy passes the data back to IBM Spectrum Discover in the form of tags so you can do a search, and do the following activities:

- Set up a filter to do a search query that finds the candidates to apply the policy.
For example, you can set an action for filtered candidates AUTOTAG: tag1: value, tag2: value
- Set a schedule to apply the policy by specifying the following methods:
 - Immediately
 - Periodically

Applications

IBM Spectrum Discover policies might contain applications in the actions parameters.

Use an application when you want to do a specific action on data or metadata on IBM Spectrum Discover.

You can define an application when you create a new DEEP-INSPECT policy. You can add parameters for an application during the process of creating a DEEP-INSPECT policy.

When you open the window for applications, you can see a view of a table with the following information:

Application

The name of the application.

Parameters

The parameters that were assigned to the application when the policy was created.

Action ID

Actions that are supported by the application for enriching metadata. For example, CONTENTSEARCH, DEEP-INSPECT.

View or Delete

Use the delete trashcan icon to remove the application from the database.

Graphical user interface

The IBM Spectrum Discover graphical user interface is a portal that is used for running data searches, report generation, policy and tag management, and user Access Management. Based on a user's role, they might have access to one or more of these areas.

The IBM Spectrum Discover environment provides access to users and groups. The role that is assigned to a user or group determines the functions that are available. Users and groups can also be associated with collections, which use policies that determine the metadata that is available to view.

User and group access can be authenticated by IBM Spectrum Discover, an LDAP server, or the IBM Cloud Object Storage System. The administrator can manage the user access functions.

Roles

Roles determine how users and groups can access records or the IBM Spectrum Discover environment.

If a user or group is assigned to multiple roles, the least restrictive role is used. For example, if a user is assigned a role of Data User, and is included in a data administrator role, the user has the privileges of a data administrator.

Dashboard

An example of the IBM Spectrum Discover dashboard is shown.

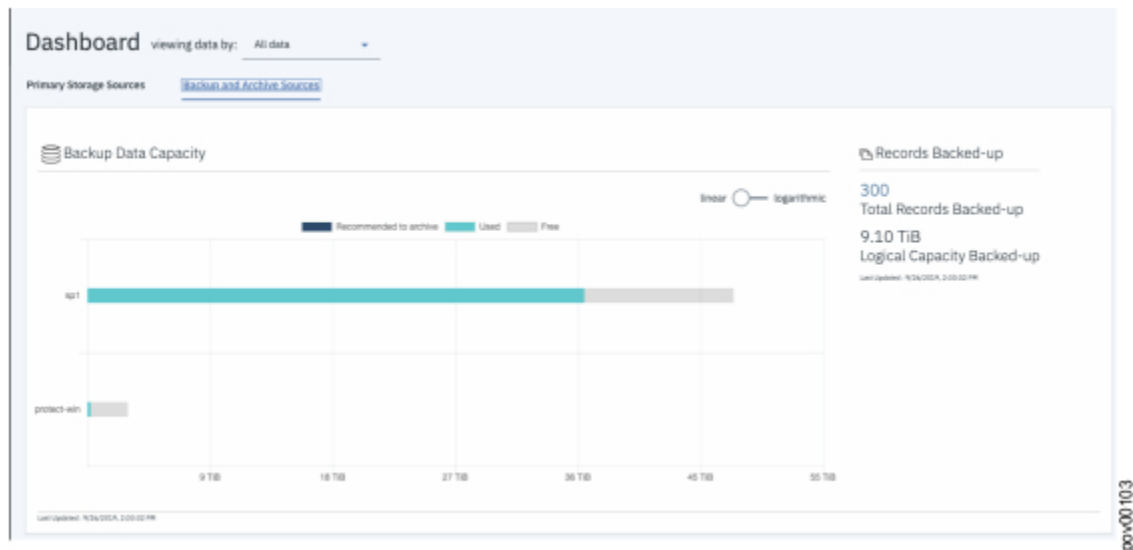
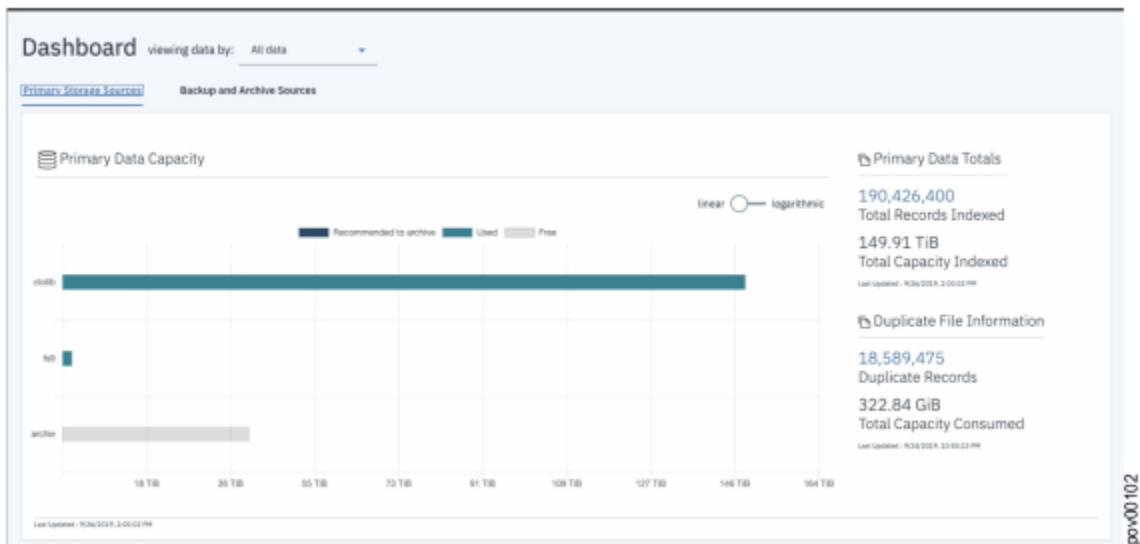


Figure 3. Example of the IBM Spectrum Discover dashboard

Data administrators and users can view the following:

- Metrics for the overall capacity used by every data source
- Total number of files
- Amount of capacity that is used by records with specific tags and facets, for example, owner, cluster, and size range
- Distribution of those records across data sources

Users can click any of the dashboard widgets to initiate a search and further explore and drill down into the data. Administrators and user can also perform the following:

- Monitor storage usage and data recommendations
- View total indexed data and capacity
- View duplicate file or object candidates. For example:
 - Number
 - Capacity used
- Preview capacity use by data facet - for example:

- Classification
- Owner
- File type
- View data capacity by group or collection - for example:
 - Customer defined
 - Lab or project

Understanding size and capacity differences

IBM Spectrum Discover collects size and capacity information. Generally:

- Size refers to the size of a file or object in bytes.
- Capacity refers to the amount of space the file or object consumes on the source storage in bytes.

For objects, size and capacity values always match. For files, size and capacity values can be different because of file system block overhead or sparsely populated files.

Note: Storage protection overhead (such as RAID values or erasure coding) and replication overhead are not captured in the capacity values.

Reports for IBM Spectrum Discover

Reports for IBM Spectrum Discover are grouped or non-grouped. Grouped reports have information for count and sum in columns and non-grouped reports have information in rows.

Data Curation Reports are a way for administrators, also known as data curators, to view the state of their storage environment in different ways. They can range from high-level grouped information to individual record level information.

For example, you can sort a report by owner, project, and department, or you can generate a list of records that meet a specific criterion. Additionally, you can create a report that lists the records in a project, not been reviewed for over a year. The owner of the data can evaluate whether to archive or delete the report.

For more information, see *Reports* in *IBM Spectrum Discover: Administration Guide*.

For more information, see the topic *REST API* in the *IBM Spectrum Discover: REST API Guide*.

IBM Spectrum Discover appliance

The virtual appliance is a virtual machine in Open Virtualization Format (OVF) format that you can download and includes the IBM Spectrum Discover.

The IBM Spectrum Discover virtual appliance is bundled as Open Virtualization Appliance (OVA) image to be deployed on VMware vSphere 6.0 or later. vSphere is VMware's hypervisor platform that is designed to manage large pools of virtualized computing infrastructure that includes software and hardware. Virtual appliance deployments use VMware's ESXi hypervisor architecture (6.5 or later).

The IBM Spectrum Discover virtual appliance cluster is automatically configured according to the input the user provides at the initial configuration console.

Each IBM Spectrum Discover virtual appliance is configured with the virtual resources.

Table 3. Virtual resources for the virtual appliance	
Component	Value
RAM (GB)	64
CPU	16

<i>Table 3. Virtual resources for the virtual appliance (continued)</i>	
Component	Value
For the ESX server, SCSI controller 0 is listed as LSI Logic. The second SCSI controller is listed as LSI Logic SAS.	VMware para virtual
Hard disk 1 Note: Three virtual disks (VMDK) are required including the disk that is created when installing the appliance from the OVA.	500 GB
Network adapter	VM network

Chapter 2. Planning

Software requirements

Virtual appliance specifications to use IBM Spectrum Discover at your site are as follows:
IBM Spectrum Discover is bundled as an Open Virtualization Appliance (OVA) image to be deployed on VMware vSphere (6.5 or later).

Table 4. Browser requirements for the IBM Spectrum Discover GUI	
Browser	Version
Google Chrome	67 and higher
Firefox	60 ESR and higher ESR releases
Microsoft Edge	All versions

IBM Spectrum Discover deployment models

IBM Spectrum Discover can be deployed by using a single node.
The diagram illustrates the single node deployment of IBM Spectrum Discover.



Figure 4. Single node deployment

CPU and memory requirements for single node trial and single node production IBM Spectrum Discover deployments

A description of the CPU and memory requirements for single node trial and single node production IBM Spectrum Discover deployment.

The following table shows the CPU and memory requirements for a single node production IBM Spectrum Discover deployment.

Table 5. CPU and memory requirements for single node production	
Specification	Value
Memory	128 GB
Logical processors	24

The following table shows the recommended CPU and memory for a single node trial IBM Spectrum Discover deployment.

Note: Single node trial deployments with less than the recommended value of memory and logical processors will not be able to scale to index two billion documents.

Table 6. CPU and memory requirements for single node trial		
Specification	Minimum value	Recommended value
Memory	64 GB	128 GB
Logical processors	[16]	24

Note: If using 64GB of RAM, no more than 25 million files can be indexed into IBM Spectrum Discover.

Networking requirements for IBM Spectrum Discover

IBM Spectrum Discover requires the following network parameters.

- Host name
- Virtual interface identifier
- IP address
- Netmask
- Gateway
- Domain Name Server (DNS) IP or host name
- Network Time Protocol (NTP) server IP or host name

Note: IBM Spectrum Discover requires a Fully Qualified Domain Name (FQDN) that is registered in a customer supplied DNS. The customer supplied FQDN must be resolvable by the customer supplied DNS from the IBM Spectrum Discover node in order for the IBM Spectrum Discover virtual appliance to operate properly.

The minimum recommended bandwidth for the network bandwidth is 1 GbE (Gigabit Ethernet) if application processing is not performed. If applications are used, the minimum recommended bandwidth is 10 GbE.

Note: The IBM Spectrum Discover nodes must be able to communicate with a customer supplied NTP server to operate properly.

Table 7. Network parameter example			
Parameter	Value Format	Recommended Value	Example
Host name	host.domain.com	FQDN of the node	node1234.example.com
Interface	ensXXX	The Ethernet interface to use for the virtual appliance networking	ens192
IP address	xxx.xxx.xxx.xxx	The IP address of the node	10.10.200.10
Netmask	xxx.xxx.xxx.xxx	Network mask for the IP range of the node	255.255.255.0
Gateway	xxx.xxx.xxx.xxx	IP address of the network gateway	10.10.200.1
DNS	xxx.xxx.xxx.xxx	The IP address of a single DNS	10.10.200.35
NTP	xxx.xxx.xxx.xxx or host.domain.com	FQDN or IP address of NTP server.	10.10.10.2 or Pool11.ntp.org

Storage requirements for single node trial and single node production IBM Spectrum Discover deployments

This topic describes the storage requirements when you are using IBM Spectrum Discover as a single node trial deployment or a single node production deployment.

The single node IBM Spectrum Discover production appliance requires a 500 GB RAID protected SSD or flash virtual machine disk (VMDK) storage device for the operating system and base software. The VMDK must be thick-provisioned and lazy-zeroed.

The single node production IBM Spectrum Discover virtual appliance requires an extra RAID protected SSD or flash VMDK storage device for the persistent message queue. The storage device for the persistent message queue can be locally attached storage or SAN-attached shared storage. The VMDK must be thick-provisioned and lazy-zeroed. If an optional application is installed in the IBM Spectrum Discover node, more storage capacity must be allocated for the VMDK storage device for the persistent message queue.

The single node production IBM Spectrum Discover virtual appliance also requires an extra RAID-protected SSD or flash virtual machine disk (VMDK) storage device for the database.

You must add the persistent message queue and database VMDK storage devices to the IBM Spectrum Discover virtual appliance as part of the configuration process.

The following table shows the storage requirements for a single node production IBM Spectrum Discover deployment that supports indexing the metadata for up to 2 billion files and objects.

Table 8. Storage requirements for single node production		
Use	Storage type	Size
Base OS and Software	Thick-provision and lazy-zero SSD or flash VMDK	500 GB
Persistent message queue	Thick-provisioned and lazy-zero SSD or flash VMDK	700 GB
Database	Thick-provision and lazy-zero SSD or flash VMDK	2.5 TB

Table 8. Storage requirements for single node production (continued)		
Use	Storage type	Size
[Database backup]	[Thick-provision and lazy-zero SSD or flash VMDK]	[2.5 TB]

For single node non-production trial versions of IBM Spectrum Discover, a 500 GB RAID-protected HDD, SSD, or flash VMDK storage device is required for the operating system and base software. The VMDK must be thick-provisioned and lazy-zeroed.

The single node non-production trial version of the IBM Spectrum Discover virtual appliance requires an extra RAID-protected HDD, SSD, or flash virtual machine disk (VMDK) storage device for the persistent message queue. If an optional application is installed in the IBM Spectrum Discover node, more storage capacity must be allocated for the VMDK storage device for the persistent message queue.

An extra RAID-protected SSD or flash VMDK storage device is required for the database.

You must add the persistent message queue and database VMDK storage devices to the IBM Spectrum Discover virtual appliance as part of the configuration process. The two extra storage devices might be smaller in size than a single node production deployment if less than 2 B records are indexed into the system. The following table shows the storage requirements.

Table 9. Storage requirements for single node trial		
Use	Storage type	Size
Base OS and Software	Thick-provision and lazy-zero HDD or SSD or flash VMDK	500 GB
Persistent message queue	Thick-provision and lazy-zero HDD or SSD or flash VMDK	50 GB minimum, 2 GB per 20 million indexed files and/or objects.
Database	Thick-provision and lazy-zero SSD or flash VMDK	100 GB minimum, [1] GB per 2 million indexed files and/or objects.
Database [backup]	Thick-provision and lazy-zero SSD or flash VMDK	100 GB minimum, 1 GB per 2 million indexed files and/or objects.

IBM Spectrum Storage software requirements

Use IBM Spectrum Discover to index metadata from other applications and to orchestrate the data management.

IBM Cloud Object Storage

IBM Spectrum Discover indexes metadata from IBM Cloud Object Storage by receiving notifications that contain metadata from IBM Cloud Object Storage. IBM Spectrum Discover also supports scanning IBM Cloud Object Storage to harvest metadata.

The following table shows the minimum required IBM Cloud Object Storage software version to enable metadata harvesting with IBM Spectrum Discover:

Table 10. IBM Cloud Object Storage software requirements	
Component	Version
IBM Cloud Object Storage	3.14.0 and higher

IBM Spectrum Scale

IBM Spectrum Discover indexes metadata from IBM Spectrum Scale by scanning IBM Spectrum Scale file systems. The IBM Spectrum Scale watch folders technology preview also enables IBM Spectrum Scale to send events that contain metadata to IBM Spectrum Discover.

The following table lists the minimum required IBM Spectrum Scale software versions:

Table 11. IBM Spectrum Scale software requirements		
Component	Feature	Version
IBM Spectrum Scale	Scanning	4.2.3.x and higher
IBM Spectrum Scale	Live events	(Advanced and Data Management Editions, only) 5.0.4.1 and higher
IBM Spectrum Scale	[Data management using ScaleAFM Application]	[5.1 and higher]

There are requirements for enabling live events, which include:

- You must use IBM Spectrum Scale Version 5.0.3.x. Due to an IBM Spectrum Scale performance issue that might result in the unexpected suspension of the IBM Spectrum Scale watch, IBM Spectrum Discover recommends the use of IBM Spectrum Scale Version 5.0.4.x.
- Watch folder must be enabled for the Scale cluster.
- A minimum of three nodes on the IBM Spectrum Scale cluster are required to act as Kafka brokers.
- The IBM Spectrum Scale nodes that act as brokers must meet a minimum local space requirement of 20 GB each to successfully enable the watch with a secondary sink.

IBM Spectrum Protect

IBM Spectrum Discover indexes metadata from IBM Spectrum Protect by scanning IBM Spectrum Protect file systems.

The following table lists the minimum required IBM Spectrum Protect software versions to enable metadata harvesting with IBM Spectrum Discover:

Table 12. IBM Spectrum Protect software requirements	
Component	Version
IBM Spectrum Protect	7.x and higher

IBM Spectrum Archive

IBM Spectrum Discover supports the advanced tiering function with the ScaleILM application.

The following table lists the minimum IBM Spectrum Archive software version that is required to control the data placement by IBM Spectrum Discover.

Table 13. IBM Spectrum Archive software requirements	
Component	Version
IBM Spectrum Archive Enterprise Edition (EE)	1.3.0.6 and higher

Backup and restore storage requirements for IBM Spectrum Discover

IBM Spectrum Discover provides a set of scripts for safely backing up and restoring the metadata database and file system.

The script integrates with the following backup targets:

- IBM Cloud Object Storage
- IBM Spectrum Protect
- External FTP server

The size of the backup pool for the backup targets is determined by taking the size of the backup staging pool and multiplying it by the number of backups that are kept as part of the retention policy.

Note: The backup that you use to restore an IBM Spectrum Discover system must be at the same code level as the IBM Spectrum Discover system that is being restored. For example, you must be restoring a 2.0.2.1 system if you want to use a 2.0.2.1 backup.

Example:

Single node backup staging pool = 2 TB

Number of backups = 7

Backup target capacity required = 2 TB x 7 = 14 TB

Single node IBM Spectrum Discover production deployment planning worksheet

Use this worksheet to plan for installing IBM Spectrum Discover for a single node production deployment.

Table 14. Single node IBM Spectrum Discover production deployment planning				
CPU and memory requirements				
Parameter		Recommended value		Record your values
Memory		128 GB		
Logical processor count		24 logical processors		
Networking requirements				
Parameter	Value format	Recommended value	Example	Record your values
<hostname>	host.domain.com	Specify the fully-qualified domain name of the node.	node.example.com	
<interface>	ensXXX	Specify the Ethernet interface to use for the virtual appliance networking.	ens192	
<ip>	xxx.xxx.xx x.xxx	Specify the IP address of the node.	10.10.200.10	

Table 14. Single node IBM Spectrum Discover production deployment planning (continued)

CPU and memory requirements				
Parameter		Recommended value		Record your values
<netmask>	xxx.xxx.xx x.xxx	Specify the network mask for the IP range of the node.	255.255.254.0	
<gateway>	xxx.xxx.xx x.xxx	Specify the IP address of the network gateway.	10.10.200.1	
<dns>	xxx.xxx.xx x.xxx	Specify the IP address of a single DNS server.	10.10.200.35	
<ntp>	xxx.xxx.xx x.xxx or host.domain.com	Specify the fully-qualified domain name or IP address of NTP server.	Pool1.ntp.org	
Storage requirements				
Parameter		Recommended value		Record your values
Base OS SW VMDK		500 GB thick provision, lazy zero SSD or flash		
Persistent message queue VMDK		Persistent message queue: 700GB thick-provision, lazy-zero SSD or flash VMDK		
		Database VMDK	2.5 TB thick provision, lazy zero SSD or flash	

Single node IBM Spectrum Discover trial deployment planning worksheet

Use this worksheet to plan for installing IBM Spectrum Discover for a single node trial deployment.

Table 15. Single node IBM Spectrum Discover trial deployment planning

CPU and memory requirements				
Parameter		Recommended value		Record your values
Memory		64 GB minimum 128 GB recommended		
Logical processor count		8 logical processors minimum 24 logical processors recommended		
Networking requirements				
Parameter	Value format	Recommended value	Example	Record your values

Table 15. Single node IBM Spectrum Discover trial deployment planning (continued)

CPU and memory requirements				
Parameter		Recommended value		Record your values
<hostname>	host.domain.com	Specify the fully-qualified domain name of the node.	node.example.com	
<interface>	ensXXX	Specify the Ethernet interface to use for the virtual appliance networking.	ens192	
<ip>	xxx.xxx.xx x.xxx	Specify the IP address of the node.	10.10.200.10	
<netmask>	xxx.xxx.xx x.xxx	Specify the network mask for the IP range of the node.	255.255.254.0	
<gateway>	xxx.xxx.xx x.xxx	Specify the IP address of the network gateway.	10.10.200.1	
<dns>	xxx.xxx.xx x.xxx	Specify the IP address of a single DNS server.	10.10.200.35	
<ntp>	xxx.xxx.xx x.xxx or host.domain.com	Specify the fully-qualified domain name or IP address of the NTP server.	Pool11.ntp.org	
Storage requirements				
Parameter		Recommended value		Record your values
Base OS SW VMDK		500 GB thick provision, lazy zero SSD or flash		

Table 15. Single node IBM Spectrum Discover trial deployment planning (continued)

CPU and memory requirements			
Parameter	Recommended value		Record your values
Persistent message queue VMDK	Persistent message queue: 50 GB minimum + 2 GB per 20 million indexed files, thick-provision, lazy-zero HDD or SSD or flash		
	Database VMDK	Database (does not include capacity for database backup): 100 GB minimum, 1 GB per 2 million indexed files, thick provision, lazy zero SSD or flash VMDK	
		Database (includes capacity for database backup): 100 GB minimum, 2 GB per 2 million indexed files, thick provision, lazy zero SSD or flash VMDK	

Note: If using 64GB of RAM, no more than 25 million files can be indexed into IBM Spectrum Discover.

Chapter 3. Deploying and configuring

This section provides information on how to deploy and configure IBM Spectrum Discover single node trial or production virtual appliance.

Deploy and configure a single node production IBM Spectrum Discover appliance cluster

The following section provides information on how to deploy and configure IBM Spectrum Discover single node trial or single node production virtual appliance.

You must consider virtual machine minimum requirements. For more information, see [“CPU and memory requirements for single node trial and single node production IBM Spectrum Discover deployments”](#) on page 12.

Deploying a single node trial or single node production IBM Spectrum Discover virtual appliance

The IBM Spectrum Discover software is available as an OVA (open virtualization appliance) file. You can deploy it on your VMware ESXi server by using the **VMware vSphere Client**.

Before you begin

- Download the IBM Spectrum Discover OVA file on the local system or obtain the URL to an IBM Spectrum Discover OVA file accessible on the internet.
- Review deployment and configuration known issues and workarounds. For more information, see [“Known issues with deploying and configuring for single node”](#) on page 43.

About this task

Deploy the IBM Spectrum Discover virtual appliance as follows by using the **Deploy OVF Template** wizard of the VMware vSphere Client.

Important: Use Firefox or Chrome to deploy the IBM Spectrum Discover virtual appliance.

Procedure

1. In the vSphere Client, right-click the ESXi server on which you want to deploy the virtual appliance and then click **Deploy OVF Template**.

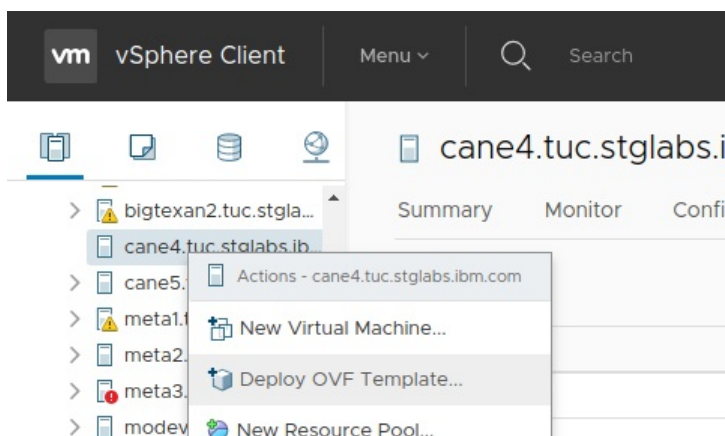


Figure 5. ESXi server menu

The **Deploy OVF Template** wizard appears.

2. Select the IBM Spectrum Discover virtual appliance that you want to deploy and click **Next**.

You can either select an OVA file that you download on the local system or you can specify a URL to the OVA file.

Deploy OVF Template

1 Select an OVF template

2 Select a name and folder

3 Select a compute resource

4 Review details

5 Select storage

6 Select networks

7 Ready to complete

Select an OVF template

Select an OVF template from remote URL or local file system

Enter a URL to download and install the OVF package from the Internet, or browse to a location accessible from your computer, such as a local hard drive, a network share, or a CD/DVD drive.

☒ URL

http://modevdump.tuc.stglabs.ibm.com/master/MetaOcean_master-3108.ova

☐ Local file

Choose Files No file chosen

CANCEL BACK NEXT

Figure 6. Deploy OVF template wizard

3. Specify the name of the virtual appliance or accept the default name and click **Next**.

Deploy OVF Template

✓ 1 Select an OVF template

2 Select a name and folder

3 Select a compute resource

4 Review details

5 Select storage

6 Select networks

7 Ready to complete

Select a name and folder

Specify a unique name and target location

Virtual machine name:

Select a location for the virtual machine.

✓ vcenter-136.tuc.stglabs.ibm.com

> Howard's Lab

> **Newies**

> Oldies

> Performance

CANCEL

BACK

NEXT

Figure 7. Virtual appliance location

4. Select the physical server on which you want to deploy the virtual appliance and click **Next**.

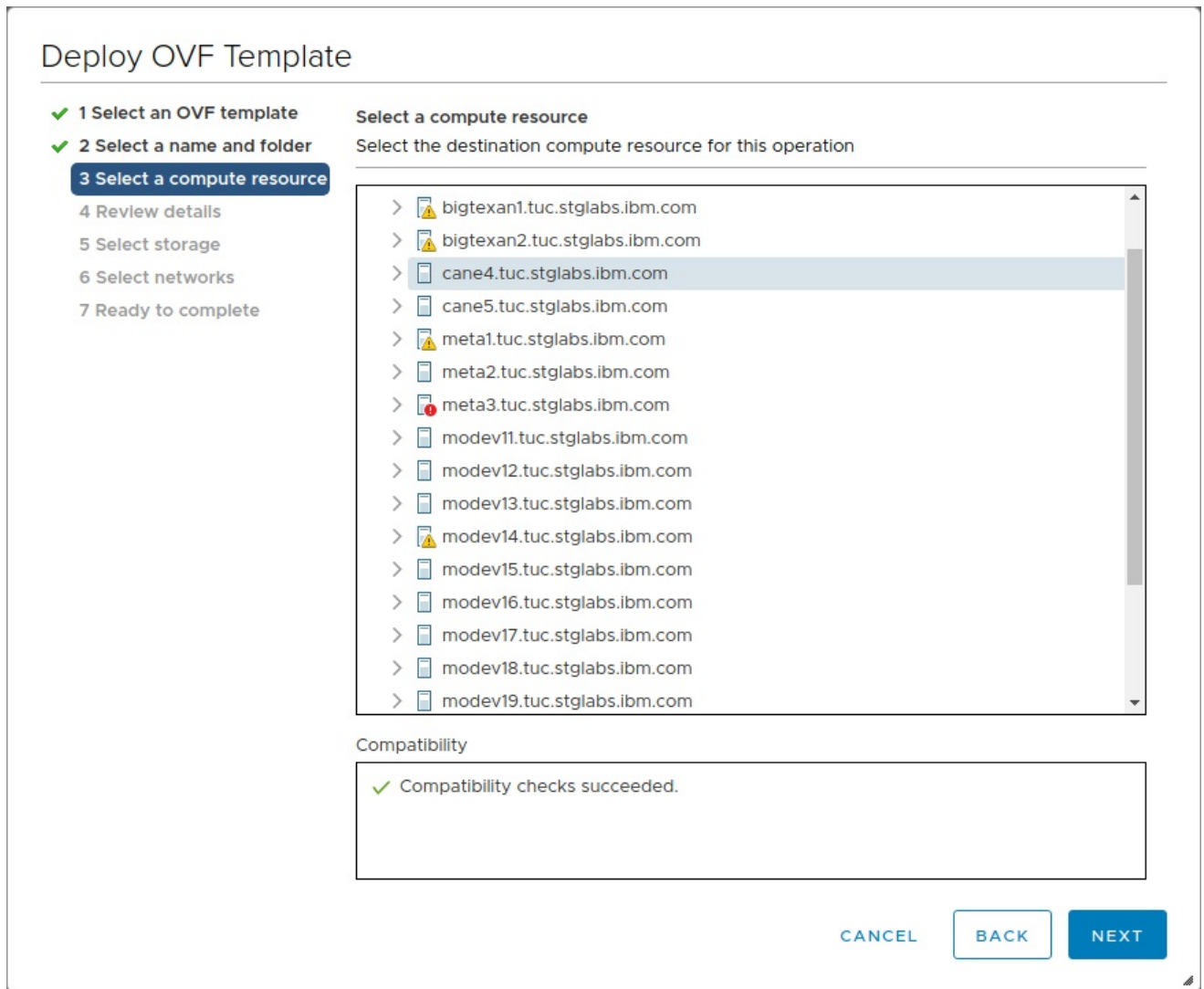


Figure 8. Physical server location

5. Review the details and click **Next**.

To validate your code integrity, see [“Validating code integrity”](#) on page 111.

6. Select the data store for the virtual appliance and the virtual disk format, and click **Next**. The recommended format is **Thick Provision Lazy Zeroed**.

Deploy OVF Template

✓ 1 Select an OVF template

✓ 2 Select a name and folder

✓ 3 Select a compute resource

✓ 4 Review details

5 Select storage

6 Select networks

7 Ready to complete

Select storage

Select the storage for the configuration and disk files

☐ Encrypt this virtual machine (Requires Key Management Server)

Select virtual disk format:

Thick Provision Lazy Zeroed

VM Storage Policy:

Datastore Default

Name	Capacity	Provisioned	Free	Type
Boot2	1,023.75 GB	1.28 TB	209.56 GB	VM
datastore4	103.25 GB	972 MB	102.3 GB	VM
MO_DATA1	1,023.75 GB	1.8 TB	82.52 GB	VM
MO_DATA2	1,023.75 GB	997.66 GB	122.33 GB	VM

Compatibility

✓ Compatibility checks succeeded.

CANCEL

BACK

NEXT

Figure 9. Select storage window

7. Select the VM network for the virtual appliance and click **Next**.

Deploy OVF Template

✓ 1 Select an OVF template

✓ 2 Select a name and folder

✓ 3 Select a compute resource

✓ 4 Review details

✓ 5 Select storage

6 Select networks

7 Ready to complete

Select networks

Select a destination network for each source network.

Source Network	Destination Network
VIS232	VM Network

1 items

IP Allocation Settings

IP allocation:

Static - Manual

IP address:

203.0.113.19

IP protocol:

IPv4

CANCEL

BACK

NEXT

Figure 10. Select virtual machine network

8. Review the settings and click **Finish**.

Deploy OVF Template

✓ 1 Select an OVF template

✓ 2 Select a name and folder

✓ 3 Select a compute resource

✓ 4 Review details

✓ 5 Select storage

✓ 6 Select networks

7 Ready to complete

Ready to complete

Click Finish to start creation.

Provisioning type	Deploy OVF From Remote URL
Name	modevvm15_master-3108
Template name	MetaOcean_master-3108
Folder	Newies
Resource	cane4.tuc.stglabs.ibm.com
Location	Boot2

CANCEL

BACK

FINISH

Figure 11. Review Settings

The IBM Spectrum Discover virtual node gets created and the storage is provisioned.

Note: Do not power on the virtual appliance until storage, CPU, and memory are configured.

Configuring storage for a single node trial or single node production of IBM Spectrum Discover virtual appliance

The IBM Spectrum Discover trial and production virtual appliance node requires three additional VMDK storage devices.

Note: The persistent message queue and database VMDK storage devices are in addition to the base OS and software VMDK that was automatically configured during the initial IBM Spectrum Discover virtual appliance deployment. A total of three virtual disks (VMDK) are required including the disk that is created when installing the appliance on the OVA. For more information, see [“Storage requirements for single node trial and single node production IBM Spectrum Discover deployments” on page 13.](#)

Procedure

1. Add virtual disk for the IBM Spectrum Discover persistent message queues.

For more information, see [“Adding virtual disk for IBM Spectrum Discover persistent message queues” on page 28.](#)

2. Add virtual disk for the IBM Spectrum Discover database.

For more information, see [“Adding a virtual disk for the IBM Spectrum Discover database” on page 31.](#)

3. [

Add virtual disk for the backups. For more information, see [“Adding a virtual disk for the IBM Spectrum Discover backups” on page 34](#)

]

Adding virtual disk for IBM Spectrum Discover persistent message queues

You can use the VMware vSphere Client to add the virtual disk that is required for IBM Spectrum Discover persistent message queues to the virtual appliance.

