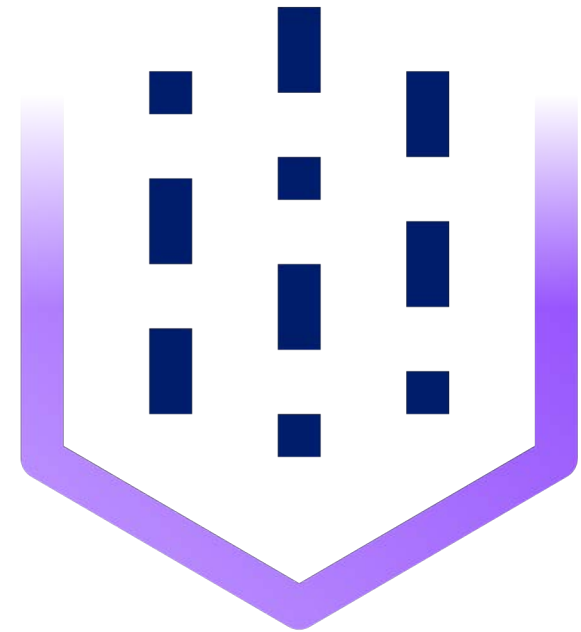# Introduction to Spectrum Discover

Spectrum Scale German User Meeting 2022
Cologne, Germany – October 19-20, 2022

Lars Lauber (IBM)
lauberla@de.ibm.com

# Disclaimer

**Spectrum Scale**

IBM's statements regarding its plans, directions, and intent are subject to change or withdrawal without notice at IBM's sole discretion. Information regarding potential future products is intended to outline our general product direction and it should not be relied on in making a purchasing decision. The information mentioned regarding potential future products is not a commitment, promise, or legal obligation to deliver any material, code, or functionality. The development, release, and timing of any future features or functionality described for our products remains at our sole discretion.

IBM reserves the right to change product specifications and offerings at any time without notice. This publication could include technical inaccuracies or typographical errors. References herein to IBM products and services do not imply that IBM intends to make them available in all countries.

# Content

- The struggle with AI and big data

- Metadata holds the key for unstructured data

- What can Spectrum Discover do to help?

- Using Spectrum Discover for data analysis

- Spectrum Discover Integration with Spectrum Scale

Those who can harness the power of their data have a competitive advantage!

**20%**

of the world's data is searchable and anybody can get to it

The other

**80%**

is like gold

"Initial data analysis for the simplest project requires brutal, excruciating data wrangling. More than 80% of a project's time spent is me trying to **parse these files** and **cross data sources** in order to build viable datasets."

– Data Scientist at a large company

# Existing data is difficult to leverage and trapped

## 50%
want to use long term data for AI/Analytics

## 84%
of digital transformation projects fail due to siloed data and unreliable integration

- Complex to manage
- Difficult to leverage data from silos
- Not designed to leverage AI
- Inefficient to search

# AI & big data *can* be challenging …

❑ Data volume and data quality

❑ Difficult to combine multiple data sources

❑ Skills gaps to analyze data

❑ Current infrastructure complexity

❑ Where and how to begin?

… but they don't have to be!

# Metadata is the Key to Data Organization & Insight

```
<!DOCTYPE html PUBLIC "-//W3C/
<html xmlns="http://www.w3.org
    <head>
        <meta http-equiv="Content-
        <meta http-equiv="Content-
        <meta http-equiv="Content-
```

## Three Types of Metadata

**1**

Metadata can come
from a system

*(owner, last modified,
size, type, etc.)*

**2**

Metadata can be custom
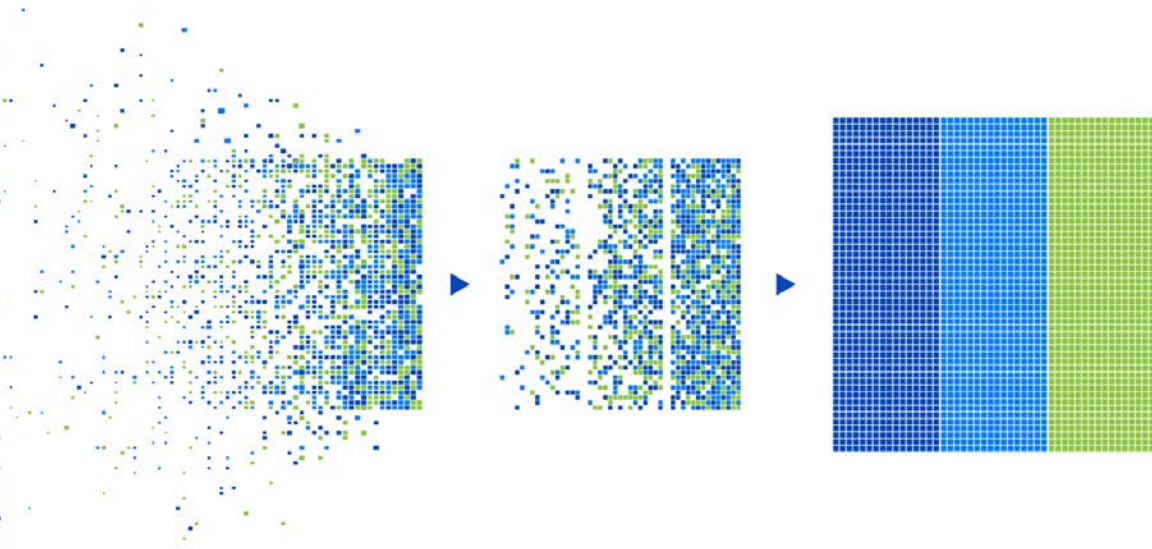
*(map to various business &
scientific aspects)*

