# Best practices for data collection - a critical foundation for obtaining valuable data insights
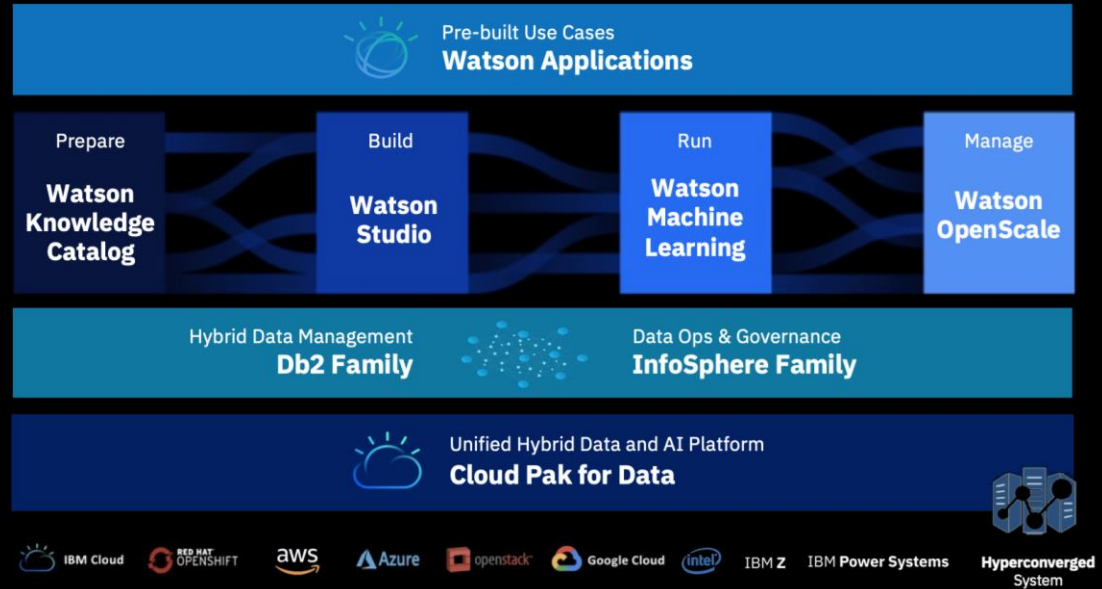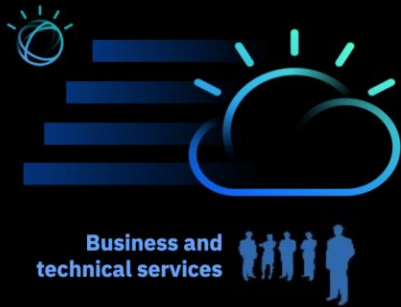
Jason Mathew
Offering Manager, Watson Knowledge Catalog

IBM Data and AI


David Wohlford
Product Marketing Manager, Spectrum Discover

IBM Systems

# Metadata unlocks data by making it visible and understandable.

# Watson™ Knowledge Catalog provides metadata management for the IBM Data and AI Portfolio

# Benefits of metadata management

## Regulatory compliance

Metadata management conducted on a unified platform that provides stewardship, data lineage, and impact analysis services is the best assurance that an organization can validate and demonstrate that the data reported is true.

## Productivity and discovery

Data is abundant. Much of it comes from existing systems and data stores for which no documentation exists, or the documentation does not reflect the changes and updates of those systems and data stores.

## Mitigating risk

Metadata management provides the measure of trust that businesses need. Through data lineage and impact analysis, businesses can know the accuracy, completeness and currency of the data used in their planning or decision-making models.