About this task

Important:

See the Chapter 2, “Planning,” on page 11 section for detailed requirements for the persistent message queue VMDK. For a single node production deployment, a 4.5 TB thick provisioned, and lazy zeroed VMDK is required. If an optional IBM Spectrum Discover applications is to be configured, an additional 1.6 TB of capacity is required.

For a trial IBM Spectrum Discover deployment, a minimum of 50 GB capacity is required for the VMDK and 1 GB per 2 million indexed files can be used as a sizing metric. These requirements are valid if IBM Spectrum Discover applications are not to be configured. If IBM Spectrum Discover applications are to be configured, a minimum of 50 GB capacity is required for the VMDK and 2 GB per 2 million index files can be used as a sizing metric.

Procedure

1. In the vSphere Client, right-click the IBM Spectrum Discover virtual appliance to which you want to add the virtual disk and click **Edit Settings**.

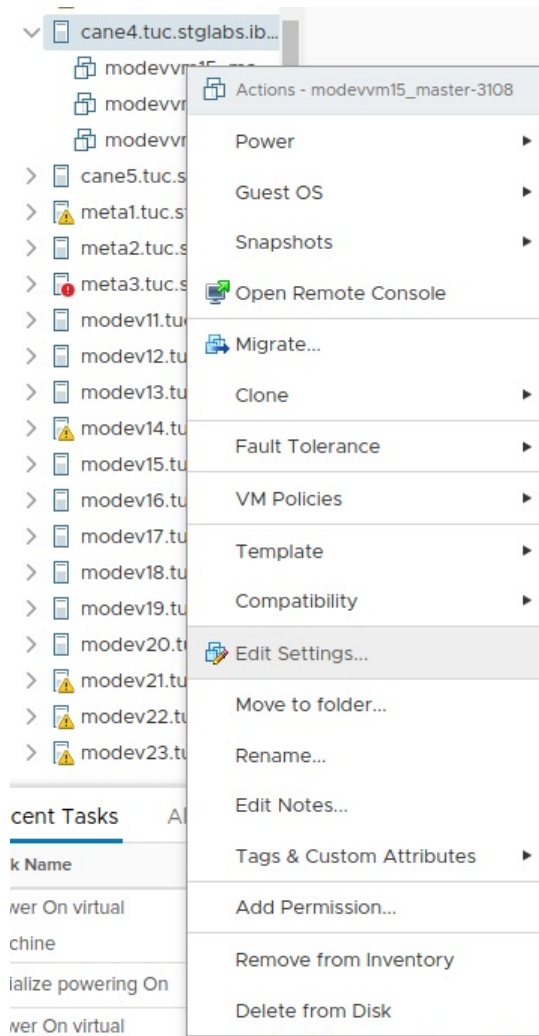


Figure 12. Edit Settings

2. From the **ADD NEW DEVICE** drop-down menu in the upper-right hand corner of the dialog box, select **Hard Disk**.

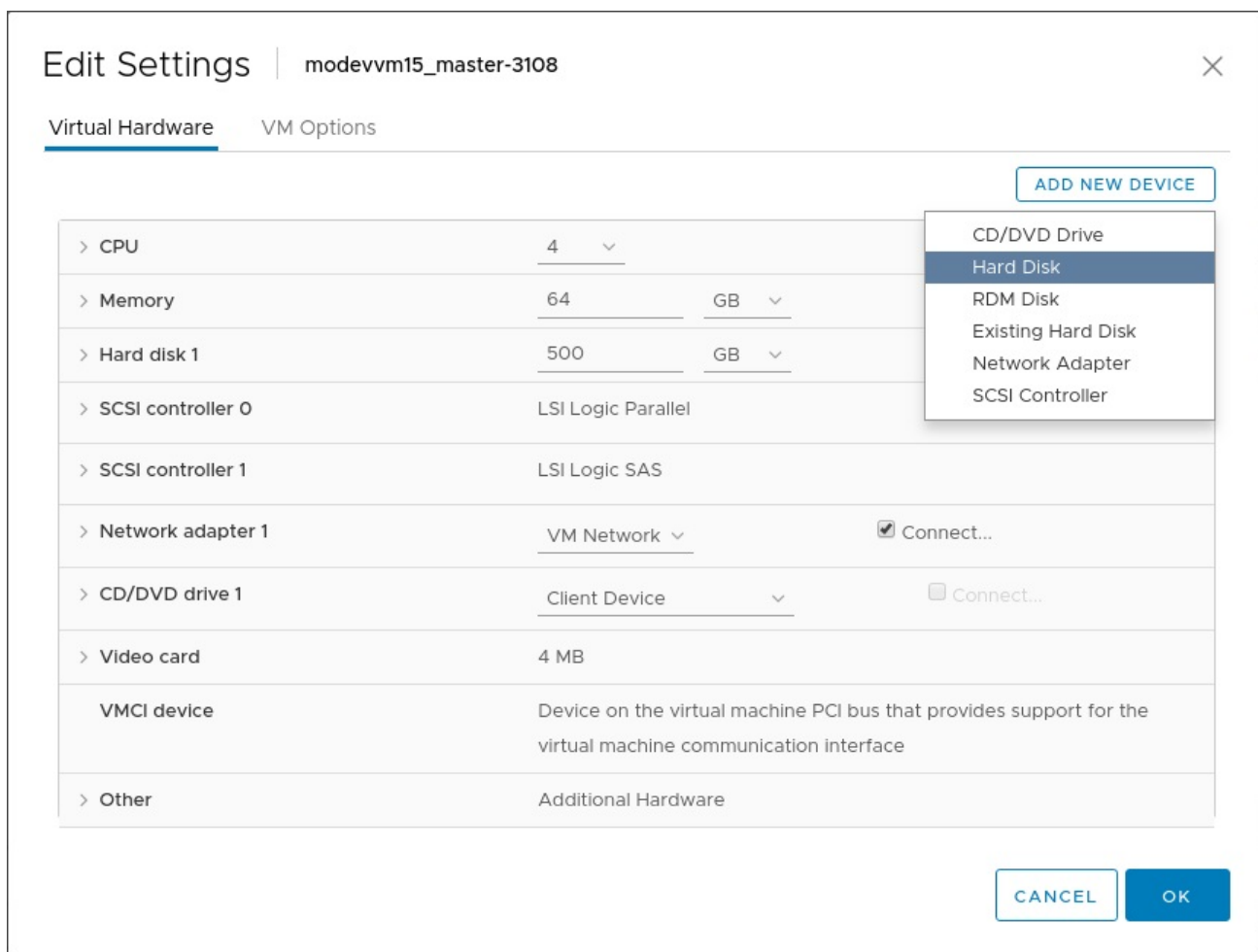


Figure 13. Edit hard disk settings

A **New Hard Disk** entry appears under **Virtual Hardware**.

3. Click **New Hard Disk** to expand the menu and select options for the disk.

Edit Settings
modevm15_master-3108

> Hard disk 1	500	GB	▼
▼ New Hard disk *	20	GB	▼
Maximum Size	827.77 GB		
VM storage policy	Datastore Default ▼		
Location	Store with the virtual machine ▼		
Disk Provisioning	Thick Provision Lazy Zeroed ▼		
Sharing	Unspecified ▼		
Shares	Normal ▼	1000	
Limit - IOPs	Unlimited ▼		
Virtual flash read cache	0	MB	▼
Disk Mode	Dependent ▼		
Virtual Device Node	SCSI controller 0 ▼	SCSI(0:1) New Hard disk ▼	
> SCSI controller 0	LSI Logic Parallel		
> SCSI controller 1	LSI Logic SAS		

CANCEL
OK

Figure 14. Hard disk options

Now you can set the size, provisioning, and location of the virtual disk. The default location is the data store where the virtual appliance resides. If needed, you can select a different data store

Note: Hard disk options show an example of a new hard disk size of 20 GB, but this number might be much larger. For production environments, it is required to allocate more space for the persistent message queue. For more information, see [Chapter 2, “Planning,” on page 11](#).

4. Click **OK** to confirm your settings and create the virtual disk.

Adding a virtual disk for the IBM Spectrum Discover database

You can use the VMware vSphere Client to add the virtual disk that is required for IBM Spectrum Discover database to the virtual appliance.

Before you begin

Important:

See the [Chapter 2, “Planning,” on page 11](#) section for detailed requirements for the database VMDK. For a single node production IBM Spectrum Discover deployment, a 2.5 TB thick provisioned and lazy zeroed VMDK are required.

For a trial IBM Spectrum Discover deployment, a minimum of 100 GB capacity is required for the VMDK and 1 GB per 2 million indexed files can be used as a sizing metric if backup and restore are not required. If backup and restore are required, a minimum of 100 GB capacity is required for the VMDK and 2 GB per 2 million indexed files can be used as a sizing metric.

You can use the VMware vSphere Client to add the virtual disk that is required for the IBM Spectrum Discover database to the virtual appliance.

Procedure

1. In the vSphere Client, right-click the IBM Spectrum Discover virtual appliance to which you want to add the virtual disk and then click **Edit Settings**.

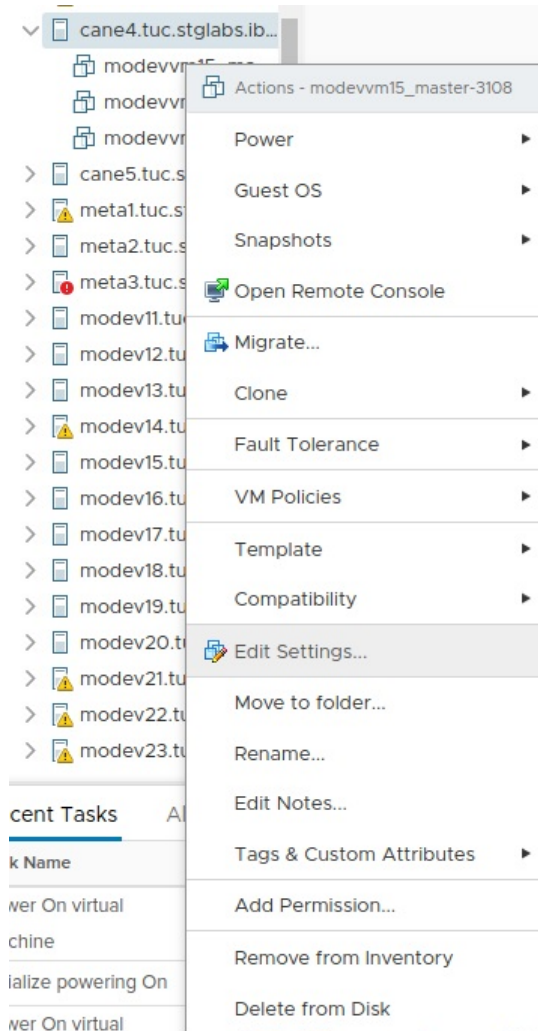


Figure 15. Edit Settings

2. From the **ADD NEW DEVICE** list, select **Hard Disk**.

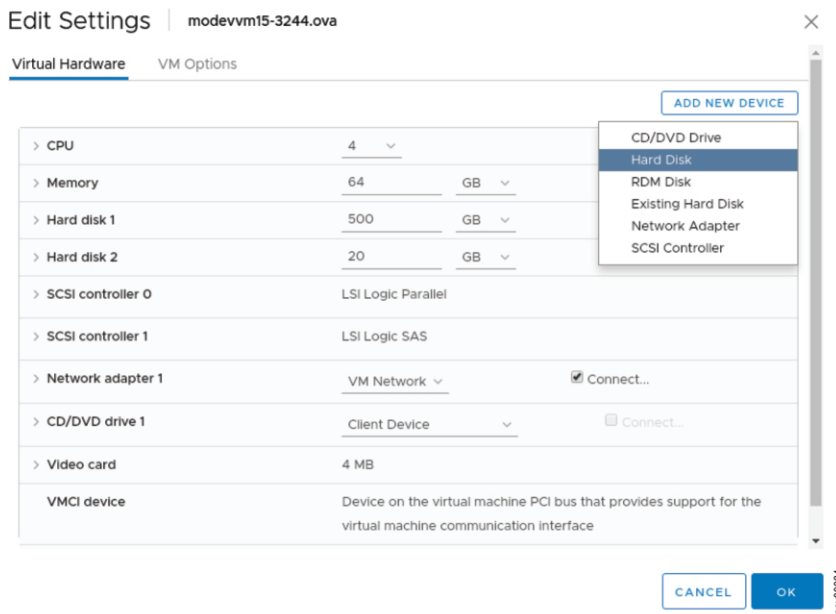


Figure 16. Edit settings dialog box

A **New Hard Disk** entry appears under **Virtual Hardware**.

3. Click **New Hard Disk** to expand the menu and select options for the disk.

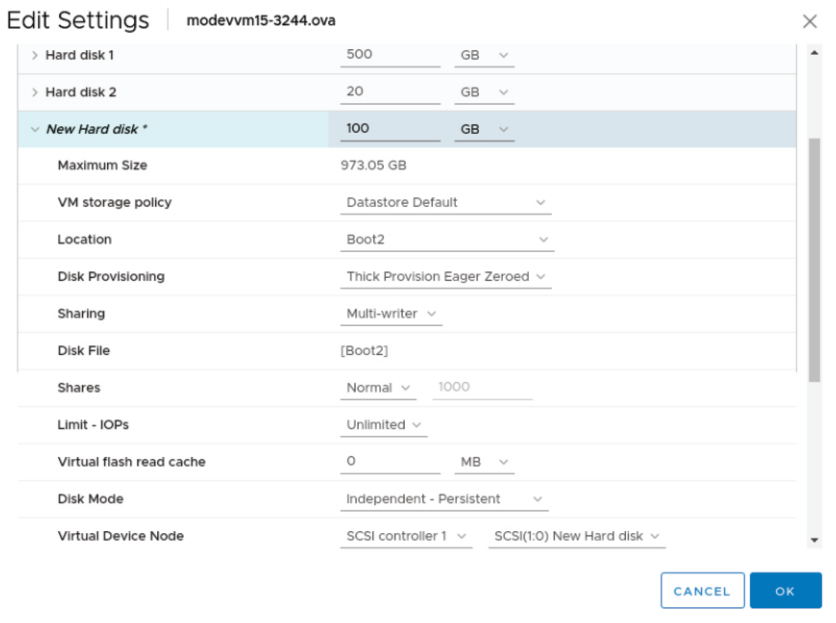


Figure 17. Hard disk options

You can set the size, provisioning, and location of the virtual disk. The default location is the data store where the virtual appliance is located. If needed, you can select a different data store.

4. Click **OK** to confirm your settings and create the virtual disk.

Adding a virtual disk for the IBM Spectrum Discover backups

You can use the VMware vSphere Client to add the virtual disk for IBM Spectrum Discover backup, to the virtual appliance.

Before you begin

Important:

See the [“Backup and restore storage requirements for IBM Spectrum Discover”](#) on page 16 section for the backup and restore storage requirements for IBM Spectrum Discover VMDK. For a single node production IBM Spectrum Discover deployment, you need a 2.5 TB thick provisioned and lazy zeroed VMDK.

For a trial IBM Spectrum Discover deployment, you need the following requirements if you do not need backup and restore.

- A minimum of 100 GB capacity for the VMDK and
- 1 GB per 2 million indexed files can be used as a sizing metric.

If you need backup and restore, then the following requirements must be fulfilled.

- A minimum of 100 GB capacity for the VMDK
- 2 GB per 2 million indexed files can be used as a sizing metric.

Procedure

1. In the vSphere Client, right-click the IBM Spectrum Discover virtual appliance to which you want to add the virtual disk and then click **Edit Settings**.

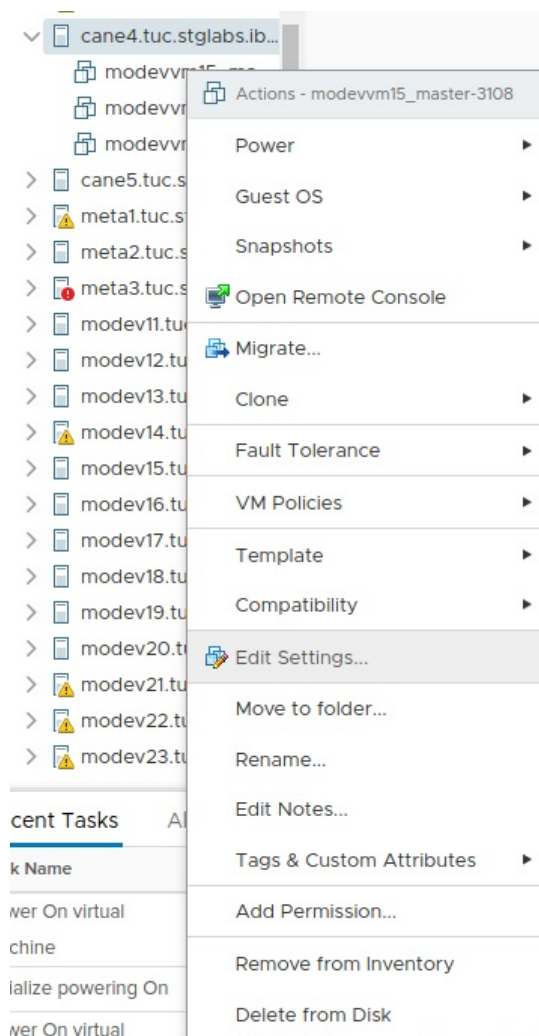


Figure 18. Selecting the **Edit Settings** menu option

2. From the **ADD NEW DEVICE** list, select **Hard Disk**.

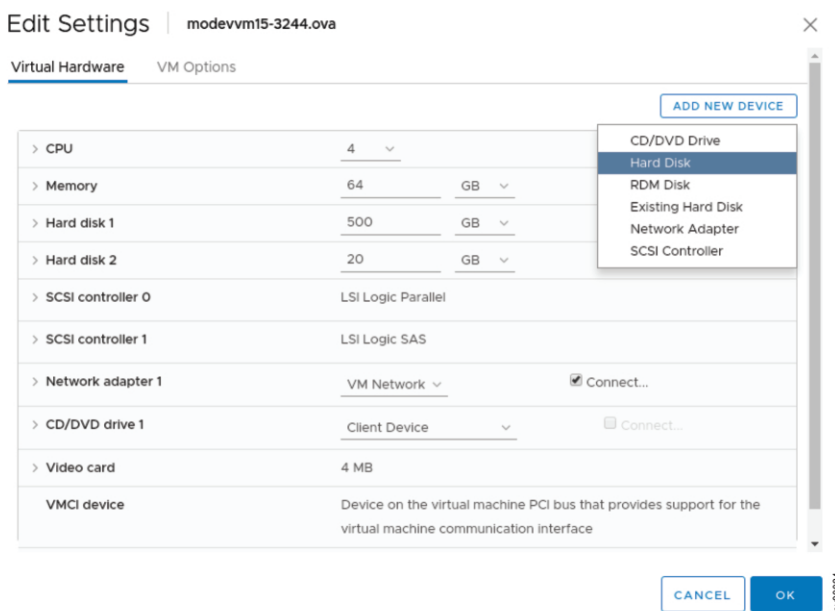


Figure 19. Selecting options in the **Edit settings** dialog box

A **New Hard disk** entry appears under **Virtual Hardware**.

3. Click **New Hard disk** to expand the menu and select options for the disk.

The screenshot shows the 'Edit Settings' dialog for a virtual machine named 'modevmm15-3244.ova'. The 'New Hard disk' entry is selected, showing a size of 100 GB. Below this, various configuration options are listed, including Maximum Size, VM storage policy, Location, Disk Provisioning, Sharing, Disk File, Shares, Limit - IOPs, Virtual flash read cache, Disk Mode, and Virtual Device Node. The 'OK' button is highlighted.

Setting	Value
Hard disk 1	500 GB
Hard disk 2	20 GB
New Hard disk *	100 GB
Maximum Size	973.05 GB
VM storage policy	Datastore Default
Location	Boot2
Disk Provisioning	Thick Provision Eager Zeroed
Sharing	Multi-writer
Disk File	[Boot2]
Shares	Normal 1000
Limit - IOPs	Unlimited
Virtual flash read cache	0 MB
Disk Mode	Independent - Persistent
Virtual Device Node	SCSI controller 1 SCSI(1:0) New Hard disk

Figure 20. Selecting hard disk options

You can set the size, provisioning, and location of the virtual disk. The default location is the data store where the virtual appliance is located. If needed, you can select a different data store.

4. Click **OK** to confirm your settings and create the virtual disk.

Configuring CPU and memory allocation for the single node IBM Spectrum Discover virtual appliance

This section lists the step to increase the default allocations of CPU and memory for each IBM Spectrum Discover virtual appliance.

About this task

It is recommended to reserve the assigned memory assigned to the IBM Spectrum Discover virtual appliance to avoid running out of physical memory and swapping, which severely impacts database performance and stability.

Important: A single node production IBM Spectrum Discover virtual appliance requires 128 GB RAM and 24 logical processors. 128 GB RAM and 24 logical processors is also recommended for the single node trial IBM Spectrum Discover virtual appliance. However, 64 GB and [16] logical processors can be configured that support indexing of up to 25 million files and objects. For more information, see [Chapter 2, "Planning,"](#) on page 11.

Procedure

1. In the vSphere client, right-click the IBM Spectrum Discover virtual appliance for which you want to change the CPU and memory allocation and click **Edit Settings**.

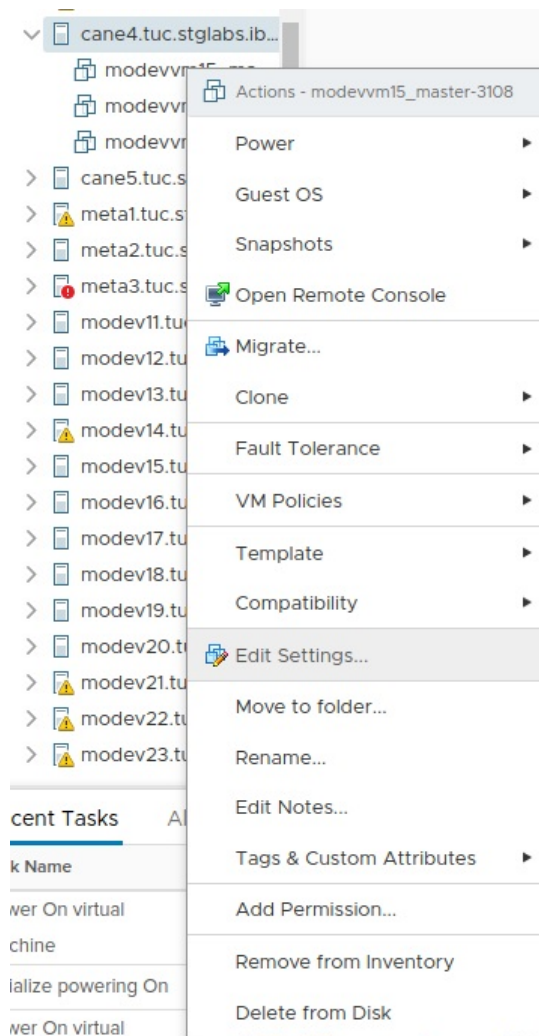


Figure 21. Edit settings menu

2. Under **Virtual Hardware**, from the **CPU** list, select the number that you want to increase the CPU allocation to.

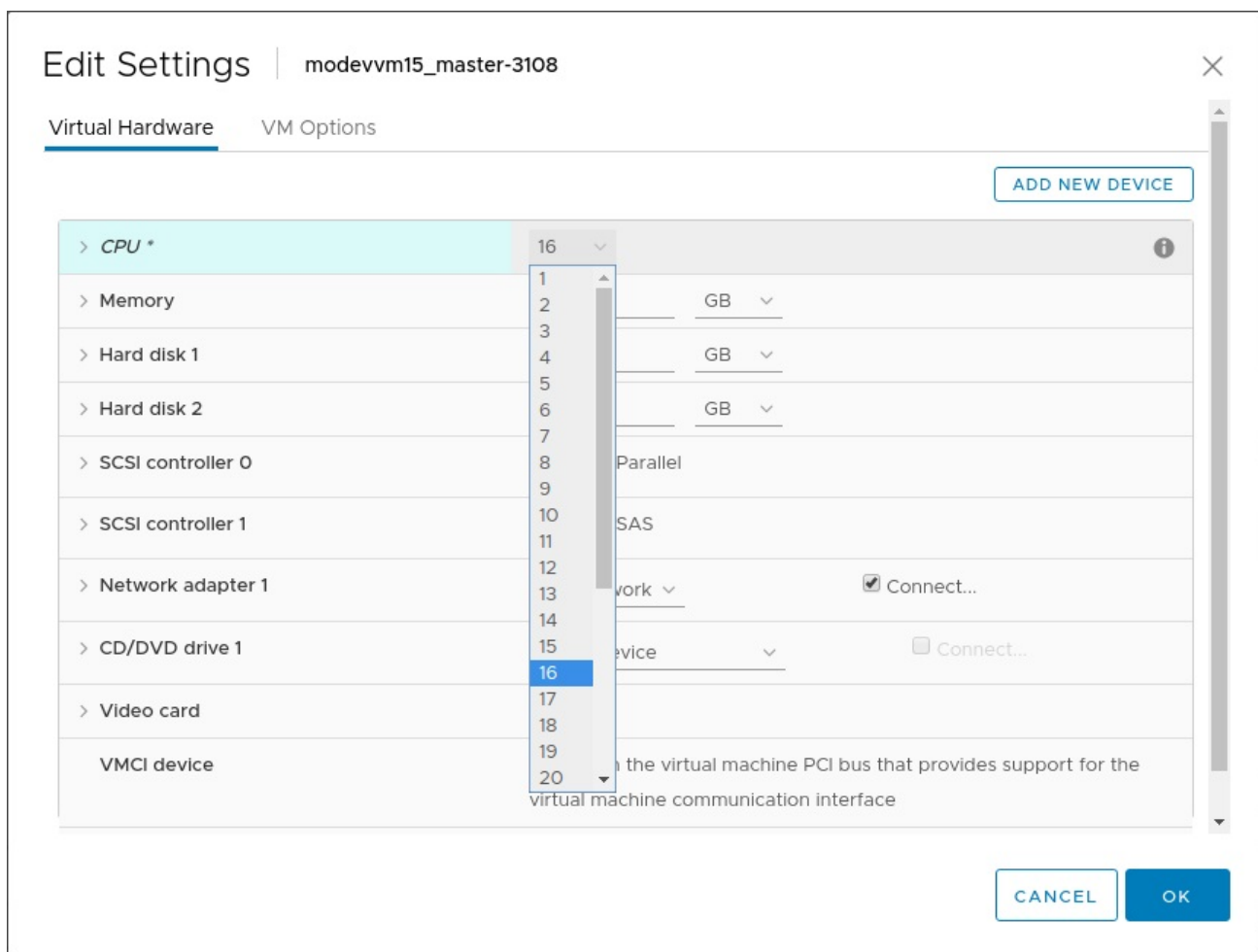


Figure 22. CPU allocation settings

3. In the **Memory** field, enter the number that you want to change the memory allocation to and select the memory unit from the drop-down list.

Edit Settings | modevbm15_master-3108

Virtual Hardware

VM Options

ADD NEW DEVICE

> CPU *	16		
> Memory *	96	GB	
> Hard disk 1	500	GB	
> Hard disk 2	20	GB	
> SCSI controller 0	LSI Logic Parallel		
> SCSI controller 1	LSI Logic SAS		
> Network adapter 1	VM Network	<input checked="" type="checkbox"/> Connect...	
> CD/DVD drive 1	Client Device	<input type="checkbox"/> Connect...	
> Video card	4 MB		
VMCI device	Device on the virtual machine PCI bus that provides support for the virtual machine communication interface		

CANCEL

OK

Figure 23. Memory allocation settings

- In the **Reservation** field under **Memory**, change the number according to the changed memory allocation and select the memory unit from the drop-down list.

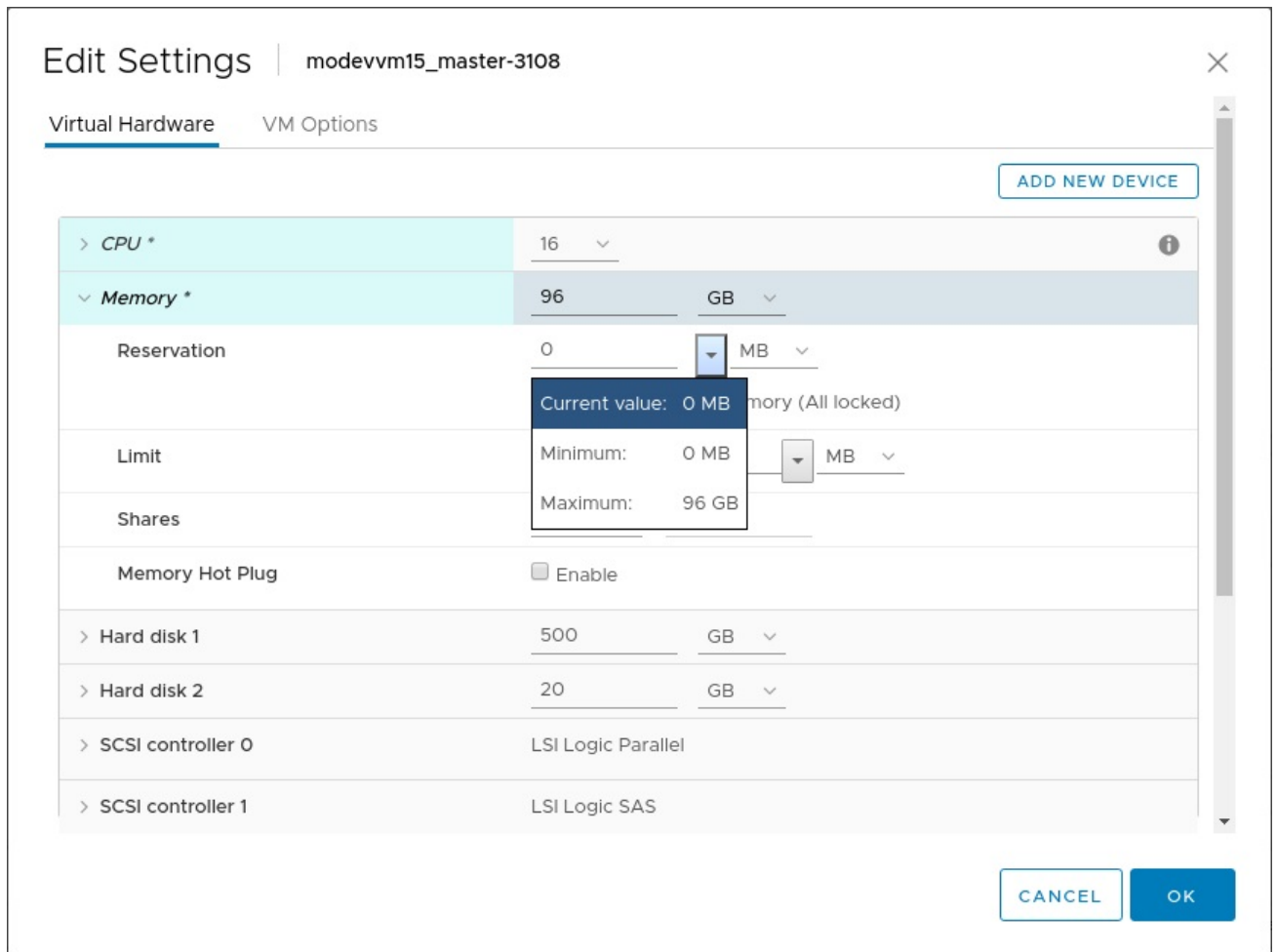


Figure 24. Reserved memory allocation settings

5. Click **OK** to confirm the changes in CPU and memory allocation.

Configuring networking and perform provisioning of a single node trial or single node production IBM Spectrum Discover virtual appliance

After you deploy a virtual appliance in the IBM Spectrum Discover (and storage, CPU, and memory are configured), you must configure networking. After you configure networking, you must provision the virtual appliances by using a provisioning tool.

Procedure

1. Power on the virtual appliance.
2. In the vSphere Client, right-click the virtual appliance and click **Open Remote Console**.