**3**

Metadata can be
derived from analytics

*(percent confident)*

7

# Metadata is the key

**Bring structure to unstructured data …**



**Why metadata?**

- Improves management
- Organizes large amounts of data
- Creates "indexes" or "tags" for data
- Provides a fast way to search and analyze

**Metadata is data about data**

- Context for data classification and management

**Types of metadata**

- **System**: information about file and object types, their sizes, when they were last modified, etc.
- **Custom**: user/organization-defined based on unique taxonomy
- **Derived**: derived from analytics and applied to your data enriching the metadata model with additional meaning

**Benefits of metadata**

- **Identify and manage** assets that add value
- **Simplify search & access** to critical data
- **Define and execute policies** based on metadata
- **Improve** time-to-value and storage economics
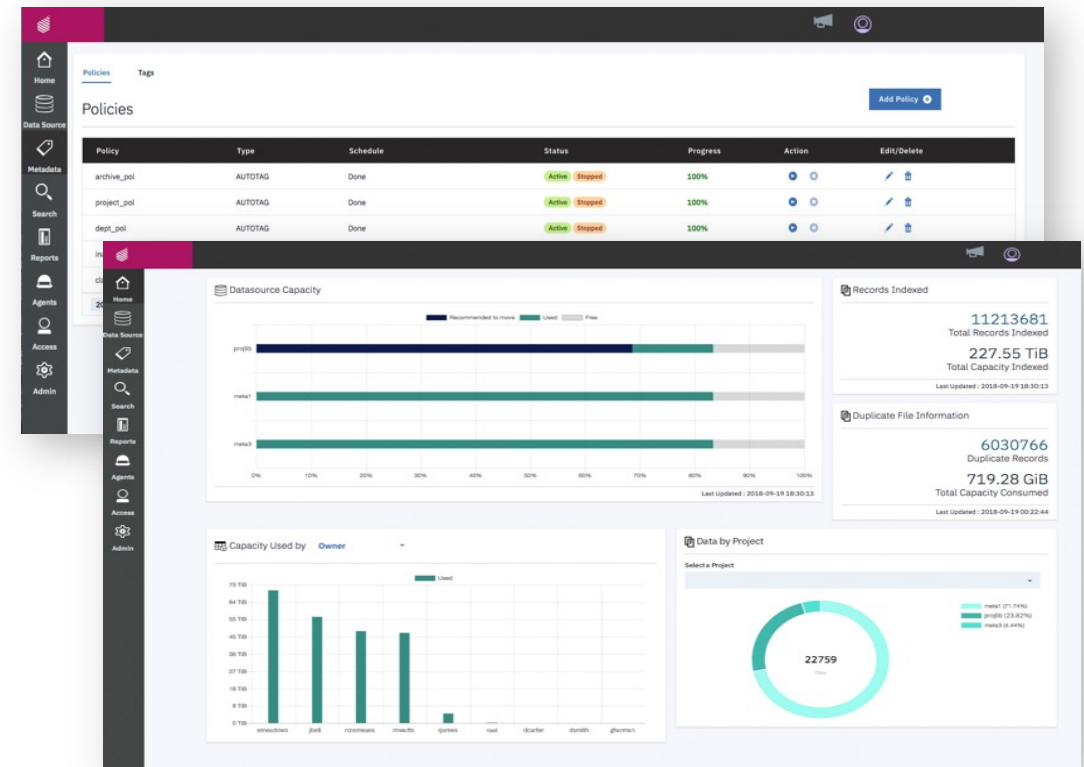- **Enrich and increase** the value of data

# IBM Spectrum Discover

- **Data insights:**
  Multi-vendor connections

- **Better AI:**
  Search billions of records < second

- **Optimize data workflows:**
  Policy based automation and auto tag data

- **Data security and compliance:**
  Discover security, compliance and governance
  issues before they become problematic

- **Business Value:**
  Link data to IBM Watson solutions and
  IBM Cloud Pak for Data

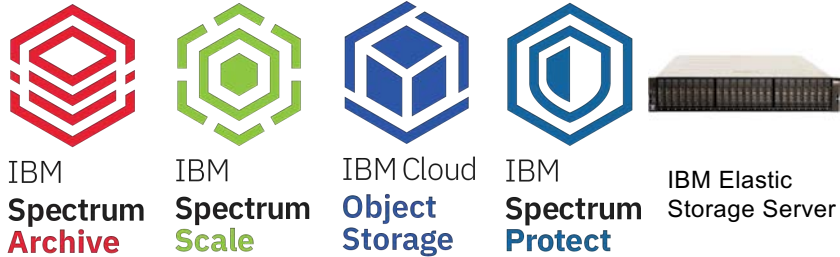Comprehensive search and AI analysis
with a data catalog and policy engine

IBM Spectrum
Discover

# IBM Spectrum Discover overview

## Where

### Backup, File & Object Storage



IBM **Spectrum Archive**

IBM **Spectrum Scale**

IBM Cloud **Object Storage**

IBM **Spectrum Protect**

IBM Elastic Storage Server

ISILON

amazon web services™ S3

NetApp™

ceph

## What

### Index & Tag Big Data



IBM **Spectrum Discover**

Search    Reporting    Dashboard

- Simple to deploy
  (VMware virtual appliance)

- Metadata curation

- Custom metadata tagging

- Automatic indexing

- Content inspection

- Policy-Engine

- Application plugin API / SDK

## Why

### High Speed Data Insight

**Large-Scale Analytics and AI/ML**
- Data discovery
- Dataset identification
- Data pipeline progression

**Data Governance**
- Data inspection and classification
- Data clean-up

**Data Optimization**
- Archive / tiering
- Duplicate data removal
- Trivial data removal

**Data Management**
- Automate Tags for custom insight
- Create reports or directly search data
- Search content for fast discovery

