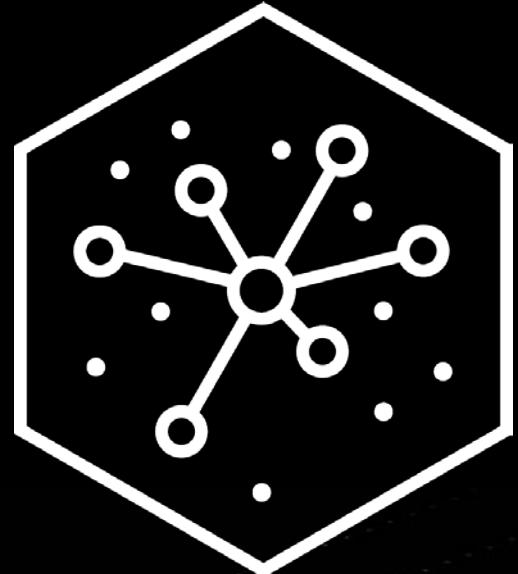


IBM Journey to Cloud and AI Analytics Modernization Workshop

Featuring: Cloud Pak for Data 3.0.1



Welcome to the IBM® Briefing Center

Location logistics

- ✓ Access restrictions
- ✓ Restrooms
- ✓ Emergency exits
- ✓ Smoking policy
- ✓ Breakfast / Lunch / Snacks
- ✓ Special meal requirements

Introductions

- ✓ IBM Speakers
- ✓ IBM Proctors
- ✓ IBM Sales Reps
- ✓ Attendees (optional)

IBM Analytics Modernization Workshop

Agenda

Part 1	<ul style="list-style-type: none">• Introduction• Business Use Case	<ul style="list-style-type: none">• Lab 01• Lab 02
Part 2	<ul style="list-style-type: none">• Collect: Connect• Organize• Collect: Virtualize	<ul style="list-style-type: none">• Lab 03• Lab 04• Lab 05
Part 3	<ul style="list-style-type: none">• Analyze• Deploy• Infuse – OpenScale• Infuse – Cognos Analytics• Wrap up	<ul style="list-style-type: none">• Lab 06• Lab 07• Lab 08• Lab 09• Lab 10

IBM Analytics Modernization Workshop

The workshop Unified Desktop

The screenshot shows a desktop interface with a dark theme. At the top, there's a toolbar with icons for Applications, Maces, and system controls. Below the toolbar, the title "IBM Journey to Cloud and AI Analytics Modernization" and subtitle "Workshop v9.0.1" are displayed. A large text "Featuring: Cloud Pak for Data" is centered. On the left, there's a vertical dock with icons for "IBM Cloud Pak for Data", "OPENSHIFT OpenShift Web Console", "Lab Solutions", and "IBM Trash". The main workspace contains several desktop icons: "Cluster Configuration.png", "Terminal", "Home", and a "Cloud Pak for Data" icon. A callout box titled "Installed Software" lists the following packages:

- Cloud Pak for Data v3.0.1 Enterprise
- OpenShift (RHOCP) v3.11.219
- Kubernetes v1.11
- Red Hat Enterprise Linux Server Release 7.8 (Maipo)
- Netezza Performance Server 11.0.3.1

Another callout box titled "Installed IBM Cloud Pak for Data Services" lists:

- Watson Knowledge Catalog
- Decision Optimization
- Watson Studio Local
- Db2 Advanced Edition
- Watson Machine Learning
- MongoDB
- Watson OpenScale
- Cognos Analytics
- Data Virtualization
- Cognos Dashboard Embedded
- DataStage Edition
- SPSS Modeler
- Metrics Server

Workshop lab exercise tips

- Use full screen with a large screen computer
- Use Chrome or Firefox browsers
- Use [Ctrl][+] and [Ctrl][-] or [Ctrl][Mouse-Scroll-Wheel] to zoom
- Use a mouse (keyboard alone is difficult)