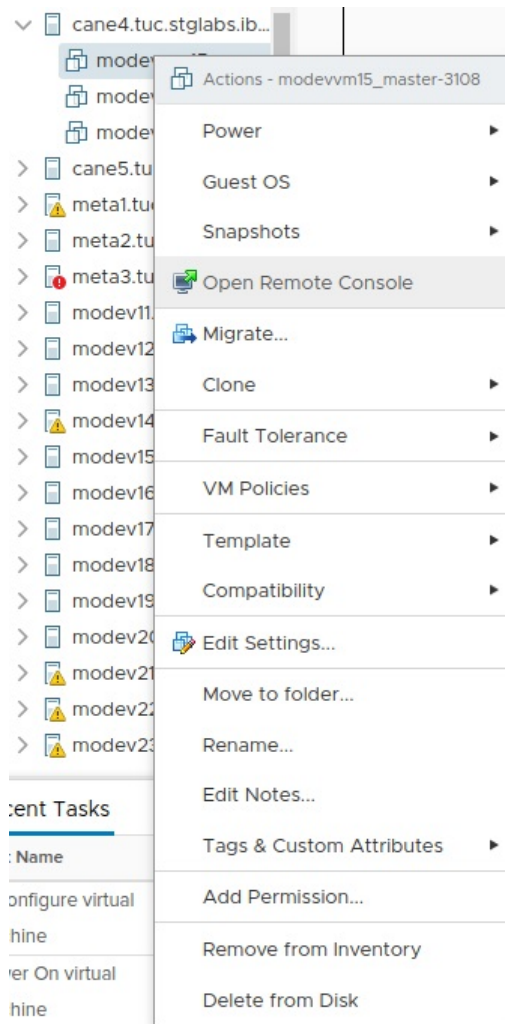


Figure 25. vSphere client settings menu

3. At the virtual appliance login prompt, enter the username and the password to log in. The default username is **moadmin** and the default password is **PasswOrd**.
4. Change the directory to `/opt/ibm/metaocean/configuration`.

```
cd /opt/ibm/metaocean/configuration
```

Note: IBM Spectrum Discover requires a fully qualified domain name (FQDN) that is registered in a customer supplied domain name server (DNS). The customer supplied FQDN must be resolvable by the customer supplied DNS from the IBM Spectrum Discover node for the IBM Spectrum Discover virtual appliance to operate properly.

5. Configure the IBM Spectrum Discover virtual appliance networking settings by running the following command:

```
sudo ./mmconfigappliance
```

Note: One or more nodes IBM Spectrum Discover must be able to communicate with a customer supplied network time protocol (NTP) server to operate properly.

The following table lists the definitions of the required settings:

Parameter	Value format	Recommended value	Example
<i>HostName</i>	host.domain.com	The fully qualified domain name of the node	node1234.example.com

Parameter	Value format	Recommended value	Example
<i>Interface</i>	ensXXX	The Ethernet interface to use for the virtual appliance networking.	ens192
<i>IPAddress</i>	xxx.xxx.xxx.xxx	The IP address of the node.	10.10.200.10
<i>NetMask</i>	xxx.xxx.xxx.xxx	The network mask for the IP range of the node.	255.255.254.0
<i>Gateway</i>	host.domain.com	The IP address of the network gateway.	10.10.200.1
<i>DNS</i>	ensXXX	The IP address of a single DNS server.	10.10.200.35
<i>NTPServer</i>	xxx.xxx.xxx.xxx or host.domain.com	The fully qualified domain or the IP address of the NTP server.	pool1.ntp.org
Mode	Single or Multi	Single	Single

Note: During this step, you are prompted to:

- Accept the IBM Spectrum Discover license agreement.
- Set the time zone. To change the timezone to something other than Coordinated Universal Time (UTC), go through the prompts to choose your continent, country, or region. (No other setup is required to set the time zone.)
- Change the moadmin password.

The **mmconfigappliance** process takes approximately 1 hour to complete for a production system.

Check for the Kubernetes `network_cidr` and `service_cluster_ip_range` values. These specified values must not conflict with the existing host network IP data. The default values are:

```
network_cidr: 10.1.0.0/16
service_cluster_ip_range: 10.0.0.0/16
```

If you are prompted for the `network_cidr` or `service_cluster_ip_range`, consider this information to avoid network conflicts. Private network ranges that might be a good choice are:

- 10.0.0.0 to 10.255.255.255
- 172.16.0.0 to 172.31.255.255

For example, you can enter:

- 172.31.0.0/16 for the `network_cidr`
- 172.30.0.0/16 for the `service_cluster_ip_range`

```
network_cidr: 172.31.0.0/16
service_cluster_ip_range: 172.30.0.0/16
```

IBM Spectrum Discover checks to see whether the system host IP overlaps with the Kubernetes default network and service IP range. If this overlap is detected, you get an error message similar to this error message:

```
Host network (10.1.10.10) is overlapped with the default Kubernetes network (10.1.0.0/16).
Please enter in a new value for the Kubernetes network.
```

Upon successful completion, an output similar to the following output is displayed.

```
PLAY RECAP *****
*****
*****
203.0.113.14      : ok=241   changed=209   unreachable=0    failed=0
canevm7.example.com : ok=9     changed=6     unreachable=0    failed=0
canevm8.example.com : ok=9     changed=6     unreachable=0    failed=0
canevm9.example.com : ok=9     changed=6     unreachable=0    failed=0
```

The process is completed successfully when you do not see any messages that say failed or when you see a message that the failed count = 0.

For more information about using Kubernetes settings, see [Customizing the cluster with the config.yaml file in the IBM Cloud Private Knowledge Center](#).

Known issues with deploying and configuring for single node

In case of any errors that you might encounter while deploying or configuring IBM Spectrum Discover, review the following information for details and possible workarounds.

Issue	Description	Resolution or workaround
Log shows error after deployment	<p>Even after successful deployment, log might show some errors with messages similar to the following:</p> <pre>2018-10-15 11:33:00,557 p=12895 u=root fatal: [203.0.113.18]: FAILED! => {"changed": false, "cmd": "awk '/ sse4_2/ {exit 42}' /proc/cpuinfo", "delta": "0:00:00.004105", "end": "2018-10-15 11:33:00.540782", "failed": true, "rc": 42, "start": "2018-10-15 11:33:00.536677", "stderr": "", "stderr_lines": [], "stdout": "", "stdout_lines": []} 2018-10-15 11:33:00,558 p=12895 u=root ...ignoring</pre>	These error messages can be ignored.

Deploying the IBM Spectrum Discover open virtualization appliance on the Kernel-based Virtual Machine virtualization module

The IBM Spectrum Discover open virtualization appliance (OVA) can be extracted and imported onto the Kernel-based Virtual Machine (KVM) virtualization module.

Before you begin

To deploy IBM Spectrum Discover, you must install KVM with a network bridge. (A network bridge connects two separate networks as if they were a single network.) A network bridge is required because IBM Spectrum Discover host names must resolve in Domain Name System (DNS).

About this task

Use the following steps to import and extract the IBM Spectrum Discover open virtualization appliance (OVA) onto the Kernel-based Virtual Machine (KVM) virtualization module. This deployment is similar to the VMware ESXi deployment.

Procedure

1. Extract the OVA.

An OVA file can be extracted as a tar file:

```
tar xvf spectrum-discover*.ova
```

These files include:

- A manifest
- An Open Virtual Machine (OVF) file
- A disk image (for example, `vmname.vmdk`)

2. Convert the disk image.

KVM can mount a disk image that is in Virtual Machine Disk (VMDK) format, but in read-only mode. To use IBM Spectrum Discover, convert the disk format to one that KVM can access in read/write mode. Run this command:

```
qemu-img convert -f vmdk -O raw MetaOcean*disk1.vmdk discover.img
```

3. Create more disk images.

IBM Spectrum Discover requires two extra disks (similar to the VMware ESXi deployment). To generate the required extra disks, run this command:

```
qemu-img create -f raw kafka.img 500G
qemu-img create -f raw scale.img 500G
```

4. Create the virtual machine.

Use the disk images with the `virt-manager` application to create a KVM. You can also use the command-line (where `${bridge}` refers to the name of the network bridge that is created as a prerequisite). To create a KVM, run this command:

```
virt-install \
  --name=discover \
  --memory 234000 \
  --vcpu 16 \
  --numatune 0,1,mode=interleave \
  --virt-type=kvm \
  --os-variant=centos7.0 \
  --boot hd \
  --network=bridge:${bridge} \
  --disk discover.img,device=disk,bus=scsi,format=raw,cache=none \
  --disk kafka.img,bus=scsi,format=raw,cache=none \
  --disk scale.img,bus=scsi,format=raw,cache=none \
  --graphics vnc \
  --console pty,target_type=serial \
  --noautoconsole
```

5. Deploy IBM Spectrum Discover.

Use a VNC viewer to access the virtual machine console and:

- Configure the IBM Spectrum Discover network.
- Configure the time.
- Use the `mmconfigappliance` command to accept the license.
- Use the `launch_ansible` command to install IBM Spectrum Discover.

Configure data source connections

Data source connections describe the source data systems for which IBM Spectrum Discover indexes metadata.

Creating data source connections in IBM Spectrum Discover identifies source storage systems that are to be indexed by IBM Spectrum Discover.

For some data source types, a network connection is (optionally) created to allow for automated scanning and indexing of the source system metadata. IBM Spectrum Discover will not index data from unknown sources, so creating a data source connection is the first step towards cataloging any source storage system.

Remember: Depending on when the scan is stopped, stopping a running scan might result in an inconsistent database state for the connection.

You can add data connections to the source storage systems from the IBM Spectrum Discover graphical user interface.

IBM Spectrum Discover discards any data that comes in from an unknown connection. Therefore, connections must be established before data ingestion. To see the list of defined connections, use the **Data Connection** tab under the **Admin** window of the GUI.

Remember: If you use a MAC, you might have to adjust the scroll bar settings in **System Settings** to see all available connection types. For example, activate the **Show scroll bars: Always** option.

Typically, a data source is equivalent to a single file system or object vault or bucket. A data source connection is an alias for the combination of a cluster name and a data source within the cluster. This allows multiple file systems or buckets or vaults with the same name to be indexed by IBM Spectrum Discover when they are in separate clusters.

Remember: IBM Spectrum Discover does not support file or file path names that use characters that are not part of the UTF-8 character set.

IBM Spectrum Scale data source connections

You can create an IBM Spectrum Scale data source connection, scan a data source, and manually initiate a scan.

Tip: If the data source connection is used by the ScaleILM application to tier data by using IBM Spectrum Archive, the host setting of the IBM Spectrum Scale connection must specify one of the IBM Spectrum Archive nodes. For more information, see the topic *Tiering data by using ScaleILM application* in the *IBM Spectrum Discover: Administration Guide*.

IBM Spectrum Discover supports the following method of scanning IBM Spectrum Scale data sources:

Automated scanning

A data source connection is defined on IBM Spectrum Discover, including details of how to connect to the IBM Spectrum Scale system. IBM Spectrum Discover connects to the IBM Spectrum Scale system by using these details, scans the file system and sends the details back to IBM Spectrum Discover. For more information, see the topic *Prerequisites for automated scanning* in the *IBM Spectrum Discover: Administration Guide*.

IBM Spectrum Scale scanning considerations

The following sections comprise considerations that you need to understand to scan IBM Spectrum Scale effectively.

Security considerations

Use this information to securely scan a system connection.

Scanning an IBM Spectrum Scale instance involves the use of the **mmapplypolicy** command on the IBM Spectrum Scale system, which requires superuser permissions. When you are creating the data source connection for the target IBM Spectrum Scale system in the IBM Spectrum Discover interface, you are prompted for a *userid* and *password* to enable automated scans. You are not required to provide these credentials if scans are run only manually on the target IBM Spectrum Scale system by an administrator. However, if you want to run automation and/or schedule scans, then the authentication credentials are required. By default, IBM Spectrum Discover uses password authentication to the Scale cluster to run commands remotely. However, you can supply your own RSA private key by selecting the shared key authentication option when you are configuring the connection if you want to avail passwordless authentication.

Rather than providing root login credentials, an administrator must create a special user ID with limited permissions on the IBM Spectrum Scale system. The administrator must also enable a password-less **sudo** for the user ID, to the binaries needed for scanning. This prevents someone from gaining root access to the target IBM Spectrum Scale system if the IBM Spectrum Discover system is somehow compromised.

Changing passwordless SSH keys

You can rotate RSA authentication key pairs for passwordless SSH on a frequency and remove old security keys from the `authorized_hosts` file on the IBM Spectrum Scale node that IBM Spectrum Discover connects to. To update the authentication keys, follow these steps:

1. Make sure that the `id_rsa.pub` contents for the new authentication key pair are in the `~/.ssh/authorized_hosts` file for the user that is specified in the IBM Spectrum Discover connection document for the IBM Spectrum Scale target file system.
2. Edit the connection and paste the contents of the new private key file (`id_rsa`) in the input form.

After you edit the connection with the new private key file, IBM Spectrum Discover uses it to connect to the IBM Spectrum Scale target system.

Performance considerations

Use this information to scan a system connection without degrading performance.

Running a scan policy on an IBM Spectrum Scale system can be resource intensive and cause noticeable performance degradation on the IBM Spectrum Scale system. Often, system administrators choose to designate certain nodes or node classes for running the scans. The IBM Spectrum Discover interface has an input field when creating IBM Spectrum Scale connections for the administrator to specify which nodes or node class(es) they would like to run the scan on. The value `all` will run the scan across all nodes in the cluster. Any other list (comma separated) will be treated as a list of nodes or node classes on which to run the scan. Scan times vary by the size of the filesystem, how many nodes are used in the scan, how many CPUs are used per node, and whether or not the IBM Spectrum Scale cluster metadata is in flash memory.

Prerequisites for automated scanning

You can use IBM Spectrum Scale automated scanning features.

IBM Spectrum Discover supports two levels of automated scanning of IBM Spectrum Scale systems. Both levels require that IBM Spectrum Discover must establish a password-less Secure Shell (SSH) connection to the IBM Spectrum Scale clustered system that is being scanned.

The difference between the two levels lies in the manner in which the output is handled. The factors that are considered inspect whether:

- The output of the IBM Spectrum Scale policy that is run to do the scan is stored in a file (which then must be automatically copied back to IBM Spectrum Discover and ingested locally); Or
- If the output is, instead, pushed to the ingest Kafka queue of IBM Spectrum Discover system directly from the IBM Spectrum Scale policy output.

The Kafka queue of IBM Spectrum Discover system is more space-efficient and time-efficient but has certain dependencies on the IBM Spectrum Scale clustered system that must be met so that it can function. The IBM Spectrum Discover automated scan code determines whether the dependencies are met on the IBM Spectrum Scale clustered system.

If the dependencies are met on the IBM Spectrum Scale clustered system, it attempts to scan the system by using the optimized path. If the dependencies are not met on the IBM Spectrum Scale clustered system, it defaults to the file copy path.

The following sections list the prerequisites for creating a connection and performing automated scanning. These prerequisites consider security and performance factors. For more information, see [IBM Spectrum Scale scanning considerations](#).

Identify the IBM Spectrum Scale cluster and node list

Identify a node to connect with the GPFS cluster.

Identify a node in the target IBM Spectrum Scale cluster to use for the IBM Spectrum Discover connection to the GPFS cluster.

You must identify the node list or node class that participates in the scanning activity.

Creating or identifying a user ID and password for scanning

Identify a user to perform scanning or create a new user ID.

About this task

You can identify an existing user to perform scanning or follow these steps on the IBM Spectrum Scale system to create a special user ID for scanning.

Procedure

1. Log in to a IBM Spectrum Scale management node as **root**.
Alternatively, you can **sudo** to **root** from another user ID.
 2. Use the following adduser steps to ensure that you are able to ssh into the cluster:
 - a) `adduser <user> -m`
 - b) `passwd <user>`
 3. Run: **visudo**
 - a) Add this line in the users section:
`<user> ALL=NOPASSWD: /usr/lpp/mmfs/bin/mmapplypolicy, /usr/lpp/mmfs/bin/mmrepquota`
 - b) Write and quit: `:wq`
 4. Create an IBM Spectrum Discover working directory and ensure that <user> has write permissions.
For example: `mkdir -p /gpfs/fs1/sd_scan -m 770; chown <user> /gpfs/fs1/sd_scan`
 5. To execute a Data Mover policy for tiering data on the IBM Spectrum Scale cluster by using the ScaleILM application on the IBM Spectrum Scale connection , do the following:
 - a. Run: **visudo**
 - b. Add or update this line in the users section:
`<user> ALL=NOPASSWD: /usr/lpp/mmfs/bin/mmapplypolicy, /usr/lpp/mmfs/bin/mmrepquota, /usr/lpp/mmfs/bin/mmlspool, /opt/ibm/ltfsee/bin/eeadm, /usr/bin/ls`
 - c. Write and quit: `:wq`
- For more information, see the topic *Tiering data by using ScaleILM application* in the *IBM Spectrum Discover: Administration Guide*.
6. [
To execute a Data Mover policy for copying data to the IBM Spectrum Scale cluster by using the ScaleAFM application on the IBM Spectrum Scale connection, do the following:
 - a. Run: **visudo**
 - b. Add or update this line in the users section:
`<user>ALL=NOPASSWD: /usr/lpp/mmfs/bin/mmapplypolicy, /usr/lpp/mmfs/bin/mmrepquota, /usr/lpp/mmfs/bin/mmlspool, /usr/lpp/mmfs/bin/mmfmctl, /usr/lpp/mmfs/bin/mmfmcosctl, /usr/lpp/mmfs/bin/mmaddcallback, /usr/lpp/mmfs/bin/mmdelcallback`
 - c. Write and quit: `:wq`

For more information, see the topic *Copying data using ScaleAFM application* in the *IBM Spectrum Discover: Administration Guide*.

]

Validate scan user permissions and configuration

Ensure that the user has the required permissions and configurations for scanning.

You need to ensure that it is possible to use Secure Shell (SSH) to log in to the IBM Spectrum Scale system with the scanning user ID and password.

Run the **sudo /usr/lpp/mmfs/bin/mmapplypolicy**.

You also need to validate the listed factors for the scanning-related working directory:

- It exists
- It is globally accessible by the scan worker nodes.
- The scan user has write permissions to the directory (ownership of the directory is preferred but not mandatory).

Place the `id_rsa` and `id_rsa.pub` files in `/opt/ibm/metaocean/data/connections/scale/` directory on the IBM Spectrum Discover instance if a specific RSA key pair for password-less SSH is wanted.

You need to validate if the listed things are installed on the IBM Spectrum Scale node. This node is identified while you are identifying a node in the target IBM Spectrum Scale cluster. The node in the target is identified for use in the IBM Spectrum Discover connection to the GPFS cluster:

- A python2 level of at least Python 2.7.5.
- A sufficient level of librarian.
- An appropriate level of confluent-Kafka.

Note: This validation is recommended for optimized scan ingestion. However, it is optional and can be skipped.

Check the dependencies for optimized scanning

Identify the dependencies that must be met on the IBM Spectrum Scale system to optimize automated scan ingest.

The dependencies that must be satisfied on the IBM Spectrum Scale system to optimize automated scan ingest from IBM Spectrum Discover are:

- A `librdkafka` library version 0.11.4 or later
- [A Python version 3.0 or later with accompanying Python package installer (pip3)]
- A `confluent-kafka` version that is greater than or equal to the installed `librdkafka` version.

If these dependencies are met, the scan output is pushed to the ingest Kafka queue of IBM Spectrum Discover system directly from the IBM Spectrum Scale policy output.

An administrator can determine whether `librdkafka` is installed on the IBM Spectrum Scale node by running the `find /usr -name "*librdkafka*"` or `ls /lib64/librdkafka*` commands. The `librdkafka` package is included with newer levels of IBM Spectrum Scale on **x86** and **ppc64le** platforms. However, it can also be built from the source code on older levels of IBM Spectrum Scale or **ppc64** platforms. If the IBM Spectrum Scale system runs on Red Hat Enterprise Linux® (RHEL) and is connected to a Red Hat Satellite, you can install it by running the Yellowdog Updater Modified (YUM) command `yum install librdkafka` as root. You can find source packages of `librdkafka` here: <https://github.com/edenhill/librdkafka>

[The user ID specified in the data source connection must be able to locate the following two binaries by using the OS shell path:]

[

1. A Python 3 binary as either python or python3
2. A Python package installer as pip3

Note: Symbolic links or aliases may be used to locate the Python executables.

]

After you install a sufficient version of Python, you can install confluent-kafka by using pip. To get pip, you must install the python-setuptools package, which provides a binary called easy_install. For more information, see <https://pypi.org/project/setuptools/#files>.

After easy_install is available, you can install pip by running easy_install-2.7 pip as root. After you install pip, you can install confluent-kafka by running pip install confluent-kafka as root.

Creating an IBM Spectrum Scale data source connection

You can use the IBM Spectrum Discover graphical user interface to create data connections from the source storage systems.

Procedure

1. Log in to the IBM Spectrum Discover web interface with a user ID that has the **Data Admin** role that is associated with it.

The **Data Admin** access role is required for creating connections. For more information, see *Managing User Access* in the *IBM Spectrum Discover: Administration Guide*.

2. Select **Admin** from the left navigation menu.

Click **Admin** to display the different types of data source connection names, platforms, clusters, data source, size, and **Add Connection** button.

3. Click **Add Connection** to display a new window that shows **Data Connections Add Data Source Connection**.

You can enter in the connection name and connection type. The connection types are:

- IBM Spectrum Scale
- IBM Cloud Object Storage
- Network File System (NFS)

4. Complete the following steps:

- a) In the field for **Connection Name**, define a **Connection Name**.
- b) Click the **Connection Type** drop-down menu and **Choose an option** to display the connection type options.

5. Set the connection type to IBM Spectrum Scale. The page displays the connection name, user, password, working directory, and scan directory information that you can enter. You can also schedule a data scan, select a collection, or enable live events.

If you click **Enable Live Events** you can enable a IBM Spectrum Scale watch folder on the specified file system.

6. Complete updating values for all the fields to add the IBM Spectrum Scale connection type, and click **Submit Connection**.

For IBM Spectrum Scale connections, you can enter the following information:

Connection name

The name of the connection, an identifier for the user, for example filesystem1.

Note: It must be a unique name within IBM Spectrum Discover.

User

A user ID that has permissions to connect to the data source system and initiate a scan.

Password

The password for the user ID specified in user.

Authentication Type

The password authentication can be done by using the password provided to authenticate with the Scale cluster. The shared RSA key authentication will perform a passwordless authentication by using a private key that is provided by the system administrator and whose public key exists in the authorized keys for the specified user on the Scale host.

Note: Release version 2.0.3.1 removes the support for self-generated RSA key pair for IBM Spectrum Discover. Any existing connections that use that method is updated to password based authentication and the self-generated key pair is removed during the upgrade to 2.0.3.1 or later. If the password for the scan user that is stored in IBM Spectrum Discover is no longer valid, this may result in scan failures after the update. To rectify this, you must edit the connection and provide a valid password for the scan user or a valid RSA private key for authentication.

Working Directory

A scratch directory on the source data system where IBM Spectrum Discover can put its temporary files.

Note: When you edit an existing connection and change the User from a root user to a non-root user, you must also change the Working Directory. This change is necessary because the non-root User cannot access the files previously created by the root user in the existing Working Directory.

Scan Directory

The root directory of the scan. All files and directories under this directory are scanned. Typically, this directory is the base directory of the file system, for example `/gpfs/fs1`.

Connection Type

The type of source storage system this connection represents.

Site

An optional physical location tag that an administrator can provide to see the physical distribution of their data.

Cluster

The IBM Spectrum Scale or GPFS cluster name. To obtain, run the following command from the IBM Spectrum Scale file system: `/usr/lpp/mmfs/bin/mmlscluster`.

Host

The hostname or IP address of an IBM Spectrum Scale node from which a scan can be initiated, for example a quorum-manager node.

File system

The short name (omit `/dev/`) of the file system to be scanned. For example, `fs1`.

Note: It is important to exactly match the file system name (data source) that IBM Spectrum Scale populates in the scan file. Run the following command on the IBM Spectrum Scale system: `/usr/lpp/mmfs/bin/mmlsmount all`

Node list

The comma-delimited list of nodes or node classes that participates in the scan of an IBM Spectrum Scale file system. For example, `scale01,scale02`.

Node	Daemon node name	IP address	Admin node name	Designation
1	msys111-10g	172.16.8.111	msys111-dmz	quorum-manager-perfmon
2	msys112-10g	172.16.8.112	msys112-dmz	quorum-manager-perfmon
3	msys113-10g	172.16.8.113	msys113-dmz	quorum-manager-perfmon

Note: When you create data source connections for IBM Spectrum Scale file systems, it is important to exactly match the cluster name and the file system name (data source) that IBM Spectrum Scale populates in the scan file.

Run the following commands on the IBM Spectrum Scale system.

- Run this command to display information about the cluster:

```
$ /usr/lpp/mmfs/bin/mmlscluster

GPFS cluster information
=====
GPFS cluster name:      modevvm19.tuc.example.com,
GPFS cluster id:        7146749509622277333
GPFS UID domain:        modevvm19.tuc.example.com
Remote shell command:   /usr/bin/ssh
Remote file copy command: /usr/bin/scp
Repository type:        CCR
Node  Daemon node name      IP address      Admin node name      Designation
-----
--
1      modevvm19.tuc.example.com  203.0.113.24   modevvm19.tuc.example.com  quorum-
manager
```

- b. Run this command to display information about file systems that are mounted:

```
$ /usr/lpp/mmfs/bin/mmlsmount all
File system gpfs0 is mounted on 1 nodes
File system Data_Science_8M is mounted on 7 nodes.
File system icp4D_data_fs_master1 is mounted on 8 nodes.
File system icp4D_data_fs_master2 is mounted on 8 nodes.
File system icp4D_data_fs_master3 is mounted on 8 nodes.
```

Automated scanning of an IBM Spectrum Scale data source

As an administrator, you can initiate an IBM Spectrum Scale scan from IBM Spectrum Discover to collect system metadata from IBM Spectrum Scale file system.

About this task

When a scan is initiated from the IBM Spectrum Discover graphical user interface, the data moves asynchronously back to the IBM Spectrum Discover.

Remember: Before you initiate a scan, see [“ IBM Spectrum Scale scanning considerations ”](#) on page 45.

Automated scanning and data ingestion relies on an established and active network connection between the IBM Spectrum Discover instance and the source IBM Spectrum Scale management node. If the connection cannot be established, the state of the data source connection displays 'unavailable' and the option for automated scanning does not appear in the IBM Spectrum Discover GUI for that connection.

Note: You cannot run scans unless you add override warnings in the configuration file.

Procedure

1. Go to the IBM Spectrum Discover GUI.
2. Go to **Admin > Data Source Connections**.
3. Select the data source connection that you want to scan. Make sure that the connection is online for your system ready to scan. (There is an indicator in the **Online** column.)
4. Select **Scan Now** to start the scan, and a small message appears to confirm that the data source connection you specify is being scanned.

You can view the status of the scan on the table in the **Scan Status** column for the target connection. After the **Scan Status** has a check mark next to it, the scan is complete.

Remember: You can also specify a time to begin the scan. Any time zones specified default to Coordinated Universal Time (UTC) time. So, if you specify your scan for 12 noon, it is 12 noon in UTC.

Automated scanning of an IBM Spectrum Scale file set

As an administrator, you can initiate an IBM Spectrum Scale scan from IBM Spectrum Discover to collect system metadata from a IBM Spectrum Scale file set or file sets.

Before you begin

This feature adds a requirement for non-root user IDs that are used for scanning IBM Spectrum Scale data source systems. This feature uses the **mmfsfileset** command to retrieve the list of available file sets from the target system when you have root-level permissions. So, if you use a non-root user ID it must have sudo access to **mmfsfileset** for this function to work.

There is already a requirement for a non-root scan user to have sudo access to **mmapplypolicy**, so this requirement adds **mmfsfileset** as an extra required command.

Note: You cannot query the available file sets on a target IBM Spectrum Scale connection or initiate a file set level scan unless you fulfill this requirement.

About this task

Scan a IBM Spectrum Scale file set or file sets to insert or update the records for the files that are found by IBM Spectrum Discover in that file set or file sets. The scan is scoped to the specified file set, which ensures a faster total scan than scanning the entire file system. Multiple file sets can be specified in a single scan operation, but the scanning of each file set is done successively.

As the scan progresses, the status message is updated to indicate the following information:

- The status message indicates which file set is being scanned.
- The status message indicates when data operations (such as transferring files or indexing data) occur.

This status message can be seen in the GUI on the data source connections table or it can be queried by using the REST API.

This feature works irrespective of whether the data is returned to IBM Spectrum Discover by using a direct Kafka connection or by using the file copy method. After a file set level scan completes, a scan generation is recorded or committed.

Additionally, an internal reclamation policy is generated to remove any deleted files that did not appear in the updated scan. The scope of this reclamation policy is limited to the file set that is scanned and does not affect other file sets or the actual file system. This limitation helps you achieve consistency with the source IBM Spectrum Scale system at file set level granularity.

Procedure

1. Go to the IBM Spectrum Discover GUI.
2. Go to **Admin > Data Source Connections**.

Select the wanted connection and click **Scan Now**, which opens the **Select Scan Type** dialog box.

You can select whether to scan the entire file system or to scan a list of file sets.

Important: Connection types other than IBM Spectrum Scale and SMB/CIFS do not open this dialog box. Additionally, **Scan Now** continues to function as it has, which means that there is an immediate initiation of a full connection scan.

3. Select either **Scan All** to scan all file sets or **Select Filesets** to scan a specific file set.

Selecting **Scan All** initiates a full scan of the file system. If you choose to scan all file sets, click **Scan** to run the scan.

Selecting **Select Filesets** initiates a specific file scan. Click **Next** to open the **Select Individual Filesets** dialog box. Use this dialog box to select the specific file sets that you want to scan. Search the table by using the table search header:

- a. You can select file sets by clicking the row of the table that represents that file set. Clicking the row highlights that row and the count under **View X selected filesets** increases by 1.
- b. You can also select file set by filtering the search criteria. The table can be filtered to show only the selected file sets by clicking **View X selected filesets**, for ease of review. For example, you can enter `fs` to display all file set with those characters in that order. Click the file set in the table row that you want to select to run the scan on that file set.

To go back to viewing all available file sets, click **View X selected filesets** again. The button changes to **View all filesets** when you view only the selected file sets.

4. After you select all wanted file sets, you can initiate the scan by clicking **Scan**. Clicking **Scan** takes you to the **Data Source Connections** table.

A notification indicates when the scan starts (or that the scan fails if there is a problem). You can view the status of the scan on the table in the **Scan Status** column for the target connection.

Remember: After the **Scan Status** has a check mark next to it, the scan is complete for all selected file sets.

Manual scanning of an IBM Spectrum Scale data source

How to configure IBM Spectrum Discover to connect to IBM Spectrum Scale. After completing these steps, data can be ingested from an IBM Spectrum Scale data source to IBM Spectrum Discover for metadata indexing.

Before you begin

Create the data source connection to IBM Spectrum Scale. For more information, see [“Configure data source connections”](#) on page 44.

The minimum connection parameters required for manual scanning are:

- Connection Name
- Connection Type
- Cluster
- Filesystem

Restriction: IBM Spectrum Discover uses a unit separator (ASCII code 0x1F) as the field delimiter for ingestion into the database. This means that data which contains this character in path/file/object names results in improper parsing of the input data and the records are rejected by IBM Spectrum Discover.

Procedure

1. Perform a file system scan to collect system metadata from IBM Spectrum Scale to be ingested into IBM Spectrum Discover. For more information, see [“Performing file system scan to collect metadata from IBM Spectrum Scale”](#) on page 54.
2. Copy the output of the file system scan to the IBM Spectrum Discover master node. For more information, see [“Copying the output of the IBM Spectrum Scale file system scan to the IBM Spectrum Discover master node”](#) on page 56.
3. Ingest data from the file system scan in IBM Spectrum Discover. For more information, see [“Ingesting metadata from IBM Spectrum Scale file system scan in IBM Spectrum Discover”](#) on page 57.
4. Ingest quota information from the file system. For more information, see [“Ingesting quota information from the file system”](#) on page 57.

Performing file system scan to collect metadata from IBM Spectrum Scale

You can use the file system scanning tool, IBM Spectrum Scale Scanner, to collect system metadata from IBM Spectrum Scale to be ingested into IBM Spectrum Discover.

About this task

The IBM Spectrum Scale Scanner tool uses the IBM Spectrum Scale information lifecycle management (ILM) policy engine to obtain the system metadata about the files stored on the file system. The system metadata is written to a file, which is then transferred to the IBM Spectrum Discover master node. The file is then ingested within the node and analytics is carried out to provide search, duplicate file detection, archive data detection, and capacity show-back report generation. The following system metadata is collected from the file system scan:

Key name	Description
Site	The site where the file or object resides.
Platform	The source storage platform that contains the file or object.
Size	The size of the file.
Owner	The owner of the file.
Path	The subdirectory where the data resides.
Name	The name of the data.
Permissions	The permissions for the file (mode).
ctime	The change time of the file metadata (inode).
mtime	The time when the data was last modified.
atime	The time when the data was last accessed.
Filesystem	The name of the IBM Spectrum Scale file system that is storing the data.
Cluster	The name of the IBM Spectrum Scale cluster.
inode	The IBM Spectrum Scale inode that is storing the data.
Group	The Linux group associated with the file.
Fileset	The file set that stores the file.
Pool	The storage pool where the file resides.
Migstatus	If applicable, indicates whether the data is migrated to tape or object.
migloc	If applicable, indicates the location of the data if migrated to tape or object.
ScanGen	Scan generation - useful to track rescans.

The IBM Spectrum Scale Scanner tool also collects quota information by calling **mmrepquota**.

The tool comprises the following files:

- `scale_scanner.py`: The tool that starts the IBM Spectrum Scale ILM policy.
- `scale_scanner.conf`: The configuration file used to customize the behavior of the `scale_scanner.py` tool.
- `createScanPolicy`: The script that is called internally by the tool.