# 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
  - ROT Detection/Disposal
  - Etc...

- **Integrate with upstream information management applications**

# Easily create a policy to "tag" items based on a filter

Data Mapping

# Search inside files/objects to find patterns and create new metadata tags

Easily create your own custom search patterns

# Discover your data with simple interface or report generation

Generate reports

Drill down

Customize view

# Discover in one screen
# duplicate records and data for archive

**Data Visualization**



Summary of capacity

Summary of duplicate records

# Create a custom "action agent" to automate a workflow

## Add new policy

Inactive ⟷ Active

**Name**

some_name

**Policy Type**

DEEP-INSPECT ▾

**Collections**

Type search collection ▾

**Filter**

datasource = 'DiscoverVault' AND filetype = 'jpg'

**Agent**

Select a value ▾

+Add tag

**Schedule**

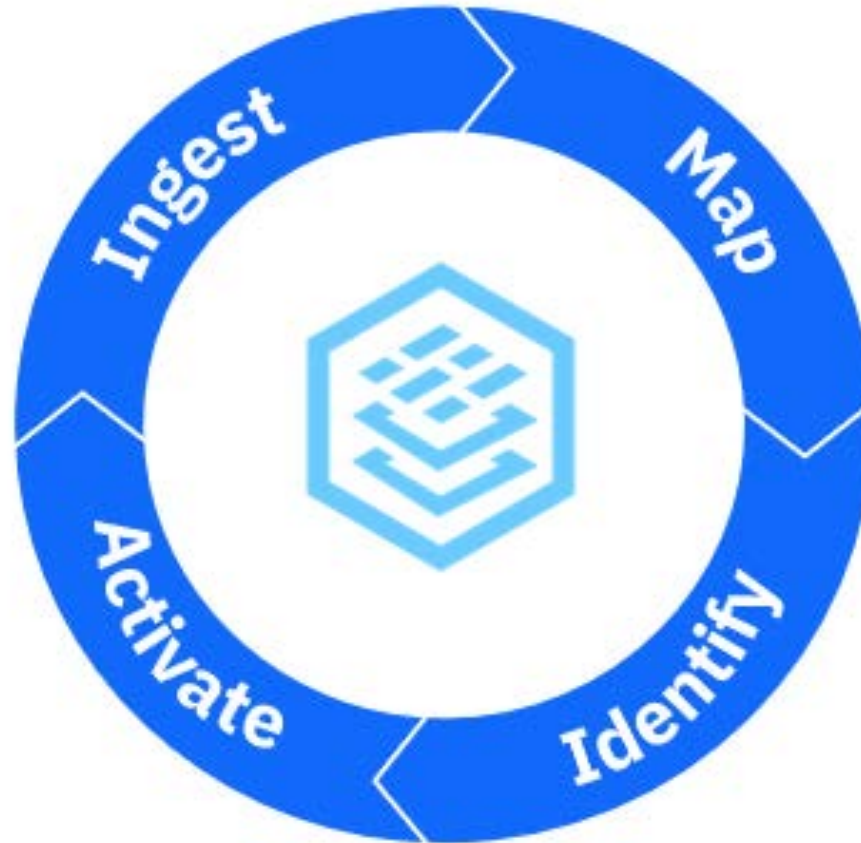◉ Now   ○ Daily   ○ Weekly   ○ Monthly

# Using Spectrum Discover for metadata-fueled data analysis

**Large Scale Data Ingest**

- Scan billions of records per day[1]
- Scale without limits
- Multi-source support

**Business-Oriented Data Mapping**

- Custom data tagging
- Content based inspection
- Automatic indexing

**Data Activation**

- Create automated actions
- Solution blueprints
- Live event notifications
- Policy-driven workflows

