

# Dell Data Lakehouse and Diskover: Creating AI Datasets from Unstructured Data

Transforming unstructured data into AI-Driven insights

February 2025

H04428

## White Paper

### Abstract

This white paper highlights how Dell Data Lakehouse and Diskover integrate to convert unstructured data into contextual, actionable datasets. By combining Diskover's metadata inventory capabilities with Dell Data Lakehouse's analytics power, the solution streamlines data indexing and contextualization.

## Copyright

The information in this publication is provided as is. Dell Inc. makes no representations or warranties of any kind with respect to the information in this publication, and specifically disclaims implied warranties of merchantability or fitness for a particular purpose.

Use, copying, and distribution of any software described in this publication requires an applicable software license.

Copyright © 2025 Dell Inc. or its subsidiaries. Published in the USA February 2025 H04428.

Dell Inc. believes the information in this document is accurate as of its publication date. The information is subject to change without notice.

# Contents

- Executive summary ..... 4
- Solution Components..... 5
- Solution Architecture..... 7
- Solution Validation..... 10
- Conclusion..... 27
- References..... 28

## Executive summary

### Overview

Unstructured data dominates enterprise environments but often remains untapped due to its complexity. The integration of Dell Data Lakehouse (DDLH) and Diskover revolutionizes how organizations transform this raw data into meaningful datasets for AI and generative AI (GenAI) workloads. Diskover's metadata inventory capabilities efficiently index and organize unstructured data, while DDLH's analytics unlock insights and contextual value. This seamless combination streamlines workflows, reduces redundancies, and enhances data-driven decision making. It enables businesses to aggregate and contextualize metadata, creating rich datasets ideal for machine learning (ML) and generative AI applications.

Scalable processing ensures the solution can handle vast datasets as businesses grow. GenAI applications benefit significantly from the contextualized data, improving training and deployment outcomes. Automation reduces manual effort, driving productivity while fostering innovation. Designed to adapt to both current and future data management needs, this solution positions organizations to innovate rapidly and respond to new challenges. By leveraging this integration, enterprises unlock the full potential of their unstructured data, fueling smarter operations and innovative advancements in AI and GenAI.

### Audience

This document is intended for enterprises with data lakes or a data lake strategy interested in empowering their organizations to act more quickly, effectively, and efficiently on their data. Audience roles include:

- Data and application administrators
- Data engineers
- Data scientists
- IT decision makers

A data lakehouse can assist more traditional analytics customers looking to modernize their data collection. It can also help analytics systems to get more value from their data or standardize their data for modern analytics workloads.

### Revisions

Date	Part number/ revision	Description
February 2025	H04428	Initial release

### We value your feedback.

Dell Technologies and the authors of this document welcome your feedback on this document. Contact the Dell Technologies team by [email](#).

**Authors:** Kirankumar Bhusanurmath, DA/AI Specialist | Dell Technologies

Chris Park, Brandon Langley | Diskover

---

**Note:** For links to other documentation for this topic, see the [Data Analytics Info Hub](#).

---

## Solution Components

### Dell Data Lakehouse

DDLH is a turnkey solution, boasting the Dell Data Analytics Engine, a potent federated data lake query engine powered by Starburst. The Dell Lakehouse System software ensures life cycle management, while tailor-made compute hardware provides seamless integration. Notably, this platform is designed to support AI solutions with its AI-ready data platform. For storing and processing large datasets in open file and table formats, Dell's leading S3 storage platforms, ECS, ObjectScale, and PowerScale, deliver exceptional performance, reliability, and security.

The DDLH core lies the Dell Data Analytics Engine (DDAE), powered by Starburst, facilitating the discovery, querying, and processing of enterprise-wide data assets regardless of their physical location. By reducing data movement requirements and enhancing query efficiency, the DDLH sets a new benchmark in data platform optimization and performance.

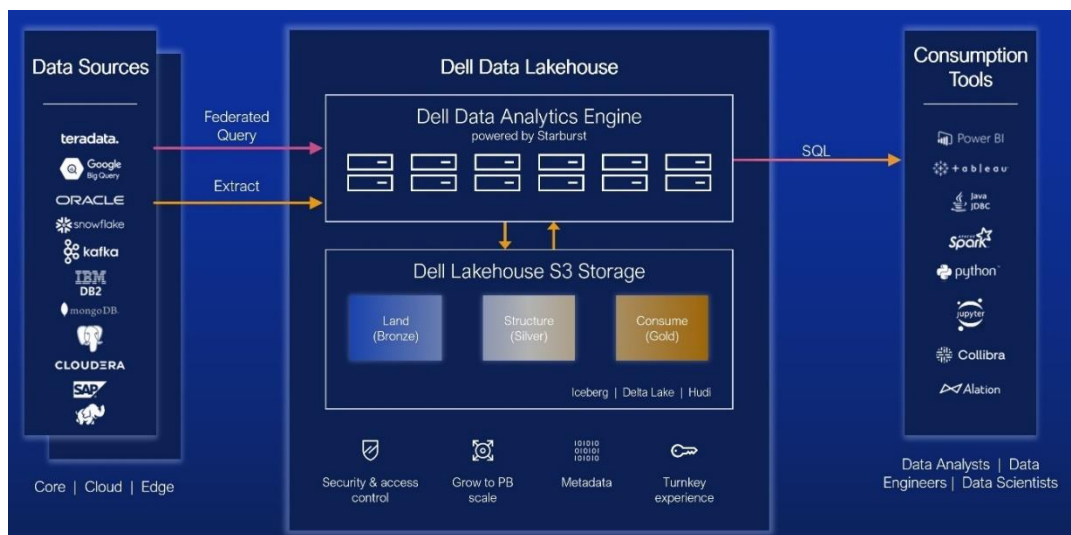


Figure 1. Dell Data Lakehouse diagram

### Dell Data Analytics Engine

DDAE contains an analytics query engine powered by Starburst. It is a fully supported enterprise-grade distributed SQL query engine designed for high-performance analytics. It allows users to query large amounts of data stored in various data sources throughout an organization using standard SQL syntax.

One of the key features of DDAE is its ability to run queries across different data sources simultaneously, in the same query. These sources include relational databases, NoSQL databases, Object Storage systems, and more. Response times are fast enough to support real-time analysis. Users can use this query engine to process data across multiple systems and data sources.

With the integrated query engine, administrators can implement a layer on top of data that abstracts away details on location, connectivity, language variations, and API. This layer of abstraction is critical to simplify data analytics over a diverse set of data sources.

### DDAE Metastore

DDAE Metastore is a dedicated Hive Metastore Server (HMS).

The Hive Metastore (HMS) is a central repository of metadata for Hive, Iceberg, Delta, and Hudi tables, and provides DDAE client access to this information using the Metastore service API. It is the building block for DDLH that uses the diverse world of open-source software, such as Apache Spark and DDAE's Trino.

## Diskover

Diskover is a powerful solution designed to provide global visibility into unstructured data, helping organizations unlock the full potential of their digital assets. With features like global indexing, search, and analytics, Diskover enables businesses to efficiently manage data spread across diverse systems and locations. Its metadata inventory capabilities allow for seamless aggregation, organization, and contextualization of data, making it an invaluable tool for addressing the challenges of data complexity.

The platform enhances security by identifying and mitigating risks associated with uncontrolled or siloed data. The Diskover index provides visibility without actual file system access enabling Data Scientist the ability to find and request access to data desired for AI pipeline. By enabling efficient search and retrieval of relevant datasets, Diskover significantly reduces the time and resources required for data management. This not only minimizes operational costs but also boosts productivity across teams. Diskover's ability to integrate with multiple systems and environments, including business intelligence (BI) and AI pipelines, positions it as a versatile solution for enterprises.

A standout feature of Diskover is its ability to enrich data with metadata context, enabling advanced curation processes such as deduplication, tagging, and purging. This ensures that data pipelines are streamlined, improving the quality and relevance of datasets for AI and generative AI applications. By fostering smarter decision making and accelerating innovation, Diskover empowers businesses to harness their unstructured data effectively, driving success in the modern data-driven landscape.

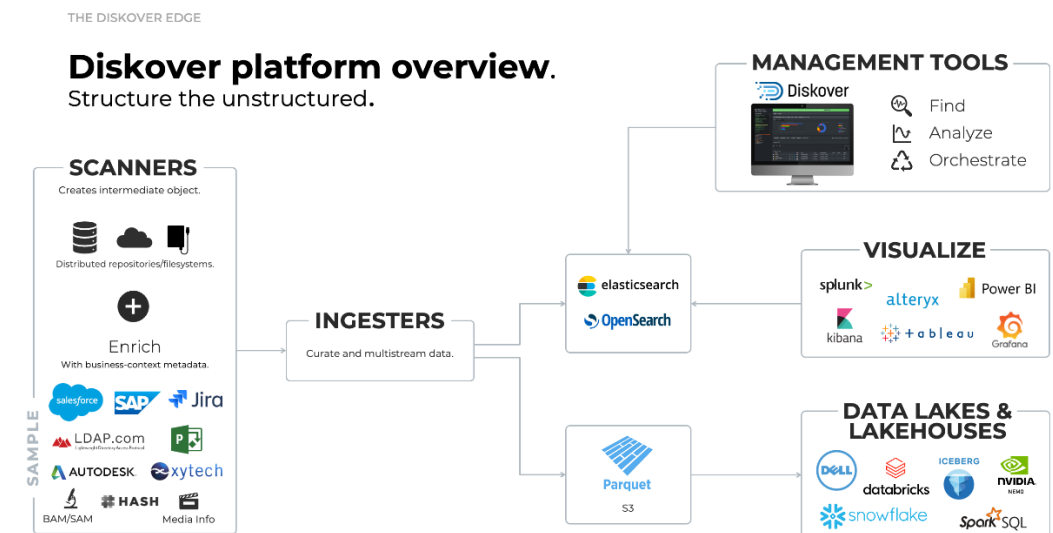
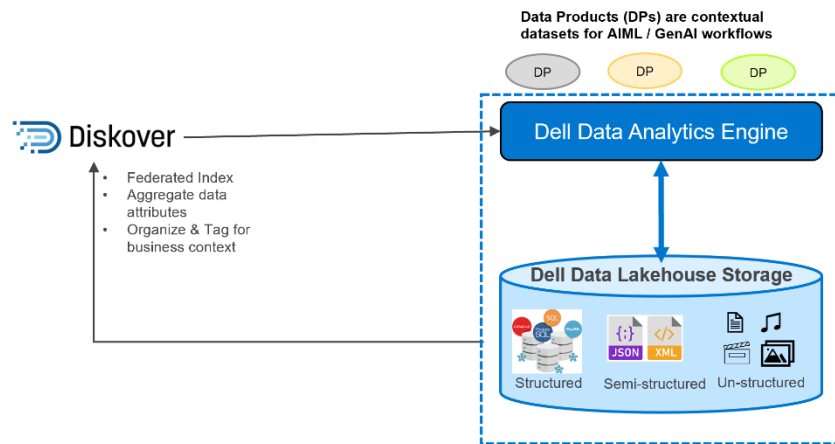