Procedure

Install the IBM Spectrum Scale Scanner tool by unpacking the utility from the IBM Spectrum Discover node to the required location on the IBM Spectrum Scale cluster node.

1. Log in to the IBM Spectrum Discover node through Secure Shell (SSH) with the moadmin username and password:

```
ssh modadmin@spectrum.discover.ibm.com
```

2. Change to the directory that contains the Spectrum Scale scanning utility:

```
[/opt/ibm/metaocean/spectrum-scale]
```

3. scp the createScanPolicy, _init_.py, scale_scanner.conf, and scale_scanner.py files to a node in the IBM Spectrum Scale cluster:

```
scp * root@spectrumscale.ibm.com:/my_scanner_directory
```

```
createScanPolicy 100% 3320 3.2KB/s 00:00
init.py 100% 427 0.4KB/s 00:00
scale_scanner.conf 100% 1595 1.6KB/s 00:00
scale_scanner.py 100% 13KB 13.2KB/s 00:00
```

4. On the IBM Spectrum Scale node where you install the scanning utility, edit the configuration file (scale_scanner.conf) as follows:
 - a) Use the IBM Spectrum Discover UI to create a connection to the SS system on which you start a manual scan for. Set the filesystem and scandir fields, and optionally set the outputdir and site fields in the [spectrumscale] stanza of the file.

```
[spectrumscale]
# Spectrum Scale Filesystem which hosts the scan directory
# example: /dev/gpfs0
filesystem=/dev/gpfs0
# The directory path on Spectrum Scale Filesystem to perform scan on
# example: /gpfs0
# specifies a global directory to be used for temporary storage during
# mmappolicy command processing. The specified directory must be
# mounted with read/write access within a shared file system
mountpoint=mount point of the gpfs filesystem
# It is unclear what the mount_point should be, but setting the mount point
# to the mount point of the scale file system on the IBM Spectrum Scale node works.
scandir=/gpfs0
# The directory to store output data from the scan in (default is
# scandir)
outputdir=
# The site tag to specify a physical location or organization identifier.
# If you use this field, remove the comment (#)
#site=
```

- b) Set the scale_connection, master_node_ip, and username fields in the [spectrumdiscover] stanza of the file.

Note: scale_connection refers to the name of the IBM Spectrum Scale file system that is scanned and ingested into IBM Spectrum Discover. The scale_connection value must match the value that is defined in the Data Source column of the **Data Connections** page in the IBM Spectrum Discover GUI.

The username must be a valid name of a IBM Spectrum Discover user who has the dataadmin role. The **username** field takes the format of <domain_name>/<username>. To determine a domain and username with the dataadmin role, go to the **Access Users** page in the IBM Spectrum Discover GUI and click the view for the defined users.

For the local domain, it is not necessary to specify the domain as part of the username field as it is the default domain. For example, to define username for user1 in the local domain that is assigned the dataadmin role, in the configuration file, enter the following value: username=user1

```
[spectrumdiscover]
# Name of the Spectrum Scale connection to scan files from
# Check using the Spectrum Discover connection manager APIs
scale_connection=fs3
# Spectrum Discover Master Node IP
master_node_ip=203.0.113.23
# Spectrum Discover user name, having 'dataadmin' role
# Use format <domain_name>/<username>
# e.g. username=Scale/scaleuser1
username=user1
```

Note: The scanner output file generates approximately 1 K of metadata for every file in the system. If there are 12 M files, the size is expected to be approximately 12 GB. By default, the output file is written to the same directory that is being scanned. The log file output location can be customized by setting the outputdir field.

5. Run the scan by using the following command:

```
./scale_scanner.py
```

Note: While you run the `./scale_scanner.py` command, you can start another scan. If you start another scan, ensure that you run the scan with another connection that is online and is not being scanned currently. When the scanner is running, the scanner hides the **scan now** automatically.

Note: As you run the `scale_scanner.py` script, you are prompted for the password for the IBM Spectrum Discover user that is configured in the `scale_scanner.conf` file with the username under the `spectrumdiscover` section. You must provide the correct password for the configured user. As described in the configuration file, this user needs to be a valid user configured in the IBM Spectrum Discover Authentication service (Access management). Also, this user must be assigned to the dataadmin role.

For example:

```
$ ./scale_scanner.py
Enter password for SD user 'user1':
Scale Scan Policy is created at: ./scanScale.policy
```

Note:

- After you see a line similar to “0 ‘skipped’ files and/or errors”, press **enter** to return to the command prompt.
- The scan takes about 2 minutes 30 seconds for every 10 M files on the following configuration:

```
x86 -based Spectrum Scale Cluster
•4 M4 NSD client nodes
•2 M4 NSD server nodes
•DCS3700 350 2TB NL SAS drives & 20 200GB SSD
•QDR InfiniBand cluster network
```

Copying the output of the IBM Spectrum Scale file system scan to the IBM Spectrum Discover master node

After you have scanned your IBM Spectrum Scale file system and have the `list.metaOcean` output file, copy it to the IBM Spectrum Discover master node.

Procedure

As an IBM Spectrum Discover administrator, use **scp** to copy `list.metaOcean` file from the scan output directory to the `/opt/ibm/metaocean/data/producer` directory on the master node.

Note: If there are multiple file systems in the same cluster that are being scanned, you can rename the `list.metaOcean` file to avoid name conflicts and to not overwrite an existing `list.metaOcean` file that is in use. For example:

```
$ mv list.metaOcean list.metaOcean.myfilesystem  
[$ scp list.metaOcean.myfilesystem moadmin@MasterNodeIP:/opt/ibm/metaocean/data/producer]
```

Ingesting metadata from IBM Spectrum Scale file system scan in IBM Spectrum Discover

Records are inserted into IBM Spectrum Discover for indexing when they are pushed to a Kafka connector topic corresponding to the type of data being ingested. In the case of IBM Spectrum Scale, the Kafka connector topic type is `scale-scan-connector-topic`.

About this task

A Kafka client producer is required to put the IBM Spectrum Scale file system scan file records onto the Kafka connector topic. The following steps show how to use the **ingest** alias command to push the records in the `list.metaOcean` file (or another named file) onto the Kafka connector topic.

Procedure

1. Run the following command to ingest the data:

```
[$ ingest /opt/ibm/metaocean/data/producer/list.metaOcean]
```

2. Replace the `list.metaOcean` path with the path of the file that you want to ingest.

Ingesting quota information from the file system

The file system scanning tool, IBM Spectrum Scale Scanner, has the ability to harvest and send quota information to IBM Spectrum Discover.

Procedure

To perform quota ingestion, run the following command on the IBM Spectrum Scale cluster node:

```
./scale_scanner.py --quota-only
```

For example:

```
$ sudo ./scale_scanner.py --quota-only  
Enter password for SD user 'user1':
```

IBM Spectrum Archive data source connections

You can define tags and policies in IBM Spectrum Discover based on values that are derived from IBM Spectrum Archive metadata to help in searching and categorizing files.

IBM Spectrum Discover integrates with IBM Spectrum Archive to display search results that include the following archive state of files:

Migration status migstatus

Search results display details for the following migration status:

migrtd

Indicates that the file is migrated to tape.

resdnt

Indicates that the file is resident in the file system.

premig

Indicates that the file is pre-migrated to tape.

Migration "location" migloc

Search results display information on the tape cartridge in the following format: "1 *tape cartridge volser@tape storage pool id@tape library serial number*". Any additional copies must be separated by colons.

Actual size of file in the associated IBM Spectrum Scale file system Size Consumed Bytes

IBM Spectrum Discover displays a zero if the file is moved to tape.

Remember: The migration state information is collected and summarized in the IBM Spectrum Discover **State** facet. You can access this facet by using IBM Spectrum Discover visual search. For more information, see the topic *Searching* in the *IBM Spectrum Discover: Administration Guide*.

The IBM Spectrum Discover interface displays the search results including the metadata information.

Important: The location information that is displayed in IBM Spectrum Discover is provided by IBM Spectrum Scale. It corresponds to the `dmap1`. `IBMTPS` attribute for the file. Run the `mmfsattr -L` command for more details.

IBM Cloud Object Storage data source connection

You can create a IBM Cloud Object Storage (COS) connection and initiate a scan.

IBM COS uses a connector residing on the storage system to push events to a Kafka topic residing in the IBM Spectrum Discover cluster. When configured, the IBM Spectrum Discover consumes the events and indexes them into the IBM Spectrum Discover database.

Restriction: IBM Spectrum Discover uses a unit separator (ASCII code 0x1F) as the field delimiter for ingestion into the database. This means that data which contains this character in path/file/object names results in improper parsing of the input data and the records are rejected by IBM Spectrum Discover.

Prerequisites

The IBM Cloud Object Storage Scanner and Replay prerequisites are listed:

IBM Cloud Object Storage Scanner prerequisite

For the Scanner, you must enable the **Get Bucket Extension** for all accesser devices.

To enable the **Get Bucket Extension**, you must set the `s3.listing-name-only-enabled` equal to `true` in the Manager System Advanced Configuration.

See [Figure 26](#) on page 58.



Figure 26. Example of the system advanced configuration

Remember: You do not need to restart the Accesser, but you might need to wait for 5 minutes before the setting takes effect if you do not restart it.

IBM Cloud Object Storage Replay prerequisite

For the Replay, `access_logs` are uploaded to the management vaults within 1 hour after rotation. Rotation can be triggered earlier by setting the **Rotation Period** to the minimum value of 15 minutes in Manager under **Maintenance/Logs/Device Log Configuration**. Refer to the [IBM Cloud Object Storage Knowledge Center](#) to make sure that this is configured, and the relevant access logs are present before you run the Replay.

Creating an IBM Cloud Object Storage data source connection

You can create an IBM Cloud Object Storage (COS) data source connection and from the storage system.

Procedure

1. Log in to the IBM Spectrum Discover web interface with a user ID that has the **Data Admin** role that is associated with it.

The **Data Admin** access role is required for creating connections. For more information, see [“Role-based access control”](#) on page 3.

2. Select **Admin** from the left navigation menu.

Clicking **Admin** displays the different types of data source connection names, platforms, clusters, data source, size, and **Add Connection**.

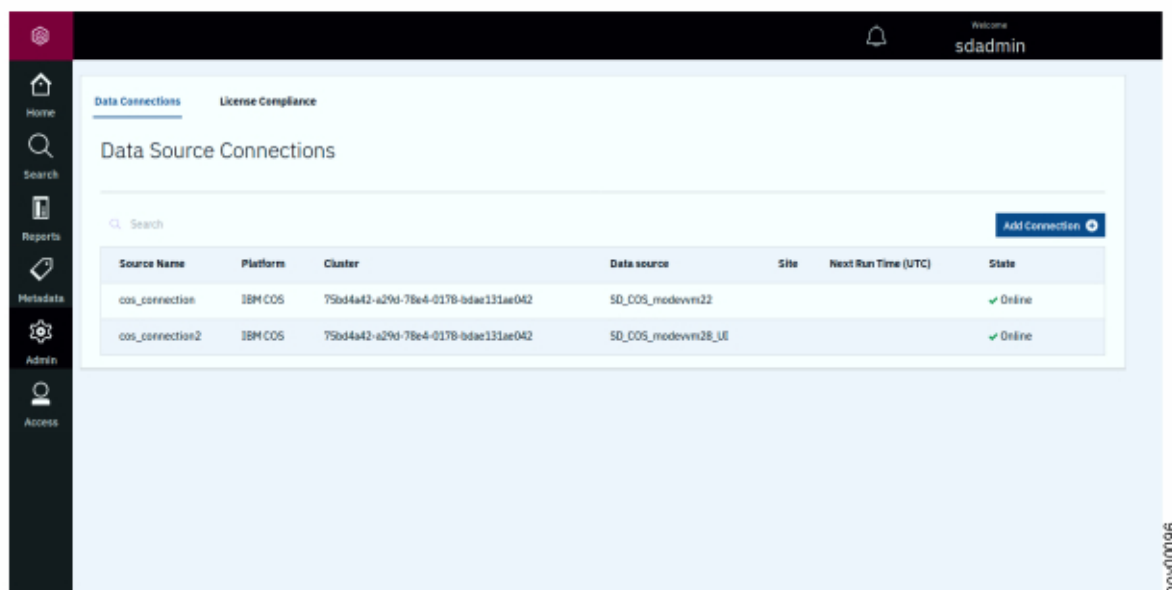


Figure 27. Displaying the source names for data source connections

3. Click **Add Connection** to display a new window that shows **Data Connections Add Data Source Connection**.

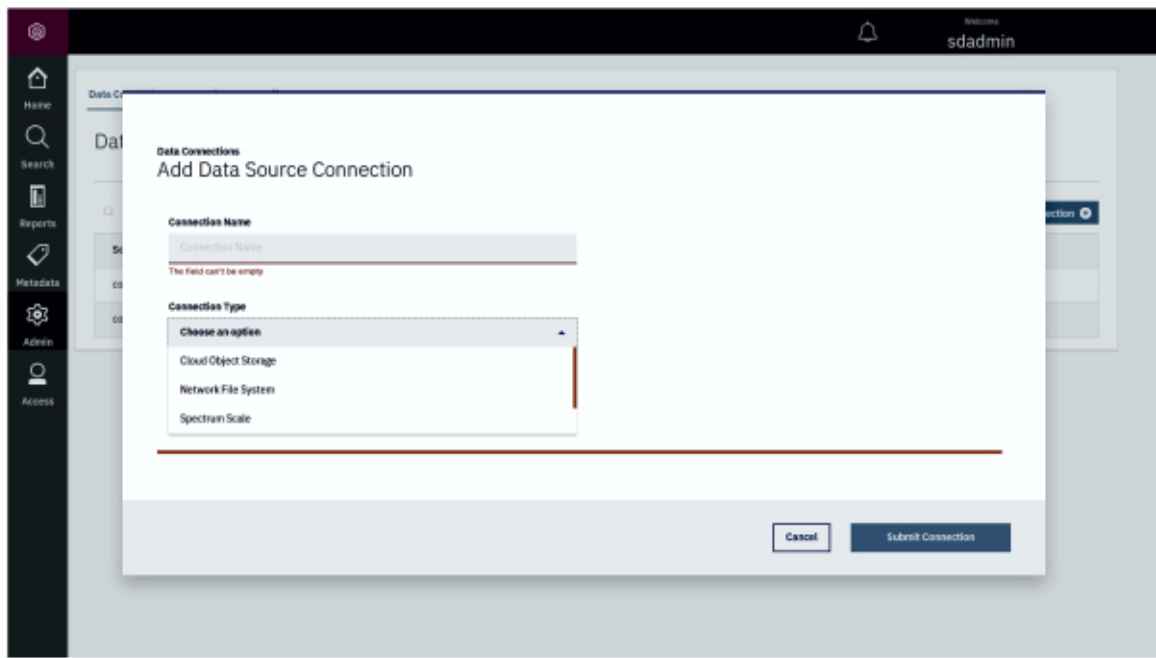


Figure 28. Example of window that shows Data Connections Add data source Connection

4. Do the following steps:
 - a) In the field for **Connection Name**, define a **Connection Name**.
 - b) Click the down arrow for **Connection Type** to display a drop-down menu for the connection type.
5. Select the connection type **Cloud Object Storage**.

[Figure 29 on page 61](#) shows an example of the screen for an IBM COS connection.

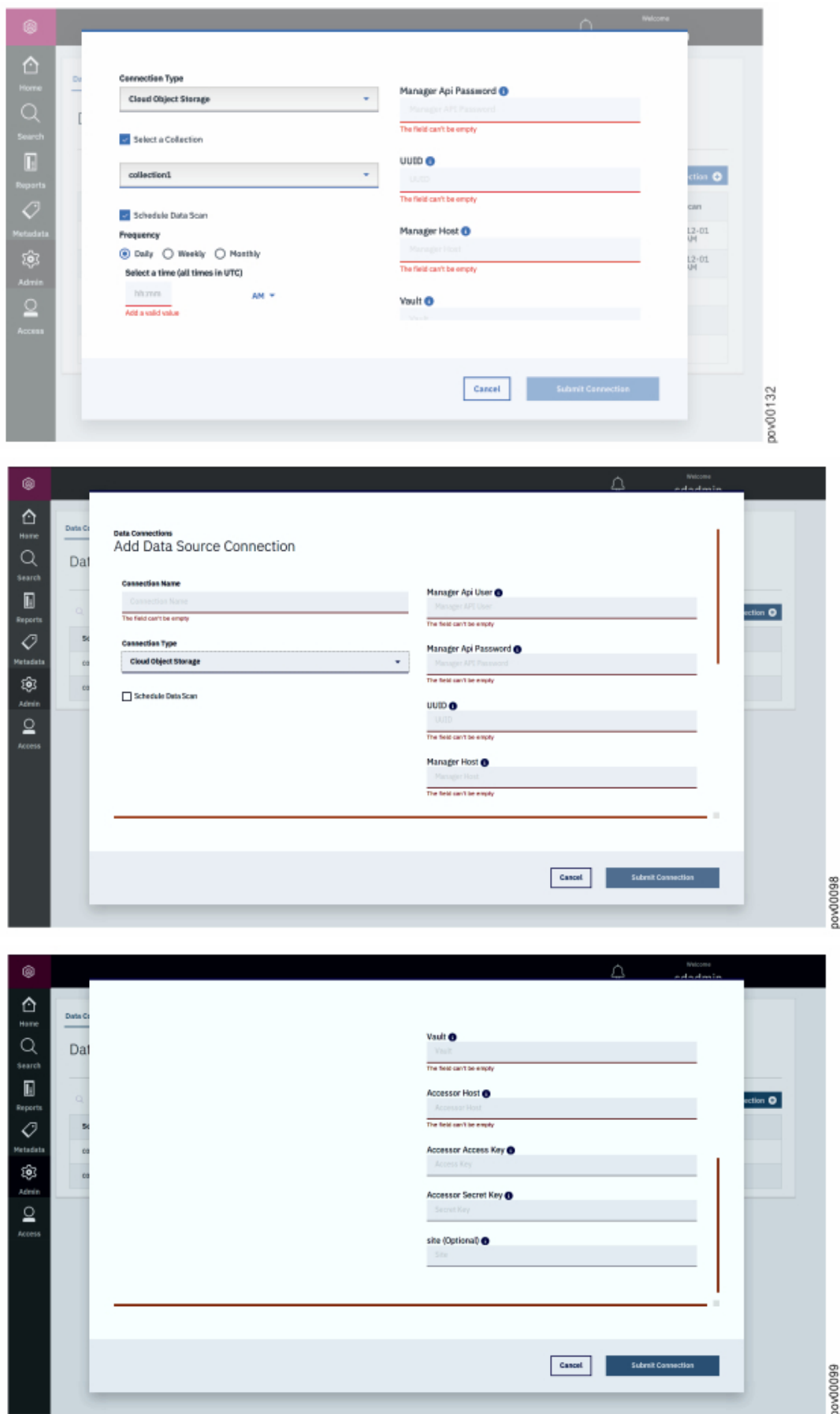


Figure 29. Example of the screen for an IBMCOS connection

6. In the screen for **Cloud Object Storage**, complete fields, and click **Submit Connection**.

For Cloud Storage Object Connections Manager

Manager API user

A user ID that has permissions to connect to the data source system.

Manager API Password

The password for the user ID specified above.

UUID

The unique ID of the DSNet cluster. To obtain the UUID, log in to the COS Manager GUI and click **Help > About this system** on the upper-right corner of the window.

Host

The IP or hostname of the manager node within the DSNet.

Vault

The specific data vault represented by this connection.

Site

An optional physical location tag that an administrator can provide to see the physical distribution of their data.

Accesser

The IP address or hostname of the Accesser® node on DSNet.

Accesser access key

The Accesser access key that has permission to access data in the data vault that is to be scanned. If the accessor access key value is blank, the value is retrieved (for the manager API user) from the manager API.

Accesser secret key

The Accesser secret key that has permission to access data in the vault that is to be scanned. If the secret access key value is blank, the value is retrieved (for the manager API user) from the manager API.

Scanning an IBM Cloud Object Storage data connection

You can initiate an IBM connection scan to collect system metadata from an IBM Cloud Object Storage system.

About this task

When you initiate a scan from the IBM Spectrum Discover graphical user interface (GUI), the metadata is transferred asynchronously back to the IBM Spectrum Discover instance.

Note:

IBM Spectrum Discover does not support scanning of vaults in a dsNet that has any of the following things:

- Proxy vault
- Mirrored vault
- Vault setup for migration

Automated scanning and data ingestion relies on an established and active network connection between the IBM Spectrum Discover instance and the IBM Cloud Object Storage storage source. If the connection cannot be established, the state of the data source connection shows as unavailable, and the option for automated scanning does not appear in the IBM Spectrum Discover GUI for that connection.

Note: If a scan does not complete successfully, check the log file for errors and warnings. If the error or warning message indicates a need to check the configurations file or the settings file, then you must modify the file as required. For example, in some cases you must update the override warnings in the settings file by adding: "override_warnings": true at the root level.

The settings or the configuration file is available in the following location: `/opt/ibm/metaocean/data/connections/cos/scan/scanner-settings.json`

Procedure

1. Log in to the IBM Spectrum Discover graphical user interface (GUI).
2. Under **Admin**, select **Data Source connections**.

The following example shows the **Admin** data connections menu page:

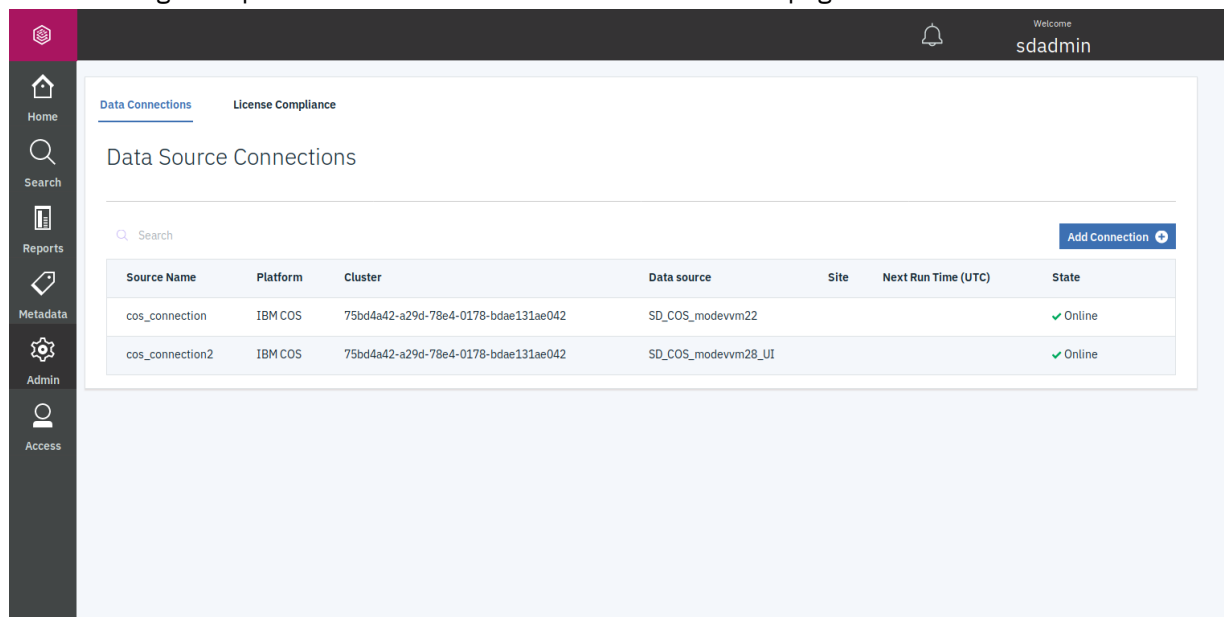


Figure 30. Data source connections

3. Select the data source connection that you want to scan. Ensure that the State is listed as **Online** to make your system scan ready.

The following example shows how to connect to the IBM Spectrum Discover library.

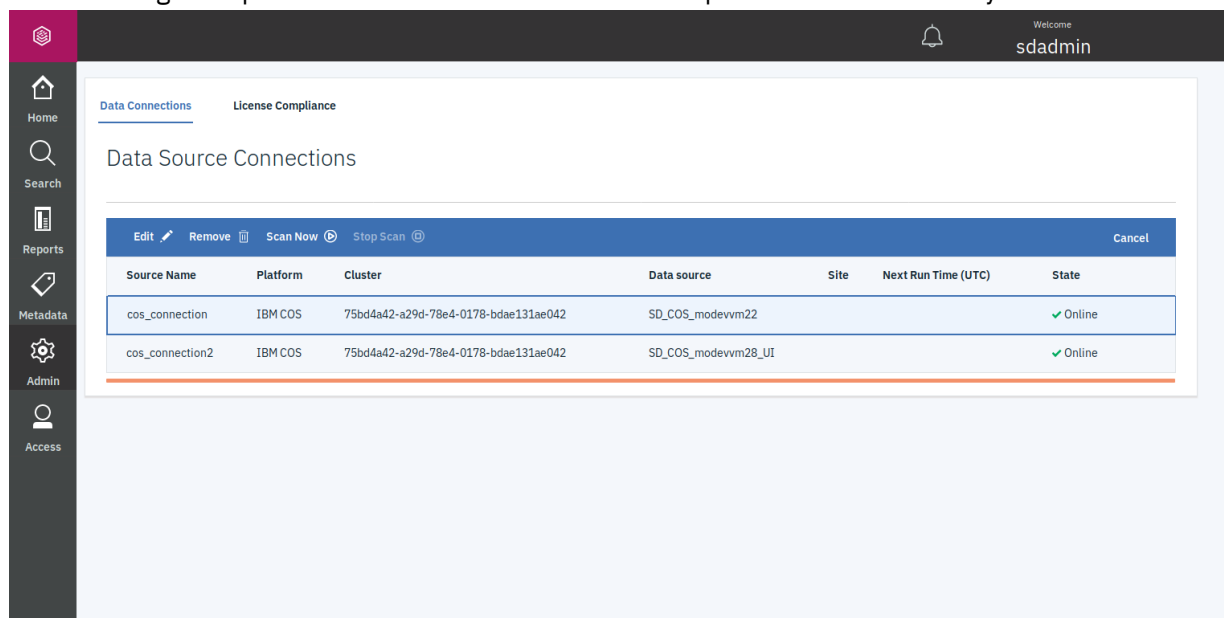


Figure 31. Selecting a data source connection to scan

4. Select **Scan Now** to change the status to **Scanning**.

The following example shows an active scan.

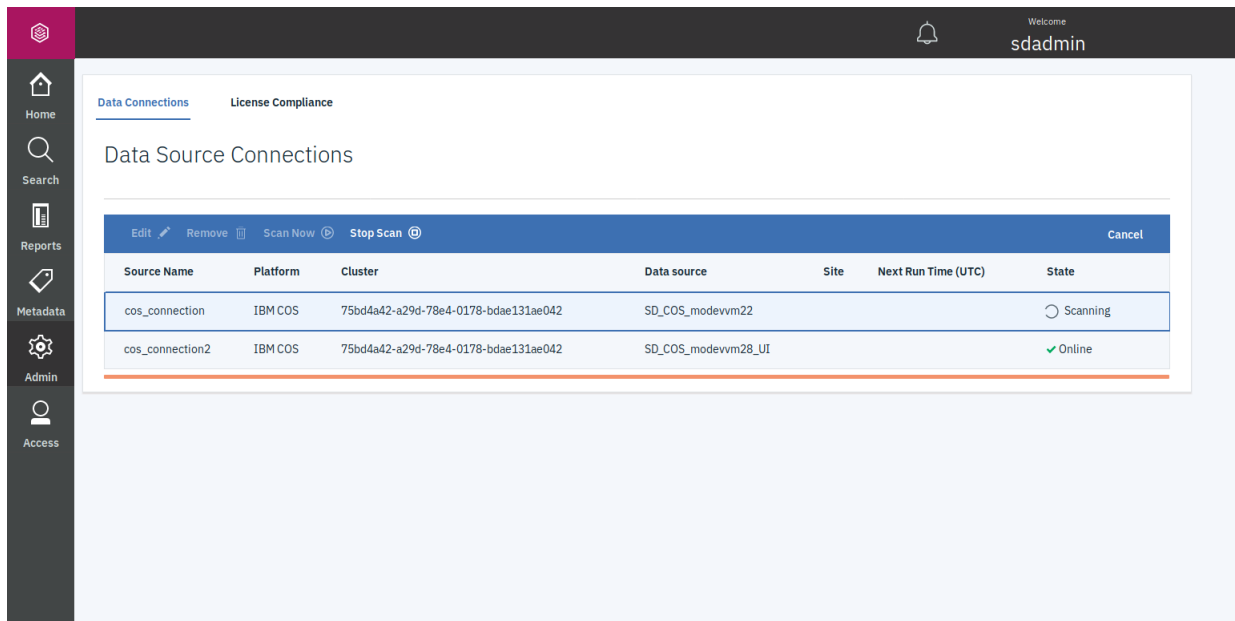


Figure 32. Active scans

5. When the scan finishes, the state field returns to a status of **Online**.

Best practices for scanning IBM Cloud Object Storage systems

Use best practices for scanning IBM Cloud Object Storage (IBM COS) systems.

It is recommended to check the log files in the following directories after each scan:

`/opt/ibm/metaocean/data/connections/cos/<connection_name>/debug/<scan_timestamp>/scanner.debug` indicates whether the scan was successful or not.

`/opt/ibm/metaocean/data/connections/cos/<connection_name>/error/<scan_timestamp>/scanner.error` contains a list of all the messages that are not delivered to IBM Spectrum Discover.

This file contains a list of all the messages that are not delivered to IBM Spectrum Discover.

`/opt/ibm/metaocean/data/connections/cos/<connection_name>/data/<scan_timestamp>/` contains a subfolder with the scanned data source name. There is a stats folder inside this folder that contains information about the number of objects in the data source or the number of objects or scanned files.

You can also compare the total size of the bucket that is reported in IBM Spectrum Discover with the total size of the IBM COS at its source (if it is available).

Enabling bucket notifications for Ceph Object Storage

Use this information to enable bucket notifications for Ceph Object Storage.

Before you begin

Make sure you have the following features set up:

- IBM Spectrum Discover 2.0.2.1
- Red Hat Ceph Storage 4.0 (available starting with version Beta 8)
- A Ceph Object Gateway node that is set up with an HTTPS endpoint

Restriction:

- Ceph Object bucket names must be unique across all data sources. You cannot use the same bucket name to reach a Ceph data source. For example, if there is a IBM Cloud Object Storage or Amazon S3

bucket with the name "my_bucket", you cannot reach a Ceph data source with the bucket name "my_bucket".

- Notifications from versioned buckets are not supported.
- Only one IBM Spectrum Discover node can be configured for push notifications from Ceph Object Storage cluster at a time.

About this task

Use the following steps to enable bucket notifications for Ceph Object Storage.

Procedure

1. Create a data source connection to the Ceph Object Storage cluster.

A Ceph Object Storage source is established as an Amazon S3 data source connection.

Remember: Each bucket must have its own data source connection entry in IBM Spectrum Discover.

2. To enable Ceph Object Storage bucket notifications:

- a. Copy the `ca.crt` file from IBM Spectrum Discover node to a directory on the Ceph Object Gateway nodes.
- b. Locate the file in the `/etc/kafka` directory on the IBM Spectrum Discover node.
- c. Give this file a unique name on the Ceph node after it is copied over.

Remember: Make sure that the file has the same name and in the same location on each Ceph Object Gateway node.

You can choose to use `/etc/ssl/certs` as the copy target directory on the Ceph Object Gateway node.

3. Create a topic entity by using Ceph bucket notification REST API. The topic contains the push endpoint on IBM Spectrum Discover where the notifications are sent to.

Remember: To enable notifications to be sent to IBM Spectrum Discover you must provide push endpoint parameters when you create the topic entity.

These parameters include the IBM Spectrum Discover Kafka topic and credentials that are required to securely produce messages to the topic. For more information about the REST API, see <https://docs.ceph.com/docs/master/radosgw/notifications/#create-a-topic>.

The following parameters must be in the POST request:

```
POST
Action=CreateTopic
&Name=ceph-le-connector-topic
&push-endpoint=<endpoint>
&Attributes.entry.5.key=use-ssl&Attributes.entry.5.value=true
&Attributes.entry.6.key=ca-location&Attributes.entry.6.value=<file path>
```

In this example:

<endpoint>

Indicates the URI of the IBM Spectrum Discover Kafka broker in this format: `kafka://cos:<password>@<discover_fqdn>:9092`

<password>

Indicates the password that can be obtained by an administrator on the IBM Spectrum Discover node from the following location: `/etc/kafka/sasl_password`

<discover_fqdn>

Indicates the fully qualified domain name of the IBM Spectrum Discover node.

<file path>

Indicates the location and file name of the Kafka certificate authority (CA) file on the Ceph Object Gateway Node.