**Data Visualization**

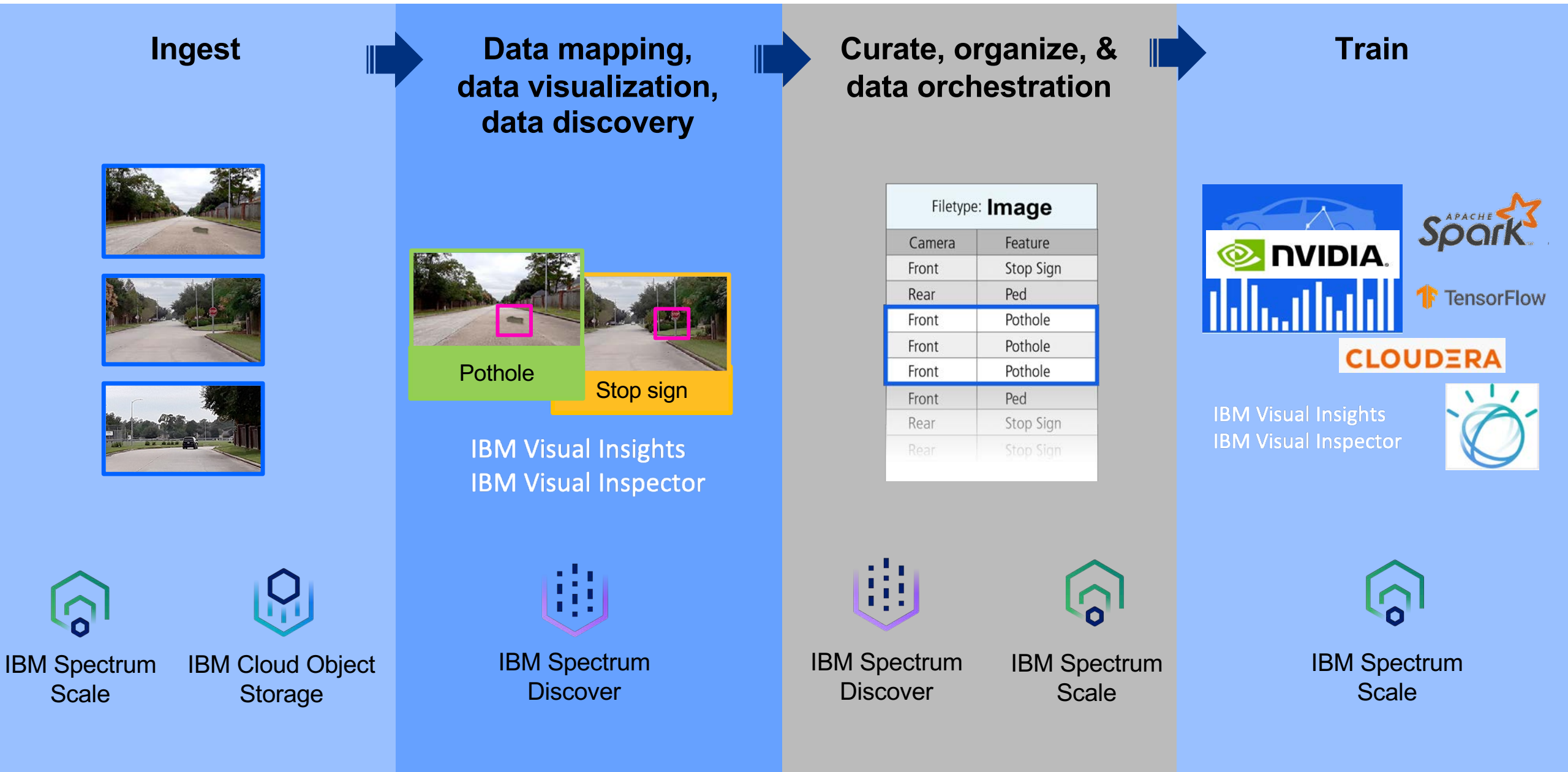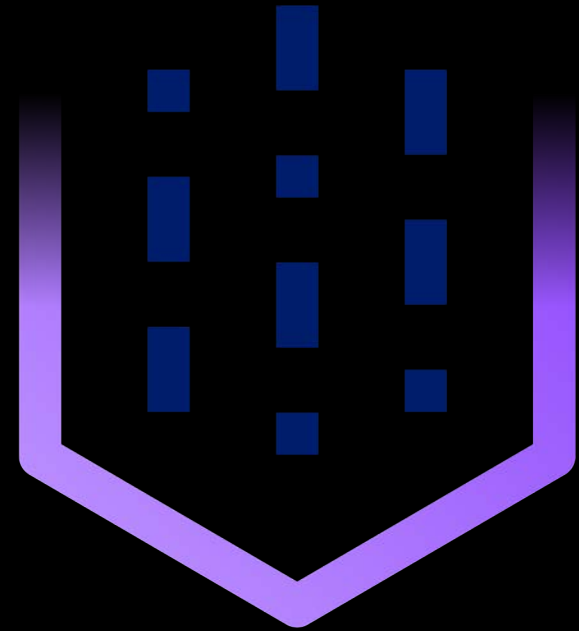- Query billions of records in seconds
- Multi-faceted SQL like search
- Drilldown dashboard
- Customizable reports

Ingest · Map · Identify · Activate

[1] Up to 30K/sec; ~2.5 billion rec/day w/ IBM Spectrum Scale; ~432 million/day w/ IBM Cloud Object Storage

# IBM Storage for AI workflows

**Ingest**



IBM Spectrum Scale

IBM Cloud Object Storage

**Data mapping, data visualization, data discovery**



Pothole

Stop sign

IBM Visual Insights
IBM Visual Inspector

IBM Spectrum Discover

**Curate, organize, & data orchestration**

| Filetype: **Image** | |
| --- | --- |
| Camera | Feature |
| Front | Stop Sign |
| Rear | Ped |
| Front | Pothole |
| Front | Pothole |
| Front | Pothole |
| Front | Ped |
| Rear | Stop Sign |
| Rear | Stop Sign |

IBM Spectrum Discover

IBM Spectrum Scale

**Train**

NVIDIA

APACHE Spark

TensorFlow

CLOUDERA

IBM Visual Insights
IBM Visual Inspector
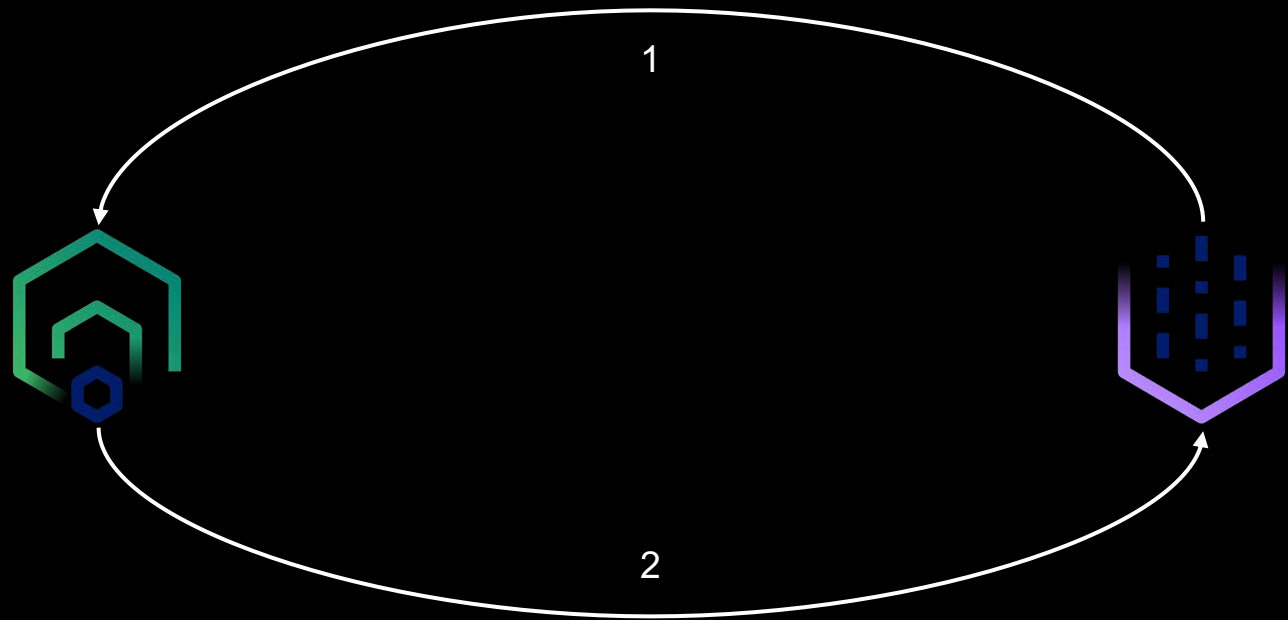
IBM Spectrum Scale