Figure 2. Diskover platform overview

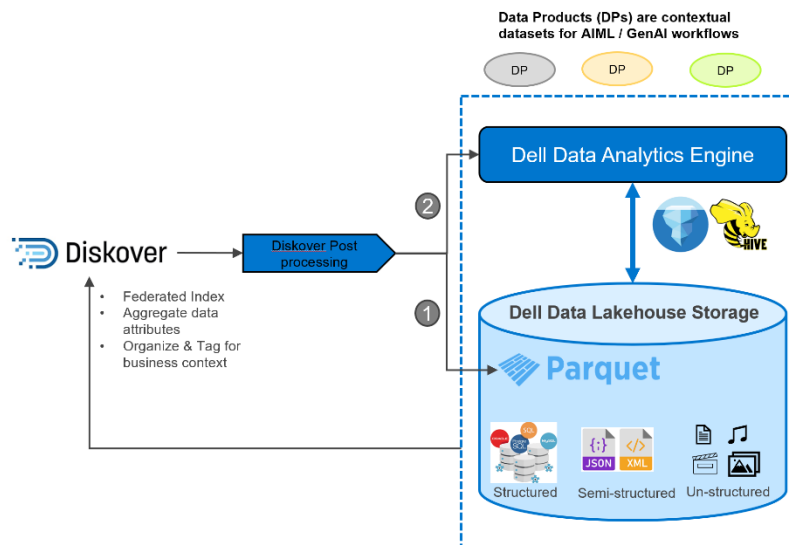
## Solution architecture

The integration of Discover and DDLH is designed to optimize metadata management, enabling enhanced organization, contextualization, and usability of unstructured data across AI and GenAI workflows. This architecture streamlines the flow of metadata for improved accessibility and scalability, ensuring seamless support for advanced analytics and decision making. While both solutions are fundamentally the same, they differ in their approach:

The first solution positions Discover as the metadata source, where DDLH actively pulls the metadata inventory files for processing:



The other utilizes DDLH as the destination, with Discover pushing metadata inventory files as partitioned Parquet files and registering the corresponding schema with DDLH.



The solutions within this framework are:

**Solution 1:** Discover Metadata Inventory as a Federated Source for DDLH

**Solution 2:** DDLH as a Destination for Discover Metadata Inventory Files

## Solutions

### Solution 1: Discover Metadata Inventory as a Federated Source for DDLH

This solution demonstrates how Discover and DDLH work together to streamline metadata management and make unstructured data more accessible and usable for AI and GenAI workflows. Here is a detailed step-by-step process:

#### **Step 1: Metadata Inventory Creation**

Discover scans and indexes unstructured data sources across various systems and generates a comprehensive metadata inventory. This inventory provides critical information about each file, including attributes like file type, size, location, timestamps, and content-based metadata.

#### **Step 2: Integration with ElasticSearch**

The metadata inventory created by Discover is stored in ElasticSearch, which acts as a central repository. This enables fast querying and retrieval of metadata, supporting high-performance integration workflows.

#### **Step 3: Metadata Registration in DDLH**

DDLH connects to the ElasticSearch repository and pulls metadata created by Discover. By leveraging this metadata as a federated source, DDLH integrates it into its system, creating a unified view of metadata records.

#### **Step 4: Exposing Metadata Inventory**

Once the metadata is registered within DDLH, it is exposed as tabular datasets. These datasets are structured in a way that makes them easily accessible and interpretable by downstream applications. This ensures compatibility with various BI tools and AI workloads, facilitating seamless collaboration between teams.

#### **Step 5: Leveraging Metadata for Advanced Workflows**

With the metadata inventory now accessible as a federated source, businesses can utilize it to build contextual datasets. These datasets serve as the foundation for advanced data analysis and AI model development, enabling more informed decision making and optimized operations.

#### **Integration and Data Flow**

The interaction between Discover and DDLH in this solution exemplifies a smooth data flow process:

- Discover acts as the upstream source, generating enriched metadata.
- ElasticSearch bridges the storage and querying of metadata between the two systems.
- DDLH pulls metadata from ElasticSearch, contextualizing and organizing it for usage in various workflows.

#### **Role of Metadata Management**

This solution highlights the pivotal role of metadata management, providing businesses with the ability to:

- Organize unstructured datasets for better accessibility
- Contextualize data to enhance usability across AI workflows



- Streamline the preparation of data for BI and AI applications

By harnessing the power of metadata management through the integration of Diskover and DDLH, organizations can unlock new efficiencies and accelerate their AI and GenAI initiatives.

## **Solution 2: DDLH as a Destination for Diskover Metadata Inventory Files**

This solution outlines how Diskover pushes its metadata inventory directly to DDLH, providing a structured and enriched foundation for BI and AI workflows. The following steps detail the process:

### ***Step 1: Metadata Inventory Creation***

Diskover scans unstructured data across various data sources and builds a detailed metadata inventory. This inventory captures file attributes such as size, type, location, timestamps, and content-based metadata, offering deep visibility into datasets.

### ***Step 2: Partitioned Parquet File Generation***

Diskover organizes the metadata inventory into partitioned Parquet files, an efficient and standardized file format. Partitioning ensures optimal data retrieval and processing performance by categorizing metadata based on logical criteria, such as time or file type.

### ***Step 3: Metadata Transfer to DDLH***

Diskover actively pushes these partitioned Parquet files to DDLH. By doing so, metadata files are directly delivered into the DDLH ecosystem, bypassing intermediary repositories like ElasticSearch.

### ***Step 4: Schema Registration in DDLH***

After transferring the metadata files, Diskover registers the corresponding schema with DDLH. This schema registration ensures that DDLH can effectively interpret and process the metadata files, enabling seamless integration with DDLH's querying and data transformation capabilities.

### ***Step 5: Structuring Metadata for BI and AI Workloads***

Within DDLH, the partitioned Parquet metadata files are made available as organized, tabular datasets. These datasets are optimized for querying and analysis, making them suitable for use in various BI tools and AI projects. This structured format allows teams to explore metadata insights and create sophisticated models tailored to their operational needs.

### ***Integration and Data Flow***

The integration process reflects the seamless handover of metadata between Diskover and DDLH:

- Diskover operates as the upstream system, processing and pushing metadata directly into DDLH.
- DDLH serves as the destination, where metadata is stored, interpreted, and made accessible for downstream workflows.

### ***Role of Metadata Management***

This solution underscores the importance of robust metadata management in enabling:

- Efficient data organization by leveraging partitioned file formats for scalability

## Solution validation

- Seamless schema alignment between source and destination systems
- Enhanced usability of metadata for diverse data-driven applications, such as decision making, AI models, and predictive analytics

By utilizing DDLH as the destination, organizations gain a centralized, scalable, and actionable metadata repository. This integration empowers businesses to maximize their unstructured data's potential, driving success in AI and GenAI initiatives.

## Key benefits

### Improved Data Accessibility

Seamless access to metadata provides teams with actionable insights directly from organized datasets.

### Enhanced Data Contextualization

Metadata integration ensures that data is not only accessible but also meaningful and relevant for analysis.

### Streamlined AI and GenAI Workflows

Easy access to prepared, structured data accelerates the development and deployment of AI models.

### Efficient Data Organization

Leveraging partitioned Parquet files and unified schemas ensures faster data retrieval and processing.

### Scalability

The solutions support growing data volumes without compromising performance or accessibility.

### Better Decision Making

Holistic and well-organized metadata lays the groundwork for better, data-driven insights and strategies.

### Optimized Operations

Reduced complexity in cataloging and accessing metadata ultimately improves operational efficiency.

### Accelerated AI Initiatives

Faster access to high-quality metadata supports rapid innovation in AI and GenAI projects.

These benefits collectively enable organizations to unlock the full potential of their unstructured data, driving success in data-intensive applications.

## Solution validation

The validation process for both Solution 1 and Solution 2 focuses on ensuring seamless integration and functionality between Discover and DDLH. Initially, the environment is set up for both Discover and DDLH, making sure all required configurations are in place. The

integration is then established to enable Diskover to scan the DDLH storage for unstructured data and build a comprehensive metadata inventory. This metadata inventory is further processed and integrated with the DDLH Engine, where it is transformed into structured tabular datasets and contextual data sets that can be utilized as data products. These steps ensure that the solutions are fully optimized for BI and AI workloads, meeting the set objectives for accuracy and efficiency.

## Solution 1: Diskover Metadata Inventory as a Federated Source for DDLH

### Set up Dell Data Lakehouse

It is assumed the DDLH appliance and Dell PowerScale storage cluster is installed and configured.

### Set up Dell Data Lakehouse storage

Configure Dell PowerScale as the primary storage cluster for the DDLH.