The following example shows topic creation by using the **s3curl** utility:

```
$ ./s3curl.pl --id=rhceph -- -k -X POST https://<ceph object gateway address>:8080/ -d
"Action=CreateTopic&Name=ceph-le-connector-topic&push-endpoint=kafka://cos:
<password>@<discover_fqdn>:9092&Attributes.entry.5.key=use-ssl&Attributes.entry.5.value=true&
Attributes.entry.6.key=ca-location&Attributes.entry.6.value=/etc/ssl/certs/ca.crt"
```

The **--id** parameter identifies the credentials to use in the **s3curl** configuration file.

4. Create a notification entity by using the Ceph bucket REST API. This associates events on a specific bucket to a topic. For more information, see: <https://docs.ceph.com/docs/master/radosgw/s3/bucketops/#create-notification>

The following example shows how to establish a bucket notification by using the **s3curl** utility:

```
$ ./s3curl.pl --id=rhceph --put=notif.xml -- -k https://<ceph object gateway address>:8080/
<bucket>?notification
```

```
Contents of notif.xml:
<NotificationConfiguration xmlns="http://s3.amazonaws.com/doc/2010-03-31/">
  <TopicConfiguration>
    <Id>id1</Id>
    <Topic>arn:aws:sns:default::ceph-le-connector-topic</Topic>
  </TopicConfiguration>
</NotificationConfiguration>
```

You can now capture events on objects within the configured buckets.

Replaying IBM Cloud Object Storage notifications

Use the IBM Cloud Object Storage (COS) Replay feature to resend notifications that failed because of an outage or loss of data.

The IBM COS Replay reads object metadata from vaults and submits the metadata to IBM Spectrum Discover by using Kafka notifications.

Overview of architecture

This topic describes a high-level overview of IBM Cloud Object Storage Scanner architecture.

The following figure shows a high-level overview of IBM Cloud Object Storage Replay architecture.

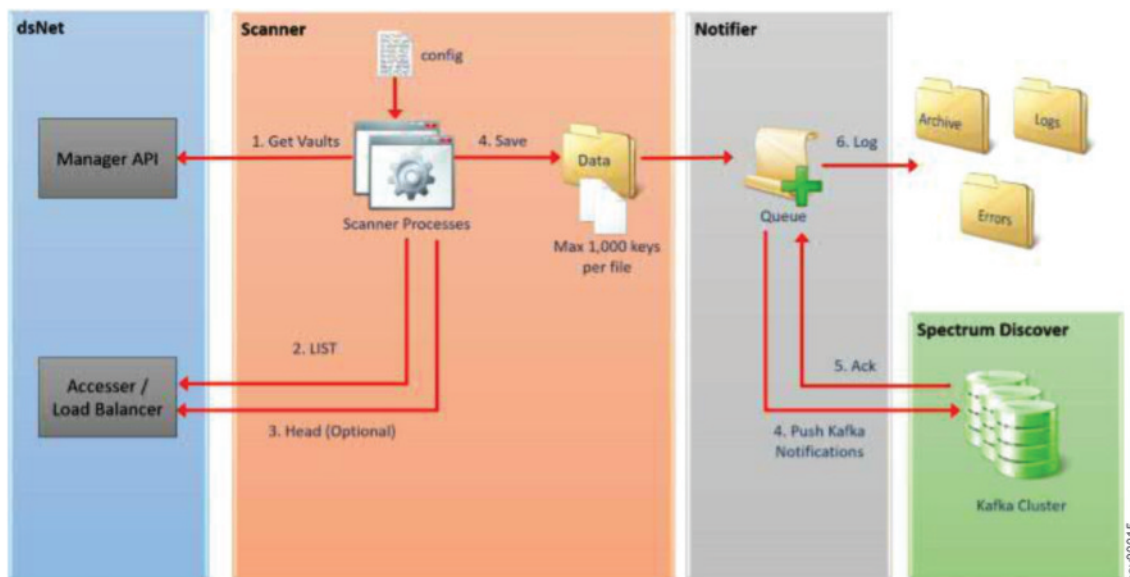


Figure 33. IBM Cloud Object Storage Replay architecture

The `/opt/ibm/metaocean/data/connections/cos/replay/output/data` folder puts the replay output and the notifier reads Kafka messages from the directory. Putting replay output onto disk means that a cold restart is possible.

The IBM Cloud Object Storage Replay consists of two major components:

Replay

Downloads system's logs and re-creates notifications that are sent during a defined time period.

Notifier

Submits the extracted information to IBM Spectrum Discover.

Configuration file

The configuration file is used by the Notifier and Replay.

The configuration file includes:

- Information regarding the net
- Runtime parameters for the Notifier and Replay
- A list of vaults to scan

The configuration file is named `scanner-settings.json` and must sit in the `/opt/ibm/metaocean/data/connections/cos/replay` directory.

The rules for IBM Cloud Object Storage Replay settings are:

- All access logs are scanned.
- All objects that are created or updated since Coordinated Universal Time (UTC) 00:00:01 from 11 April 11 2018 to Coordinated Universal Time (UTC) 10:01:53 on 21 September 2018 are scanned in batches of 1000.
- Custom metadata is retrieved for each object or version.
- Ten vaults are processed in parallel.
- Each vault has a single process LIST that issues requests and 15 processes that issue HEAD requests.

The following example shows every setting. Most settings have default values and can be omitted, but these screens show a typical example by using default values.

Example of the Cloud Object Storage Replay settings

```

{
  "system": {
    "name": "Test dsnet",
    "uuid": "00000000-0000-0000-0000-000000000000",
    "manager_ip": "172.1.1.1",
    "accesser_ip": "172.1.1.2",
    "accesser_supports_https": false,
    "manager_username": "admin",
    "manager_password": "password",
    "is_ibm_cos": true
  },
  "timestamps": {
    "min_utc": "2018-01-01T00:00:00Z",
    "max_utc": "2018-09-21T10:01:53Z"
  },
  "policy_engine": {
    "spectrum_discover_host": "modevvm32.tuc.stglabs.ibm.com"
    "user": "sdadmin",
    "password": "password"
  },
  "scanner": {
    "max_requests_per_second": 5000,
    "max_parallel_list": 10,
    "parallel_head_per_list": 5,
    "list_objects_size": 100
  },
  "notifier": {
    "kafka_format": 1,
    "kafka_endpoint": "192.168.1.1:9092",
    "kafka_topic": "cos-le-connector-topic",
    "kafka_username": "cos",
    "kafka_password": "password",
    "kafka_pem": "-----BEGIN CERTIFICATE-----...\n-----END CERTIFICATE-----\n"
  },
  "logging": {
    "debug_log_max_bytes": 10000000,
    "debug_log_backup_count": 10000,
    "notification_log_max_bytes": 10000000,
    "notification_log_backup_count": 10000,
    "notification_log_all": true
  },
  "include_all_vaults": false,
  "has_custom_metadata": true,
  "override_warnings": true,
  "exclude_vaults": ["Manager"],
  "vaults": [
    {
      "vault_name": "Vault-1"
    },
    {
      "vault_name": "Vault-2",
      "has_custom_metadata": false
    },
    {
      "vault_name": "Vault-3",
      "has_custom_metadata": false,
      "prefix": "customers/live"
    }
  ]
}

```

Typical Cloud Object Storage configuration settings

```

{
  "dsnet": {
    "name": "Test dsnet",
    "uuid": "00000000-0000-0000-0000-000000000000",
    "manager_ip": "172.1.1.1",
    "accesser_ip": "172.1.1.2",
    "accesser_supports_https": false,
    "manager_username": "admin",
    "manager_password": "password",
    "is_ibm_cos": true
  },
  "timestamps": {
    "min_utc": "2018-01-01T00:00:00Z",
    "max_utc": "2018-09-21T10:01:53Z"
  },
  "policy_engine": {
    "spectrum_discover_host": "modevm32.tuc.stglabs.ibm.com"
    "user": "sdadmin",
    "password": "password"
  },
  "scanner": {
    "max_requests_per_second": 5000,
    "max_parallel_list": 10,
    "parallel_head_per_list": 5,
    "list_objects_size": 100
  },
  "notifier": {
    "kafka_format": 1,
    "kafka_endpoint": "192.168.1.1:9092",
    "kafka_topic": "cos-le-connector-topic",
    "kafka_username": "cos",
    "kafka_password": "password",
    "kafka_pem": "-----BEGIN CERTIFICATE-----...\n-----END CERTIFICATE-----\n"
  },
  "logging": {
    "debug_log_max_bytes": 10000000,
    "debug_log_backup_count": 10000,
    "notification_log_max_bytes": 10000000,
    "notification_log_backup_count": 10000,
    "notification_log_all": true
  },
  "include_all_vaults": false,
  "has_custom_metadata": true,
  "override_warnings": true,
  "exclude_vaults": ["Manager"],
  "vaults": [
    {
      "vault_name": "Vault-1"
    },
    {
      "vault_name": "Vault-2",
      "has_custom_metadata": false
    },
    {
      "vault_name": "Vault-3",
      "has_custom_metadata": false,
      "prefix": "customers/live"
    }
  ]
}

```

```

{
  "dsnet": {
    "manager_ip": "192.168.2.106",
    "accesser_ip": "192.168.2.111"
  },
  "timestamps": {
    "min_utc": "2018-04-11T00:00:01.000Z",
    "max_utc": "2018-09-21T10:01:53Z"
  },
  "scanner": {
    "max_requests_per_second": 5000
  },
  "include_all_vaults": true
}

```

```

{
  "system": {
    "manager_ip": "192.168.2.106",
    "accesser_ip": "192.168.2.111"
  },
  "policy_engine" : {
    "spectrum_discover_host": "modevwm32.tuc.stglabs.ibm.com"
  },
  "timestamps": {
    "min_utc": "2018-04-11T00:00:01.000Z",
    "max_utc": "2018-09-21T10:01:53Z"
  },
  "scanner":{
    "max_requests_per_second": 5000
  },
  "include_all_vaults": true
}

```

IBM Cloud Object Storage Scanner is highly configurable. Each element in the file is described in [Table 16](#) on page 70.

Remember: IBM Spectrum Discover does not support file or file path names that use characters that are not part of the UTF-8 character set.

Table 16. Explanation of the configuration file					
Element	Description	Optional	Default value	Restart scanner if changed	Restart notifier if changed
System section					
name	Free-text name of the dsNet. Appears in the 'system_name' element in all Kafka messages.	✓	Retrieved from Manager API if configured. If not, the name does not appear in Kafka messages.	✓	✗
uuid	UUID of the dsNet. Appears in the 'system_uuid' element in all Kafka messages.	✓	Retrieved from Manager API.	✓	✗
manager_ip	Single IP address or host name of the manager device.	✗	Not applicable	✓	✗
accesser_ip	Single IP address or host name of an accesser device or load balancer that routes to the accessers.	✗	Not applicable	✓	✗
accesser_supports_https	Boolean value that indicates whether http or https can be used when you send requests to the accesser or load balancer.	✓	true	✓	✗
manager_username	Username for accessing the manager API. For testing only. Not to be used in production.	✓	Supplied by user at prompt	✓	✗
manager_password	Password for accessing the Manager API. For testing only. Not to be used in production.	✓	Supplied by user at prompt	✓	✗

Table 16. Explanation of the configuration file (continued)

Element	Description	Optional	Default value	Restart scanner if changed	Restart notifier if changed
System section					
is_ibm_cos	Boolean value that indicates whether the system is an IBM Cloud Object Storage or another s3 compliant system. If true, the IBM Get Bucket Extension is used to retrieve object keys from the vaults. Note: Setting the value to false is not currently supported by the Scanner and Notifier.	✓	True	✓	✗
accesser_access_key	Access key ID for S3 calls to the accesses or load balancer. For testing only. Not to be used in production.	✓	Supplied by user at prompt if you cannot retrieve it from Manager API for the user account that is specified in dsNet/manager_username.	✓	✗
accesser_secret_key	Secret key for S3 calls to the accesser or load balancer. For testing only. Not to be used in production.	✓	Supplied by user at prompt if you cannot retrieve from Manager API.	✓	✗
Time stamps section					
min_utc	Only objects or version in the vaults that have a LastModified datetime on or after this timestamp is submitted to IBM Spectrum Discover. Needs to be less than the max_utc value. Note: Changing min_utc and restarting scanner applies only to objects not yet scanned. Objects scanned before restart might have a LastModifiedDate value that is earlier than the min_utc value.	✗		✓ See note.	✗

Table 16. Explanation of the configuration file (continued)

Element	Description	Optional	Default value	Restart scanner if changed	Restart notifier if changed
System section					
max_utc	Only objects or version in the vaults that have a LastModified datetime on or before this time stamp is submitted to IBM Spectrum Discover. Needs to be more than min_utc and less than current time. Note: Changing max_utc to a more recent time and restarting does not mean that new objects written since the old max_utc is scanned. The scanner continues from the last object's key that is scanned in lexicographic order. This means that new objects with names smaller than the last object scanned are not scanned.	✓		✓ See note.	✗
Policy engine section					
		(Only required for IBM Spectrum Discover release 2.0.0.3 and later)			
spectrum_discover_host	Host name or IP address of the policy engine service from which the Kafka certificate is retrieved.	✗	none	✓	✓
user	Username for authorization on policy engine.	✗	none	✓	✓
password	Password for authorization on policy engine.	✗	none	✓	✓
Replay section					
access_log_directory	The access_log_directory is where the dsNet access log files are stored after download. Access logs must be in the root input folder. Files in subdirectories are not processed.	✓	[IBM Cloud Object Storage Replay]/ access_logs	Restart Replay if changed	Restart Replay if changed

Table 16. Explanation of the configuration file (continued)

Element	Description	Optional	Default value	Restart scanner if changed	Restart notifier if changed
System section					
download	If download is set to false, access logs are not downloaded and are assumed to already be present in access_log_directory.	✓	true	Restart Replay if changed	Restart Replay if changed
Notifier section					
kafka_format	Format of the Kafka message.	✓	1	✗	✓
kafka_endpoint	IP address and port of the Kafka endpoint.	✓	Retrieved from Manager API	✗	✓
kafka_topic	Name of the Kafka topic.	✓	Retrieved from Manager API	✗	✓
kafka_username	The username for authentication with Kafka. Note: For testing only. Not to be used in production.	✓	Supplied by user at prompt if you cannot retrieve from Manager API.	✗	✓
kafka_password	The password for authentication with Kafka. Note: For testing only. Not to be used in production.	✓	Supplied by user at prompt if cannot be retrieved from Manager API.	✗	✓
kafka_pem	The certificate PEM for authentication with Kafka. Must include '\n' characters to ensure correct formatting. Note: For testing only. Not to be used in production.	✓	Supplied by user at prompt if it cannot be retrieved from the system	✗	✓
Logging section					
debug_log_max_bytes	The scanner.debug and notifier.debug roll over when this size is reached.	✓	1,000,000	✓	✓
debug_log_backup_count	The number of scanner.debug and notifier.debug files to retain.	✓	10	✓	✓
notification_log_max_bytes	The notification.log rolls over when this size is reached.	✓	1,000,000	✓	✓
notification_log_backup_count	The number of notification.log files to retain.	✓	10	✓	✓

Table 16. Explanation of the configuration file (continued)

Element	Description	Optional	Default value	Restart scanner if changed	Restart notifier if changed
System section					
notification_log_all	Boolean value that controls the level of Notifier logging. When true: an entry is written to notification.log for message you send to the Kafka cluster. When false: only failed sends are written to notification.log.	✓	False	x	✓
Root-level items					
include_all_vaults	Boolean value that determines whether all vaults in the dsNet are scanned. If false, the details of the vaults to be scanned must be specified in the 'vaults' element. Boolean value that determines whether custom metadata and content type are retrieved for each object by using individual HEAD requests.	✓	False	✓	x
has_custom_metadata	This value is only relevant when a versioned vault is scanned. For IBM Cloud Object Storage systems, non-versioned vaults always require a HEAD request for every object. Can be overridden for each vault in the 'vaults' element.	✓	True	✓	x
override_warnings	Boolean value that allows the scanner to run and ignore any warnings that are generated on start-up. For example, a warning is raised on start-up if versioning is suspended on a vault.	✓	False	✓	x
exclude_vaults	Comma-separated list of vault names to be excluded from scanning, such as: "exclude-vaults": ["COSVault", "COSVault-V"]	✓	[] Empty list	✓	x

Table 16. Explanation of the configuration file (continued)

Element	Description	Optional	Default value	Restart scanner if changed	Restart notifier if changed
System section					
vaults	<p>List of vaults to be scanned. If <code>include_all_vaults</code> is true, the vaults list can be left empty.</p> <p>This list can be used to define more detailed scanning parameters for individual vaults. Any settings that are defined here take precedence over the settings that are described.</p> <p>Each element in the list contains:</p> <p>The <code>vault_name</code> is the name of the vault.</p> <p>The <code>has_custom_metadata</code> is an optional Boolean that overrides the <code>has_custom_metadata</code> that is described.</p> <p>The <code>prefix</code> is an optional string that is used to filter the objects or versions that are retrieved from the vault.</p>	✓	Dependent on settings <code>include_all_vaults</code> and <code>exclude_vaults</code>	✓	✗

Replay performance

The number of requests that are issued by IBM Cloud Object Storage Replay is throttled to ensure that overall dsNet performance remains at the agreed level.

You can control throttling by the number of settings in a configuration file. All settings are optional. The following screen shows an example of the default values.

```
"replay": {
  "max_requests_per_second": 1000,
  "max_parallel_list": 10,
  "parallel_head_per_list": 15,
  "list_objects_size": 1000
}
```

Process count

The following list shows an example of how 161 processes are divided. [Figure 34 on page 76](#) shows a caution message of how the number of processes should not exceed 161.

- One main process
- 10 List worker processes
- 150 HEAD worker processes

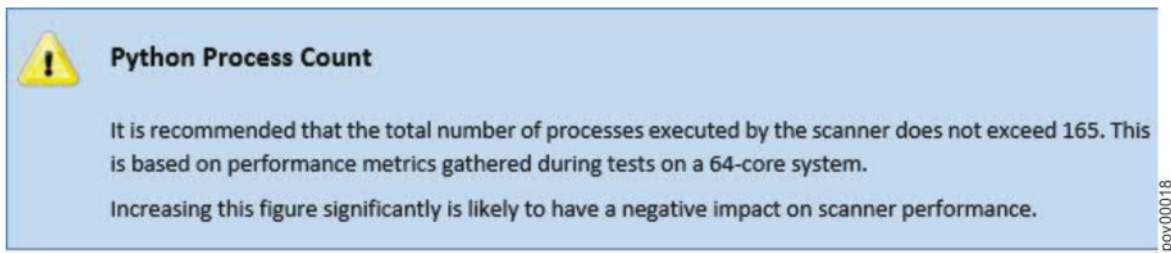


Figure 34. Python process count

Maximum Replay performance

Replay and Notifier maximize performance on a 64 core Intel(R) Xeon(R) CPU E5-2690 v3 @ 2.60GHz server is 2300 objects that are scanned and notified per second with a dsNet with 6 accessers and 12 slicestors under customer load at 50 percent capacity.

```
"replay": {
  "max_requests_per_second": 2300,
  "max_parallel_list": 10,
  "parallel_head_per_list": 15,
  "list_objects_size": 1000
}
```

The recommendation is to start the replay at a rate of 1000 objects scanned per second. Measure the latency degradation of customer traffic and increase the scanning rate until the maximum acceptable degradation is reached.

One thousand objects per second on the net, which is a 5 - 27 percent increase of write operations, the latency (larger increase for smaller size files) and around 10 percent for read operations latency were measured.

At 2000 objects a second, a 10 - 50 percent increase of write operations latency and in the range 18 - 28 percent, and 10 percent for read operations latency were measured.

Replay tasks and vault settings

A few scenarios exist that prevent the Replay from operating correctly.

Certain combinations of the following IBM Cloud Object Storage vault settings prevent the Replay from running a full scan:

- Vault versioning
- Name index
- Recover listing

Figure 35 on page 77 shows settings for the first three items on the vault configuration page in the DsNet Manager user interface.

Figure 35. Settings for three items on the vault configuration page in the net Manager user interface

The scenarios that are invalid are reported at startup.

Remember: You must correct the scenarios before you can run the Replay.

If a scan does not complete successfully, make sure that you check the log file for errors and warnings. In some cases, you must modify the settings file as detailed in the errors and warning messages. The settings file is located at: /opt/ibm/metaocean/data/connections/cos/scan/scanner-settings.json

Table 17 on page 77 shows the behavior for the Replay for different combinations of the four variables.

Table 17. Behaviors for Replay for four variables				
ID	Name index	Recovery listing	Versioned	Cloud Object Storage Replay behavior
0	X	X	X	<p>❑ Stop start-up and report error in config file:</p> <p>Error: Objects cannot be listed because Name Index and Recovery Index are both disabled. You might enable Recovery Listing on the vault or add this vault to the "exclude_vaults" list in the configuration file.</p> <p>For example:</p> <pre>"exclude-vaults": ["vault-name"]"</pre>
1	X	X	✓	<p>❑ Stop start-up and report error in config file:</p> <p>Error: Objects cannot be listed because Name Index and Recovery Index are both unavailable. You might enable Recovery Listing on the vault or add this vault to the "exclude_vaults" list in the configuration file.</p> <p>For example:</p> <pre>"exclude-vaults": ["vault-name"]"</pre>
2	X	✓	X	<p>☑ Object Listing is run.</p>

Table 17. Behaviors for Replay for four variables (continued)

ID	Name index	Recovery listing	Versioned	Cloud Object Storage Replay behavior
3	x	✓	✓	<p>☑ Object Listing runs. Only the most recent version of each object is listed. A warning is logged:</p> <p>Warning: Versions cannot be listed as Name Index is unavailable. An object scan is run and only the most recent version of each object is listed. You must add <code>override_warnings: true</code> in the configuration file to ignore this warning.</p> <p>Switching Name Index on does not enable scanning of a full version history. Objects created while Name Index is off is not present when it is enabled.</p>
4	✓	x	x	☑ Object Listing is run.
5	✓	x	✓	☑ Object Listing is run.
6	✓	✓	x	☑ Object Listing is run.
7	✓	✓	✓	<p>! Stop start-up and report warning:</p> <p>This is a versioned vault but version scanning is not possible as Recovery Listing is enabled. You might either disable Recovery Listing on the vault to allow version scanning, or rerun the Replay with the argument <code>override-warnings: true</code> to allow object scanning.</p>

Important: You might receive system errors about records not being scanned to the database if you scan a IBM Cloud Object Storage vault with **Name Index** disabled and **Recovery Listing** enabled. You cannot specify a prefix listing on a vault that has **Recovery Listing** enabled (you cannot specify a blank prefix either).

Including and excluding vaults

You can set the vaults that you scan with various settings in the configuration file.

Use the following settings in the configuration file to scan the vaults:

- `include_all_vaults` (Boolean)
- `exclude_vaults` (List)
- `vaults` (Dictionary)

When `include_all_vaults` is true, all vaults in the system are scanned except for any vaults specified in the `exclude_vaults` list.

You might consider `exclude_vaults` a list of vaults to ignore and `vaults` is a list that specifies details of individual vaults to be scanned.

If `include_all_vaults` is true and the `vaults` list is populated, the list of vaults that are scanned is the superset of all vaults that are returned by the Manager that are merged with the `vaults` list from the config file.

An error is raised and the Scanner aborts on start-up if the same vault appears in both `vaults` and `exclude_vaults`.

Mirror, Proxy, Data Migration

IBM Cloud Object Storage Scanner does not support scanning of the following:

- Mirrored vaults
- Proxy vaults
- Vaults that are set up for migration

Any vaults of these types are ignored by the scanner and a warning logged in the debug log.

Examples for including and excluding vaults

To summarize the rules for including and excluding vaults, following are some examples:

Example 1

- The system contains 1000 vaults.
- Five of the 1000 vaults are management vaults (named mgmt-1 to mgmt-5).
- The scan includes all vaults except the management vaults.

```
"include_all_vaults": true,  
"exclude-vaults": ["mgmt-1", "mgmt-2", "mgmt-3", "mgmt-4", "mgmt-5"]
```

Example 2

- The system contains 1000 vaults.
- 5 of the 1000 vaults are management vaults (named mgmt-1 to mgmt-5).
- The scan includes all vaults except the management vaults.
- The scan includes a filter for scanning a vault that is named vault-x.
- The scan includes only a scan of the objects whose key starts with **production/finance**.

```
"include_all_vaults": true,  
"exclude-vaults": ["mgmt-1", "mgmt-2", "mgmt-3", "mgmt-4", "mgmt-5"],  
"vaults": [  
  {"vault_name": "vault-x", "prefix": "production/finance"}  
]
```

Example 3

- The system contains 1000 vaults.
- 5 of the 1000 vaults are management vaults (named mgmt-1 to mgmt-5).
- The scan includes all vaults except the management vaults.
- The scan includes a filter for scanning a vault that is named vault-x.
- The scan includes only a scan of the objects whose key starts with **production/finance** or **production/marketing**.

```
"include_all_vaults": true,  
"exclude-vaults": ["mgmt-1", "mgmt-2", "mgmt-3", "mgmt-4", "mgmt-5"],  
"vaults": [  
  {"vault_name": "vault-x", "prefix": "production/finance"},  
  {"vault_name": "vault-x", "prefix": "production/marketing"}  
]
```

Example 4

- The system contains 1000 vaults.
- Run a test on three vaults named vault-a, vault-b, and versioned-vault-c.

- Run a scan on versioned-vault-c and issue LIST requests. Do not issue HEAD requests because the objects do not have custom amz headers.

```
"include_all_vaults": false,
"vaults": [
  {"vault_name": "vault-a"
  {"vault_name": "vault-b"
  {"vault_name": "vault-c",  "has_custom_metadata": false}
]
```

Stats files

The IBM Cloud Object Storage Scanner tracks each LIST process status to a stats file.

During a scan, the Scanner runs multiple processes. Each LIST processes and tracks the progress, saves the next_key, and optionally the next_version to a stats file named task.stats that is stored with the log files in the /opt/ibm/metaocean/data/connections/cos/replay/output/data directory.

```
{
  "estimated_object_count": 1000,
  "list_objects_size": 100,
  "next_key": "",
  "next_version": "",
  "prefix": "",
  "scan_type": "Object Scan",
  "status": "Complete",
  "total_bytes_output": 1126809,
  "total_bytes_scanned": 1126809,
  "total_objects_output": 47,
  "total_objects_scanned": 47,
  "vault_name": "dsmgmt-sp1",
  "vault_uuid": "868daa21-9e56-4c41-b6fd-845a4c85cea9"
}
```

From the Scanner, you can start, stop, recover files from a crash, and restart at the point where the scan was interrupted.

When you start the scanner:

1. Processing of the Scanner continues from next_key and next_version.
2. Queue of the Notifier is optimized by reloading from the files in the data folder instead of requerying the dsNet.
3. Batches that were processed partially are reprocessed. Duplicate Kafka notifications might occur, but are handled safely by the IBM Spectrum Discover system.

Replay

When a severe outage occurs and causes the loss of notifications sent by the system to IBM Spectrum Discover, the IBM Cloud Object Storage Scanner Replay feature can be used to recover lost notifications.

Replay parses the access logs of a system and reconstitutes the notifications. Also, the Notifier can resend the notifications.

Initialization for Replay

During the startup, Replay reads the configuration file and issues requests to the Manager of the dsNet device similar to the Scanner.

Data from the configuration file is validated to ensure that appropriate permissions are granted in the dsNet. This allows access to management vaults and regular vaults. Startup errors or warnings are logged and printed to the console.

After initialization, Replay extracts accesser log files from the management vaults of dsNet and enables Replay to process and write notifications to the output directory.

Error conditions

Sometimes Replay does not have enough information to replay the original notification. If this occurs, you must fix the problems manually.

For example, if vault versioning was suspended when you made the request and you receive an s3 DeleteObject for an object or delete marker, the following error is logged:

```
error_code=True, error_description="Delete operation with [no version_id|null
version_id|version_id] for vault with versioning = [suspended/enabled]"
```

The error message displays because Replay cannot distinguish when a notification with s3:CreateDeleteMarker or s3:CreateDeleteMarker:NullVersionDeleted is sent.

If vault versioning is disabled, and an s3 PutObject request is received for an object that is deleted, the following error is logged:

```
error_code=404, error_message="Not Found"
```

The error message displays because Replay cannot determine the tag of the object that was deleted.

Output

Messages are batched by 1,000 or to the Scanner list objects size configuration setting, if specified.

The messages are written to the output folder with the same Notification format used by the Scanner.

```
{
  "system_name": "Test",
  "object_etag": "\"de37d2cee49596916f62a233dfc790a4\"",
  "request_time": "2018-09-24T18:49:29.383Z",
  "format": 1,
  "bucket_uuid": "ac89915b-d4ec-7ff1-00be-9c32b2aca580",
  "system_uuid": "f7d033c2-9066-499a-a883-829860d4d865",
  "object_length": "12319",
  "object_name": "test_version",
  "bucket_name": "vault3",
  "content_type": "binary/octet-stream",
  "request_id": "17451c3d-e81e-40ed-939a-4534780daaa8",
  "operation": "s3:PutObject"
}
```

If an error occurs, the error messages are written to the /opt/ibm/metaocean/data/connections/cos/replay/output/data/access_log_error/ directory. Take note of the extra error_code and error_description elements.

```
{
  "system_name": "Test",
  "object_version": "null",
  "request_time": "2018-09-24T17:07:59.471Z",
  "format": 1,
  "bucket_uuid": "ac89915b-d4ec-7ff1-00be-9c32b2aca580",
  "system_uuid": "f7d033c2-9066-499a-a883-829860d4d865",
  "object_length": "12319",
  "object_name": "object5.2",
  "bucket_name": "vault3",
  "request_id": "ebc472a1-f955-4605-895b-840867b12e01",
  "operation": "s3:PutObject",
  "error_description": "Not Found",
  "error_code": 404
}
```

Renaming a vault for Replay

When you rename a vault, it is possible that Replay can abort.

Replay aborts when you:

- Delete the vault.
- Rename the vault.

- Discover that the read permission is revoked for the credentials that are supplied by the operator or manager API.

All other scans of a vault continue scanning until complete.

You can find the details of the errors that include stack trace in the `replay.debug` file in the `/opt/ibm/metaocean/data/connections/cos/replay/debug/replay/[timestamp]` directory.

Starting the Replay

The guidelines and rules for using Replay are documented in this topic.

To start Replay, run the following command:

```
cos-replay
```

The following rules apply for Replay:

- Configure Replay according to the guidelines in [Table 16 on page 70](#).
- Replay component requires `min_utc` and `max_utc` time stamps defined in the “[Configuration file](#)” on [page 67](#).
- Only notifications sent between `min_utc` and `max_utc` are parsed and replayed.
- Replay automatically shuts down when all accessor logs are downloaded and processed. The message **Complete Replay Process** appears in the console.

This is an example of how to start Replay:

```
Starting COS Replay - Version 0.1 Log file and config file are in directory /Users/weebrew/
Documents/Development/ibmworkspace/cosscanner/output/debug/scanner/20180925-125528-232131
Starting Accessor Log Extraction Downloading files...
('Downloaded', 10, 'of', 36)
('Downloaded', 20, 'of', 36)
('Downloaded', 30, 'of', 36)
Download complete.
Total files: 36
Complete Accessor Log Extraction
Starting Replay process
Complete Replay process
```

Debug mode for Replay

Run Replay in debug mode to troubleshoot problems.

To start debug mode, run the following command:

```
cos-replay --log=DEBUG
```

Running debug mode creates large log files and creates a significant drop in performance. Do not run debug mode for long periods especially when you are in production mode.

Notifier

The Notifier is the component that reads the JSON notifications that are written by Scanner or Replay and sends notifications to the Kafka cluster.

When notifications are acknowledged by Kafka, the Notifier moves the file to the archive folder.

On start-up, Notifier calls the Manager API and retrieves details of any Notification Service Configurations (NSC) configured in the dsNet for IBM Spectrum Discover. If more than one is found, the first one is used.

Retrieval of NSCs is overridden by defining the details of the Kafka configuration in the config file.

```
"notifier":{
  "kafka_format": 1,
  "kafka_endpoint": "192.168.1.34:9092",
  "kafka_topic": "cos-le-connector-topic"
}
```

Limitations

Limitations apply when the Notifier uses a Kafka configuration retrieved from the Manager API.

- If more than one NSC exists, the first one is used for all vaults.
- If more than one host name is defined in the NSC, the first one is used for all vaults.

Starting the Notifier

Running the Notifier has rules and limitations.

To start the Notifier, run the following command:

```
cos-notify
```

After you start the Notifier, you are prompted for security credentials for the manager API and Kafka cluster.

```
Starting COS Notifier - Version 0.1
Enter the Manager API username: admin
Enter the Manager API password:
Enter the Kafka username: cos
Enter the Kafka password:
Enter the Kafka pem:
Creating Kafka producer...
Done
Notifier is running
Log file and config file are in directory C:\dev\cos-scanner\output\debug\notifier
\20180912-121641-283000
Checking for files in \data
- 11 files found Checking for files in output\data
- 256 files found
```

Rules and limitations

The following rules and limitations apply to the Notifier:

- You cannot start the Notifier in the background because the Notifier requires user input at the terminal window.
- You can stop the Notifier and force the Notifier to run in the background.
- The passwords and pem do not display when you type and paste the passwords in the console.
- The certificate pem is approximately 1,600 characters. If you use an SSH connection, the certificate pem might be truncated to 1,000 characters.
- If the number is truncated to 1,000 characters, include the certificate pem in the config file.

