IBM **Spectrum Discover**

# Data Insight for Petabyte-Scale Unstructured Data Storage

**Indulis Bernsteins**
Systems & Storage Architect

CIUK, Manchester
December 12, 2018

IBM

# Please note

IBM's statements regarding its plans, directions, and intent are subject to change or withdrawal without notice and 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. Information about potential future products may not be incorporated into any contract.

The development, release, and timing of any future features or functionality described for our products remains at our sole discretion.

Performance is based on measurements and projections using standard IBM benchmarks in a controlled environment. The actual throughput or performance that any user will experience will vary depending upon many factors, including considerations such as the amount of multiprogramming in the user's job stream, the I/O configuration, the storage configuration, and the workload processed. Therefore, no assurance can be given that an individual user will achieve results similar to those stated here.

# Harnessing the Value of Data

"The world's most valuable resource is no longer oil, but **data**."

*The Economist, May, 2017*

*…how to **harness the value?***
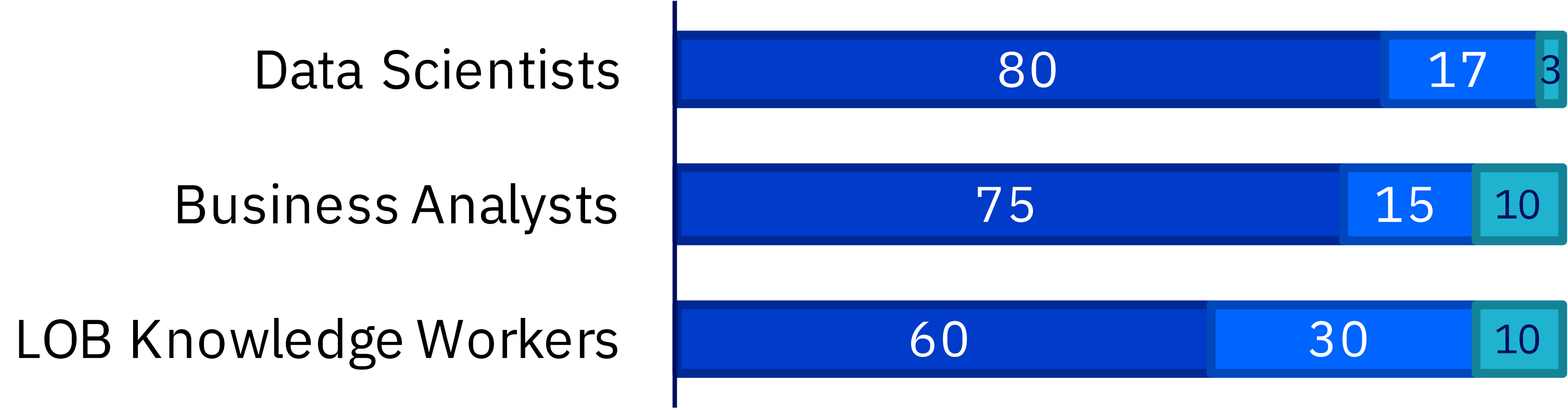
- **Identify**
- **Categorize**
- **Utilize**

IBM

# Unstructured Data Challenges

Enterprises **with**
**1,000 TB+**
**unstructured data**

**grew 3X**

**2016 ⇨ 2017**

**39%**

of firms see sourcing, gathering, managing & **governing data** as their biggest **challenges** when using systems of insight

IBM

# Unstructured Data Challenges

**Data Scientists** spend **2/3** **of their time finding data**

| | | | |
|---|---|---|---|
| Data Scientists | 80 | 17 | 3 |
| Business Analysts | 75 | 15 | 10 |
| LOB Knowledge Workers | 60 | 30 | 10 |

**Challenges for exabyte-scale data: how to…**

- Pinpoint & activate relevant data for AI, analytics etc.