1. Log in to the DDLH system software.
2. Under **Storage** configure the Dell PowerScale S3 endpoint.

The screenshot displays the Dell Data Lakehouse System Software interface. On the left, a navigation sidebar includes options for Cluster, Catalogs, Storage (highlighted with a red box), Alerts, Logs, Infrastructure, and Licenses. The main content area is titled 'Edit connection' and shows configuration for 'default-s3storage'. Below the title, it instructs the user to 'Edit the parameters required to connect to a Dell ECS.' Two input fields are visible: 'Host name' with the value '172.17.1.22' and 'Port' with the value '9020'. Both fields are enclosed in a red rectangular box. At the bottom, there is an unchecked checkbox labeled 'Use SSL'.

Figure 3. Add S3 storage endpoint

## Setup Hive Catalog

Configure the Hive Catalog to storage Hive tables on the Dell PowerScale of the DDLH.

1. Log in to The DDLH system software
2. Under **Catalogs** select **Connect Catalog**.
3. Select **Properties**. For **Type** choose **Hive**. Enter the configuration parameters in the **Configuration** field.

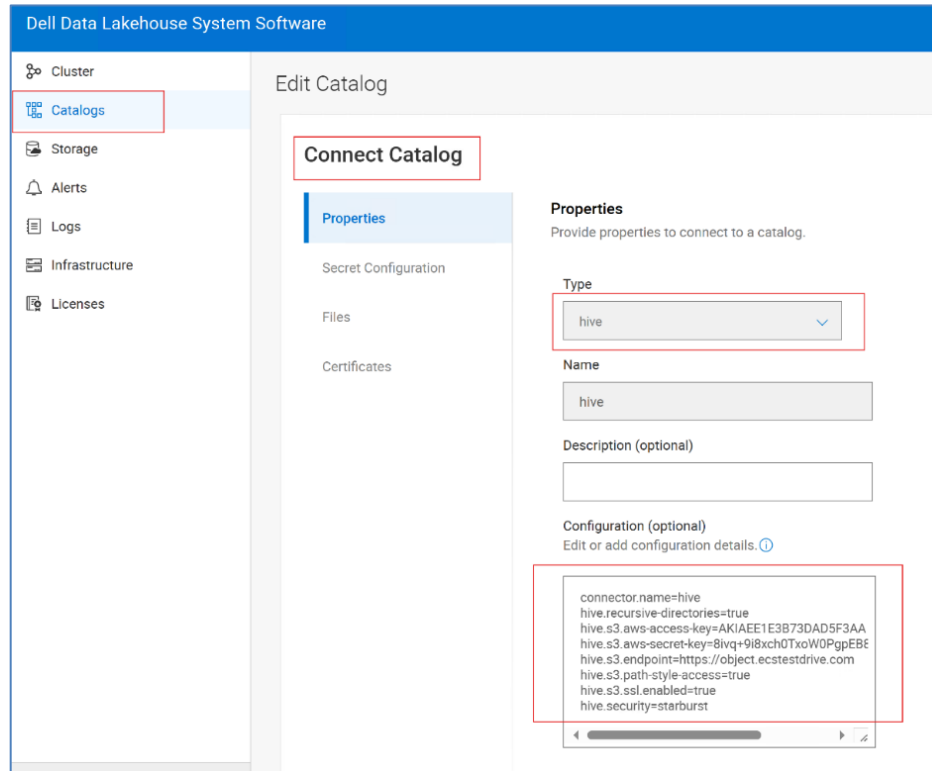


Figure 4. Add Hive Catalog

## Setup Diskover

### Diskover installation

On the utility node install Diskover. See [here](#) for further information. For more information, see [Diskover documentation](#).

### Diskover and DDL Storage S3 setup

On the utility node where Diskover is installed, S3fs either as python library or Linux package. Mount the S3 bucket of DDLH storage to the Diskover client Linux file system as an S3fs mount point. This enables the Diskover plugin to run additional content metadata extraction of each unstructured object it scans.

```
Last login: Fri Dec 20 14:46:31 2024 from 172.16.13.36
[diskover@localhost ~]$ mount | grep s3fs
s3fs on /mnt/ddlh_images type fuse.s3fs (rw,nosuid,nodev,relatime,user_id=0,group_id=0)
s3fs on /mnt/ddlh_pdfs type fuse.s3fs (rw,nosuid,nodev,relatime,user_id=0,group_id=0)
s3fs on /mnt/ddlh_videos type fuse.s3fs (rw,nosuid,nodev,relatime,user_id=0,group_id=0)
[diskover@localhost ~]$
```

If Dell PowerScale is used as the DDLH Storage, it is recommended to mount the PowerScale to Diskover as an NFS mount for better performance.

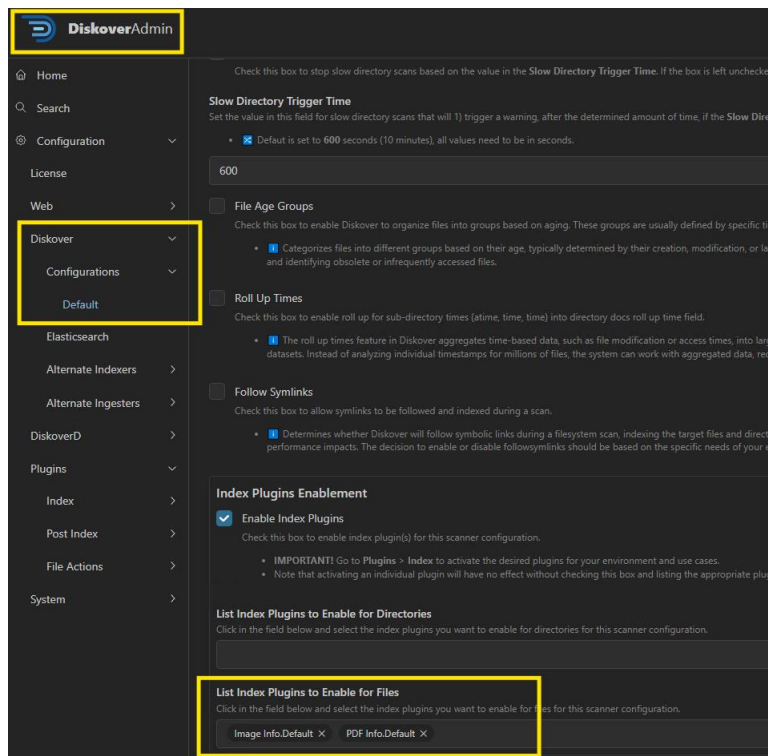
### Diskover plugin configuration

Diskover's powerful extensibility enables the development of custom plugins by Diskover, third parties, or even end users to enrich metadata catalogs. For this solution validation, we focus on two specific plugins, PDFs and images. These plugins are enhanced to extract additional system metadata, including keywords, subject, and the absolute path within the S3 bucket. They also capture detailed content metadata, providing insights into unstructured files and contextual objects contained within them.

```
[diskover@localhost plugins]$ pwd
/opt/diskover/plugins
[diskover@localhost plugins]$ ls -ltr | grep -E "imageinfo|pdfinfo"
drwxr-xr-x 3 root root 61 Nov 7 09:38 imageinfo
drwxr-xr-x 3 root root 61 Nov 7 09:51 pdfinfo
[diskover@localhost plugins]$
```

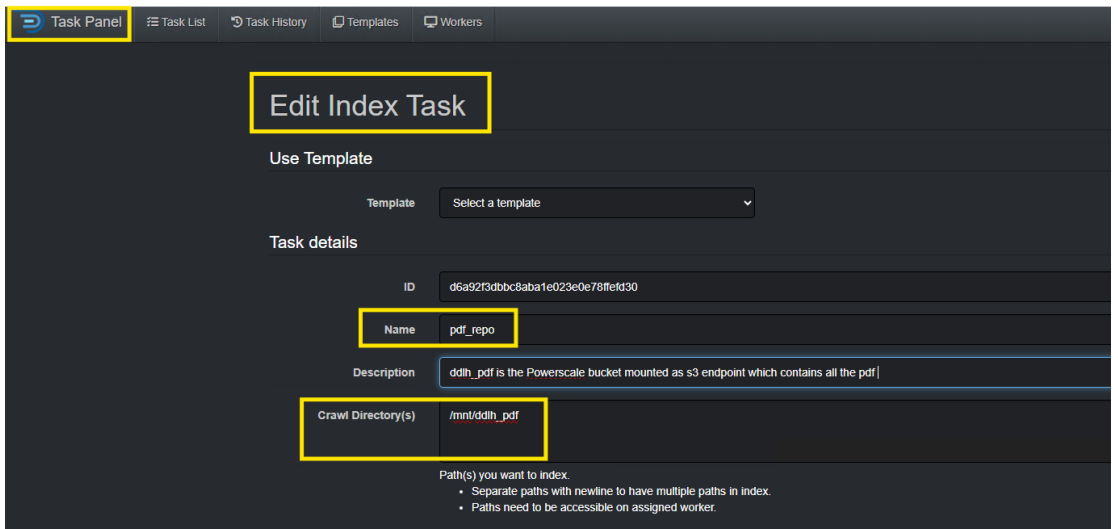
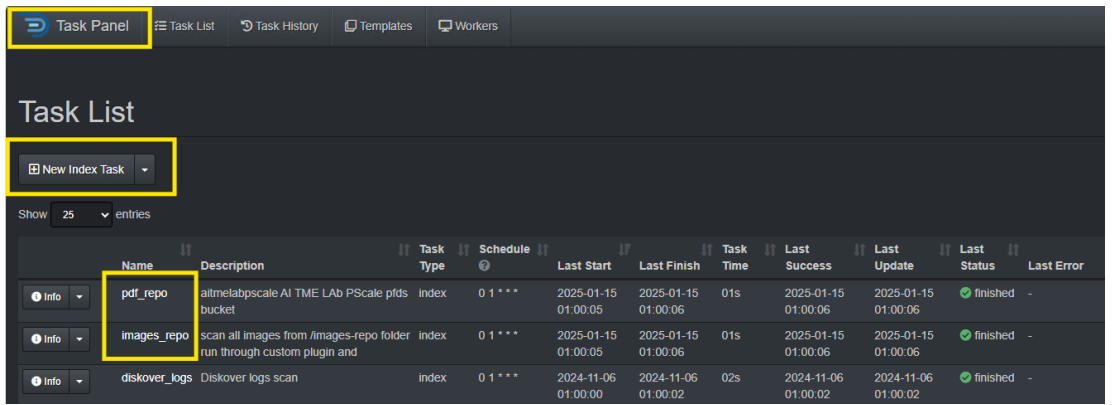
### Diskover enable plugins

Navigate to **Diskover > Configurations** and enable the enriched plugins.