Spectrum Discover Integration in
Spectrum Scale and ESS

# Connecting Spectrum Scale and Spectrum Discover

1

2

1. Scan Files System / Single Fileset
   - Spectrum Discover scans actively

2. Live Events
   - Spectrum Scale sends Watchfolder events to Spectrum Discover

## Data source capacity

SpectrumScale

## Records indexed

### 159,380
Total records indexed

### 10.49 TiB
Total capacity indexed

# System metadata collected by Spectrum Discover

**IBM Spectrum Scale**

- Filesystem
- Site
- Platform
- Cluster
- Inode
- Owner
- Group
- uid
- gid
- Mode

- Fileset
- Path
- mtime
- atime
- ctime
- Pool
- Size
- migstatus
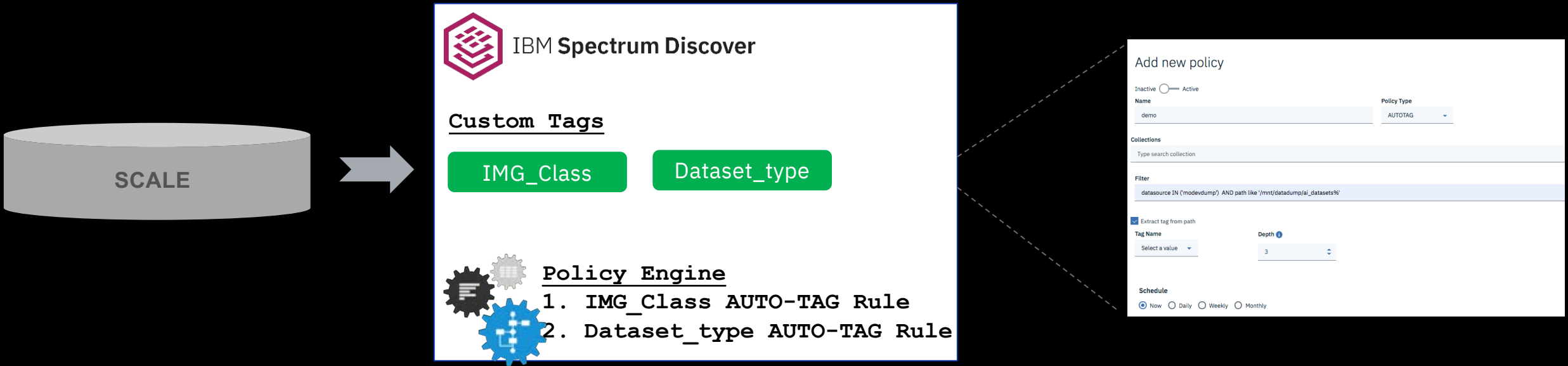- migloc

**IBM Cloud Object Storage**

- Operation
- Bucket Name
- Object Name
- Object Length
- Object etag

- Content Type
- Bucket UUID
- System UUID

# Leverage filesystem patterns to automatically create custom tags for dataset management

- Automatically parse file system sub-directories and insert as custom tags into Spectrum Discover
- Track and manage data across storage pools by custom tags

`/mnt/datadump/ai_datasets/POWER-AI-VISION/EPRI/Training Set Final/Conductor Damaged`

Platform  Dataset Name  Dataset Type  Image Category

**SCALE**

**IBM Spectrum Discover**

**Custom Tags**

IMG_Class    Dataset_type

**Policy Engine**
1. IMG_Class AUTO-TAG Rule
2. Dataset_type AUTO-TAG Rule

Add new policy

Inactive ⚪— Active

Name
demo

Policy Type
AUTOTAG

Collections
Type search collection

Filter
datasource IN ('modevdump') AND path like '/mnt/datadump/ai_datasets%'