- Get fine-grained visibility of value of data

- Remove redundant, trivial & obsolete ("ROT") data

- Identify & classify sensitive data

# Metadata is data about data

Provides **context** to **classify** & **manage** unstructured data

- Last time accessed
- Author/user
- Type of Data
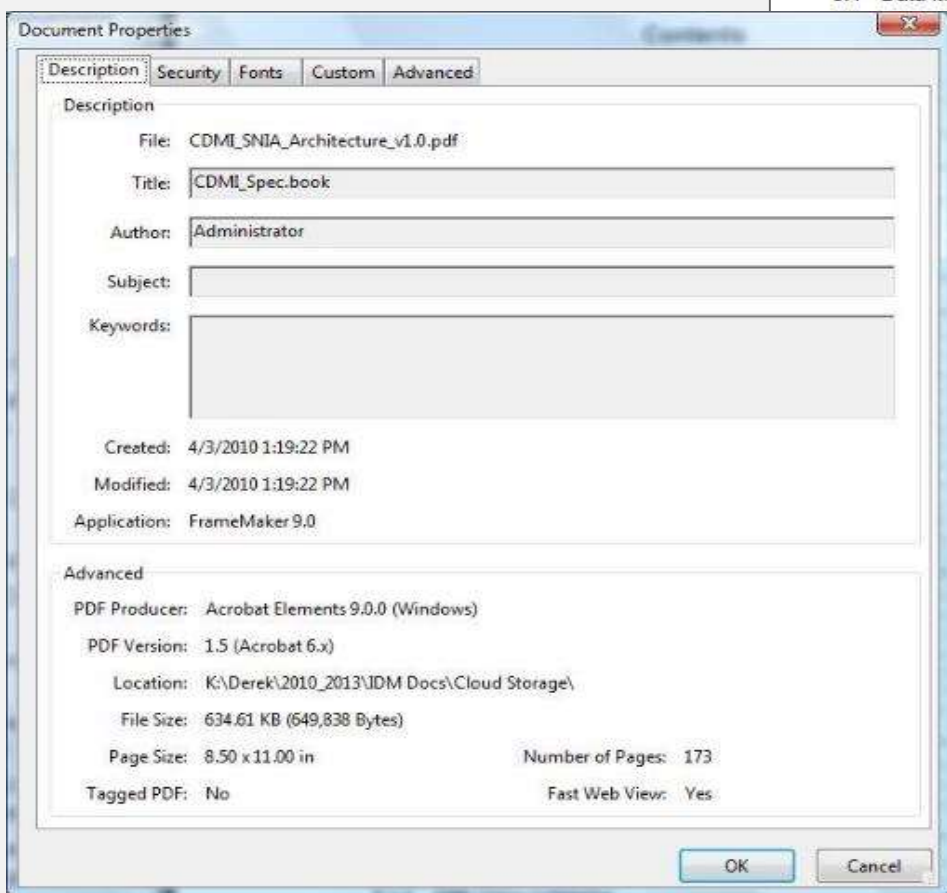- Key words
- Custom tags
  - e.g. project, department, etc.

**Files/Objects**

**PDF**

**System Metadata**

# Types of metadata

**System:** Typical file system metadata, date created, owner, last modified, size, file type, etc.

**Custom:** User-defined or based on unique organizational schema/taxonomy



**Files/Objects**

**PDF**

**System Metadata**

# 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

- **Install policy scan agent on node in target Spectrum Scale cluster**
- **No EAs/XATTRs**
- **No user MD from Objects**