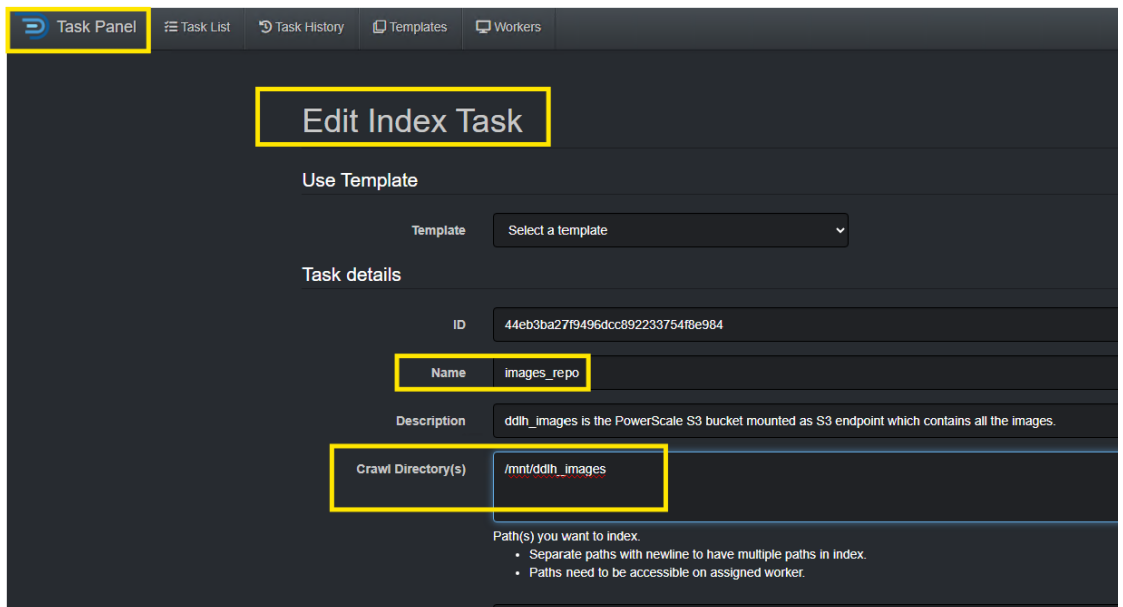


### Diskover Add New Index Task

Select **Task Panel** and add two **New Index Task** items, one for *pdf\_repo* and other for *images\_repo*.



New Index Task for pdf\_repo.

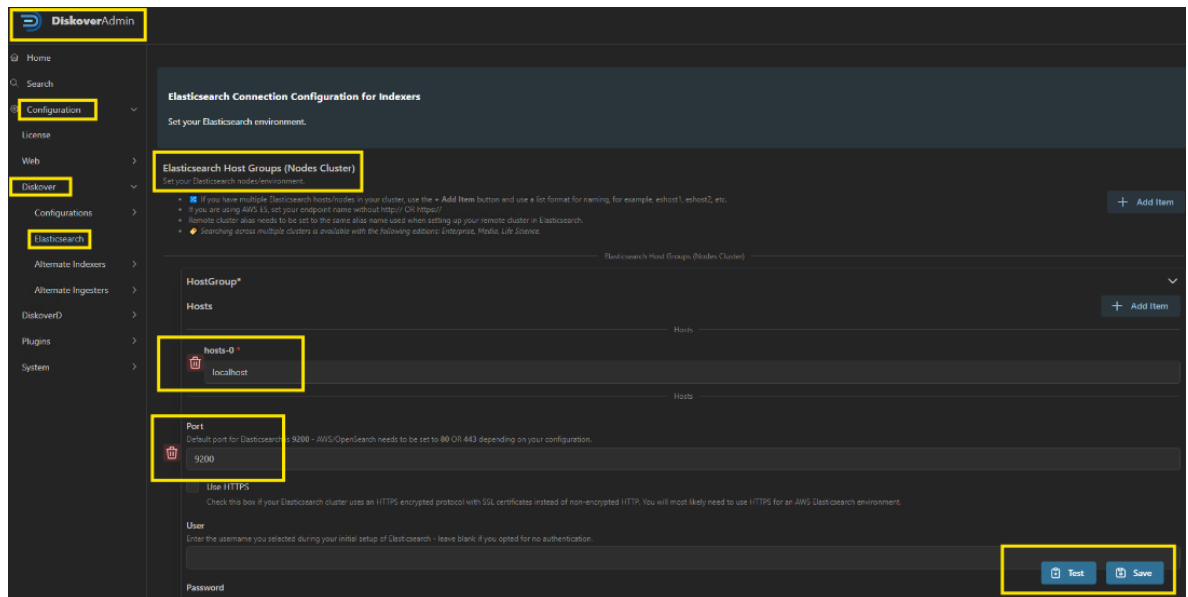


New Index Task for image\_repo.

## Set up ElasticSearch

To configure ElasticSearch in the Diskover Admin UI:

1. Log in and navigate to the Configuration section.
2. Locate the **ElasticSearch** configuration options and ensure the settings point to the default ElasticSearch instance included with Diskover. Verify that all configuration fields, such as host, port, and credentials match the default instance.
3. Save the settings and test the connection to confirm integration, ensuring seamless metadata indexing and search functionality.



## Set up DDLH and Diskover ElasticSearch Integration

To establish a connection between DDLH and the Diskover ElasticSearch instance, ensure network accessibility between the two systems and configure authentication credentials for secure communication. Verify proper metadata field mapping to maintain consistency during integration.

On the DDLH software UI, navigate to **Catalog** and create a new ElasticSearch catalog for Diskover Metadata Inventory to complete the setup.

## Solution Validation

### Validate metadata index tasks

Indices

Max indices to load: 250  Total 3 indices, indices are loaded in order by creation date. Max index is not used when use latest indices is checked. Always use latest indices (auto select)

Show indices newer than: All  Index name contains:

Select all  Unselect all  3 index(s) selected

Show 25  entries

Index	Index 2	Index Name	Top Path(s)	Start Time	Finish Time	Crawl Time	Files	Folders	Inodes/sec
<input checked="" type="checkbox"/>		diskover-opt_diskover-241210181500	/opt/diskover	2024-12-10 17:15:01	2024-12-10 17:15:01	00s	320	105	7,636.1
<input checked="" type="checkbox"/>		diskover-images_repo	/images_repo	2024-12-10 01:00:08	2024-12-10 01:00:20	12s	1,803	1	148.5
<input checked="" type="checkbox"/>		diskover-pdf_repo	/aitmelabpscale	2024-12-10 01:00:08	2024-12-10 01:00:09	01s	65	1	52.4

Showing 1 to 3 of 3 entries

To verify that ElasticSearch is populated with metadata indices for the configured PDF and image paths, check that the metadata inventory has been successfully scanned, extracted, and stored as indices by Diskover. Confirm that the indices corresponding to these file paths are present in ElasticSearch, ensuring they are correctly populated. Once verified, these indices are ready for immediate use.

Dell Data Lakehouse System Software

Cluster

Catalogs

Storage

Alerts

Logs

Infrastructure

Licenses

Edit Catalog

Connect Catalog

Properties

Secret Configuration

Files

Certificates

Summary

Properties

Provide properties to connect to a catalog.

Type

elasticsearch

Name

diskover-es

Description

dv-admin

Configuration

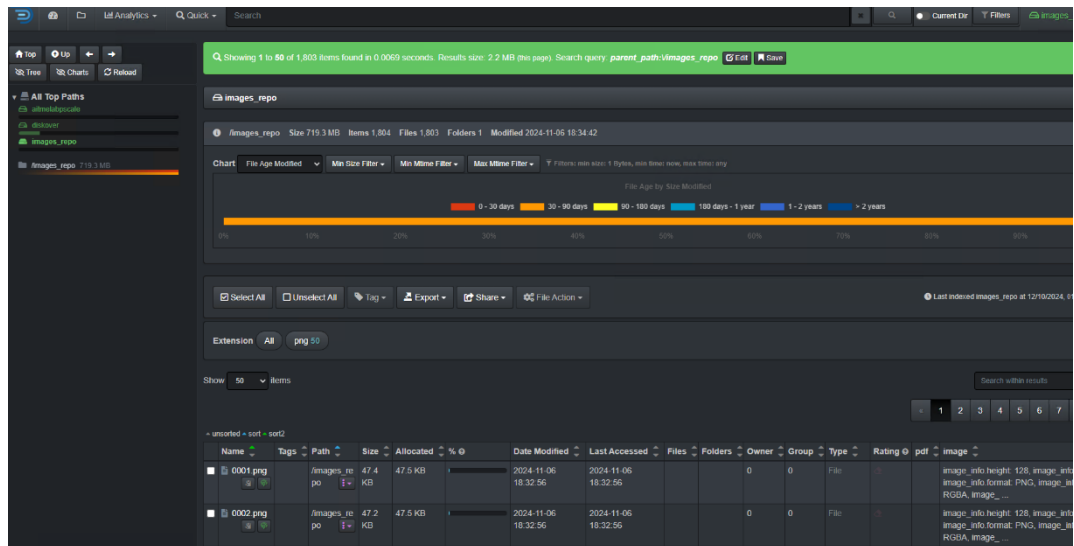
Edit or add configuration details.