Notifier operation

The Notifier enumerates and processes all .log files in the Scanner data directory.

After all files are processed, the Notifier repeats the process so that new .log files that are generated by the Scanner are processed. The Notifier sleeps repeatedly in 1-second intervals if no new files are found in the /opt/ibm/metaocean/data/connections/cos/replay/output/data directory.

The Notifier does not automatically shut down. The Notifier continues to monitor the Scanner data directory for new .log files. Monitor the progress of the Notifier by using the status report. When the operator or administrator determines that all scanned objects are submitted successfully to the IBM Spectrum Discover, shut down the Notifier by using the kill switch.

Stopping the Notifier

You might need to stop and restart the Notifier.

Before you stop and restart the Notifier:

1. Create a file named `kill.notifier` in the `/opt/ibm/metaocean/data/connections/cos/replay/output/command` directory.

2. Ensure that the processing of any batches is complete before you stop the Notifier.

Stopping the Notifier displays the following output:

The shutdown is complete when the "**Shutdown is complete**" message displays.

```
Starting COS Notifier - Version 0.1
Enter the Manager API username: admin
Enter the Manager API password:
Enter the Kafka username: cos Enter the Kafka password:
Enter the Kafka pem: Creating Kafka producer...
Done
Notifier is running
Log file and config file are in directory C:\dev\cos-scanner\output\
debug\notifier\20180912-121641-283000
Checking for files in output\data
- 11 files found Checking for files in output\data
- 256 files found Detected the kill trigger file. Shutting down...
Shutdown is complete
```

Restarting the Notifier

When you stop the Notifier following a shutdown with the kill.notifier file, you must rename the file manually or delete the file before you do a restart.

If you do not rename or delete the kill.notifier file, the system finds the file and displays the following message:

```
C:\dev\cos-scanner>python main_notifier.py Starting COS Notifier - Version 0.1
The file 'kill.notifier' is preventing the notifier from running.
You should delete or rename the file and re-start the notifier.
File location: 'C:\dev\cos-scanner\command\kill.notifier'
```

The Scanner and Notifier are separate solutions that share the config file. The files operate independently, so you can start and stop either file independently any time.

Progress report

The Progress Report provides an instant snapshot of status for the Scanner and Notifier.

To create a progress report, run the following command:

```
cos-report
```

The progress report displays in plain text format to the console in a static HTML file that is named: /opt/ibm/metaocean/data/connections/cos/replay/output/cos-scanner-report.html

If a progress report exists, the new progress report overwrites the existing progress report. See [Figure 36 on page 85](#).

IBM COS Scanner Progress Report

Scans in progress: 6
Scans complete: 17
Scan progress: 12.00%

Scan Type	Vault Name	Vault UUID	Est. Object Count	Scan Status	Last Scan Activity	Scanned	Output	Queued	Notified	Error	Approx % Scanned	Approx % Notified
Object	COSVault-N	e0e08245-53b9-7bf0-0024-c116dc33fa80	21,001	In progress	2018-08-01 11:03:48	13,000	13,000	13,000	0	0	62%	0%
Object	COSVault-NR	ec3e0eb6-799d-7062-0143-3f59b4118180	3	Complete	2018-08-01 11:02:42	2	2	2	2	0	100%	100%
Object	COSVault-NRV	88c53420-884f-749c-11f0-f27fe2875980	25,931	In progress	2018-08-01 11:03:49	13,000	13,000	7,000	6,000	0	50%	46%
Version	COSVault-NV	6bd9011d-0d28-76f8-1132-dcc0540bc380	69	Complete	2018-08-01 11:02:43	68	68	0	0	0	100%	0%
Version	COSVault-NV-Suspended	687a13f7-a287-7bb8-1069-f90847582080	19	Complete	2018-08-01 11:02:42	18	18	18	0	0	100%	0%
Object	COSVault-R	854bcc22-b657-7a29-0029-0c4760284280	5,000	Complete	2018-08-01 11:03:13	5,000	5,000	2,000	3,000	0	100%	60%
Object	COSVault-RV	b8cf437c-65b4-726a-10d8-8c29a406c980	20,001	In progress	2018-08-01 11:03:47	12,000	12,000	3,000	7,000	194	99%	58%
Object	EV01_AWSV4_PERF1	43d01b33-a84e-7d4b-1080-326024c3d880	101	Complete	2018-08-01 11:02:43	100	0	0	0	0	100%	100%
Object	EV01_AWSV4_PLUGIN4	f355f8a7-461a-7aa4-10b0-5ef47c519e80	1	Complete	2018-08-01 11:02:42	0	0	0	0	0	100%	100%
Object	EV01_SMALL_VAULT1	af7db680-57c9-7359-1091-b5b763913b80	201	Complete	2018-08-01 11:02:43	200	0	0	0	0	100%	100%
Object	mega_vault?prefix=test1	97ecd99e-b78b-7f1f-0198-aab16d346080	977,200	In progress	2018-08-01 10:51:33	107,000	107,000	80,000	27,000	0	10%	25%
Object	mega_vault?prefix=test2	97ecd99e-b78b-7f1f-0198-aab16d346080	977,200	Complete	2018-08-01 10:51:13	98,863	98,863	71,863	27,000	0	100%	27%
Object	mega_vault?prefix=test3	97ecd99e-b78b-7f1f-0198-aab16d346080	977,200	In progress	2018-08-01 10:51:32	109,000	109,000	61,000	48,000	0	11%	44%
Object	MGMTO01	a914abe-4dcf-73bd-10d4-f87f1c3f380	14,187	Aborted	2018-08-01 11:03:46	13,000	9,701	9,701	0	0	91%	0%
Version	SuspendedVersioningTest-Vault	ab5fccc0-9018-7c0c-108e-e506470bd280	10	Complete	2018-08-01 11:02:42	9	9	9	0	0	100%	0%
Object	Threading-Test-2.15	91f6b76-0ed6-767c-0153-c29101a74b80	0	Complete	2018-08-01 11:02:42	0	0	0	0	0	100%	100%
Object	Threading-Test-2.15.2	781e51b0-cf6e-78c7-1156-2f4ce53eb80	0	Complete	2018-08-01 11:02:42	0	0	0	0	0	100%	100%
Object	TimestampTest	65e502c4-477e-7c99-11f7-6e91a409ce80	3	Complete	2018-08-01 11:02:42	2	2	2	0	0	100%	0%
Object	Vault-Empty	2892a6de-7ed6-736f-00da-a20deb42080	0	Complete	2018-08-01 11:02:42	0	0	0	0	0	100%	100%
Version	Vault-N	2344f296-a5af-7de0-10d0-75acbcfd180	9	Complete	2018-08-01 11:02:42	8	8	8	0	0	100%	0%
Version	Vault-V	15ad56da-f304-77c3-00c5-f31ddff13480	12	Complete	2018-08-01 11:02:42	11	11	11	0	0	100%	0%
Object	vault1	c7b18026-87bd-7e19-10ae-5b7b350ff880	0	Complete	2018-08-01 11:02:42	0	0	0	0	0	100%	100%
Version	version-delete-test	2db39e94-f498-7dd7-0135-1dae35772280	20	Complete	2018-08-01 11:02:42	19	19	19	0	0	100%	0%

Notes

- Est. Object Count may not accurately reflect the number of objects in the vault. The discrepancy is typically very small.

prov00021

Figure 36. IBM Cloud Object Storage Scanner progress report

See Table 18 on page 85 for a description of information in the progress report.

Table 18. Description for IBM Cloud Object Storage Scanner progress report	
Column name	Description
Scan Type	<p>Either Object or Version.</p> <p>Non-Versioned vaults show Object.</p> <p>Versioned vaults show Version. However, there are some exceptions. If the Name Index for a versioned vault is unavailable but Recovery Listing is enabled, an object scan might be run. The user is alerted that an object scan can be done, but this object scan requires changes to the configuration file.</p>
Vault Name	<p>The name of the vault. Any prefix that is defined in the configuration file is also shown.</p> <p>Example: mega_vault?prefix=test</p>
Vault UUID	The UUID of the vault.
Estimated Object Count	<p>The estimated number of objects in the vault, as reported by the Manager API.</p> <p>This value is refreshed from the Manager API each time the scanner is started, regardless of the status of each scan. Given that the number of objects in each vault might be constantly changing, the number of objects that are reported in this column becomes out of date during long running scans.</p> <p>Note: This issue affects only the status report but does not affect the data integrity of the Scanner.</p>

Table 18. Description for IBM Cloud Object Storage Scanner progress report (continued)

Column name	Description
Scan Status	<p>Shows the status of the scanner.</p> <p>Not started The task is queued but not started.</p> <p>In progress The task is running.</p> <p>Complete The task finished.</p> <p>Aborted The task encountered an unrecoverable error and aborted. Shut down the Scanner and the debug file, and inspect the file to investigate the problems. After you resolve the problems, restart the Scanner.</p> <p>The debug file is in the /opt/ibm/metaocean/data/connections/cos/replay/output/data/<vault-name>/<prefix> directory. For each vault, see the /opt/ibm/metaocean/data/connections/cos/replay/output/data/<vault-name>/<prefix> directories.</p>
Last scan activity	The last time data was retrieved from the vault.
Scanned	Number of objects or versions that are scanned. For a versioned vault, this value shows a figure that is higher than the Estimated Object Count .
Output	<p>Number of objects or versions that scanned AND whose LastModified time stamp is inside the time window that is defined in the configuration file.</p> <p>The figure in the column is Queued + Notified + Error.</p>
Queued	Number of objects or versions that are Output and are waiting to be sent to the Kafka cluster.
Notified	Number of objects or versions that are submitted successfully to the Kafka cluster.
Error	Number of objects or versions that failed to send to the Kafka cluster. Details of all errors are logged to notifier.debug.
Approximate percentage scanned	<p>Scanned as a percentage of Est. Object Count.</p> <p>The cell background shows a progress bar.</p>
Approximate percentage scanned	<p>Notified as a percentage of Output.</p> <p>The cell background shows a progress bar.</p>

Table 19. What is reported beneath the report title

Measure	Description
Scans in progress	Number of scans with the status "In Progress". Applies to Scanner only.
Scans complete	Number of scans with the status "Complete". Applies to Scanner only.
Scan progress	Sum (number of objects scanned) as a percentage of sum (estimated object count).

Logging

You can view the list of directories that are generated by scanner, notifier, and replay.

Table 20 on page 87 lists the directories that are generated on start-up by the Scanner, Notifier, and Replay.

Directory	Description
For IBM Spectrum Discover: /opt/ibm/metaocean/ data/connections/cos/replay/ output/data	Contains .log files (Kafka messages), stats files and debug information for each scanned vault. 
For IBM Spectrum Discover: /opt/ibm/metaocean/ data/connections/cos/replay/ output/debug/[scanner replay]/	Contains Scanner and Replay debug or troubleshooting information. A new subdirectory is created each time the Scanner and Replay starts. Each subdirectory contains a copy of the configuration file and scanner.debug (replay.debug). 
For IBM Spectrum Discover: /opt/ibm/metaocean/ data/connections/cos/replay/ output/debug/notifier	Contains Notifier debug or troubleshooting information. A new subdirectory is created each time the Notifier starts. It contains a copy of the configuration file and notifier.debug. Same directory-naming convention as shown for the scanner. The notifier.debug rolls over when it reaches a predefined size as defined in the configuration file. See “Configuration file” on page 67 .

Table 20. List of directories generated by scanner, notifier, and replay (continued)

Directory	Description
For IBM Spectrum Discover: /opt/ibm/metaocean/data/connections/cos/replay/output/archive/	Contains all .log files that are successfully processed by the Notifier, and their contents are successfully submitted to the Kafka cluster. Remember: Files in this directory are truncated – they contain only the object key and they can also contain the version.
For IBM Spectrum Discover: /opt/ibm/metaocean/data/connections/cos/replay/output/error/	Contains any .log files that failed to submit to the Kafka cluster.
For IBM Spectrum Discover: /opt/ibm/metaocean/data/connections/cos/replay/output/notification_log/	The notification.log file contains details of any errors (including stack trace) that occur when it attempts to send notifications to the Kafka cluster. If logging/notification_log_all is true in the config file, all successful sends are also logged. The notification.log file rolls over when it reaches a predefined size as defined in the configuration file. See “Configuration file” on page 67 .

IBM Cloud Object Storage Scanner output data

The Scanner generates a directory beneath the output data directory for each vault or vault prefix as defined in the configuration file.

The /opt/ibm/metaocean/data/connections/cos/replay/output/data directory is the Scanner output data directory.

The following screen shows an example of a configuration file and also shows that all vaults are scanned, but mega_vault has four separate prefixes that are defined which means the four scans of the vault occurred.

```
"include_all_vaults": true,
  "vaults": [
    {"vault_name": "mega_vault", "prefix": "main/production/finance"},
    {"vault_name": "mega_vault", "prefix": "main/production/sales"},
    {"vault_name": "mega_vault", "prefix": "main/production/marketing"},
    {"vault_name": "mega_vault", "prefix": "main/production/hr"}
  ]
```

[Figure 37 on page 89](#) shows the directory structure.

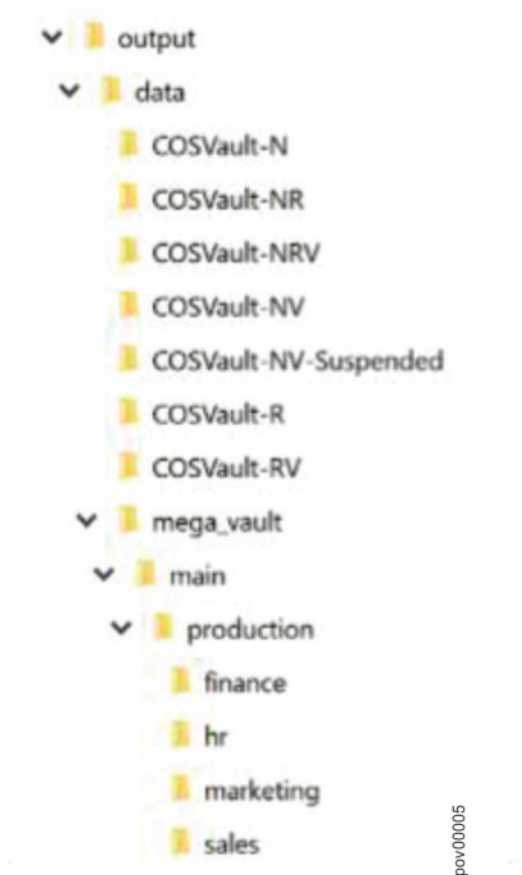
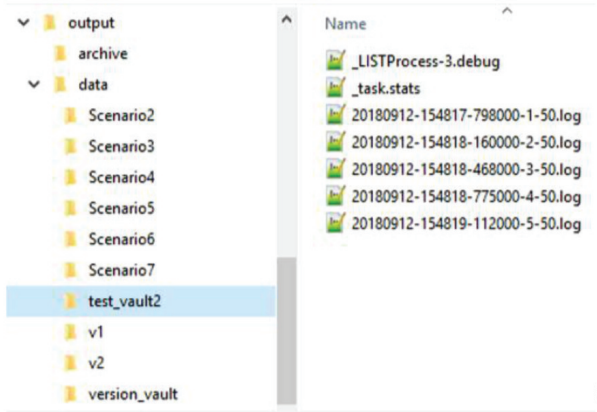


Figure 37. Directory structure from the configuration file

The status and progress of each scan must be maintained so a separate directory structure is created for each scan. Table 21 on page 89 shows the leaf directories that contain the file names and description.

Table 21. Leaf directory file names	
File name	Description
_LISTProcessN.debug	<p>The N in the file name is different for each process (0 - 9 if there are 10 processes).</p> <p>Contains detailed debug information and details of any errors that are encountered when you scan the vault. Figure 38 on page 89 shows an example of running in debug mode.</p> <pre> 12-Sep-2018 15:50:13 LISTProcess-2 test_vault2 +++ Adding batch 16 12-Sep-2018 15:50:13 LISTProcess-2 test_vault2 16 Working batches: 1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16 12-Sep-2018 15:50:13 LISTProcess-2 test_vault2 16 Working batches: 1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16 12-Sep-2018 15:50:13 LISTProcess-2 test_vault2 --- Removing batch 1 12-Sep-2018 15:50:13 LISTProcess-2 test_vault2 Buffer is full. Sleeping for 1 second... 12-Sep-2018 15:50:14 LISTProcess-2 test_vault2 +++ Adding batch 17 12-Sep-2018 15:50:14 LISTProcess-2 test_vault2 16 Working batches: 2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17 12-Sep-2018 15:50:15 LISTProcess-2 test_vault2 16 Working batches: 2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17 12-Sep-2018 15:50:15 LISTProcess-2 test_vault2 --- Removing batch 2 12-Sep-2018 15:50:15 LISTProcess-2 test_vault2 +++ Adding batch 18 12-Sep-2018 15:50:15 LISTProcess-2 test_vault2 16 Working batches: 3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18 12-Sep-2018 15:50:15 LISTProcess-2 test_vault2 16 Working batches: 3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18 12-Sep-2018 15:50:15 LISTProcess-2 test_vault2 --- Removing batch 3 12-Sep-2018 15:50:15 LISTProcess-2 test_vault2 +++ Adding batch 19 12-Sep-2018 15:50:15 LISTProcess-2 test_vault2 16 Working batches: 4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19 12-Sep-2018 15:50:15 LISTProcess-2 test_vault2 16 Working batches: 4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19 </pre> <p>Figure 38. Example of running in debug mode</p>

Table 21. Leaf directory file names (continued)	
File name	Description
task.stats	<p>Scanner starts in JSON format for a single vault. Updated following successful processing of each batch of objects.</p> <pre> "estimated_object_count": 1718, "list_objects_size": 50, "next_key": "", "next_version": "", "prefix": "", "scan_type": "Object Scan", "status": "Complete", "total_bytes_output": 45257, "total_bytes_scanned": 45257, "total_objects_output": 1717, "total_objects_scanned": 1717, "vault_name": "test_vault2", "vault_uuid": "06c1641d-082f-7ba2-011b-c7550651a780" </pre>
*.log	<p>The Scanner creates multiple .log files for each vault. Each .log file contains up to 1000 Kafka messages, ready to be submitted to the Kafka cluster by the Notifier.</p> <p>The naming convention for the log files is</p> <pre> <date>-<time>-<milliseconds>-<batch number>- <number of messages in file>.log </pre> 

Appendix

The appendix shows an example of a log file and examples of Scanner debug data.

Log file

Figure 39 on page 91 shows an example with extra line breaks.

```
{
  "system_name": "Test",
  "object_version": "38d811e6-dba1-4830-859d-6275f2016bc3",
  "object_etag": "\73b8c5dcca9cc6aab928396a2d98a340\\",
  "request_time": "2018-08-24T18:39:16Z",
  "format": 1,
  "bucket_uuid": "cbc03649-d218-727d-10ec-df8c22873280",
  "system_uuid": "f7d033c2-9066-499a-a883-829860d4d865",
  "meta_headers": [],
  "object_length": 12,
  "object_name": "my-object-black",
  "bucket_name": "version_vault",
  "content_type": "text/plain",
  "request_id": "7ed39768-1184-4d5f-8c0c-7912515bf8fe",
  "operation": "s3:PutObject"
}

{
  "system_name": "Test",
  "object_version": "8a68208b-9b5f-4107-9eae-fa08200c7913",
  "object_etag": "\73b8c5dcca9cc6aab928396a2d98a340\\",
  "request_time": "2018-08-24T18:39:15Z",
  "format": 1,
  "bucket_uuid": "cbc03649-d218-727d-10ec-df8c22873280",
  "system_uuid": "f7d033c2-9066-499a-a883-829860d4d865",
  "meta_headers": [],
  "object_length": 12,
  "object_name": "my-object-black",
  "bucket_name": "version_vault",
  "content_type": "text/plain",
  "request_id": "6ed2a774-f384-4cba-96fd-81dfbb681482",
  "operation": "s3:PutObject"
}

{
  "system_name": "Test",
  "object_version": "322d9ed2-ca86-4efd-b95a-a2467ded9202",
  "object_etag": "\73b8c5dcca9cc6aab928396a2d98a340\\",
  "request_time": "2018-08-24T18:39:15Z",
  "format": 1,
  "bucket_uuid": "cbc03649-d218-727d-10ec-df8c22873280",
  "system_uuid": "f7d033c2-9066-499a-a883-829860d4d865",
  "meta_headers": [],
  "object_length": 12,
  "object_name": "my-object-black",
  "bucket_name": "version_vault",
  "content_type": "text/plain",
  "request_id": "ced1de9e-e62f-4966-b803-7aff3ddab245",
  "operation": "s3:PutObject"
}

{
  "system_name": "Test",
  "object_version": "0e2b0369-a506-4ce3-8336-ea42eae16489",
  "object_etag": "\73b8c5dcca9cc6aab928396a2d98a340\\",
  "request_time": "2018-08-24T18:39:15Z",
  "format": 1,
  "bucket_uuid": "cbc03649-d218-727d-10ec-df8c22873280",
  "system_uuid": "f7d033c2-9066-499a-a883-829860d4d865",
  "meta_headers": [],
  "object_length": 12,
  "object_name": "my-object-black",
  "bucket_name": "version_vault",
  "content_type": "text/plain",
  "request_id": "ee5b8ded-86af-4e9e-9054-8902f7a7a5c0",
  "operation": "s3:PutObject"
}

{
  "system_name": "Test",
  "object_version": "41518cb8-632f-45c4-989b-44b4d2b19b2e",
  "object_etag": "\73b8c5dcca9cc6aab928396a2d98a340\\",
  "request_time": "2018-08-24T18:39:15Z",
  "format": 1,
  "bucket_uuid": "cbc03649-d218-727d-10ec-df8c22873280",
  "system_uuid": "f7d033c2-9066-499a-a883-829860d4d865",
  "meta_headers": [],
  "object_length": 12,
  "object_name": "my-object-black",
  "bucket_name": "version_vault",
  "content_type": "text/plain",
  "request_id": "5adaa4d9-7aa6-4c46-bcd1-1260c074a972",
  "operation": "s3:PutObject"
}

{
  "system_name": "Test",
  "object_version": "c75c6faf-e7a9-41ce-a576-6bfc9d374dfc",
  "object_etag": "\73b8c5dcca9cc6aab928396a2d98a340\\",
  "request_time": "2018-08-24T18:39:14Z",
  "format": 1,
  "bucket_uuid": "cbc03649-d218-727d-10ec-df8c22873280",
  "system_uuid": "f7d033c2-9066-499a-a883-829860d4d865",
  "meta_headers": [],
  "object_length": 12,
  "object_name": "my-object-black",
  "bucket_name": "version_vault",
  "content_type": "text/plain",
  "request_id": "35aa2dee-8700-4cfb-a62e-dc7f7ec7b8e1",
  "operation": "s3:PutObject"
}
```

pov00009

Figure 39. Example of a log file

Scanner debug data

Figure 40 on page 92, Figure 41 on page 93, Figure 42 on page 94, and Figure 43 on page 95 show that throttling settings are logged, and multiple HEAD processes are started for each LIST process.


```

12-Sep-2018 15:57:25 |
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | Python package setup
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | """Setup"""
12-Sep-2018 15:57:25 | from setuptools import setup, find_packages
12-Sep-2018 15:57:25 | setup(
12-Sep-2018 15:57:25 |     name='ibm_cos_scanner',
12-Sep-2018 15:57:25 |     version='2.0.0',
12-Sep-2018 15:57:25 |     packages=find_packages(),
12-Sep-2018 15:57:25 |     include_package_data=True,
12-Sep-2018 15:57:25 |     zip_safe=True,
12-Sep-2018 15:57:25 |     url='www.ibm.com',
12-Sep-2018 15:57:25 |     license='See LICENSE folder',
12-Sep-2018 15:57:25 |     author='IBM',
12-Sep-2018 15:57:25 |     description='IBM COS Scanner / Spectrum Discover Notifier'
12-Sep-2018 15:57:25 | )
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | Initialising IBM COS Scanner. Reading config
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | Retrieving System Advanced Configuration
12-Sep-2018 15:57:25 | Calling https://172.19.17.38/manager/api/json/1.0/viewSystemConfiguration.adm
12-Sep-2018 15:57:26 | | - OK
12-Sep-2018 15:57:26 | | - dsNet Name: est
12-Sep-2018 15:57:26 | | - dsNet UUID: f7d033c2-9066-499a-a883-829860d4d865
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Retrieving user's access keys
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Calling https://172.19.17.38/manager/api/json/1.0/listMyAccessKeys.adm
12-Sep-2018 15:57:26 | | - OK
12-Sep-2018 15:57:26 | | - Accesser credentials successfully retrieved from Manager API
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Retrieving vault information from dsNet
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Calling https://172.19.17.38/manager/api/json/1.0/viewSystemConfiguration.adm
12-Sep-2018 15:57:27 | | - OK
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Retrieving vault size information from dsNet
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Calling https://172.19.17.38/manager/api/json/1.0/listVaults.adm
12-Sep-2018 15:57:27 | | - OK
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Retrieving dsNet device information
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Calling https://172.19.17.38/manager/api/json/1.0/listDevices.adm
12-Sep-2018 15:57:28 | | - OK
12-Sep-2018 15:57:28 | 3 devices found
12-Sep-2018 15:57:28 | | - device_id:1    manager    172.19.17.38    0e547a7f-849a-79e8-102c-2026af755443vm-eqx204201-
mgr-03
12-Sep-2018 15:57:28 | | - device_id:2    accesser    172.19.17.39    39a14ec8-090c-79d8-10f9-255249197f43vm-eqx204201-

```

Figure 40. Scanner debug

```

12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | Python package setup
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | """Setup"""
12-Sep-2018 15:57:25 | from setuptools import setup, find_packages
12-Sep-2018 15:57:25 | setup(
12-Sep-2018 15:57:25 |     name='ibm_cos_scanner',
12-Sep-2018 15:57:25 |     version='2.0.0',
12-Sep-2018 15:57:25 |     packages=find_packages(),
12-Sep-2018 15:57:25 |     include_package_data=True,
12-Sep-2018 15:57:25 |     zip_safe=True,
12-Sep-2018 15:57:25 |     url='www.ibm.com',
12-Sep-2018 15:57:25 |     license='See LICENSE folder',
12-Sep-2018 15:57:25 |     author='IBM',
12-Sep-2018 15:57:25 |     description='IBM COS Scanner / Spectrum Discover Notifier'
12-Sep-2018 15:57:25 | )
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | Initialising IBM COS Scanner. Reading config
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | Retrieving System Advanced Configuration
12-Sep-2018 15:57:25 | Calling https://172.19.17.38/manager/api/json/1.0/viewSystemConfiguration.adm
12-Sep-2018 15:57:26 | | - OK
12-Sep-2018 15:57:26 | | - dsNet Name: est
12-Sep-2018 15:57:26 | | - dsNet UUID: f7d033c2-9066-499a-a883-829860d4d865
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Retrieving user's access keys
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Calling https://172.19.17.38/manager/api/json/1.0/listMyAccessKeys.adm
12-Sep-2018 15:57:26 | | - OK
12-Sep-2018 15:57:26 | | - Accesser credentials successfully retrieved from Manager API
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Retrieving vault information from dsNet
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Calling https://172.19.17.38/manager/api/json/1.0/viewSystemConfiguration.adm
12-Sep-2018 15:57:27 | | - OK
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Retrieving vault size information from dsNet
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Calling https://172.19.17.38/manager/api/json/1.0/listVaults.adm
12-Sep-2018 15:57:27 | | - OK
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Retrieving dsNet device information
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Calling https://172.19.17.38/manager/api/json/1.0/listDevices.adm
12-Sep-2018 15:57:28 | | - OK
12-Sep-2018 15:57:28 | 3 devices found
12-Sep-2018 15:57:28 | | - device_id:1    manager    172.19.17.38    0e547a7f-849a-79e8-102c-2026af755443vm-eqx204201-
mgr-03
12-Sep-2018 15:57:28 | | - device_id:2    accesser    172.19.17.39    39a14ec8-090c-79d8-10f9-255249197f43vm-eqx204201-

```

Figure 41. Scanner debug (continued)

```

12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | Python package setup
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | """Setup"""
12-Sep-2018 15:57:25 | from setuptools import setup, find_packages
12-Sep-2018 15:57:25 | setup(
12-Sep-2018 15:57:25 |     name='ibm_cos_scanner',
12-Sep-2018 15:57:25 |     version='2.0.0',
12-Sep-2018 15:57:25 |     packages=find_packages(),
12-Sep-2018 15:57:25 |     include_package_data=True,
12-Sep-2018 15:57:25 |     zip_safe=True,
12-Sep-2018 15:57:25 |     url='www.ibm.com',
12-Sep-2018 15:57:25 |     license='See LICENSE folder',
12-Sep-2018 15:57:25 |     author='IBM',
12-Sep-2018 15:57:25 |     description='IBM COS Scanner / Spectrum Discover Notifier'
12-Sep-2018 15:57:25 | )
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | Initialising IBM COS Scanner. Reading config
12-Sep-2018 15:57:25 | -----
12-Sep-2018 15:57:25 | Retrieving System Advanced Configuration
12-Sep-2018 15:57:25 | Calling https://172.19.17.38/manager/api/json/1.0/viewSystemConfiguration.adm
12-Sep-2018 15:57:26 | | OK
12-Sep-2018 15:57:26 | |- dsNet Name: est
12-Sep-2018 15:57:26 | |- dsNet UUID: f7d033c2-9066-499a-a883-829860d4d865
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Retrieving user's access keys
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Calling https://172.19.17.38/manager/api/json/1.0/listMyAccessKeys.adm
12-Sep-2018 15:57:26 | | OK
12-Sep-2018 15:57:26 | |- Accesser credentials successfully retrieved from Manager API
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Retrieving vault information from dsNet
12-Sep-2018 15:57:26 | -----
12-Sep-2018 15:57:26 | Calling https://172.19.17.38/manager/api/json/1.0/viewSystemConfiguration.adm
12-Sep-2018 15:57:27 | | OK
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Retrieving vault size information from dsNet
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Calling https://172.19.17.38/manager/api/json/1.0/listVaults.adm
12-Sep-2018 15:57:27 | | OK
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Retrieving dsNet device information
12-Sep-2018 15:57:27 | -----
12-Sep-2018 15:57:27 | Calling https://172.19.17.38/manager/api/json/1.0/listDevices.adm
12-Sep-2018 15:57:28 | | OK
12-Sep-2018 15:57:28 | 3 devices found
12-Sep-2018 15:57:28 | |- device_id:1    manager    172.19.17.38    0e547a7f-849a-79e8-102c-2026af755443vm-eqx204201-
mgr-03
12-Sep-2018 15:57:28 | |- device_id:2    accesser    172.19.17.39    39a14ec8-090c-79d8-10f9-255249197f43vm-eqx204201-

```

pov00012

Figure 42. Scanner debug (continued)


```

12-Sep-2018 15:57:29 | 10 tasks
12-Sep-2018 15:57:29 | -----
12-Sep-2018 15:57:29 | |- Object Scan of v1
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0500 seconds, Head: 0.0050 seconds
12-Sep-2018 15:57:29 | |- Object Scan of Scenario3
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0500 seconds, Head: 0.0050 seconds
12-Sep-2018 15:57:29 | |- Object Scan of Scenario2
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0500 seconds, Head: 0.0050 seconds
12-Sep-2018 15:57:29 | |- Object Scan of v2
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0500 seconds, Head: 0.0050 seconds
12-Sep-2018 15:57:29 | |- Version Scan of version_vault
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0010 seconds, Head: n/a
12-Sep-2018 15:57:29 | |- Object Scan of Scenario6
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0500 seconds, Head: 0.0050 seconds
12-Sep-2018 15:57:29 | |- Version Scan of Scenarios
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0010 seconds, Head: n/a
12-Sep-2018 15:57:29 | |- Object Scan of Scenario4
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0500 seconds, Head: 0.0050 seconds
12-Sep-2018 15:57:29 | |- Object Scan of Scenario7
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0500 seconds, Head: 0.0050 seconds
12-Sep-2018 15:57:29 | |- Object Scan of test_vault2
12-Sep-2018 15:57:29 | |   |- Throttling list: 0.0500 seconds, Head: 0.0050 seconds
12-Sep-2018 15:57:29 | -----
12-Sep-2018 15:57:29 | Ignoring vaults
12-Sep-2018 15:57:29 | -----
12-Sep-2018 15:57:29 | | Scenario1
12-Sep-2018 15:57:29 | | Scenario0
12-Sep-2018 15:57:29 | -----
12-Sep-2018 15:57:29 | Queuing scanner tasks
12-Sep-2018 15:57:29 | -----
12-Sep-2018 15:57:29 | | Queuing task 'Object Scan of v1'
12-Sep-2018 15:57:29 | | Queuing task 'Object Scan of Scenario3'
12-Sep-2018 15:57:29 | | Queuing task 'Object Scan of Scenario2'
12-Sep-2018 15:57:29 | | Queuing task 'Object Scan of v2'
12-Sep-2018 15:57:29 | | Queuing task 'Version Scan of version_vault'
12-Sep-2018 15:57:29 | | Queuing task 'Object Scan of Scenario6'
12-Sep-2018 15:57:29 | | Queuing task 'Version Scan of Scenarios'
12-Sep-2018 15:57:29 | | Queuing task 'Object Scan of Scenario4'
12-Sep-2018 15:57:29 | | Queuing task 'Object Scan of Scenario7'
12-Sep-2018 15:57:29 | | Queuing task 'Object Scan of test_vault2'
12-Sep-2018 15:57:29 | -----
12-Sep-2018 15:57:29 | Creating 10 list processes, each with 5 head processes
12-Sep-2018 15:57:29 | -----
12-Sep-2018 15:57:31 | | Started LISTProcess-0
12-Sep-2018 15:57:32 | | | Started HEADProcess-0-0
12-Sep-2018 15:57:33 | | | Started HEADProcess-0-1
12-Sep-2018 15:57:34 | | | Started HEADProcess-0-2
12-Sep-2018 15:57:35 | | | Started HEADProcess-0-3
12-Sep-2018 15:57:36 | | | Started HEADProcess-0-4
12-Sep-2018 15:57:38 | | Started LISTProcess-1
12-Sep-2018 15:57:39 | | | Started HEADProcess-1-0
12-Sep-2018 15:57:40 | | | Started HEADProcess-1-1
12-Sep-2018 15:57:41 | | | Started HEADProcess-1-2
12-Sep-2018 15:57:42 | | | Started HEADProcess-1-3
12-Sep-2018 15:57:43 | | | Started HEADProcess-1-4
12-Sep-2018 15:57:44 | | Started LISTProcess-2
12-Sep-2018 15:57:46 | | | Started HEADProcess-2-0
12-Sep-2018 15:57:47 | | | Started HEADProcess-2-1
12-Sep-2018 15:57:48 | | | Started HEADProcess-2-2
12-Sep-2018 15:57:50 | | | Started HEADProcess-2-3
12-Sep-2018 15:57:52 | | | Started HEADProcess-2-4
12-Sep-2018 15:57:54 | | Started LISTProcess-3

```

Figure 43. Scanner debug (continued)

Configure IBM Cloud Object Storage notifications for IBM Spectrum Discover

Ingesting IBM Cloud Object Storage event records into IBM Spectrum Discover requires the user to enable the Notification service on the IBM Cloud Object Storage system. Thereafter, the user must connect the IBM Cloud Object Storage system to the IBM Cloud Object Storage connector Kafka topic on the IBM Spectrum Discover cluster. The name of this connector topic is `cos-1e-connector-topic`.