## IBM Cloud Object Storage
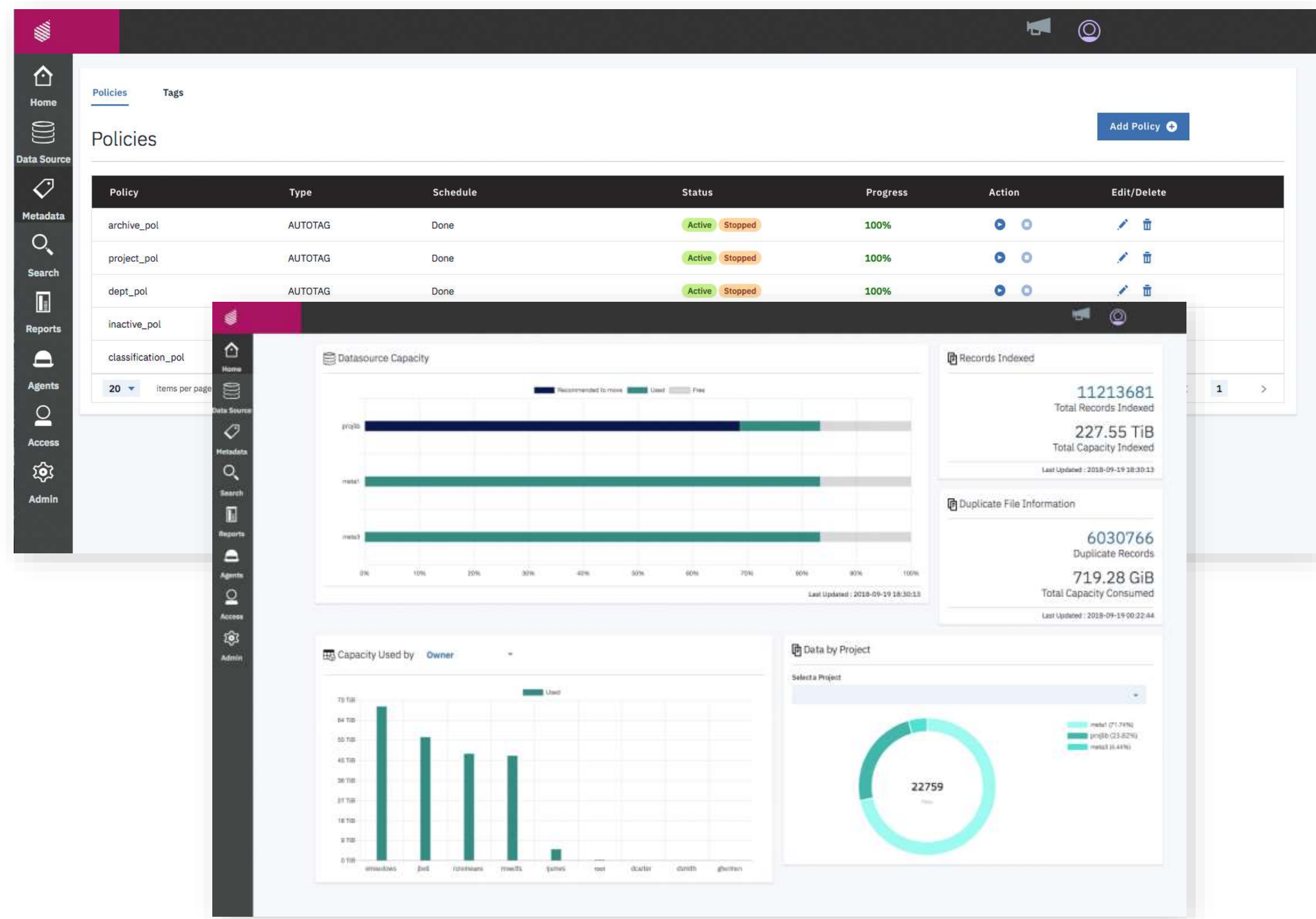
- Operation
- Bucket Name
- Object Name
- Object Length
- Object etag
- Content Type

- Bucket UUID
- System UUID

IBM **Spectrum Discover**

# Data Insight for Analytics, Governance & Optimization

- Automated cataloging

- Comprehensive insights: combine system metadata with custom tags

- Extend usefulness using the API, custom tags, content inspection and policy-based workflows

# Spectrum Discover **is**

- A scalable & simple to use virtualised s/w "appliance"
  - Capable of dealing with System & Custom metadata
  - Works with Billions of files & objects
  - Designed to be easy to use & support admins **and** end users

- ...Uses System and Custom Metadata to primarily support **storage management**
  - Extensible: add MD tags from "deep inspections"
  - For file and object curation: movement, deletion, migration
  - For data selection for processing: ad-hoc selection

IBM **Spectrum Discover**

# Spectrum Discover **is not**

- Not a replacement for all Metadata engines

  - e.g. DICOM scanner to populate hospital patient record system

- Not a replacement for all ILM/HSM features in Spectrum Scale

  - There is some overlap, as Data Migration is planned (SoD)