connector name=elasticsearch  
elasticsearch.default-schema-name=default  
elasticsearch.host=172.16.10.202  
elasticsearch.port=9200



### Validate ElasticSearch populated with Metadata inventory

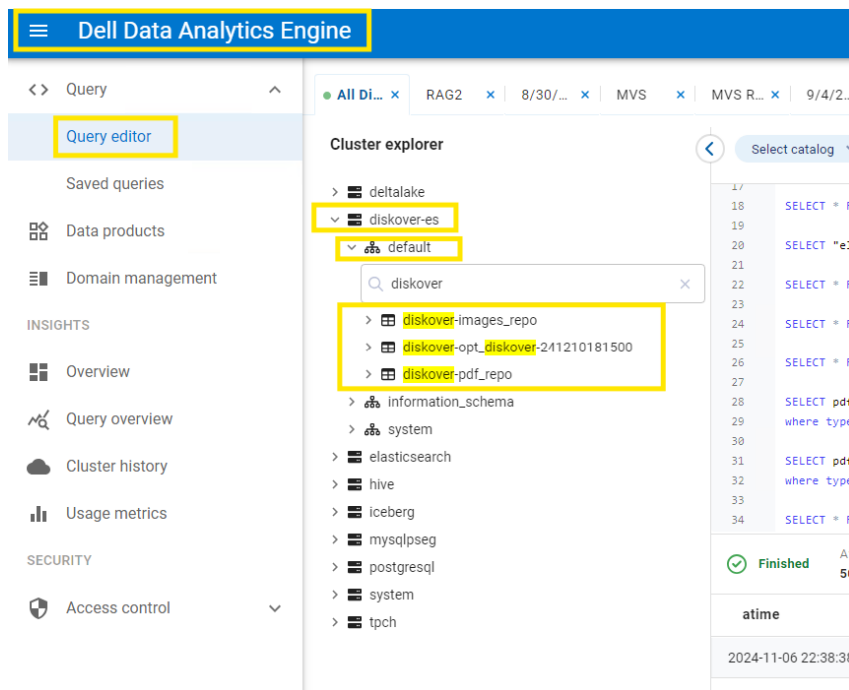
From the Diskover Web UI we can browse the metadata index, or Diskover can provide an analytics UI which populates the ElasticSearch Indices in more consumable dashboards and reports.



### Validate DDLH Metadata with ElasticSearch Indices

#### Validate ElasticSearch Catalog schema

From the DDAE UI, browse the Diskover catalog previously set up. Under the default schema you will see the Diskover metadata indices the PDF and image repos.



## Sample SQL on Diskover Metadata Indices

On the DDAE Query editor workspace, run a simple select SQL query to pull metadata information from Diskover ElasticSearch over the catalog connector that was established.

## Create sample data product as contextual datasets

The screenshot shows the Dell Data Analytics Engine interface. On the left, the 'Cluster explorer' shows a tree view with 'diskover-es' selected. The main area displays a SQL query in the 'Query editor' and its results in a table. The query is: `SELECT * FROM "diskover-es"."default"."diskover-pdf_repo" where type in ('JPG') LIMIT 10;` The results table has columns: 'atime', 'available', 'available\_percent', 'costpergb', 'crawl\_time', 'ctime', and 'dir\_co'. The table shows 10 rows of data, all with NULL values for most columns.

As the metadata is available in the DDLH, we can create contextual datasets as Data products. For more information on data products. See [Starburst Data Product Documentation](#).

From the DDAE UI, select **Data Product** and create new data product. For this validation we will create a separate image data product from the metadata inventory file from the Diskover Elasticsearch catalog. We will create two datasets under the data product: one for JPG files, and the other for PNG files.

The screenshot shows the 'Data products' page in the Dell Data Analytics Engine. A 'Create data product' button is visible at the top. Below it, two data products are listed: 'Amazon Reviews' and 'Image Data Products'. The 'Image Data Products' card is highlighted with a yellow box and has a 'PENDING CHANGES' badge. Both cards show a star rating and 'CREATED BY: DV-ADMIN'.

The screenshot displays the Dell Data Analytics Engine interface. The left sidebar contains navigation options: Query, Query editor, Saved queries, Data products (highlighted), Domain management, INSIGHTS (Overview, Query overview, Cluster history, Usage metrics), and SECURITY (Access control). The main panel shows 'Data product details' for 'Image Data Products'. It includes a warning about pending changes, an overview section with a summary of 'All Data Products consolidated', and a table showing 'Number of queries' and 'Number of users' for different time periods (7 and 30 days). Below this is a 'Description' section with a 'Datasets' list showing 'all\_jpg' and 'all\_pngs'.

## Solution 2: DDLH as a Destination for Diskover Metadata Inventory Files

### Setup Dell Data Lakehouse

These steps are same as described above in [Setup Dell Data Lakehouse](#).

### Setup Diskover

#### Diskover installation

On the utility node install Diskover. See [here](#) for further information. For more information, see [Diskover documentation](#).

#### Diskover and DDLH Storage S3 setup

On the utility node where Diskover is installed, S3fs either as python library or Linux package. Mount the S3 bucket of DDLH storage to the Diskover client Linux file system as an S3fs mount point. This enables the Diskover plugin to run additional content metadata extraction of each unstructured object it scans.

```
Last login: Fri Dec 20 14:46:31 2024 from 172.16.13.36
[diskover@localhost ~]$ mount | grep s3fs
s3fs on /mnt/diskover-meta-parquet type fuse.s3fs (rw,nosuid,nodev,relatime,user_id=0,group_id=0)
s3fs on /mnt/ddlh_pdfs type fuse.s3fs (rw,nosuid,nodev,relatime,user_id=0,group_id=0)
s3fs on /mnt/ddlh_images type fuse.s3fs (rw,nosuid,nodev,relatime,user_id=0,group_id=0)
[diskover@localhost ~]$
```

#### PowerScale S3bucket as the destination to Diskover Metadata inventory

Now we will mount the Dell PowerScale S3 bucket as the destination for Diskover to ingest Metadata inventory in the Parquet file format (open file format).

```
[diskover@localhost ~]$ mount | grep s3fs
s3fs on /mnt/diskover-meta-parquet type fuse.s3fs (rw,nosuid,nodev,relatime,user_id=0,group_id=0)
s3fs on /mnt/ddlh_pdfs type fuse.s3fs (rw,nosuid,nodev,relatime,user_id=0,group_id=0)
s3fs on /mnt/ddlh_images type fuse.s3fs (rw,nosuid,nodev,relatime,user_id=0,group_id=0)
[diskover@localhost ~]$
```

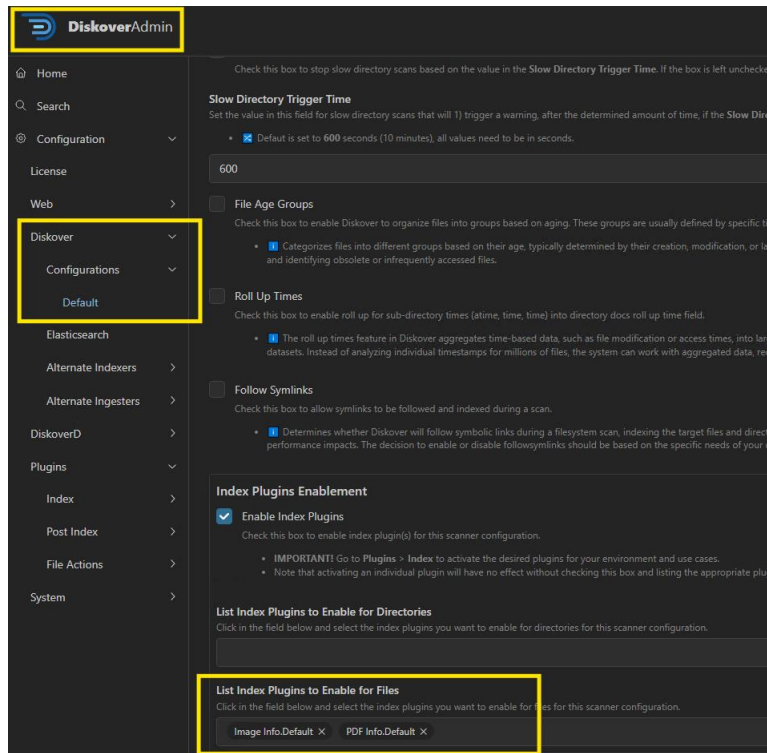
## Diskover plugin configuration

Diskover's powerful extensibility enables the development of custom plugins by Diskover, third parties, or even end users to enrich metadata catalogs. For this solution validation, we focus on two specific plugins, PDFs and images. These plugins are enhanced to extract additional system metadata, including keywords, subject, and the absolute path within the S3 bucket. They also capture detailed content metadata, providing insights into unstructured files and contextual objects contained within them.

```
[diskover@localhost plugins]$ pwd
/opt/diskover/plugins
[diskover@localhost plugins]$ ls -ltr | grep -E "imageinfo|pdfinfo"
drwxr-xr-x 3 root root 61 Nov 7 09:38 imageinfo
drwxr-xr-x 3 root root 61 Nov 7 09:51 pdfinfo
```