☑ Extract tag from path

Tag Name
Select a value

Depth ⓘ
3

Schedule
◉ Now ⚪ Daily ⚪ Weekly ⚪ Monthly

# Spectrum Scale ILM Policies

Files written in file system are placed storage pools

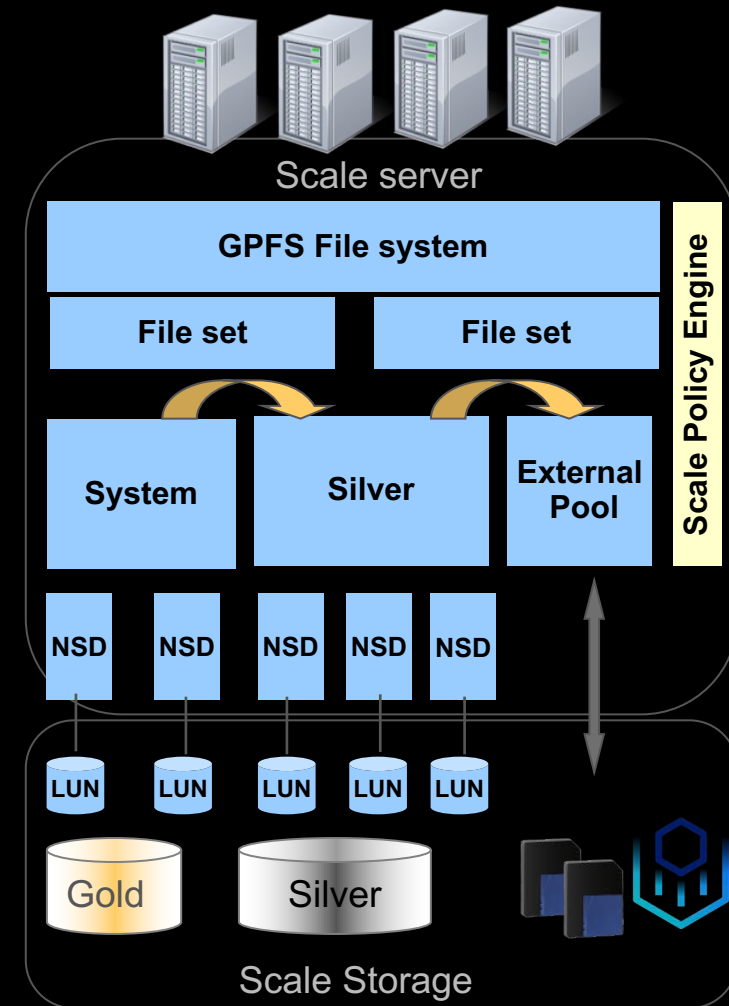– Storage poola are a collection of disks of the same type

Placement policy controls where files are placed

– Applied during file creation

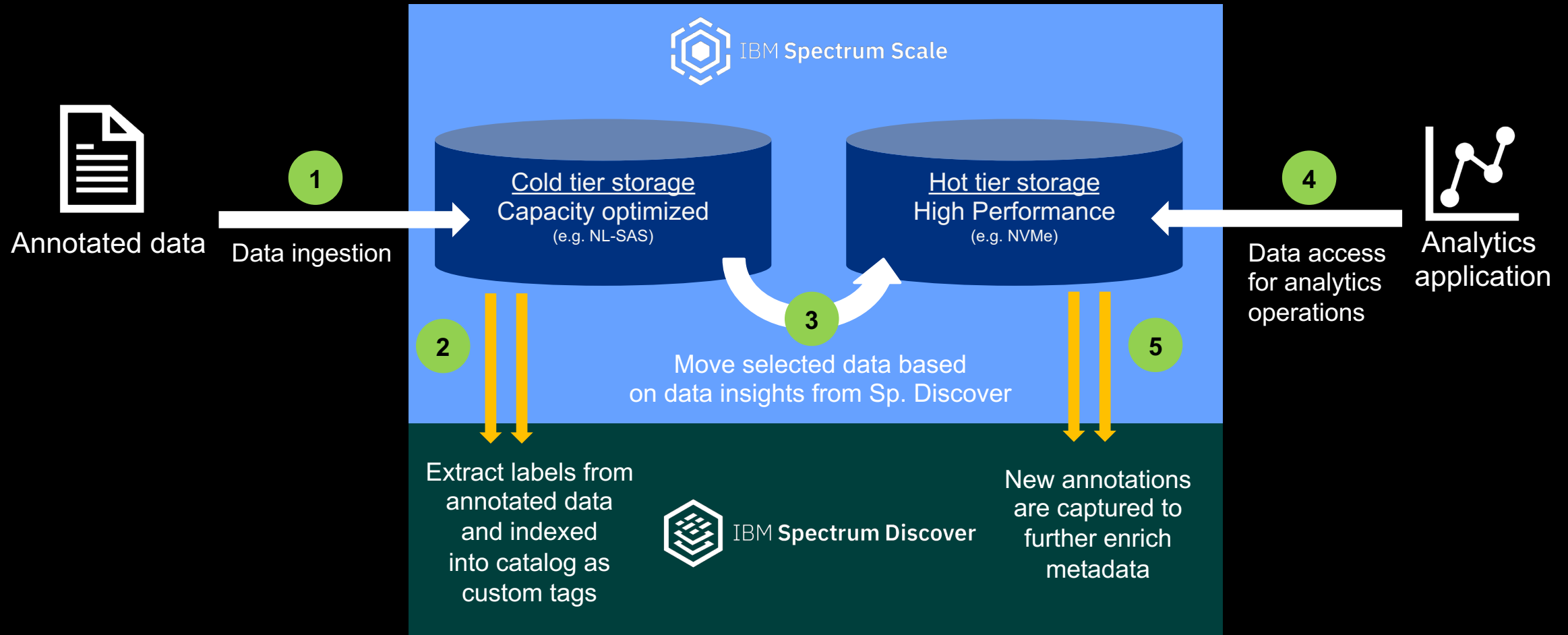Migration policies control transparent migration of files from one pool to another