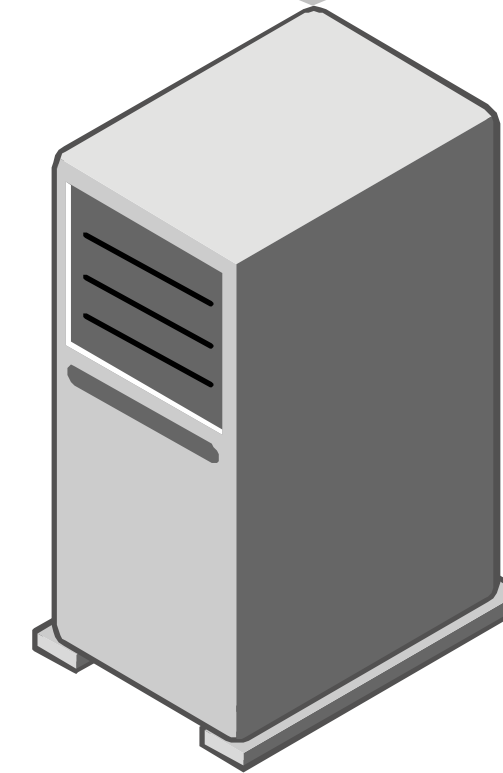    - Use cases for Spectrum Discover are ***probably*** more ad-hoc, user driven

IBM **Spectrum Discover**

# Single Node Spectrum Discover Deployment Example (2 billion files)

**Spectrum Discover**

128GB RAM
8 Physical CPU
5TB VMDK primary
2TB VMDK backup
1/10GbE Virtual Network

**VMware ESXi 6.5**

**Local attached
SSD or shared
storage**

# Multi-Node Spectrum Discover Deployment Example (10 billion files)

VMware ESXi 6.5

VMware ESXi 6.5

VMware ESXi 6.5

Server 1

Server 2

Server 3

**Min 2 x 8/16GB FC per Node**

**Optional FC Switch**

**Shared storage (Flash/SSD/NVMe) e.g. IBM Flashsystem**

**Spectrum Discover**

256GB RAM
24 Physical CPU
15TB VMDK primary
7TB VMDK backup
10GbE Virtual Network

# Use Cases

# Optimization Use Case:
## *Optimize Data Placement on Storage*

**IBM Spectrum Discover**

| Business Leaders | IT Infrastructure Team | **Problem** | **Solution** |
|---|---|---|---|
| **Business Leaders**<br>(Dir / VP of IT / CIO / CFO)<br><br>Cost<br><br>• Unnecessary purchases due to inefficient management of storage<br>• Avoid unplanned Operational Expense due to business processes not being followed | **IT Infrastructure Team**<br>(Storage Administrator) | **Managing Capacity Utilization**<br><br>• Monitor multiple vendor technologies<br>• Wait for the signs of quota threshold violations<br>• Track down Data Owners & ask them to audit their data placement, retention, & take corrective action<br>• Wait for reply from Data Owner regarding what data can be archived | **Shared Management Responsibility**<br><br>• Visualize/monitor utilization of multiple storage sources with drill-down analytics.<br>• Custom, policy-driven data tagging identifies candidate files for action.<br>• Intelligently manage data placement on appropriate storage tier (archive) |
| Risk<br><br>• Ensuring new Business Strategies are not delayed by lack of capacity<br>• Ensuring existing applications have capacity for expansion | **Data Owners**<br>(Researcher / Data Scientist / Lab Manager / RIM) | **Manually Managing the Data**<br><br>• Track data ownership in a spreadsheet<br>• Relies on peers following best-practices in order to track data<br>• Tries to figure out what data can be archived in order to save on storage costs<br>• Open ticket to have data archived & migrated | **Shared Management Responsibility**<br><br>• Custom, policy-driven data tagging identifies candidate files for action<br>• Multi-faceted search to identify other candidates for archiving.<br>• Improved understanding of data use profile |

# Product Features
# & Architecture

# IBM Spectrum Discover Overview

## File and Object Storage

IBM Cloud **Object Storage**

IBM **Spectrum Scale**

**OTHERS**

**Planned for 2019**

Scanning and
Event Notifications

## Data Insight

**IBM Spectrum Discover**

Search    Reporting    Dashboard

- Simple to deploy
  (VMware virtual appliance)
- Metadata curation
- Custom metadata tagging
- Automatic indexing
- Policy-Engine
- Action Agent API