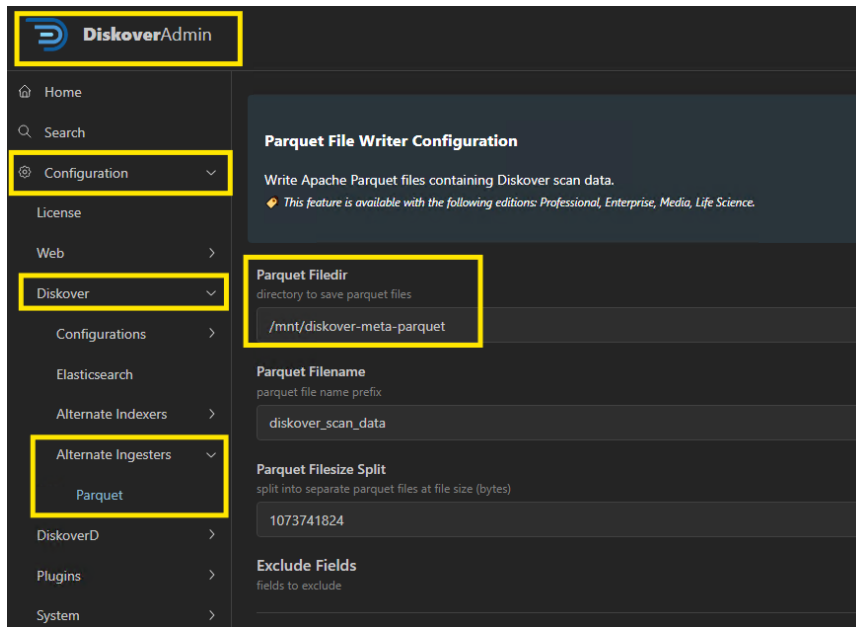
## Diskover enable plugins

Under Diskover Admin UI, enable the enriched plugins:



## Set up Diskover alternate ingesters to point to S3 bucket on Powerscale

This is a crucial step where we enable alternate ingesters in Diskover and choose the Parquet file format. Diskover converts the metadata inventory into Parquet file format and pushes it into the S3 bucket mounted on the utility node running Diskover.



## Solution Validation

### IMAGE Metadata Extraction and schema registration

1. Run the Diskover scanner on image repo.
2. On the utility node where Diskover is running we will run the Diskover scanner manually for the image repository. This is to scan the image repo S3 bucket on Dell Data Lakehouse storage that is mounted as S3 mount point to the utility node. The destination bucket will be again on the lakehouse storage S3 bucket. The metadata extracted needs to be partitioned as stored as parquet files.

```
[root@localhost diskover]# export PARQUETDIR=/mnt/diskover-meta-parquet/images-meta/year=$(date +%Y)/month=$(date +%m)/day=$(date +%d)/hour=$(date +%H)
[root@localhost diskover]# echo $PARQUETDIR
/mnt/diskover-meta-parquet/images-meta/year=2024/month=12/day=11/hour=11
[root@localhost diskover]# python3 diskover.py --alt.ingester parquet /mnt/ddlh_images/
```

```

diskover (1.0)
(v)(v)

"Crawling all your stuff."
v2.3.1
https://diskoverdata.com

2024-12-11 11:49:35,651 - diskover - INFO - Logging output to /var/log/diskover/diskover_mnt_ddlh_images_2024_12_11_11_49_35.log
2024-12-11 11:49:35,651 - diskover - INFO - Logging warnings to /var/log/diskover/diskover_mnt_ddlh_images_2024_12_11_11_49_35_warnings.log
2024-12-11 11:49:35,729 - diskover - INFO - Using alternate ingester module 'ingesters-parquet' from /opt/diskover/ingesters/parquet.py>
2024-12-11 11:49:36,085 - parquet_ingester - INFO - Found env var PARQUETDIR: /mnt/diskover-meta-parquet/images-meta/year=2024/month=12/day=11/hour=11
2024-12-11 11:49:36,102 - parquet_ingester - INFO - Making directory /mnt/diskover-meta-parquet/images-meta/year=2024/month=12/day=11/hour=11
2024-12-11 11:49:36,711 - diskover - INFO - configuration Default
2024-12-11 11:49:36,776 - diskover - INFO - Plugins loaded: PDF,Info,Default Image,Info,Default
2024-12-11 11:49:36,777 - diskover - INFO - maxwalkthreads set to 4
2024-12-11 11:49:36,777 - diskover - INFO - maxthreads set to 4
2024-12-11 11:49:36,777 - diskover - INFO - Enqueuing crawl /mnt/ddlh_images...
2024-12-11 11:49:36,778 - diskover - INFO - [Thread-1 (crawl_tree_thread)] crawling dir tree /mnt/ddlh_images...
2024-12-11 11:49:36,779 - diskover - INFO - searching for sub-dirs in /mnt/ddlh_images for config().threaddirdepth...
2024-12-11 11:49:36,959 - diskover - INFO - config().threaddirdepth set to 1
2024-12-11 11:49:36,989 - diskover - INFO - finding all sub-dirs for /mnt/ddlh_images up to level 1 directory depth...
2024-12-11 11:49:37,137 - diskover - INFO - found 0 subdirs in /mnt/ddlh_images (level 1)
2024-12-11 11:49:37,137 - diskover - INFO - starting scanning threads for /mnt/ddlh_images and 0 subdirs...
/mnt/ddlh_images/0001.png
{"height": 128, "width": 128, "format": "PNG", "mode": "RGBA", "is_animated": False, "n_frames": 1, "exif_orientation": "1"}
*****
/mnt/ddlh_images/0002.png
{"height": 128, "width": 128, "format": "PNG", "mode": "RGBA", "is_animated": False, "n_frames": 1, "exif_orientation": "1"}
*****
/mnt/ddlh_images/0003.png
{"height": 128, "width": 128, "format": "PNG", "mode": "RGBA", "is_animated": False, "n_frames": 1, "exif_orientation": "1"}
*****
/mnt/ddlh_images/0004.png
{"height": 128, "width": 128, "format": "PNG", "mode": "RGBA", "is_animated": False, "n_frames": 1, "exif_orientation": "1"}
/mnt/ddlh_images/0005.png
{"height": 128, "width": 128, "format": "PNG", "mode": "RGBA", "is_animated": False, "n_frames": 1, "exif_orientation": "1"}

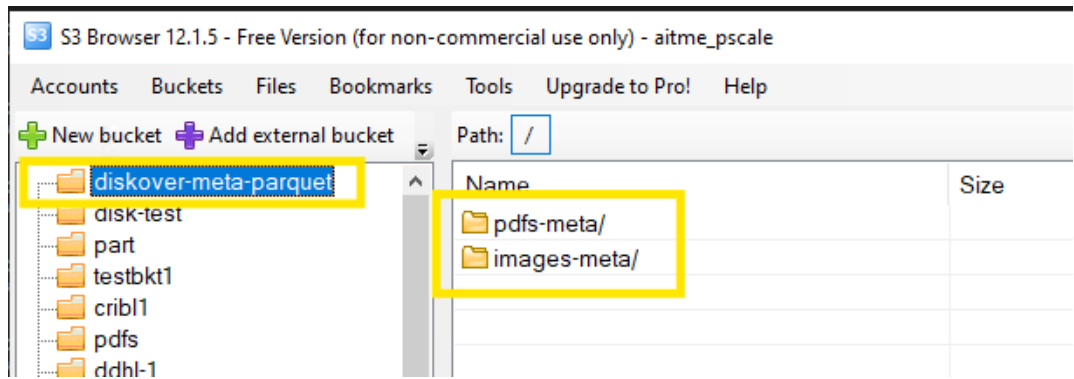
```

```

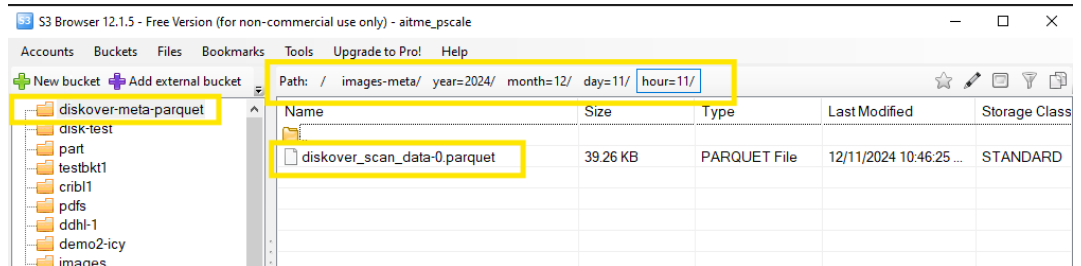
=====
/mnt/ddlh_images/r1de-97.jpg
{ "height": 275, "width": 183, "format": "JPEG", "mode": "RGB", "is_animated": false, "n_frames": 1 }
=====
/mnt/ddlh_images/r1de-98.jpg
{ "height": 303, "width": 166, "format": "JPEG", "mode": "RGB", "is_animated": false, "n_frames": 1 }
=====
/mnt/ddlh_images/r1de-99.jpg
{ "height": 159, "width": 109, "format": "JPEG", "mode": "RGB", "is_animated": false, "n_frames": 1 }
=====
2024-12-11 11:49:46.799 - diskover - INFO - [ThreadPoolExecutor-0_0] finished crawling /mnt/ddlh_images (0 dirs, 1803 files, 719.25 MB) in 0d:0h:00m:09s
2024-12-11 11:49:46.799 - diskover - INFO - *** finished walking /mnt/ddlh_images ***
2024-12-11 11:49:46.799 - diskover - INFO - *** walk files 1803, skipped 0 ***
=====
2024-12-11 11:49:46.799 - diskover - INFO - *** walk du size 719.84 MB ***
2024-12-11 11:49:46.799 - diskover - INFO - *** walk dirs 1, skipped 0 ***
2024-12-11 11:49:46.799 - diskover - INFO - *** walk took 0d:0h:00m:09s ***
2024-12-11 11:49:46.799 - diskover - INFO - *** walk perf 196.800 tnodes/s (max 0.000, min 0.000, avg 0.000) ***
2024-12-11 11:49:46.800 - diskover - INFO - *** docs indexed 1804 ***
2024-12-11 11:49:46.800 - diskover - INFO - *** indexing perf 196.712 docs/s (max 0.000, min 0.000, avg 0.000) ***
2024-12-11 11:49:46.800 - diskover - INFO - *** indexing took 0d:0h:00m:09s ***
=====
2024-12-11 11:49:46.800 - diskover - INFO - [Thread-1 [crawl_tree_thread]] crawling dir tree /mnt/ddlh_images completed in 0d:0h:00m:10s
2024-12-11 11:49:46.814 - parquet_ingester - INFO - writing data to parquet file /mnt/diskover-meta-parquet/images-meta/year=2024/month=12/day=11/hour=11/diskover_scan_data-0.parquet...
2024-12-11 11:49:46.874 - parquet_ingester - INFO - Wrote 1804 rows to 1 parquet file(s)
=====