A combination of SASL and TLS is used to authenticate and encrypt the connection between the IBM Cloud Object Storage source system and the Kafka brokers which reside in the IBM Spectrum Discover cluster. The certificate and credentials required to establish this connection might be obtained directly from the IBM Spectrum Discover cluster by the IBM Spectrum Discover storage administrator.

For information on how to enable and configure the IBM Cloud Object Storage Notification service with the IBM Spectrum Discover provided credentials, see [IBM Cloud Object Storage Administration Documentation](#).

The following information is required to establish a secure connection between IBM Cloud Object Storage and IBM Spectrum Discover:

Hosts

One or more of the IBM Spectrum Discover Kafka brokers is in the format: *host1:port,host2:port*. The Kafka producers on the IBM Cloud Object Storage system will retrieve the full list of IBM Spectrum Discover Kafka brokers from the first host that is alive and responding. The broker's host and port (the list configured might contain more than one broker) for SASL SSL can be obtained by the IBM Spectrum Discover storage administrator from the following location on the IBM Spectrum Discover master node: `/etc/kafka/server.properties`.

Authentication credentials

The username is `cos` and the password can be obtained by the IBM Spectrum Discover storage administrator from the following location on the IBM Spectrum Discover master node: `/etc/kafka/sasl_password`.

Certificate PEM for TLS encryption

This is the CA certificate that is used to sign the Kafka server and client certificates for the IBM Spectrum Discover cluster. It might be obtained by the IBM Spectrum Discover storage administrator from the following location on the IBM Spectrum Discover master node: `/etc/kafka/ca.crt`.

This file is in the PEM format and the entirety of its contents must be pasted into the **Certificate PEM** field of the **COS Notifications** configuration panel.

Enabling IBM Cloud Object Storage notification services

The IBM Cloud Object Storage notification service can be enabled with the information that follows:

Procedure

1. Log in to the IBM Cloud Object Storage Manager Admin console https://manager_host/manager/login.adm with a username of **admin** and a password of **password**.

If you defined your own password, use your pre-defined password. If you do not have a pre-defined password use the default password.

2. Select the **Administration** tab.
3. Scroll to the end of the page and select **Configure the Notification Service**.



Figure 44. Configurations


```
NAME: <NAME>
Topic: cos-le-connector-topic
Hosts: <SD hostname> :9093
Type: IBM Spectrum Discover
```

Enabling authentication

1. Check **Enable authentication**.

```
Username: cos
Password: <PASSWORD>
```

Enabling encryption

1. Check **Enable TLS for Apache Kafka network connections**.
2. Add the certificate PEM file from the IBM Spectrum Discover platform. See [Figure 45 on page 98](#).

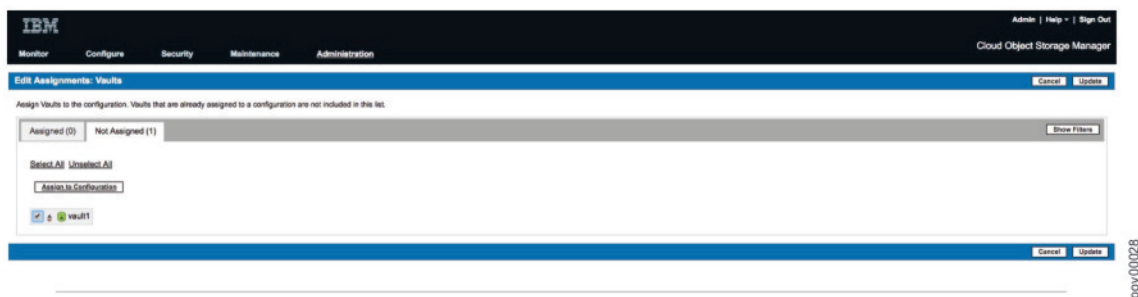


Figure 45. Add a storage vault to the configuration

Testing the IBM Cloud Object Storage notification service

To test the IBM Cloud Object Storage notification service, the tester can populate the IBM Cloud Object Storage vault with test data.

About this task

You can use a number of methods to write files to an IBM Cloud Object Storage vault, but you can use cURL directly on IBM Spectrum Discover platform. cURL is a computer software project that provides a library and command-line tool for transferring data that uses various protocols.

Procedure

1. Create a test file, for example, object_1.txt.
The test file can be any file that contains data.
2. Write a file to the IBM Cloud Object Storage vault by using cURL.

Requirements

IBM COS Vault Name (vault1) [anonymous access enabled]
IBM COS Accesser IP address

Example

```
[moadmin@spectrum_discover~]$ curl -X PUT -i -T object_1.txt http://9.11.200.208/vault1/object_1.txt
HTTP/1.1 100 ContinueHTTP/1.1 200 OK
Date: Fri, 04 Jan 2019 13:21:14 Greenwich mean time
X-Clv-Request-Id: a9ad657a-a919-4b13-9b72-961ae8c57e3c
```

```
Server: 3.14.0.23
X-Clv-S3-Version: 2.5
x-amz-request-id: a9ad657a-a919-4b13-9b72-961ae8c57e3c
ETag: "7c517c7108f7180377e7b37db2e39261"
Content-Length: 0
```

Monitoring the IBM Cloud Object Storage accesser logs

To determine whether a file is successfully written to the IBM Cloud Object Storage vault and a notification is successfully sent to the IBM Spectrum Discover server, the accesser logs can be monitored on the IBM Cloud Object Storage Accesser server.

In the following example, an object that is written to vault1 results in the sending of one notification to the IBM Spectrum Discover server. The user must have access privileges to log on to the IBM Cloud Object Storage Accesser host to check the log files.

Confirm that an object is stored in the IBM Cloud Object Storage vault.

```
root@ibm_accesser:/var/log/dsnet-core# tail -f http.log
9.11.201.78 - "" - [04/Jan/2019:13:21:14 +0000] "PUT /vault1/object_1.txt HTTP/1.1" 200 0 "-"
"curl/7.29.0" 22
```

Confirm that a notification is sent to the IBM Spectrum Discover server.

```
root@ibm_accesser:/var/log/dsnet-core# tail -f notification.log
{"time":"2019-01-04T13:21:14.668Z","request_id":"a9ad657a-a919-4b13-9b72-961ae8c57e3c","retried":true,"success":true,
"request_time":"2019-01-04T13:21:14.567Z","kafka_config_uuid":"d842c7a0-9c36-412e-8908-8ad5120a261e",
"topic":"cos-le-connector-topic"}
```

Monitoring the IBM Spectrum Discover producer IBM Cloud Object Storage logs

When the IBM Spectrum Discover server receives a notification from the IBM Cloud Object Storage platform, the IBM Spectrum Discover producer IBM Cloud Object Storage records a transaction.

A successful notification is recorded as an offset value of one, when a notification is received from IBM Cloud Object Storage platform.

```
[moadmin@spectrum_discover]$ kubectl logs -f -n producercos kindled-alligator-producer-cos-
producer-9f6966b4-8jsg7
break time. waiting for work...
2019-01-04 13:21:19.187 > offset_commit_cb: success, offsets:[{part: 0, offset: 1, err: none}]
```

Monitoring the IBM Spectrum Discover dashboard for IBM Cloud Object Storage ingestion

You can monitor the IBM Spectrum Discover dashboard for IBM Cloud Object Storage ingestion.

After IBM Cloud Object Storage notifications are ingested from the IBM Cloud Object Storage platform, the IBM Spectrum Discover dashboard displays the total number of indexed records.

Note: The IBM Spectrum Discover dashboard can take approximately 30 minutes to display the total number of indexed records.

See [Figure 46 on page 100](#) .

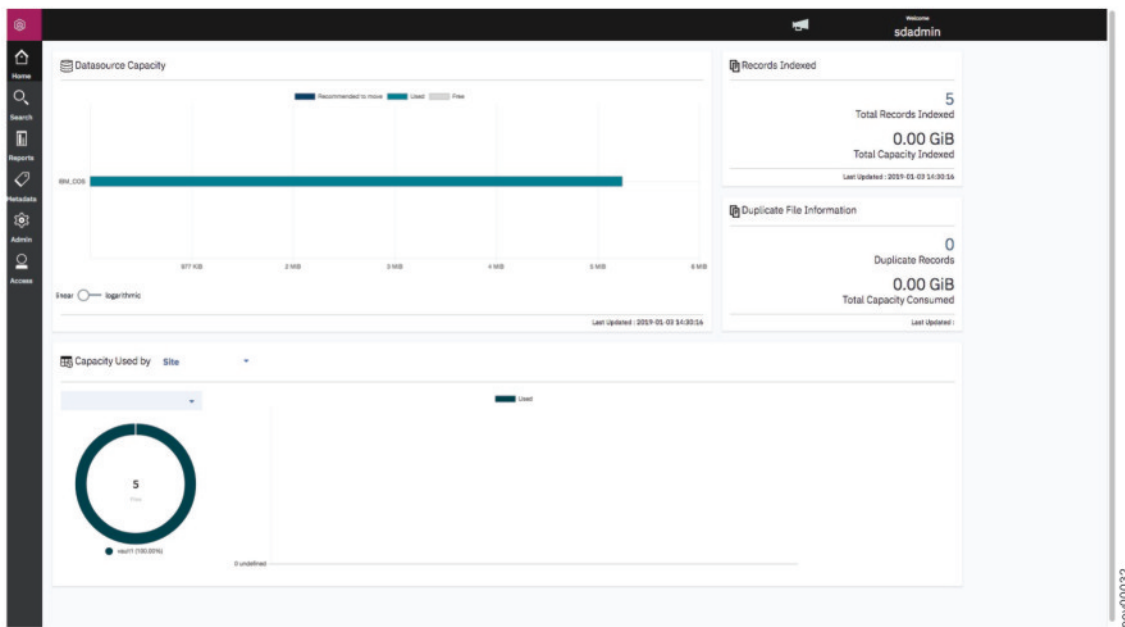


Figure 46. Total number of indexed records

IBM Spectrum Discover and S3 object storage data source connections

Use this information to understand how IBM Spectrum Discover works with S3-compliant object store.

Creating an S3 object storage data connection

Use this information to create a connection to an S3-compliant object store.

About this task

To create the S3-compliant connection:

Procedure

1. Log in to the IBM Spectrum Discover graphical user interface (GUI) with a user ID that is associated with data administration role.

The data administration access role is required for creating connections. For more information about role-based access control, go to [“Role-based access control”](#) on page 3.

2. Select **Admin** from the left navigation menu.

Clicking **Admin** displays the different types of data source connection names, platforms, clusters, data source, size, and **Add Connection**.

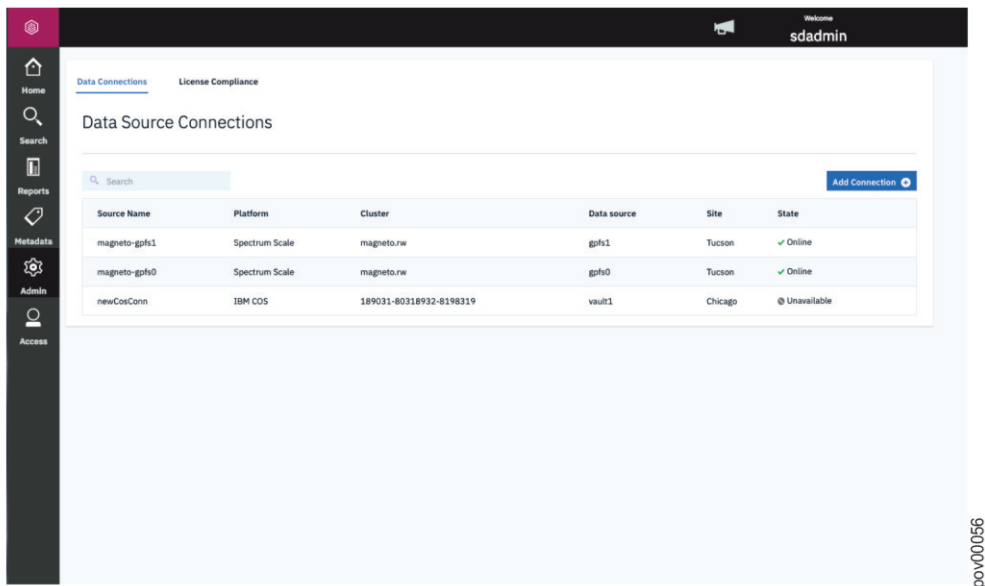


Figure 47. Displaying the source names for **Data Source Connections**

- Click **Add Connection** to display a new window that shows **Add Data Source Connection**.

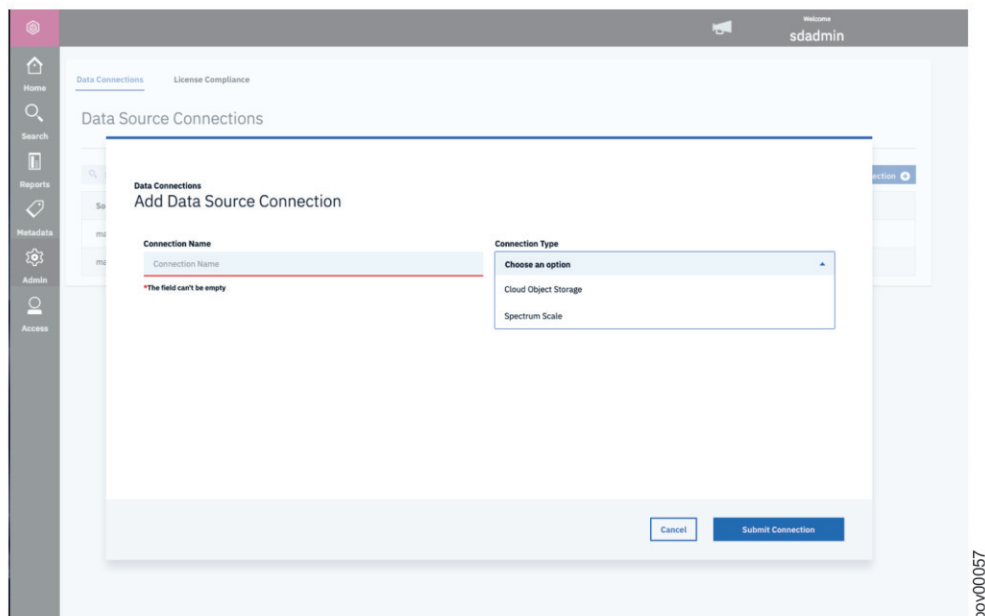


Figure 48. Example of the **Add Data Source Connection** GUI window

- Complete the following steps:
 - In the field for **Connection Name**, define a **Connection Name**.
 - Click the down arrow for **Connection Type** to display a drop-down menu for **Choose an option**.
- Select the connection type **S3 Object Storage**.

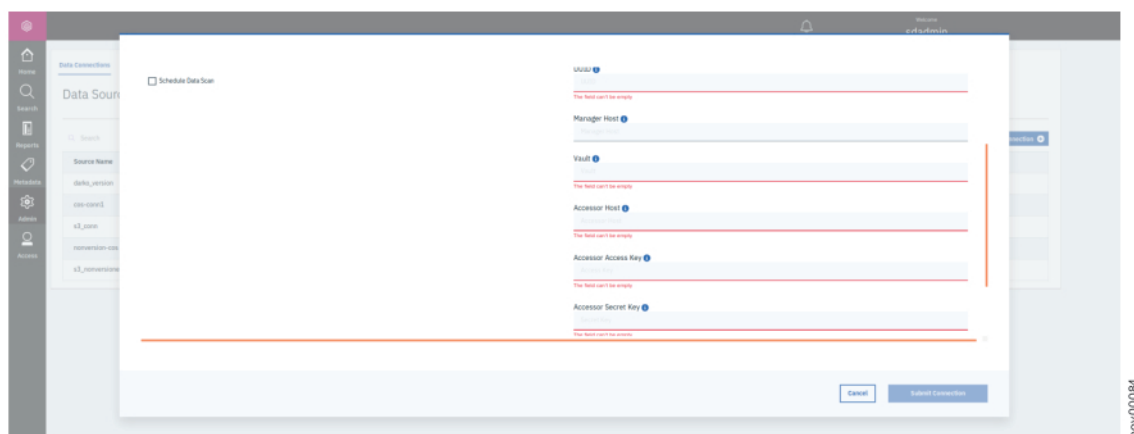


Figure 49. Selecting **S3 Object Storage**

6. In the screen for S3, complete the fields and then click **Submit Connection** for the S3 connections manager.

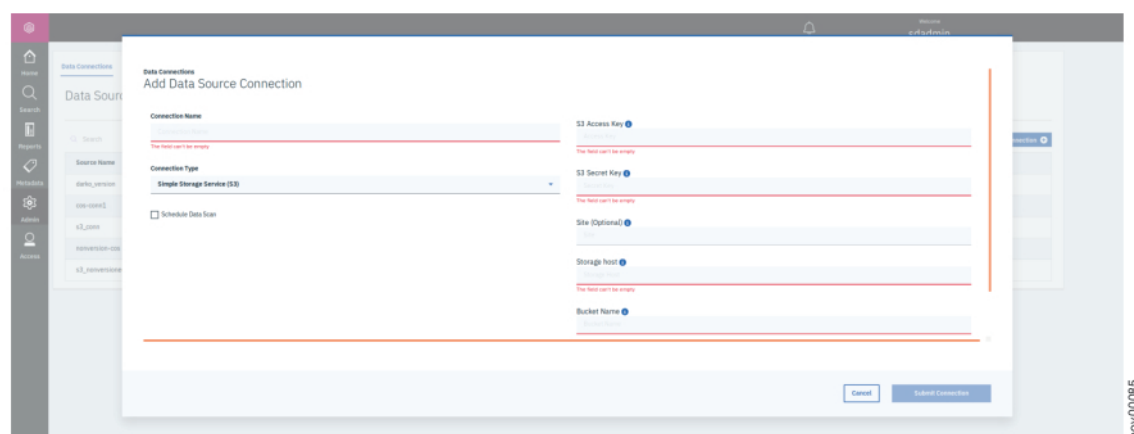


Figure 50. Completing the S3 information fields

S3 access key

Indicates the access key ID for the S3 object store.

S3 secret key

Indicates the secret key ID for the S3 object store.

Site

Indicates the physical location of the records.

Storage host

Indicates the IP address or host of the storage system.

Bucket name

Indicates the name of the bucket that you are going to scan.

Scanning an S3 object storage data connection

Use this information to scan a connection for an S3-compliant object store.

About this task

When you initiate a scan from the IBM Spectrum Discover graphical user interface (GUI), the metadata is transferred asynchronously back to the IBM Spectrum Discover instance.

Automated scanning and data ingestion relies on an established and active network connection between the IBM Spectrum Discover instance and the S3 storage source. If the connection cannot be established,

the state of the data source connection shows as unavailable, and the option for automated scanning does not appear in the IBM Spectrum Discover GUI for that connection.

Procedure

1. Go to the IBM Spectrum Discover graphical user interface (GUI).
2. Under **Admin** select **Data Source connections**.

The following example shows the **Admin** data connections menu page:

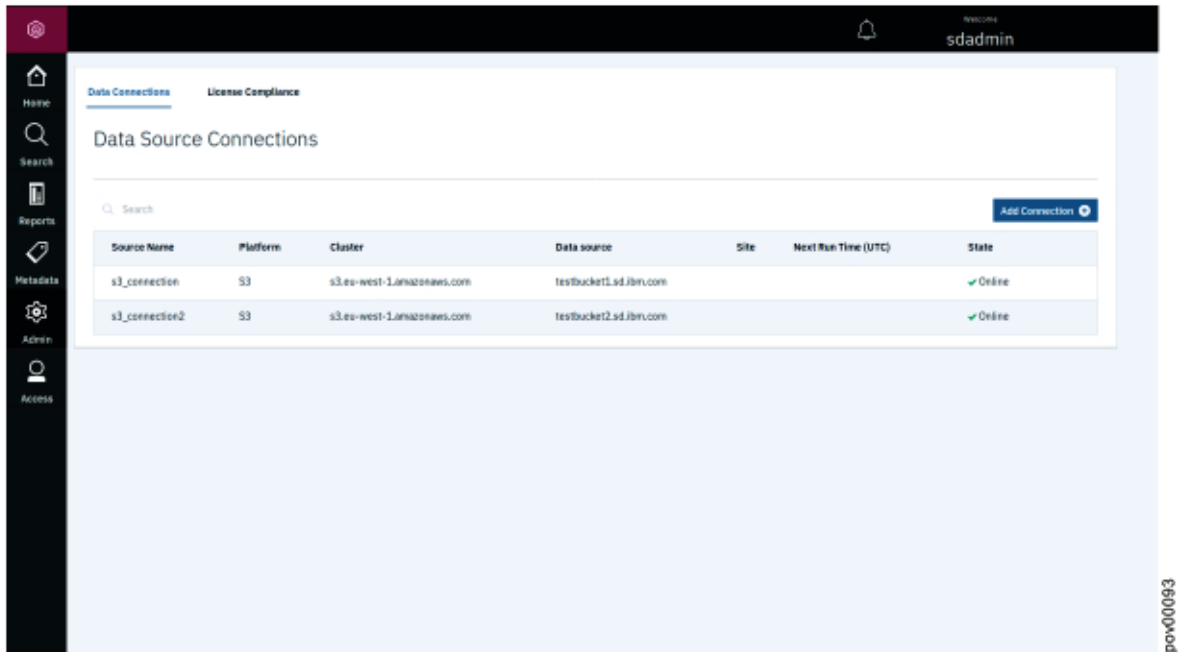


Figure 51. Data source connections

3. Select the data source connection that you want to scan. Make sure that the State is listed as Online to make your system scan ready.

The following example shows how to select a data source connection to scan.

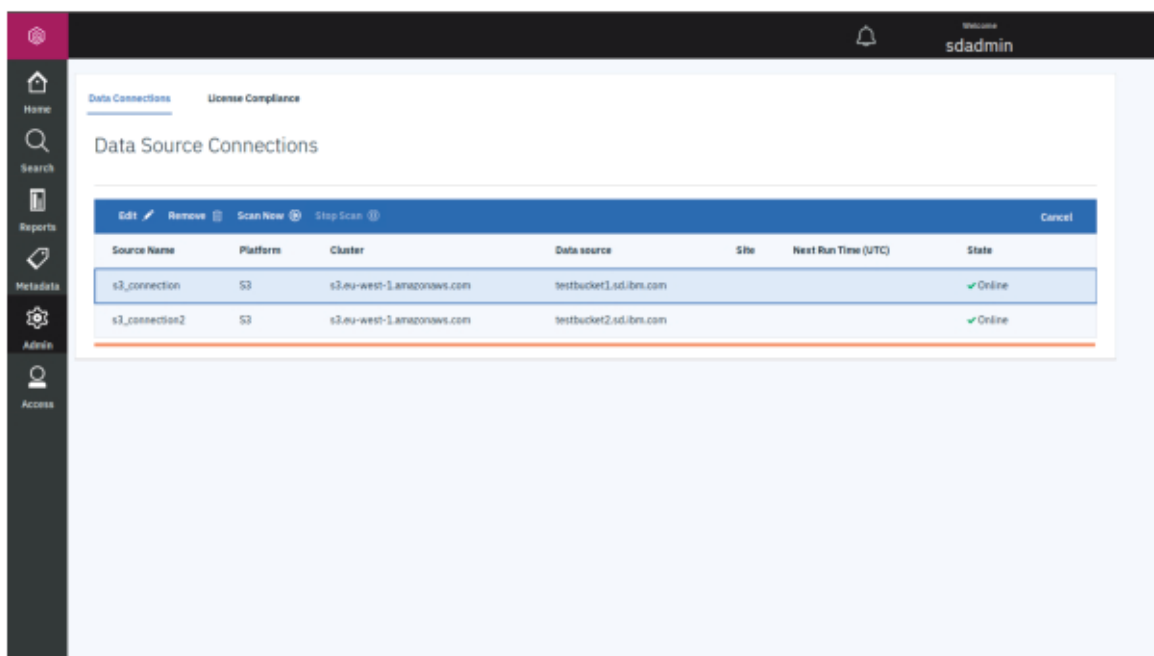


Figure 52. Selecting a data source connection to scan

4. Select **Scan Now** to change the status to **Scanning**.

The following example shows an active scan.

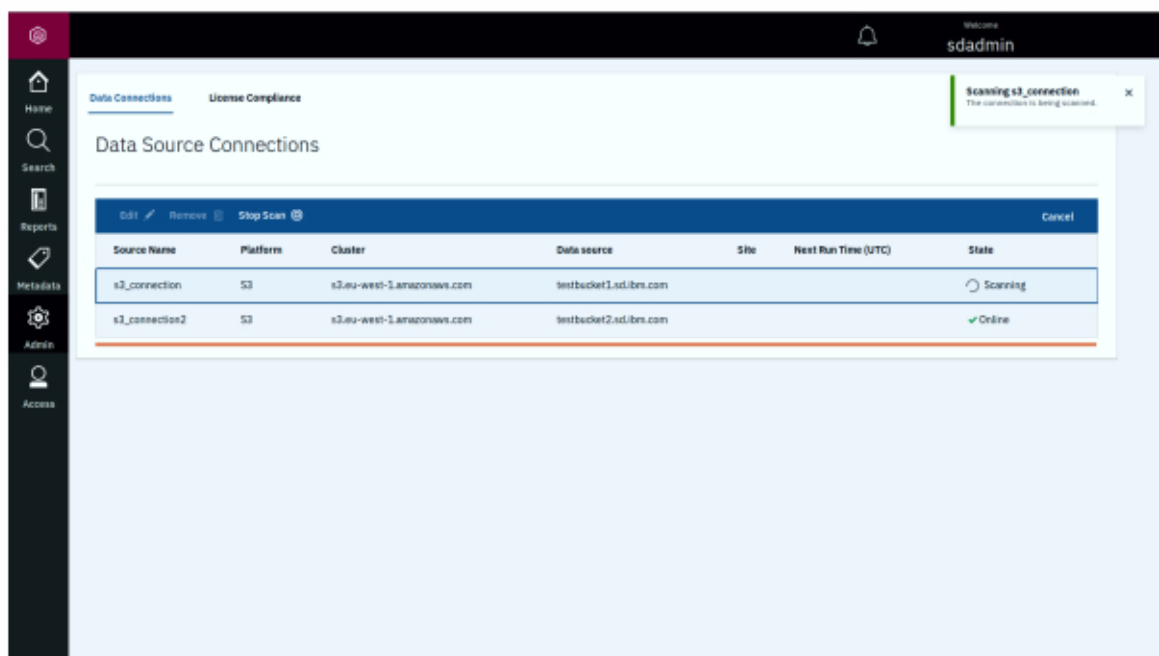


Figure 53. Active scans

5. When the scan finishes, the state field returns to a status of **Online**.

Best practices for scanning an S3 object storage system

Use best practices for scanning S3-compliant object storage systems.

It is recommended to check the log files in the following directories after each scan:

/opt/ibm/metaocean/data/connections/s3/<connection_name>/debug/
<scan_timestamp>/scanner.debug indicates whether the scan was successful or not.

/opt/ibm/metaocean/data/connections/s3/<connection_name>/error/<scan_timestamp>/scanner.error contains a list of all the messages that are not delivered to IBM Spectrum Discover.

/opt/ibm/metaocean/data/connections/s3/<connection_name>/data/<scan_timestamp>/ contains a subfolder with the scanned data source name. There is a stats folder inside this folder that contains information about the number of objects in the data source or the number of objects or scanned files.

You can also compare the total size of the bucket that is reported in IBM Spectrum Discover with the total size of the S3 object store at source (if it is available).

Scanning an Elastic Storage Server data source connection

Use the IBM Spectrum Discover GUI to scan an Elastic Storage Server (ESS) data connection.

Procedure

1. Go to the IBM Spectrum Discover GUI and log in with **Data Admin** privileges.
2. Open an internet browser to the IP address or host name of your IBM Spectrum Discover cluster.
The default credentials to put into the **IBM Spectrum Discover** dialog box are:

```
Default user: sdadmin  
Default Password: Passw0rd
```

3. Go to **Admin > Add Datasource Connection**.
4. Enter the EMS node as the host system name:

Connection Name

Indicates the connection name, such as ESS_test.

Connection Type

Indicates the connection type, which is Spectrum Scale.

User

Indicates the user that initiates the scan, such as root.

Password

Indicates the password for the user that initiates the scan.

Working Directory

Indicates the working directory for the scan, such as /gpfd/icp4D_data_fs_master1/sd_scan.

Scan Directory

Indicates the scan directory, such as /gpfd/icp4D_data_fs_master1.

Site (optional)

Indicates the scan site. This field is optional.

Cluster

Indicates the clustered system name that is being scanned, such as heliumess.tuc.stglabs.ibm.com.

Host

Indicates the host system name that is being scanned, such as 9.11.103.45.

Filesystem

Indicates the file system name that is being scanned.

Node List

Indicates the node list of nodes that are being scanned. Enter gss_ppc64 as the node class to perform the scan.

5. After you complete all fields, select **Submit Connection**.
6. Select your data source from the table (click the row) and then select **Scan Now**.

7. You can monitor the status of the scans from **Data Connections Overview**.
8. You can also search for documents from that data source by using the search navigation icon in the left side of the GUI.
In the search icon you must:
 - a. Enter a query.
 - b. Click **Search**.
 - c. Indicate the files that are found in the search results.

Creating a Network File System data source connection

You can use the IBM Spectrum Discover graphical user interface to create a Network File System (NFS) data connection.

About this task

Note: For NFS scanning that uses Data ONTAP version 8.1.2 or later, export the file system with the following configuration:

Protocol

NFSv3

Security type

UNIX

Client permissions

A minimum of read-only access is necessary.

Anonymous users

Root access must be granted.

Setuid and Setgid executable routines are not necessary.

Creating data source connections in IBM Spectrum Discover identifies source storage systems that are to be indexed by IBM Spectrum Discover.

For some data source types, a network connection can be created to allow for automated scanning and indexing of the source system metadata. IBM Spectrum Discover will not index data from unknown sources, so creating a data source connection is the first step towards cataloging any source storage system.

Procedure

1. Log in to the IBM Spectrum Discover web interface with a user ID that has the **Data Admin** role that is associated with it.
Note: The **Data Admin** access role is necessary for creating connections.
2. Select **Admin** from the left navigation menu to display the source name, platform, cluster, data source, site, next run time, and state of existing data source connections.
3. Click **Add connection** to display the **Add Data Source Connection** window.
4. Enter a name in the **Connection name** field.
5. Click **Connection type** to expand the **Choose an option** menu.
6. Select the connection type **Network File System**.
7. Complete the fields for the NFS parameters, and click **Submit Connection**.

Parameters for NFS connections

Connection name

The name of the connection, an identifier for the user, for example `filesystem1`.