Use
Cases

## Data Activation/Optimization

**Large-Scale Analytics**

- **Data discovery**
- Dataset identification
- Data pipeline progression

**Risk Mitigation**

- Data inspection
- Data classification
- Data clean-up

**Data Optimization**

- Archive / tiering
- Duplicate data removal
- Trivial data removal

# Extensible Foundation for Data Insight

**IBM Spectrum Discover**

| | |
|---|---|
| **Connectors** | • Scanner for IBM Cloud Object Storage (COS) and IBM Spectrum Scale<br>• Notifications for IBM Cloud Object Storage<br>• Tech preview of Notifications (Live Events) for IBM Spectrum Scale |
| **Platform** | • **Support for Single and Multi-Node (3-Node HA Cluster) (x86 only)**  • **Support for Role Based Access Control (RBAC)**  • Audit Trail – track and log user actions on Spectrum Discover<br>• Code upgrade for Single and Multi-node  • Remote Support – tools to collect logs and upload to IBM support  • Dashboards to monitor Spectrum Discover health<br>• **Encryption of Metadata Database and Notification Logs (Kafka)**  • **Backup / restore / DR of Metadata Database** |
| **Action Agent SDK Ecosystem** | • SDK to help extend the platform capabilities to perform custom actions around data – data migration, archiving, content-based search & tagging, etc. |
| **GUI** | • Basic Search<br>• **Advanced drill-down search**  • Support for Role Based Access Control (RBAC)<br>• **Dashboard** to visualize storage consumption on a wide range of system and custom metadata<br>• **Create and schedule policies** |
| **Scalability** | • Up to 100 billion indexed documents |
| **Performance** | • Ingest up over 1 Billion records per day |
| **Quality** | • Net Promoter Score (NPS) widget instrumented into Discover to gather NPS scores. |