```

Verify the destination S3 bucket (using S3 Browser in this example) is populated with the image metadata inventory file as Parquet files:



Verify the proper partitioning is retained for the Image metadata repository:



On DDLH, run the schema extraction script on S3 bucket where Diskover saves image metadata in Parquet file format.

```

49 ----- Infer DDL schema from the image meta parquet file from s3 bucket
50 select sql from hive.schema_discovery.discovery
51 where
52 uri='s3a://discover-meta-parquet/images-meta';
54

```

Finished	Avg. read speed	Elapsed time	Rows	Results from cache
✓	16.9 rows/s	0.06s	1	No

```

sql
CREATE SCHEMA IF NOT EXISTS "discovered" WITH (location = 's3a://discover-meta-parquet/images-meta/');

```

```

USE "discovered";
CREATE TABLE "images-meta" (
  "name" varchar,
  "extension" varchar,
  "parent_path" varchar,
  "size" bigint,
  "size_du" bigint,
  "owner" bigint,
  "group" bigint,
  "mtime" varchar,
  "atime" varchar,
  "ctime" varchar,
  "type" varchar,
  "image_info" row("exif_orientation" varchar,"format" varchar,"height" bigint,"is_animated" boolean,"mode" varchar,"n_frames" bigint,"width" bigint),
  "size_norecurs" double,
  "size_du_norecurs" double,
  "file_count" double,
  "file_count_norecurs" double,
  "dir_count" double,
  "dir_count_norecurs" double,
  "dir_depth" double,
  "year" varchar,
  "month" varchar,
  "day" varchar,
  "hour" varchar,
  "year" int,
  "month" int,
  "day" int,
  "hour" int
)
WITH (
  format = 'PARQUET',
  external_location = 's3a://discover-meta-parquet/images-meta/',
  partitioned_by = ARRAY['year', 'month', 'day', 'hour']
);

CALL system.sync_partition_metadata("discovered", 'images-meta', 'ADD');

```

Copy the schema discovery output, (DDL schema), and register it as a table under the schema *discover\_meta\_parquet*. The new table created is called *images\_meta*.

The screenshot shows the Dell Data Analytics Engine interface. On the left, the 'Cluster explorer' pane shows a tree view of schemas, with 'discover\_meta\_parquet' expanded to show the 'images\_meta' table. The main pane displays the SQL execution results, including the DDL schema for the 'images\_meta' table. The schema includes fields like 'name', 'extension', 'parent\_path', 'size', 'size\_du', 'owner', 'group', 'mtime', 'atime', 'ctime', 'type', 'image\_info', 'size\_norecurs', 'size\_du\_norecurs', 'file\_count', 'file\_count\_norecurs', 'dir\_count', 'dir\_count\_norecurs', 'dir\_depth', 'year', 'month', 'day', 'hour', and their respective data types. The table is registered with the format 'PARQUET', external location 's3a://discover-meta-parquet/images-meta/', and partitioned by 'year', 'month', 'day', and 'hour'.

### Run sample SQL queries on images metadata table:

The screenshot shows the Dell Data Analytics Engine interface. On the left, the 'Cluster explorer' shows a tree view with 'hive' selected, and 'diskover' > 'diskover\_meta\_parquet' > 'images\_meta' highlighted. The main area displays a SQL query:

```
SELECT * FROM "hive"."diskover_meta_parquet"."images_meta"
where extension='jpg' LIMIT 10;
```

The results table shows the following data:

name	extension	parent_path	size	size_du
0001.png	png	/mnt/ddlh_images	48554	48640
0002.png	png	/mnt/ddlh_images	48290	48640
0003.png	png	/mnt/ddlh_images	43858	44032
0004.png	png	/mnt/ddlh_images	47382	47616
0005.png	png	/mnt/ddlh_images	45703	46080
0006.png	png	/mnt/ddlh_images	46484	46592

### PDFs Metadata Extraction and schema registration

Run the Diskover scanner on PDF repo.

The terminal screenshot shows the following commands and output:

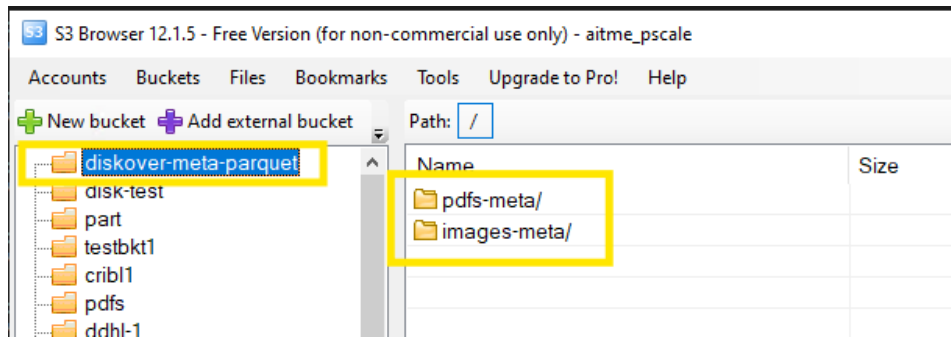
```
root@localhost diskover# export PARQUETDIR=/mnt/diskover-meta-parquet/pdfs-meta/year=${date +%Y}/month=${date +%m}/day=${date +%d}/hour=${date +%H}
root@localhost diskover# echo $PARQUETDIR
/mnt/diskover-meta-parquet/pdfs-meta/year=2024/month=12/day=11/hour=08
root@localhost diskover# python3 diskover.py --altIngestor parquet /mnt/ddlh_pdfs/
```

The output shows the diskover scanner running and logging information:

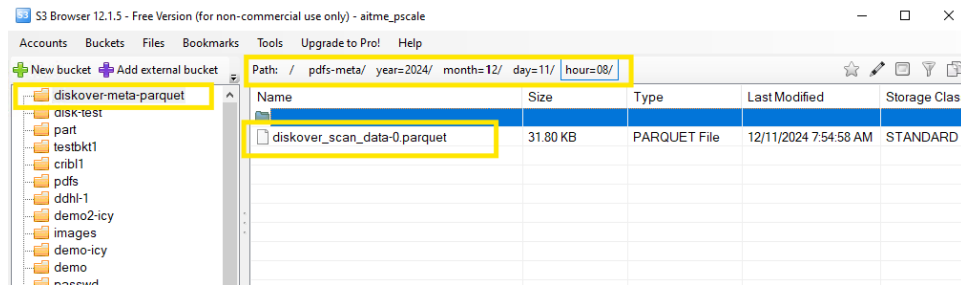
```
2024-12-11 08:58:17.100 - diskover - INFO - Logging output to /var/log/diskover/diskover_mnt_ddlh_pdfs_2024_12_11_08_58_17.log
2024-12-11 08:58:17.100 - diskover - INFO - Logging warnings to /var/log/diskover/diskover_mnt_ddlh_pdfs_2024_12_11_08_58_17_warnings.log
2024-12-11 08:58:17.181 - diskover - INFO - Using alternate ingestor module 'ingesters.parquet' from '/opt/diskover/ingesters/parquet.py' =>
2024-12-11 08:58:17.542 - parquet_ingester - INFO - Found env var PARQUETDIR: /mnt/diskover-meta-parquet/pdfs-meta/year=2024/month=12/day=11/hour=08
2024-12-11 08:58:17.542 - parquet_ingester - INFO - Making directory /mnt/diskover-meta-parquet/pdfs-meta/year=2024/month=12/day=11/hour=08
2024-12-11 08:58:17.771 - diskover - INFO - configuration default
2024-12-11 08:58:17.834 - diskover - INFO - Plugins loaded: Image Info.Default PDF Info.Default
2024-12-11 08:58:17.834 - diskover - INFO - maxwalkthreads set to 4
2024-12-11 08:58:17.834 - diskover - INFO - maxthreads set to 4
2024-12-11 08:58:17.835 - diskover - INFO - Enqueuing crawl /mnt/ddlh_pdfs...
2024-12-11 08:58:17.836 - diskover - INFO - [Thread-1 (crawl_tree_thread)] crawling dir tree /mnt/ddlh_pdfs...
2024-12-11 08:58:17.836 - diskover - INFO - searching for sub-dirs in /mnt/ddlh_pdfs for config().threaddirdepth...
2024-12-11 08:58:17.844 - diskover - INFO - config().threaddirdepth set to 1
2024-12-11 08:58:17.844 - diskover - INFO - finding all sub-dirs for /mnt/ddlh_pdfs up to level 1 directory depth...
2024-12-11 08:58:17.851 - diskover - INFO - found 0 subdirs in /mnt/ddlh_pdfs (level 1)
2024-12-11 08:58:17.851 - diskover - INFO - starting scanning threads for /mnt/ddlh_pdfs and 0 subdirs...
2024-12-11 08:58:18.057 - diskover - INFO - [ThreadPoolExecutor-0_0] finished crawling /mnt/ddlh_pdfs (0 dirs, 55 files, 54.77 MB) in 0d:0h:00m:00s
2024-12-11 08:58:18.057 - diskover - INFO - *** finished walking /mnt/ddlh_pdfs ***
2024-12-11 08:58:18.057 - diskover - INFO - *** walk files 55, skipped 0 ***
2024-12-11 08:58:18.057 - diskover - INFO - *** walk size 54.77 MB ***
2024-12-11 08:58:18.057 - diskover - INFO - *** walk du size 54.78 MB ***
2024-12-11 08:58:18.057 - diskover - INFO - *** walk perf 68.250 inodes/s (max 68.250, min 68.250, avg 68.250) ***
2024-12-11 08:58:18.057 - diskover - INFO - *** docs indexed 56 ***
2024-12-11 08:58:18.057 - diskover - INFO - *** indexing perf 69.491 docs/s (max 69.491, min 69.491, avg 69.491) ***
2024-12-11 08:58:18.057 - diskover - INFO - *** indexing took 0d:0h:00m:00s ***
2024-12-11 08:58:18.057 - diskover - INFO - *** warnlog/errlogs 0 ***
2024-12-11 08:58:18.057 - diskover - INFO - [Thread-1 (crawl_tree_thread)] crawling dir tree /mnt/ddlh_pdfs completed in 0d:0h:00m:00s
2024-12-11 08:58:18.665 - parquet_ingester - INFO - Writing data to parquet file /mnt/diskover-meta-parquet/pdfs-meta/year=2024/month=12/day=11/hour=08/diskover_scan_data-0.parquet...
2024-12-11 08:58:18.723 - parquet_ingester - INFO - Wrote 56 rows to 1 parquet file(s)
```

Verify the destination S3 bucket is populated with the PDF metadata inventory file as Parquet files.