Note: It must be a unique name within IBM Spectrum Discover.

Connection type

The type of source storage system this connection represents.

Data source

The full name of the data source.

Export path

The data path from which data is to be exported.

Host

The host name of the node from which a scan can be initiated.

Site

The location in which the data source facility is located.

Creating an SMB data source connection

Creating an SMB data connection by using the IBM Spectrum Discover graphical user interface.

About this task

Create data source connections in IBM Spectrum Discover to identify and index source storage systems.

For some data source types, a network connection is (optionally) created to allow for automated scanning and indexing of the source system metadata. IBM Spectrum Discover does not index data from unknown sources, so creating a data source connection is the first step towards cataloging any source storage system.

Procedure

1. Log in to the IBM Spectrum Discover web interface with a user ID that is associated with the **Data Admin** role.
Note: The **Data Admin** access role is required for creating connections.
2. Select **Admin** from the left navigation menu to display the source name, platform, cluster, data source, site, next run time, and state of existing data source connections.
3. Click **Add Connection** to display the **Add Data Source Connection** window.
4. Enter a name in the **Connection Name** field.
5. Click **Connection Type** to expand the **Choose an option** menu.
6. Select the connection type **Server Message Block (SMB) / CIFS**.
7. Complete the fields for the SMB parameters, and click **Submit Connection**.

Parameters for SMB connections

Connection Name

Indicates the name of the connection, an identifier for the user, for example Share1.

Note: It must be a unique name within IBM Spectrum Discover.

Share

Indicates the path name of the SMB/CIFS file share, for example //server/share.

User

Indicates the user ID that has access permissions to the SMB share, for example DOMAIN1\user1

Password

Indicates the authentication password for the user ID provided.

Site (Optional)

Indicates the location in which the data source facility is located.

Creating an IBM Spectrum Protect data source connection

Creating a IBM Spectrum Protect data connection by using the IBM Spectrum Discover graphical user interface.

About this task

Creating data source connections in IBM Spectrum Discover identifies source storage systems that are to be indexed by IBM Spectrum Discover.

For some data source types, a network connection is (optionally) created to allow for automated scanning and indexing of the source system metadata. IBM Spectrum Discover will not index data from unknown sources, so creating a data source connection is the first step towards cataloging any source storage system.

Procedure

1. Log in to the IBM Spectrum Discover web interface with a user ID that has the **Data Admin** role that is associated with it.
Note: The **Data Admin** access role is required for creating connections.
2. Select **Admin** from the left navigation menu to display the source name, platform, cluster, data source, site, next run time, and state of existing data source connections.
3. Click **Add Connection** to display the **Add Data Source Connection** window.
4. Enter a name in the **Connection Name** field.
5. Click **Connection Type** to expand the **Choose an option** menu.
6. Select the connection type **IBM Spectrum Protect**.
7. Complete the fields for the IBM Spectrum Protect connection type, and click **Submit Connection**.

Parameters for IBM Spectrum Protect connections

Connection Name

The name of the connection, an identifier for the user, for example `filesystem1`.

Note: It must be a unique name within IBM Spectrum Discover.

Spectrum Protect Server IP

The IP address or host name of the IBM Spectrum Protect server.

ODBC Port

The ODBC (open database connector) port for the Protect server (default is 51500).

Instance User

The name of the IBM Spectrum Protect database instance user. The default is `tsminst1`.

Instance user password

The password for the database instance user.

Site (Optional)

The location in which the data source facility is located.

Editing and using the TimeSinceAccess and Size Range buckets

Users can group or aggregate data into three user-defined bucket ranges. The three user-defined bucket ranges are TimeSinceAccess, Size Range and FileGroup.

The TimeSinceAccess bucket groups files and objects based on the time they were last accessed. The SizeRange bucket groups files and objects based on their size. The File Group bucket groups files based on their file type or extension. All three buckets can be customized to better align with the user's requirements. Both the TimeSinceAccess and the SizeRange buckets have up to five custom ranges with user-defined labels.

Note: The **FileGroup** bucket cannot be edited through the user interface and must be modified using REST APIs. For more information, see the topic `/db2whrest/v1/buckets/<bucket>: PUT` in the *IBM Spectrum Discover: REST API Guide*

[To access the SizeRange bucket groups, select **Metadata > Tags > edit icon for the SizeRange tag**.] For example, SizeRange can be broken up into 'T-shirt size' ranges where the ranges and labels are:

Table 22. Examples of size ranges and sizes of buckets with user-defined labels	
Size range	Size
0 - 4 K	XS
4 K - 1 M	S
1 M - 1 G	M
1 G - 1T B	L
1 TB+	XL

[After you change or update a bucket definition, IBM Spectrum Discover summarizes the current set of files and objects into their respective bucket ranges. The changes are updated periodically every half an hour; thus, it may take a half an hour or more before the changes are reflected in the Spectrum Discover GUI.]

Note: Ensure that the maximum value for each bucket is greater than the value assigned to the previous bucket.

See the [Figure 54 on page 109](#)

Figure 54. Example of how to define the settings for a SizeRange bucket

[To open the menu for the TimeSinceAccess buckets select **Metadata > Tags > edit icon for the TimeSinceAccess tag**. See the [Figure 54 on page 109](#) for an example.

[Figure 55 on page 110](#) shows an example of how to modify and define the settings of a bucket that is older than one year.

Figure 55. Example of how to modify and define the settings of a bucket that is older than one year old

Using custom TLS certificate

You can change the TLS certificate that is used by IBM Spectrum Discover for serving web pages and the REST API endpoints.

About this task

Follow the procedure to use a custom TLS certificate:

Procedure

1. Create a secret for your TLS certificate within the same namespace as the one used for deploying IBM Spectrum Discover on OpenShift® that is *"spectrum-discover"*.

Note: You can use any name for the secret. The following example uses "my-tls-secret" as the secret name.

```
kubectl create secret tls my-tls-secret --key ${KEY_FILE} --cert ${CERT_FILE} -n spectrum-discover
```

2. Modify the IBM Spectrum Discover custom resource and specify the following ingress settings:

```
kubectl edit SpectrumDiscover spectrumdiscover -n spectrum-discover
```

3. Update the *"host"* and *"tls_secret_name"* in the relevant ingress section.

```
ingress:
  host: spectrum-discover.ibm.com
  tls_secret_name: my-tls-secret
```

Note: The "ingress.host" setting must match the fully qualified domain name as specified in the TLS certificate. This domain name is the hostname that the ingress binds to.

4. Save the custom resource.

Note: The operator takes a while to go through all components and update them with the new settings. Issue the following command to check the operator log for monitoring its progress.


```
kubect1 logs $(kubect1 get po -l name=spectrum-discover-operator -n spectrum-discover -o name) -n spectrum-discover -c operator --follow
```

- The log displays "PLAY RECAP" on completing the update.
- Enter `ctrl+c` to stop following the log.

Validating code integrity

You can validate code integrity to improve your system security.

You can verify code integrity and verify the secure origin of the IBM Spectrum Discover OVA (open virtualization appliance) and IBM Spectrum Discover upgrader. Make sure that both OVA and upgrader are signed with an official IBM code-signing certificate.

For OVA deployments, you can confirm that you have a valid certificate when you use the **Deploy OVF Template wizard** in the VMware vSphere Client. When you review the details after you specify your input, you can confirm that the publisher displays something similar to the following output:

```
"DigiCert SHA2 Assured ID Code Signing CA (Trusted certificate)"
```

For more information, see the VMware vSphere Client documentation.

For upgrade deployments, download the upgrader binary along with the `.pem` and `.sig` files. To verify the upgrade integrity, run the following command:

```
openssl dgst -sha256 -verify <download pem file>.pem -signature <downloaded signature file>.sig <upgrader file>
```

After several minutes, the system returns with the following status:

```
"Verified OK"
```

Applying the license file

As a prerequisite, you must have the `ibm-spectrum-discover-unrestricted.lic` license file - or an alternative file provided by IBM.

Procedure

1. Use Secure Shell (SSH) to log in to the machine.
2. Copy the license file (`ibm-spectrum-discover-unrestricted.lic`) to the IBM Spectrum Discover machine.
3. Get an authentication token for the API.

```
TOKEN=$(curl -ks -u sdadmin:<password> https://localhost/auth/v1/token -I | awk '/X-Auth-Token/ {print $2}')
```

4. Run the following command to load the license from the license file:

```
LICENSE=$(cat ibm-spectrum-discover-unrestricted.lic)
```

5. Run the following command to push the license to the server:

```
curl -k -H "Authorization: Bearer ${TOKEN}" -H "Content-Type: application/json" -X PUT --data "{\"license\": \"${LICENSE}\"}" https://localhost/api/license/
```

6. To verify the license, go to the administration section of the Web UI, or check the license with the API. You can also upload or install the unrestricted license by using the **License Compliance** page of the IBM Spectrum Discover graphical user interface (GUI).

Chapter 4. Upgrading

Preparing to run the upgrade tool for IBM Spectrum Discover

You can upgrade IBM Spectrum Discover from the master node.

Before you begin

Remember: IBM Spectrum Discover does not support file or file path names that use characters that are not part of the UTF-8 character set.

1. The IBM Spectrum Discover upgrade tool can be downloaded from IBM Fix Central. <https://www-945.ibm.com/support/fixcentral/>
2. To validate your code integrity, see “Validating code integrity” on page 111.
3. Run the upgrade tool from the master node in the IBM Spectrum Discover cluster.

Note: After the upgrade is completed, any previous edits to deployment configuration that are still required, must be reapplied. For example, The API-KEY for the Watson Knowledge Catalog application must be reapplied. For more information, see the topic *Exporting metadata* in the *IBM Spectrum Discover: Administration Guide*.

[If you are upgrading from a 2.0.2.x or 2.0.3.x release, you must add 2 extra disks for 2.0.4 for the following purposes:

- A new disk for backups
- A new disk for database storage

During the upgrade, all data from the previous IBM Spectrum Scale disk will be migrated to the new database storage disk. On completion of the migration process, IBM Spectrum Scale is removed. After the upgrade is complete, the previous IBM Spectrum Scale disk can also be removed from the system.]

About this task

Follow the given procedure for preparing to run the upgrade tool.

Procedure

1. Stop data ingest before you run the upgrade tool.
2. Make sure that you have the most current and up-to-date authentication certificates within IBM Spectrum Discover:
 - a) Log in to the IBM Spectrum Discover console with a user ID of `<moadmin>` and a password of `<Passw0rd>`.
 - b) Run the following command:

```
sudo /etc/cron.hourly/icp_login.sh
```

Note:

When you log in to IBM Spectrum Discover or restart it and try to run the `icp_login.sh` script immediately after that, it might report a failure. Wait for some time before you rerun the script. Before you continue with the upgrade, you must ensure that the script does not report any error.

Running the upgrade tool for IBM Spectrum Discover

You can run the upgrade tool by extracting a tarball on a directory in the master node.

Before you begin

Before you run the upgrade tool, you must ensure that it is prepared and ready to run. For more information, see [“Preparing to run the upgrade tool for IBM Spectrum Discover ” on page 113.](#)

About this task

Follow the procedure to run the upgrade tool.

Procedure

1. The upgrade tool is a self-extracting archive.

To run the upgrade tool, use this command:

```
sudo bash < upgrader >
```

2. If you need to run the upgrade again, you can run it without re-extracting the media.

To run the upgrade again, use this command:

```
cd /opt/ibm/metaocean
sudo ./upgrade
```

Note: Before you run the upgrade, you must disable any automated backup that is scheduled or is running. Follow the procedure shown to manually disable any cron jobs that are created for running automated backups:

- a. Edit the cron job:

```
$ sudo crontab -e
```

- b. Delete the following backup entry.

```
@daily python3 /opt/ibm/metaocean/backup-restore/backup.py
```

- c. Save the file.

When the upgrade is completed, edit the crontab to reenter the line that you deleted in **Step b.**

The upgrade process can restart several times based on the version of IBM Spectrum Discover that is being upgraded.

A restart disconnects the Secure Shell (SSH) session. However, the upgrade process is displayed immediately upon reconnection after the restart. The message-of-the-day (MOTD) displays a message that an upgrade is in progress:

```
# BEGIN ANSIBLE MANAGED BLOCK
*****
* UPGRADE IN PROGRESS *
*****
# END ANSIBLE MANAGED BLOCK
```

This statement precedes the running upgrader status, whose output comprises the end (or tail) of the upgrader log file data. Use **CTRL + C** to escape without affecting the running upgrader process.

When the upgrade is complete, the MOTD is replaced with a message that displays the new or upgraded IBM Spectrum Discover version. You do not have to immediately reconnect after a restart because the upgrader runs automatically as soon as the system returns online. An upgrade can be initiated, and when a new SSH session displays the MOTD without upgrader progress, the upgrade is completed.

Upgrading IBM Spectrum Discover to versions 2.0.4 and higher

You can upgrade to IBM Spectrum Discover version 2.0.4. and higher by adding extra disks for purpose of backups and storage.

About this task

Follow the procedure to upgrade to IBM Spectrum Discover version 2.0.4 and higher by including extra disks.

Procedure

1. Issue the following command to enter the maintenance mode.

```
sudo /opt/ibm/metaocean/helpers/maintenance.sh on
```

2. Issue the following command to shut down the Db2® Warehouse.

```
docker exec Db2wh /mnt/blumeta0/home/db2inst1/sqllib/bin/db2fmcu -d
docker exec Db2wh stop
docker stop Db2wh
```

3. Shut down the VM.
4. Add the extra disks for backup and data storage with adequate storage size. For more information, see [“Storage requirements for single node trial and single node production IBM Spectrum Discover deployments” on page 13.](#)
5. Start the VM.
6. Ensure that all pods in the kube-system namespace are running indicates that the VM is fully running before you proceed. The sample log displays the status as shown.

```
while [[ $(kubectl get po -n kube-system --no-headers | egrep -v 'Running|Completed' | wc -l) != 0 ]]; do echo 'Waiting on pods'; sleep 20; done
```

7. Issue the following command to log in to the ICp.

```
sudo /etc/cron.hourly/icp_login.sh
```

8. Run the upgrade.

Accessibility features for IBM Spectrum Discover

Accessibility features help users who have a disability, such as restricted mobility or limited vision, to use information technology products successfully.

Accessibility features

The following list includes the major accessibility features in IBM Spectrum Discover:

- Keyboard-only operation
- Interfaces that are commonly used by screen readers
- Keys that are discernible by touch but do not activate just by touching them
- Industry-standard devices for ports and connectors
- The attachment of alternative input and output devices

IBM Knowledge Center, and its related publications, are accessibility-enabled. The accessibility features are described in [IBM Knowledge Center \(www.ibm.com/support/knowledgecenter\)](http://www.ibm.com/support/knowledgecenter).

Keyboard navigation

This product uses standard Microsoft Windows navigation keys.

IBM and accessibility

See the [IBM Human Ability and Accessibility Center \(www.ibm.com/able\)](http://www.ibm.com/able) for more information about the commitment that IBM has to accessibility.

Notices

This information was developed for products and services offered in the US. This material might be available from IBM in other languages. However, you may be required to own a copy of the product or product version in that language in order to access it.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing IBM Corporation North Castle Drive, MD-NC119 Armonk, NY 10504-1785 US

For license inquiries regarding double-byte character set (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

Intellectual Property Licensing Legal and Intellectual Property Law IBM Japan Ltd. 19-21, Nihonbashi-Hakozakicho, Chuo-ku Tokyo 103-8510, Japan

INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some jurisdictions do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM websites are provided for convenience only and do not in any manner serve as an endorsement of those websites. The materials at those websites are not part of the materials for this IBM product and use of those websites is at your own risk.

IBM may use or distribute any of the information you provide in any way it believes appropriate without incurring any obligation to you.

Licensees of this program who wish to have information about it for the purpose of enabling: (i) the exchange of information between independently created programs and other programs (including this one) and (ii) the mutual use of the information which has been exchanged, should contact:

IBM Director of Licensing IBM Corporation North Castle Drive, MD-NC119 Armonk, NY 10504-1785 US

Such information may be available, subject to appropriate terms and conditions, including in some cases, payment of a fee.

The licensed program described in this document and all licensed material available for it are provided by IBM under terms of the IBM Customer Agreement, IBM International Program License Agreement or any equivalent agreement between us.

The performance data discussed herein is presented as derived under specific operating conditions. Actual results may vary.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM

products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

Statements regarding IBM's future direction or intent are subject to change or withdrawal without notice, and represent goals and objectives only.

All IBM prices shown are IBM's suggested retail prices, are current and are subject to change without notice. Dealer prices may vary.

This information is for planning purposes only. The information herein is subject to change before the products described become available.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to actual people or business enterprises is entirely coincidental.

COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs. The sample programs are provided "AS IS", without warranty of any kind. IBM shall not be liable for any damages arising out of your use of the sample programs.

Each copy or any portion of these sample programs or any derivative work must include a copyright notice as follows:

© (your company name) (year).

Portions of this code are derived from IBM Corp.

Sample Programs. © Copyright IBM Corp. _enter the year or years_.

If you are viewing this information softcopy, the photographs and color illustrations may not appear.

Trademarks

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at [Copyright and trademark information](http://www.ibm.com/legal/copytrade.shtml) at www.ibm.com/legal/copytrade.shtml.

Intel is a trademark of Intel Corporation or its subsidiaries in the United States and other countries.

Java™ and all Java-based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates.

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both.

Microsoft and Windows are trademarks of Microsoft Corporation in the United States, other countries, or both.

UNIX is a registered trademark of the Open Group in the United States and other countries.

Terms and conditions for product documentation

Permissions for the use of these publications are granted subject to the following terms and conditions.

Applicability

These terms and conditions are in addition to any terms of use for the IBM website.

Personal use

You may reproduce these publications for your personal, noncommercial use provided that all proprietary notices are preserved. You may not distribute, display or make derivative work of these publications, or any portion thereof, without the express consent of IBM.

Commercial use

You may reproduce, distribute and display these publications solely within your enterprise provided that all proprietary notices are preserved. You may not make derivative works of these publications, or reproduce, distribute or display these publications or any portion thereof outside your enterprise, without the express consent of IBM.

Rights

Except as expressly granted in this permission, no other permissions, licenses or rights are granted, either express or implied, to the publications or any information, data, software or other intellectual property contained therein.

IBM reserves the right to withdraw the permissions granted herein whenever, in its discretion, the use of the publications is detrimental to its interest or, as determined by IBM, the above instructions are not being properly followed.

You may not download, export or re-export this information except in full compliance with all applicable laws and regulations, including all United States export laws and regulations.

IBM MAKES NO GUARANTEE ABOUT THE CONTENT OF THESE PUBLICATIONS. THE PUBLICATIONS ARE PROVIDED "AS-IS" AND WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESSED OR IMPLIED, INCLUDING BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, NON-INFRINGEMENT, AND FITNESS FOR A PARTICULAR PURPOSE.

IBM Online Privacy Statement

IBM Software products, including software as a service solutions, ("Software Offerings") may use cookies or other technologies to collect product usage information, to help improve the end user experience, to tailor interactions with the end user or for other purposes. In many cases no personally identifiable information is collected by the Software Offerings. Some of our Software Offerings can help enable you to collect personally identifiable information. If this Software Offering uses cookies to collect personally identifiable information, specific information about this offering's use of cookies is set forth below.

This Software Offering does not use cookies or other technologies to collect personally identifiable information.

If the configurations deployed for this Software Offering provide you as customer the ability to collect personally identifiable information from end users via cookies and other technologies, you should seek your own legal advice about any laws applicable to such data collection, including any requirements for notice and consent.

For more information about the use of various technologies, including cookies, for these purposes, See IBM's Privacy Policy at <http://www.ibm.com/privacy> and IBM's Online Privacy Statement at <http://www.ibm.com/privacy/details> the section entitled "Cookies, Web Beacons and Other Technologies" and the "IBM Software Products and Software-as-a-Service Privacy Statement" at <http://www.ibm.com/software/info/product-privacy>.

Index

A

- accesser [58](#)
- accessibility features for IBM Spectrum Discover [117](#)
- adding
 - certificate PEM file [97](#)
 - notification service [96](#)
- adding data source connections
 - from graphical user interface [49](#), [59](#), [106–108](#)
- Additional disks
 - Upgrade IBM Spectrum Discover [115](#)
- API commands
 - DELETE [4](#)
 - descriptions [4](#)
 - GET [4](#)
 - POST [4](#)
 - PUT [4](#)
- application
 - action ID [6](#)
 - delete [6](#)
 - parameters [6](#)
 - view [6](#)
- Applications
 - DEEPINSPECT [5](#)
- Applying the license file
 - Deploying and configuring [111](#)
- architecture
 - IBM Cloud Object Storage [66](#)
- Architecture
 - Application SDK [2](#)
 - example diagram [2](#)
 - IBM Spectrum Discover [2](#)
- authentication
 - example [97](#)
- Automated scanning of an IBM Spectrum Scale data source
 - IBM Spectrum Scale data source connection [51](#)
- Automated scanning of an IBM Spectrum Scale file set
 - IBM Spectrum Scale data source connection [52](#)
- AUTOTAG
 - example [6](#)

B

- backup
 - configuring VMDK [34](#)
- backup requirements
 - IBM Spectrum Discover [16](#)
- buckets
 - using SizeRange [108](#)
 - using TimeSinceAccess [108](#)
- Business-oriented data mapping
 - features [2](#)

C

- Cataloging metadata [4](#)
- certificate authority

- certificate authority (*continued*)
 - copying [96](#)
 - example [97](#)
 - PEM file [97](#)
 - PEM format [96](#)
- components
 - IBM Cloud Object Storage
 - Notifier [66](#)
 - Replay [66](#)
 - Scanner [66](#)
- configuration file
 - explanation of settings [67](#)
 - Notifier [67](#)
 - Replay [67](#)
 - Scanner [67](#)
 - what file includes [67](#)
- configuration setup
 - notification service [97](#)
- connect to IBM Spectrum Scale
 - copying file system scan output [56](#)
 - file system scan [54](#)
 - start consumer [57](#)
 - start producer [57](#)
- cURL
 - description [98](#)
 - writing files [98](#)
- custom TLS certificate
 - using [110](#)

D

- Dashboard
 - example [7](#)
 - purpose [7](#)
 - what is viewable [7](#)
- Data activation
 - features [2](#)
- data source connections
 - creating [44](#)
 - displaying the source names [49](#), [59](#), [106–108](#)
 - editing [49](#), [59](#), [106–108](#)
 - removing [49](#), [59](#), [106–108](#)
 - scanning [62](#)
 - scanning now [49](#), [59](#), [106–108](#)
- Data source connections
 - cluster [4](#)
 - data source name [4](#)
 - platform [4](#)
 - site [4](#)
- Data visualization
 - features [2](#)
- database
 - configuring VMDK [31](#)
- debug mode
 - creating log files [82](#)
 - Replay [82](#)
 - running to troubleshoot problems [82](#)

- debug mode (*continued*)
 - starting [82](#)
- DEEPINSPECT
 - example [6](#)
- deleting data source connections
 - from graphical user interface [49](#), [59](#), [106–108](#)
- Deploying the IBM Spectrum Discover open virtualization
 - appliance on the Kernel-based Virtual Machine
 - virtualization module
 - Deploying and configuring [43](#)
- deployment models
 - IBM Spectrum Discover [11](#)
- determining
 - file written to IBM Cloud Object Storage
 - successfully [99](#)

E

- editing data source connections
 - from graphical user interface [49](#), [59](#), [106–108](#)
- enabling
 - notification service [97](#)
 - storage vault [97](#)
 - TLS [97](#)
- Enabling bucket notifications for Ceph Object Storage
 - IBM Cloud Object Storage data source connection [64](#)
- encryption [97](#)
- Enriching
 - metadata [5](#)
- error conditions
 - Replay
 - scenarios [81](#)
- Exabyte-scale data ingest
 - features [2](#)
- example
 - authentication [97](#)
 - system advanced configuration [58](#)
- Extensible foundation for data insight
 - features [2](#)

F

- features
 - IBM Spectrum Discover
 - Business-oriented data mapping [2](#)
 - Data activation [2](#)
 - Data visualization [2](#)
 - Exabyte-scale data ingest [2](#)
 - Extensible foundation for data insight [2](#)
- file system scan [54](#)

G

- Get Bucket Extension [58](#)
- graphical user interface
 - adding data source connections [49](#), [59](#), [106–108](#)
 - description [7](#)
 - scanning data source connections [62](#)

I

- IBM Cloud Object Scanner
 - configuration file [88](#)

- IBM Cloud Object Scanner (*continued*)
 - directory structure from the configuration file [88](#)
 - output data [88](#)
- IBM Cloud Object Storage
 - architecture [66](#)
 - components
 - Notifier [66](#)
 - Replay [66](#)
 - Scanner [66](#)
 - introduction [66](#)
 - notification service [96](#)
 - notifications ingested [99](#)
 - overview [66](#)
 - platform [96](#)
 - prerequisites [58](#)
 - rules settings
 - Scanner [67](#)
 - Scanner
 - rules settings [67](#)
 - settings example [67](#)
 - writing files [98](#)
- IBM COS
 - best practices scanning [64](#)
- IBM Fix Central
 - upgrading
 - IBM Spectrum Discover [113](#), [114](#)
- IBM Spectrum Archive data source connections
 - Configure data source connections [57](#)
- IBM Spectrum Discover
 - appliance
 - definition [9](#)
 - resources [9](#)
 - architecture [2](#)
 - architecture diagram [2](#)
 - backup deployment [34](#)
 - benefits [1](#)
 - connect to COS [57](#)
 - connect to IBM Spectrum Scale [53](#), [54](#), [56](#), [57](#)
 - connecting to data sources [44](#)
 - data sheet [1](#)
 - data source connections [44](#)
 - deployment [21](#), [27](#), [28](#), [31](#), [36](#), [40](#), [44](#), [53](#), [54](#), [56](#), [57](#)
 - deployment issues [43](#)
 - deployment models [11](#)
 - displaying number of indexed records [99](#)
 - Extra disks
 - Upgrade procedure [115](#)
 - graphical user interface [7](#)
 - introduction [1](#)
 - known issues [43](#)
 - OVA file deployment [21](#)
 - overview [7](#)
 - planning [12](#)
 - producer logs
 - monitoring [99](#)
 - reports
 - DELETE [9](#)
 - examples [9](#)
 - GET [9](#)
 - POST [9](#)
 - PUT [9](#)
 - single management and data path [12](#)
 - sizing requirements for backup [16](#)
 - software requirements [11](#)

IBM Spectrum Discover *(continued)*

- upgrade [113](#), [114](#)
- upgrading [113](#), [114](#)
- virtual appliance deployment [21](#), [27](#), [28](#), [31](#), [34](#), [36](#), [40](#)
- virtual node [21](#), [27](#), [28](#), [31](#), [34](#), [36](#), [40](#)

IBM Spectrum Discover information units [xi](#)

IBM Spectrum Scale data source connections

- Configure data source connections [45](#)

L

limitations

- Notifier [83](#)

loading

- upgrade tool [113](#), [114](#)

log file

- example of [90](#)
- how a log file is created [90](#)

Logging

- Notifier [87](#)
- Replay [87](#)

M

messages

- error_code [81](#)
- error_description [81](#)
- examples [81](#)

metadata

- cataloging [4](#)
- description [4](#)
- enriching [5](#)

monitoring

- IBM Cloud Object Storage
 - accessor logs [99](#)
- IBM Spectrum Discover
 - producer logs [99](#)

multi-node deployments

- IBM Spectrum Discover [11](#)

N

networking requirements

- IBM Spectrum Discover [12](#)

notification service

- adding [96](#)
- configuration setup [97](#)
- enabled for storage vault [97](#)
- enabling [97](#)

Notifier

- acknowledgment [82](#)
- Archie folder [82](#)
- before restarting [83](#)
- before you stop [83](#)
- description [82](#)
- generating log files [83](#)
- how Notifier operates [83](#)
- how Notifier works [82](#)
- IBM Cloud Object Storage [66](#)
- limitations [83](#)
- monitoring the progress [83](#)
- restarting [84](#)
- shutdown

Notifier *(continued)*

shutdown *(continued)*

- kill.notifier file [84](#)
- stopping [83](#)
- using a Kafka configuration [83](#)

O

output

- messages [81](#)

overview

- IBM Cloud Object Storage [66](#)

P

persistent message queue

- configuring VMDK [28](#)

Policy

- action [6](#)
- AUTOTAG
 - example [6](#)
- DEEPIINSPECT
 - example [6](#)
- description [6](#)
- filter [6](#)
- policy id [6](#)
- purpose [5](#)

preparing

- upgrade tool [113](#), [114](#)

prerequisites

- IBM Cloud Object Storage [58](#)

process count

- Replay
 - default values for settings [75](#)

Progress report

- creating [84](#)
- description [84](#)
- example [84](#)
- format [84](#)

R

Replay

- behaviors [76](#)
- debug mode [82](#)
- default values for settings [75](#)
- directories generated by Notifier [87](#)
- directories generated by Replay [87](#)
- directories generated by Scanner [87](#)
- error conditions [81](#)
- example of how to start [82](#)
- guidelines [82](#)
- IBM Cloud Object Storage [66](#)
- maximum performance [75](#)
- parsing the access logs [80](#)
- performance [75](#)
- process count [75](#)
- purpose of [80](#)
- reasons for abort
 - deleting [81](#)
 - read permission revoked [81](#)
 - renaming [81](#)
- rules [82](#)

- Replay (*continued*)
 - running debug mode for troubleshooting [82](#)
 - starting [82](#)
 - throttling
 - how to control [75](#)
 - variables [76](#)

- Reports
 - endpoints
 - DELETE [9](#)
 - GET [9](#)
 - POST [9](#)
 - PUT [9](#)

- restarting
 - Notifier [84](#)

- Role-based access control
 - definition [3](#)
 - roles

- Admin [3](#)
 - CollectionAdmin [3](#)
 - Data Admin [3](#)
 - Data User [3](#)
 - Service User [3](#)

- roles
 - admin [3](#)
 - consumer [4](#)
 - data admin [3](#)
 - data user [3](#)
 - IBM Spectrum Discover [7](#)
 - producer [4](#)
 - service user [3](#)
- running
 - upgrade tool [113](#), [114](#)

S

- S3
 - best practices scanning [104](#)
 - Creating an S3 object storage connection [100](#)
 - Scanning an S3 object storage connection [102](#)
- S3 object storage data source connections [100](#)

- Scanner
 - IBM Cloud Object Storage [66](#)
 - messages [81](#)
 - tracking LIST process
 - next_key [80](#)
 - next_version [80](#)
 - stats file [80](#)
 - task.stats file [80](#)

- Scanning an Elastic Storage Server data source connection
 - Configure data source connections [105](#)

- scanning data source connections
 - from graphical user interface [62](#)

- single-node deployments
 - IBM Spectrum Discover [11](#)

- SizeRange buckets
 - defining settings [108](#)
 - definition [108](#)
 - editing [108](#)
 - modifying [108](#)
 - sizes [108](#)
 - using [108](#)

- source data types
 - IBM Cloud Object Storage Live Event [4](#)
 - IBM Cloud Object Storage Scan [4](#)

- source data types (*continued*)
 - IBM Spectrum Scale Live Event [4](#)
 - IBM Spectrum Scale Scan [4](#)
- source name [4](#)
- starting
 - Replay [82](#)
- stats files
 - Scanner
 - tracking LIST process [80](#)

T

- Tags
 - custom [5](#)
- TimeSinceAccess buckets
 - definition [108](#)
 - editing [108](#)
 - modifying [108](#)
 - using [108](#)
- TLS certificate
 - using custom certificate [110](#)
- transport layer service (TLS)
 - enabling [96](#)

U

- Upgrade
 - IBM Spectrum Discover
 - extra disks [115](#)
- upgrade tool
 - loading [113](#), [114](#)
 - preparing [113](#), [114](#)
 - running [113](#), [114](#)
- upgrading
 - IBM Spectrum Discover [113](#), [114](#)
- Using custom TLS certificate [110](#)

V

- Validating code integrity
 - Deploying and configuring [111](#)
- vault
 - renaming causes Replay to abort [81](#)
- vaults
 - excluding
 - example [78](#)
 - including
 - example [78](#)
 - invalid configuration [76](#)
 - invalid settings [76](#)
 - settings
 - exclude_all_vaults (list) [78](#)
 - include_all_vaults (boolean) [78](#)
 - vaults (dictionary) [78](#)
- virtual appliance
 - automated configuration [40](#)
 - configuring backup [34](#)
 - configuring CPU [36](#)
 - configuring memory [36](#)
 - configuring networking [40](#)
 - configuring storage [27](#)
 - configuring VMDK [28](#), [31](#)
 - deployment [21](#)

- virtual appliance deployment
 - automated configuration [40](#)
 - configuring CPU [36](#)
 - configuring memory [36](#)
 - configuring networking [40](#)
 - configuring storage [27](#)
 - configuring VMDK [27](#), [28](#), [31](#), [34](#)
 - connect to IBM Spectrum Scale [53](#), [54](#), [56](#), [57](#)

W

- writing files
 - cURL [98](#)
 - IBM Cloud Object Storage vault [98](#)



Product Number: 5737-I32
5737-SG1

SC27-9601-08