**Live events\* and scans**
**With System Metadata**

IBM Cloud Object Storage

GET, PUT, POST, DELETE

OPEN, CLOSE, DELETE, RENAME

Spectrum Scale

**Custom Agents**

Apache Tika

**Header Extraction**
(example)

kafka

# IBM Spectrum Discover

**Search Visualize**

Search    Reporting    Dashboard

**Event Consumers**

**Database**

**Policy Engine Act on Event metadata**

**Data Movement\***
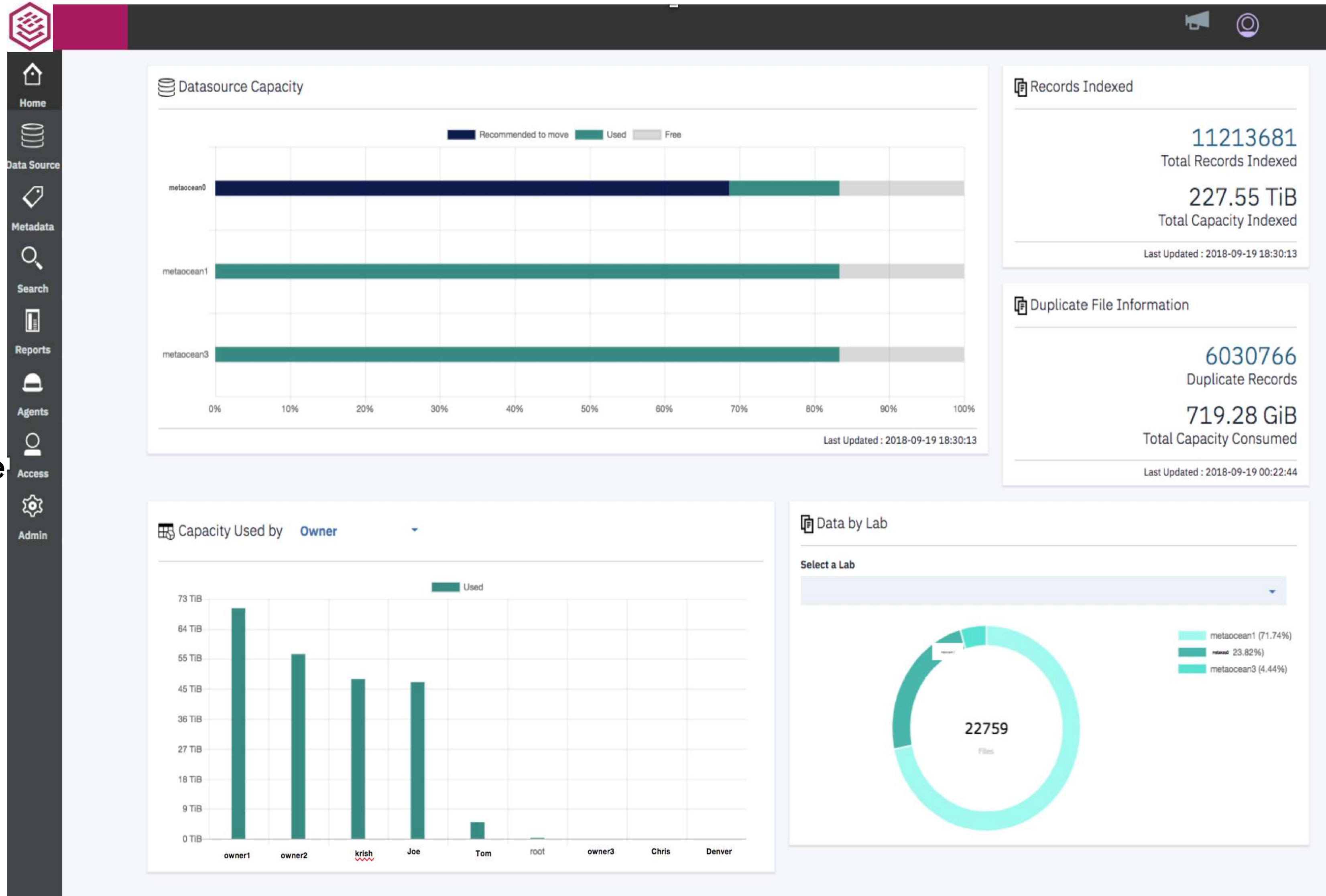
*\*Planned for 2019*

# Spectrum Discover Admin Dashboard

**Monitor storage utilization and data recommendations (Move/Archive)**

**Preview capacity use by data facet**
- **Classification**
  - **Owner**
    - **File Type**
      - **Etc.**



**Total indexed data and capacity**

**Duplicate file or object candidates**
- **Number**
- **Capacity used**

**Data capacity by group/collection**
- **Customer defined**
- **Lab/Project/etc.**

# Data Curation for AI Workloads

# Spectrum Discover for Data Preparation & Data Curation

NEW DATA

## DATA SOURCE

Traditional Business Data

Sensor Data

Data from collaboration partners

Data from mobile app & social media

Legacy Data

**Years of Data**

## DATA PREPARATION

IBM **Spectrum Discover**

Heavy IO

Pre-Processing

Training Dataset

Testing Dataset

**Weeks & months**

## MODEL TRAINING

Parallel Hyper-Parameter Search & Optimization

Monitor & Advise

Network Models

Hyper-Parameters

Iterate

Instrumentation

AI Deep Learning Frameworks *(Tensorflow & IBM Caffe)*

Distributed & Elastic Deep Learning *(Fabric)*

**Days & weeks**

## INFERENCE

Deploy in Production using Trained Model

Trained Model

**Seconds to results**

# ADAS-AD Data Pipeline

INGEST

CLASSIFY

ANALYZE / TRAIN

INSIGHTS

# Autonomous Vehicle Use Case and Demo with Spectrum Discover

| INGEST | SORT & EXTRACT | CURATE | TRAIN |
|---|---|---|---|

**EDGE**



Pothole

Stop sign

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

**IBM Cloud Object Storage**

**IBM Spectrum Discover**

**Data Scientist**

**Data Scientist**

- Global ingest of IOT data from vehicles
- Geo-dispersed COS

- Ingests & indexes system metadata via Action Agent SDK
- Extracts labels from images
- Adds as custom tags