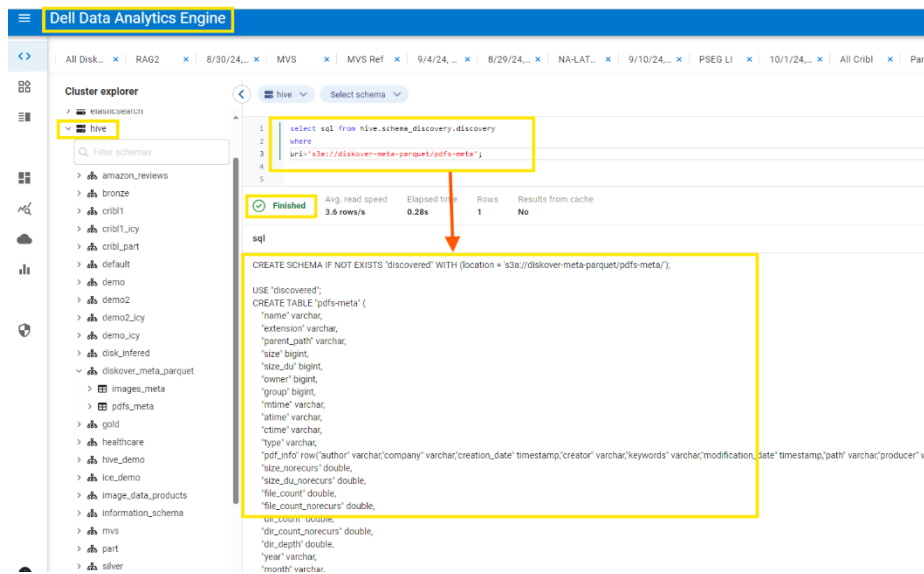




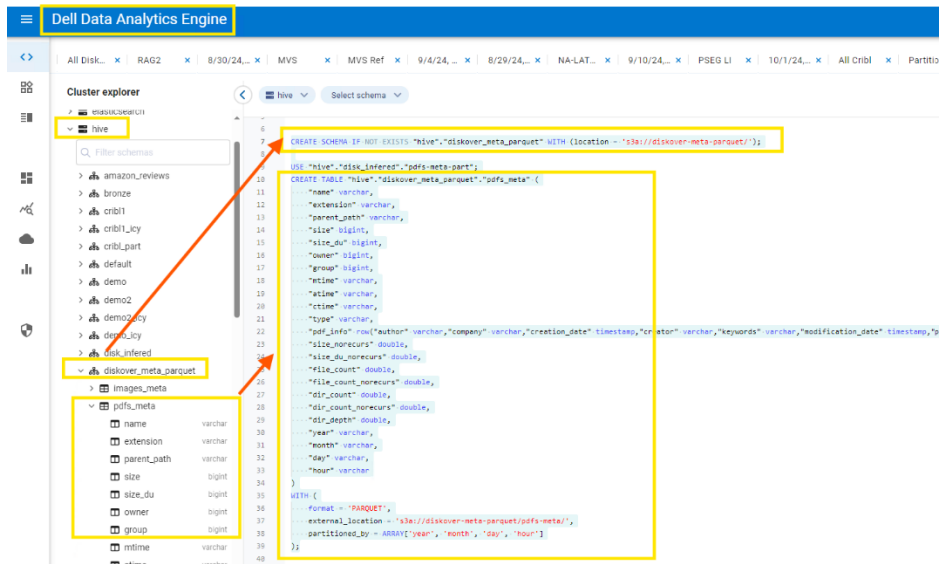
Verify the proper partitioning is retained for the PDF metadata repository.



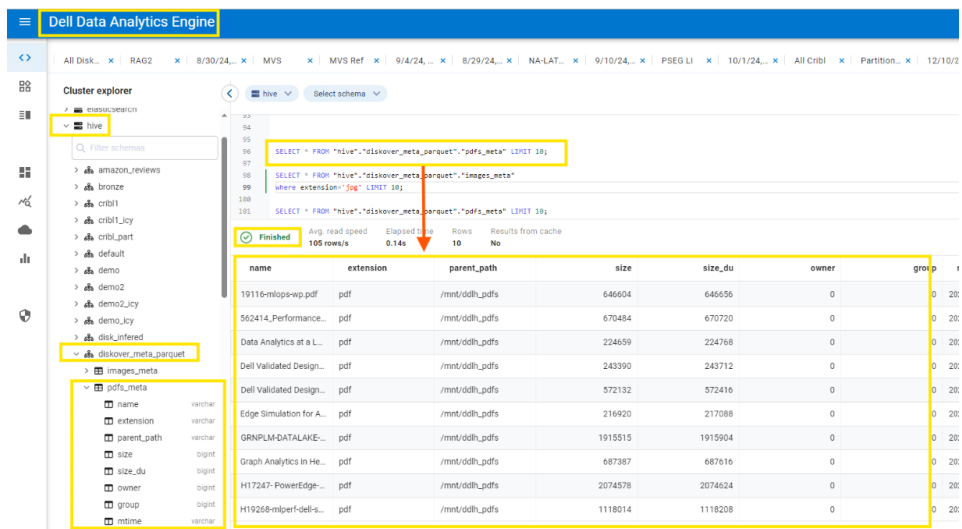
On DDLH, run the schema extraction script on the S3 bucket where Diskover saves PDF metadata in Parquet file format.



Register PDF metadata inventory file's schema as structured data in the DDLH catalog.



Run sample SQL queries on the PDF metadata table.



## Conclusion

The integration of DDLH and Diskover offers a robust approach to optimizing metadata management and addressing the challenges of unstructured data. By validating two distinct solutions, this architecture proves its ability to enhance data usability and drive efficiency for AI and GenAI workflows.

**Solution 1** positions Diskover as a federated metadata source, enabling DDLH to pull metadata inventory files from Elasticsearch, perform contextualization, and deliver actionable insights with its advanced analytics.

**Solution 2** utilizes DDLH as the destination, where Diskover exports partitioned Parquet files and registers schemas for seamless ingestion and processing.

Both solutions highlight the synergy between Diskover's powerful metadata inventory capabilities and DDLH's scalable high-performance analytics. Diskover effectively captures, organizes, and structures metadata, creating a foundation for transforming raw data into meaningful, contextualized datasets. Meanwhile, DDLH supports large-scale data processing, essential for building AI pipelines and supporting demanding workloads. This collaboration enables faster model training, improved deployment outcomes, and reduced redundancies across workflows.

The benefits of this integration extend beyond efficiency. Scalable processing ensures the ability to handle vast data volumes as businesses grow, while automation reduces manual intervention, fostering productivity and innovation. By bridging the divide between data complexity and actionable insights, the combination of DDLH and Diskover unlocks unprecedented opportunities for AI and GenAI applications, positioning organizations to drive smarter operations and achieve technological advancements in an increasingly data-centric world.

## References

### Dell Technologies documentation

The following Dell Technologies documentation provides other information related to this document. Access to these documents depends on your login credentials. If you do not have access to a document, contact your Dell Technologies representative.

Additional information can be found on the [Dell Technologies Info Hub for Data Analytics](#). If you need additional services or implementation help, contact your Dell Technologies sales representative.

Document type	Location
Dell Data Lakehouse	<a href="#">Dell Data Lakehouse Technical Solution Guide</a>
	<a href="#">Dell Data Lakehouse Specification Sheet</a>
	<a href="#">Dell Data Lakehouse Sizing and Configuration Guide</a>
	<a href="#">Dell Data Lakehouse Solution Brief</a>
Resilient Data Pipelines	<a href="#">Resilient Data pipelines on Dell Data Lakehouse</a>
Data Governance solution	<a href="#">Privacera Platform with Dell Data Lakehouse</a>
RAG Chatbot on DDLH	<a href="#">Multimodal RAG Chatbot Powered by Dell Data Lakehouse</a>
Mainframe Data Analytics	<a href="#">Unlock mainframe data using Dell Data Lakehouse</a>
Cribl Streaming with DDLG	<a href="#">Dell Data Lakehouse with Cribl Stream for Scalable Real-Time Data Processing and Analytics</a>

### Diskover documentation

The following [Diskover](#) documentation provides additional and relevant information.

Document type	Location
Diskover	<a href="#">Diskover documentation</a>

### Dell Technologies Info Hub

The [Dell Technologies Info Hub](#) is your one-stop destination for the latest information about Dell Solutions products. New material is frequently added, so browse often to keep up to date on the expanding Dell portfolio of cutting-edge products and solutions.

### More information

For more information, including sizing guidance, technical questions, or sales assistance, email [Analytics.Assist@dell.com](mailto:Analytics.Assist@dell.com), or contact your Dell Technologies or authorized partner sales representative.