– Applied during life cycle

– Migrated files can be transparently accessed

– Files can also be migrated to external pools such as

• Tape

• Object Storage

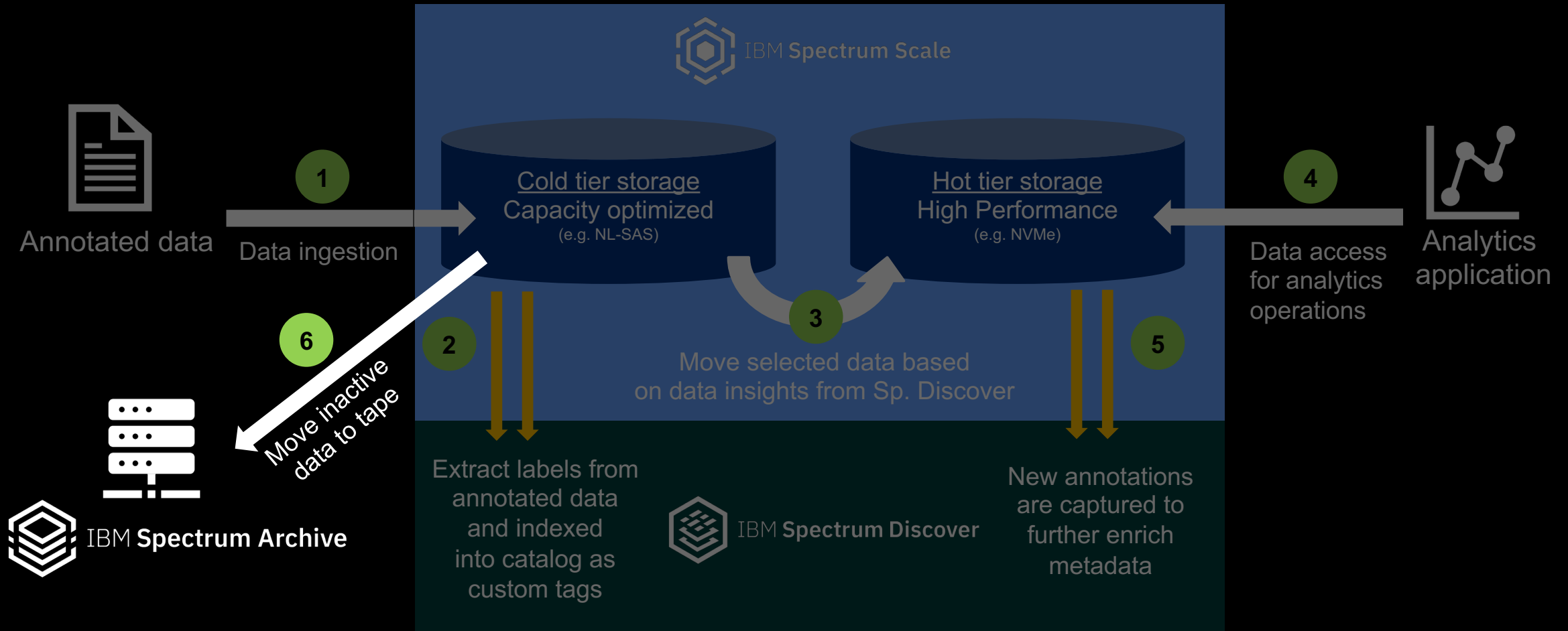Policy engine applies and executes policies based on metadata or manually.

# Data placement optimization for Spectrum Scale and ESS

- Leverage the insights from Spectrum Discover to drive better data management in Spectrum Scale
- Goes beyond the system metadata and extended attributes used by the current Spectrum Scale policy engine.

# Further data placement optimization for Spectrum Archive Enterprise Edition

- Move data to/from Spectrum Archive Tape
- A searchable catalog is maintained that includes the tape

# Trigger Scale Data tiering

Spectrum Discover triggers ILM Policies

- Tier data to different storage pool based on system metadata

**Management Policies**

# Add Policy

| ⊘ Define | ⦿ Configure | ○ Schedule | ○ Review |

Application

ScaleILM ⌄

Source connection ⓘ

SpectrumScale ⌄

Filter

tier IN ('system') AND filegroup IN ('image','compressed','video') AND atime > (NOW-60)

Destination tier ⓘ

bronze

# Trigger Scale Data tiering

Spectrum Discover
triggers ILM Policies

- Tier data to
  different storage
  pool based on
  system metadata

- Tier data to
  different storage
  pool based on
  custom metadata

**Management Policies**

# Add Policy

| ✓ Define | ● Configure | ○ Schedule |

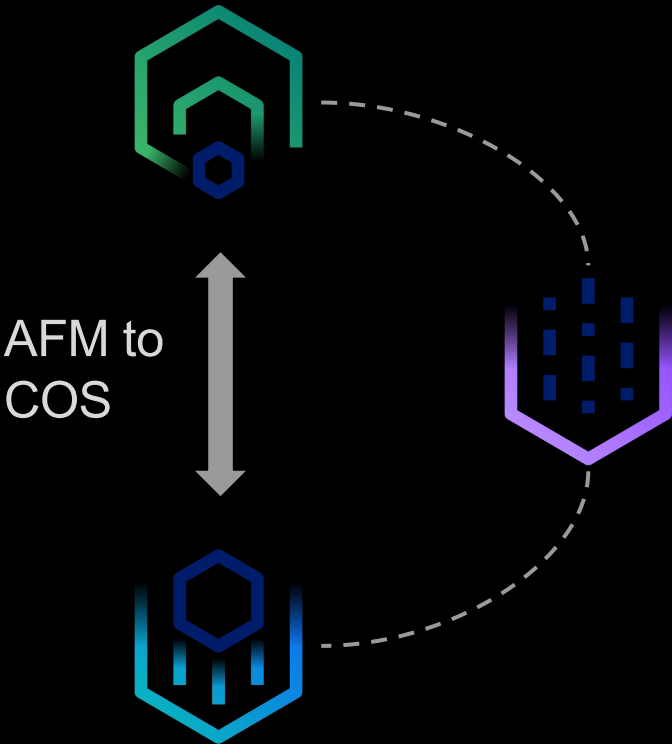Application

ScaleILM ⌄

Source connection ⓘ

SpectrumScale ⌄

Filter

project_status like "finished"

Destination tier ⓘ

archive:tapepool@lib1

# Trigger Scale Data tiering

Spectrum Discover triggers
ILM Policies

- Tier data to different
  storage pool based on
  system metadata

- Tier data to different
  storage pool based on
  custom metadata

- Recall data to hot tier
  based on custom
  metadata

# Discover and Scale AFM

AFM to COS

**Management Policies**
## Add Policy

Application

ScaleAFM ⌄

Filter

next_run like "True" AND
Datasource IN ('COS_DiscoverTest')
AND State IN ('migr')

Download objects from from cloud object store to spectrum scale ⓘ

Download objects from from Cloud Object Store to Spectrum Scale
This is used when, data is to be downloaded from the configured Cloud Object Store target vault to the Spectrum Scale AFM fileset.