- Searches for images with labeled as having 'Pothole' feature

- Trains a model

**IBM Spectrum Discover**
*via Action Agent SDK*
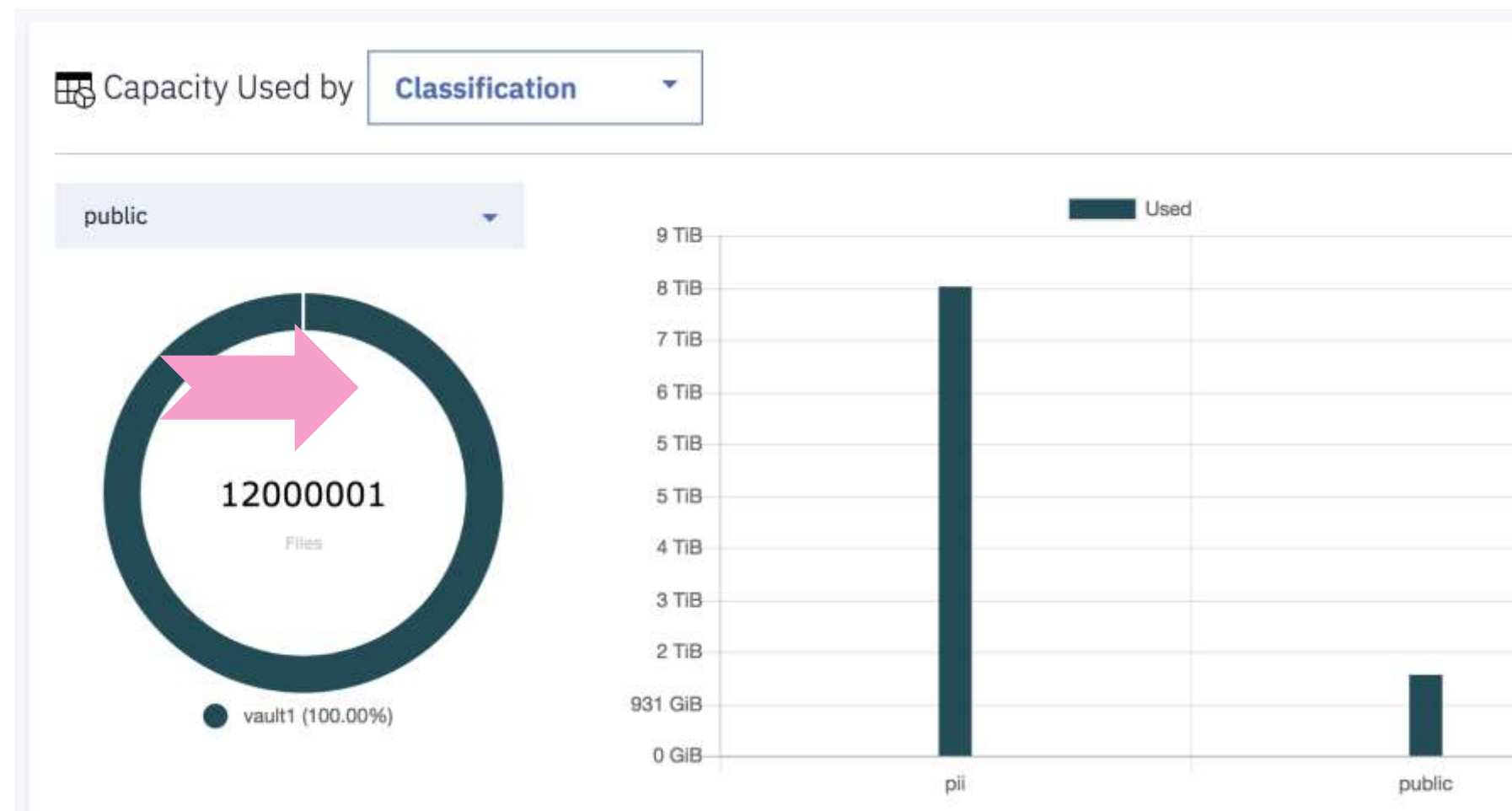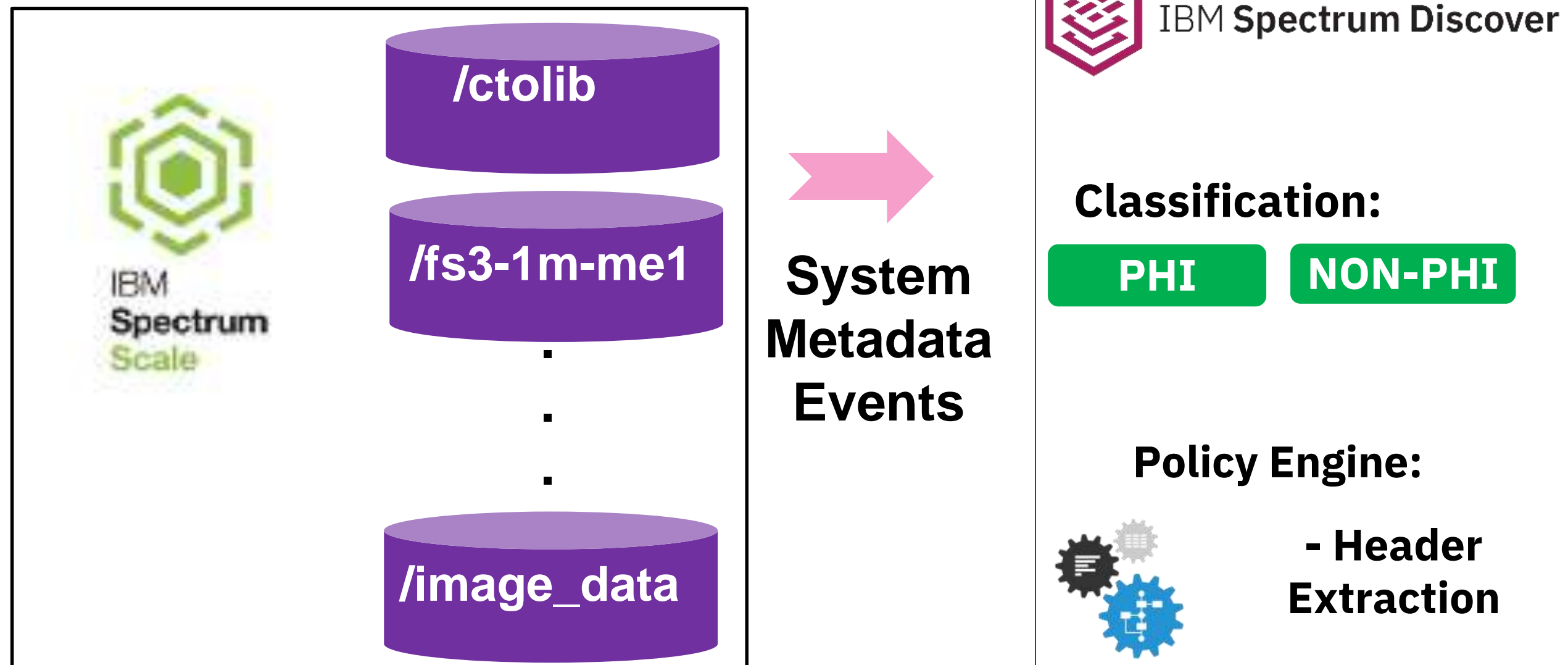
- Enriches data catalog with new tags derived from analysis