# DataOps Impact

Data Inventory Example

**85%**
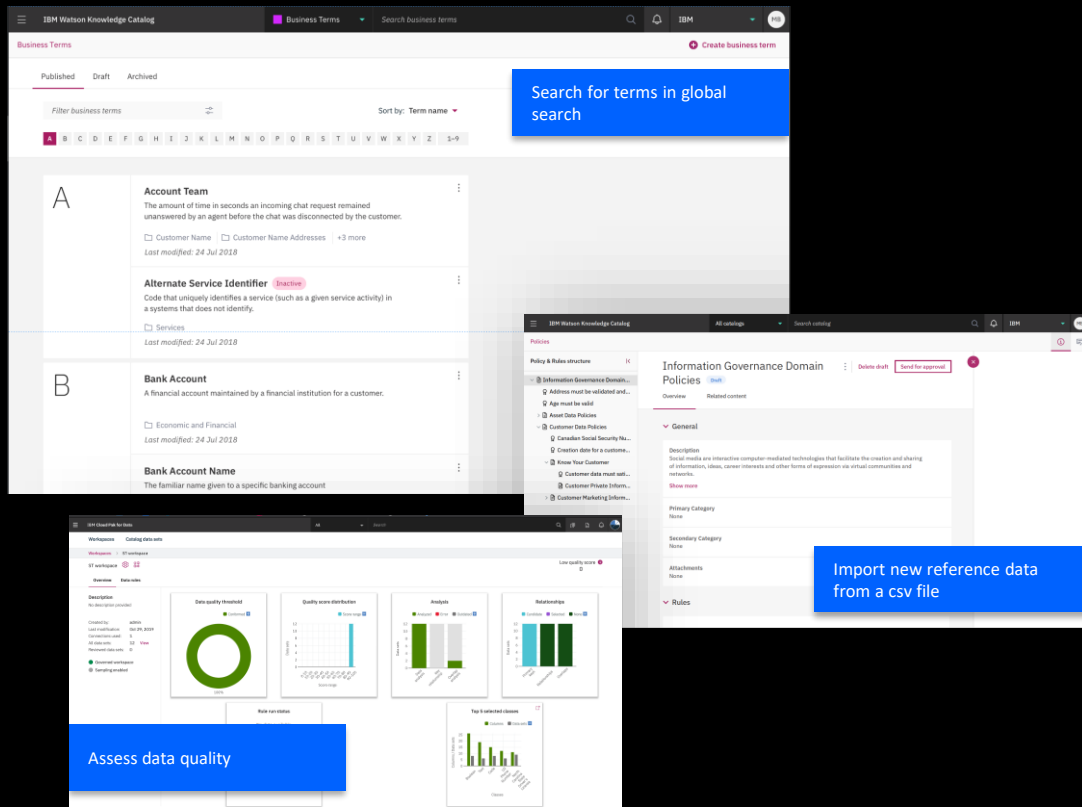
Reduction in business glossary creation time

**90%**

Reduction in time to discover metadata and assign terms

**200,000**

Number of technical assets across multiple clouds discovered in less than 5 mins

# Watson™ Knowledge Catalog

Open and intelligent data catalog for data and AI model governance, quality, and collaboration.



Search for terms in global search

Import new reference data from a csv file

Assess data quality

## WKC Use Cases

- Metadata Management

- Regulatory compliance

- Improve Data Quality

- Enterprise Data & AI Governance

- Data monetization and self-service

- Enterprise Data Cataloging

- Data lake consumption

- Reference Data Management

- DataOps Maturity

An unstructed data catalog and policy engine to organize the AI infrastructure and
help solve the data and AI puzzle faster

AI
Workflows

Data
Curation

Data Analysis

IBM Spectrum Discover

# IBM Spectrum Discover

Catalog your data and search billions of files/objects in 0.5 sec
Manage AI workflows, data security analysis and data governance

Heterogeneous data organization


IBM Spectrum Scale


IBM Cloud Object Storage

Ingest data continuously
real-time for live data updates

IBM
Spectrum
Discover

Netapp Dell/EMC   AWS
Windows   COS FA   NFS

Multi-vendor scan from edge,
core and public cloud

Policy engine and data catalog

Analyze and identify data anywhere

One click integration with IBM Watson solutions

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

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

Summary of capacity

Summary of duplicate records

# A simple one click export to IBM Watson Knowledge Catalog

- Automatically register assets with Watson Knowledge Catalog

- Leverage assets in IBM Cloud Pak for Data

- Import custom tags and create new and expanded insights from data

# Smarter data with integrated IBM solutions



Video Files

EXIF          VCF

Leverage more data and classify data in real time

Collect

Index          Organize          Analyze

Real-time ingest → IBM Spectrum Discover → Register Assets

IBM Spectrum Scale and ESS

IBM Cloud Object Storage

Extract content from source data

IBM Watson Knowledge Catalog

Build — Watson Studio
Run — Watson Machine Learning
Manage — Watson OpenScale

**IBM Cloud Pak™ for Data**

Red Hat OpenShift

# Learn more

Watson Knowledge Catalog at
ibm.com/watson-knowledge/catalog

IBM Spectrum Discover at
https://www.ibm.com/products/spectrum-discover

# Introducing DataOps

"DataOps is a collaborative data management practice focused on improving the communication, integration and automation of data flows between data managers and data consumers across an organization."

— Gartner

# DataOps Methodology Begins with Automating Metadata Management Best Practices

## Business objectives



Inventory and categorize data

Publish data and use

Deliver quality and governance

## DataOps Methodology

- Align data pipelines with business objective and success criteria.

- Automatically measures accuracy and speed of data capture, quality and use.

- Automates data and metadata ingestion and classification.

- Automatically assesses data quality issues and alerts when anomalies are detected.

- Automatically initiates remediation via workflow.

- Automatically ensures authorized use of published data assets by enforcing data privacy and governance policies.

# DataOps capabilities

**Data Sources**

Systems of record

IoT

Systems of insights

Cloud

Hadoop

Social media

Unstructured

Other external

Logs

Data access services

Automated data integration, and replication services

Automated data curation services

**Metadata management services**

Open metadata integration (Egeria)

Self-service interaction

Industry knowledge

Automated data governance, data quality and entity services

Governed data access services (virtual)

(virtual)

**Users**

Chief data officer

Governance officers

Data quality analyst

Data steward

Data scientist

Business users

Data engineer

Application developer

Application tester

DataOps ToolChain

On-prem    Private    IBM **Cloud**    Microsoft Azure    openstack.    aws    Google Cloud

# Thank you

# Legal notices

# Information and trademarks