Source connection ⓘ

COS_DiscoverTest ⌄

Destination connection ⓘ

SpectrumScale ⌄
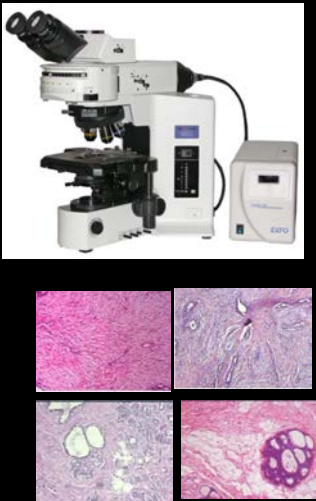
Spectrum scale afm fileset ⓘ

AFM_test

# Use Case: Tumor Classification

Event driven architecture to automatically classify and catalog biopsies of breast cancer tumors using PowerAI Vision inference model, Spectrum Discover, and Spectrum Scale / ESS

1. New imaging data ingested into Spectrum Scale / ESS

2. Storage sends Spectrum Discover system metadata events when new imaging data is ingested and Spectrum Discover builds catalog

3. Spectrum Discover policy automatically reads new imaging data from source storage, passes to the image classification model, captures results and indexes into Spectrum Discover

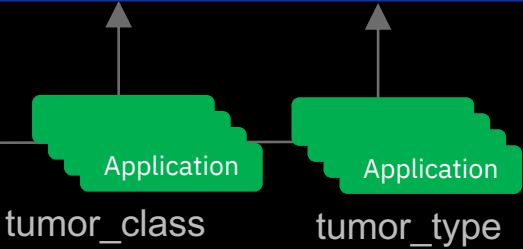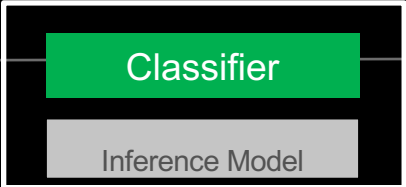Breast Cancer Tumor Biopsies

IBM ESS

2. System Metadata Events

| Tumor_class | tumor _type | score |
|---|---|---|
| malignant | ductal_carcinoma | 100% |
| benign | adenosis | 90% |

Policy Engine:

- Image Classifier

Connector Application Plugin

3. Image classification workflow

Classifier

Inference Model

Application

Application

tumor_class

tumor_type

# Free 90-day trial

Experience for yourself the game-changing insights possible with IBM Spectrum Discover
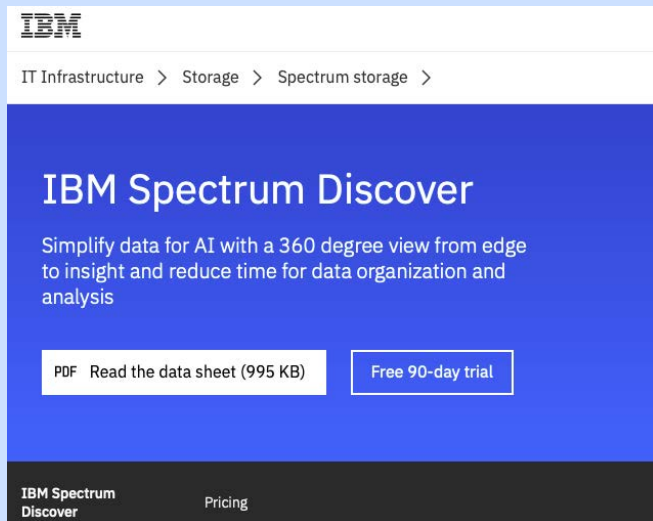
IBM Spectrum Discover
Free Trial
_____

Unleash metadata-fueled insights for your unstructured data -- free for 90 days.

→ Free trial

www.ibm.com/marketplace/spectrum-discover

# Learn more about Spectrum Discover



### Web Page and Resources

https://www.ibm.com/products/spectrum-discover

http://www.redbooks.ibm.com/abstracts/redp5603.html?Open

http://www.redbooks.ibm.com/redpapers/pdfs/redp5550.pdf

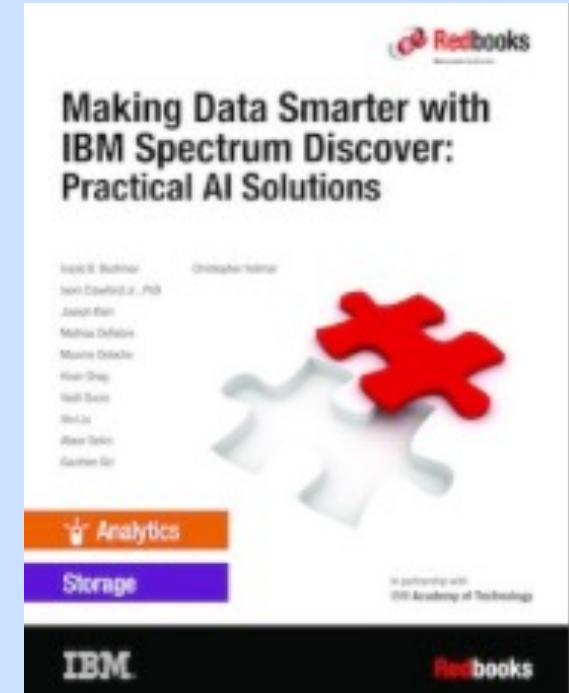http://www.redbooks.ibm.com/abstracts/sg248488.html?Open

# Thank you for using
# IBM Spectrum Scale with Spectrum Discover!