# Healthcare / Life Sciences Integration

# Use Case: PHI and non-PHI Classification of DICOM Images



Show me where my PHI data resides

IBM Spectrum Discover

**Classification:**

PHI    NON-PHI

**Policy Engine:**

- Header Extraction

IBM Spectrum Scale

/ctolib

/fs3-1m-me1

.
.
.

/image_data

System Metadata Events

**Header Extraction (pydicom)**

DICOM

PatientName: Smith

Capacity Used by  Classification

public

12000001
Files

vault1 (100.00%)

Used

9 TiB
8 TiB
7 TiB
6 TiB
5 TiB
5 TiB
4 TiB
3 TiB
2 TiB
931 GiB
0 GiB

pii    public

# Spectrum Discover Trial

# Free Trial Software Download

**IBM Spectrum Discover**

## 90 Day Free Trial

- At end of 90 days, code no longer accessible by client w/o approved extension or purchase of full license

## Full Function Version of Code

- Not limited scale or function set
- At termination of trial, access terminates

## Restriction(s)

- Cannot upgrade from single node trial to multi-node production

Support for trial: spdiscov@us.ibm.com

https://www.ibm.com/us-en/marketplace/spectrum-discover

# THANK YOU!

IBM Global Financing offerings are provided through IBM subsidiaries and divisions worldwide to qualified commercial and government clients. IBM Global Financing lease and financing offerings are provided in the United States through IBM Credit LLC. Rates and availability are based on a client's credit rating, financing terms, offering type, equipment and product type and options, and may vary by country. Non-hardware items must be one-time, non-recurring charges and are financed by means of loans. Other restrictions may apply. Rates and offerings are subject to change, extension or withdrawal without notice and may not be available in all countries. IBM and IBM Global Financing do not, nor intend to, offer or provide accounting, tax or legal advice to clients. Clients should consult with their own financial, tax and legal advisors. Any tax or accounting treatment decisions made by or on behalf of the client are the sole responsibility of the client. For IBM Credit LLC in California: Loans made or arranged pursuant to a California Financing Law license.

For more information, visit: ibm.com/financing