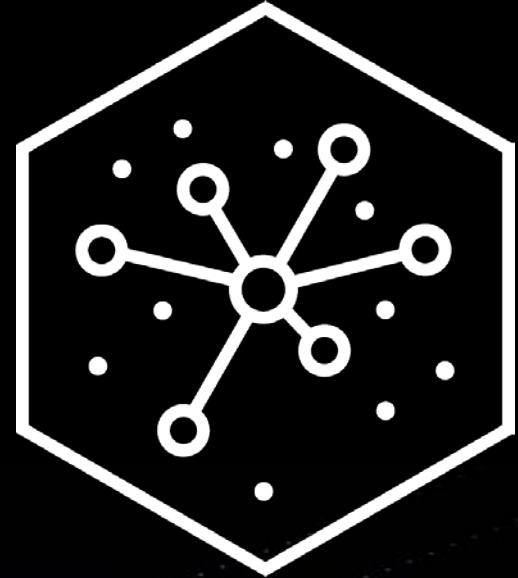
IBM Analytics Modernization Workshop

Part 1

<ul style="list-style-type: none">• Introduction• Business Use Case	<ul style="list-style-type: none">• Lab 01• Lab 02
<ul style="list-style-type: none">• Collect: Connect• Organize• Collect: Virtualize	<ul style="list-style-type: none">• Lab 03• Lab 04• Lab 05
<ul style="list-style-type: none">• Analyze• Deploy• Infuse – OpenScale• Infuse – Cognos Analytics• Wrap up	<ul style="list-style-type: none">• Lab 06• Lab 07• Lab 08• Lab 09• Lab 10

Introduction

Lab 01 – Getting Started



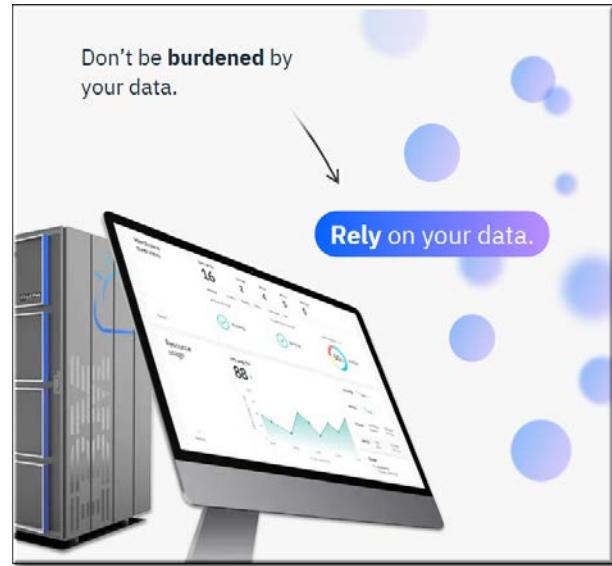
IBM Cloud Pak for Data



IBM Cloud Pak for Data is a single unified, integrated platform which helps to simplify the collection, organization and analysis of data.

With it, enterprises can turn data into insights through an integrated cloud-native architecture.

IBM Cloud Pak for Data is extensible and easily customized to unique client data and AI landscapes through an integrated catalog of IBM, open source, and third-party microservices.



Why Cloud Pak for Data?

The business case



Enterprises are not leveraging the data that is created

- **2.5 quintillion bytes** of data is created around the world **EVERY DAY!**
- **~70%** of a company's data remains unused, yet every company must be data driven in today's ecosystem
- **Only 7.3%** of leading company executives are confident in their own strategy for leveraging data



Most enterprises are challenged to be truly data driven

- **Most large enterprises** have legacy systems making it difficult to connect to and utilize all their data
- **Most industries** have to deal with increasing regulatory constraints to protect their data
- Highly paid teams spend **50-80%** of their time to find, prepare, and govern data assets with increasing silos
- **97.2%** of leading enterprise executives launch data and AI initiatives due to the knowledge that three-quarters of Fortune 1,000 companies have been replaced by data driven competitors



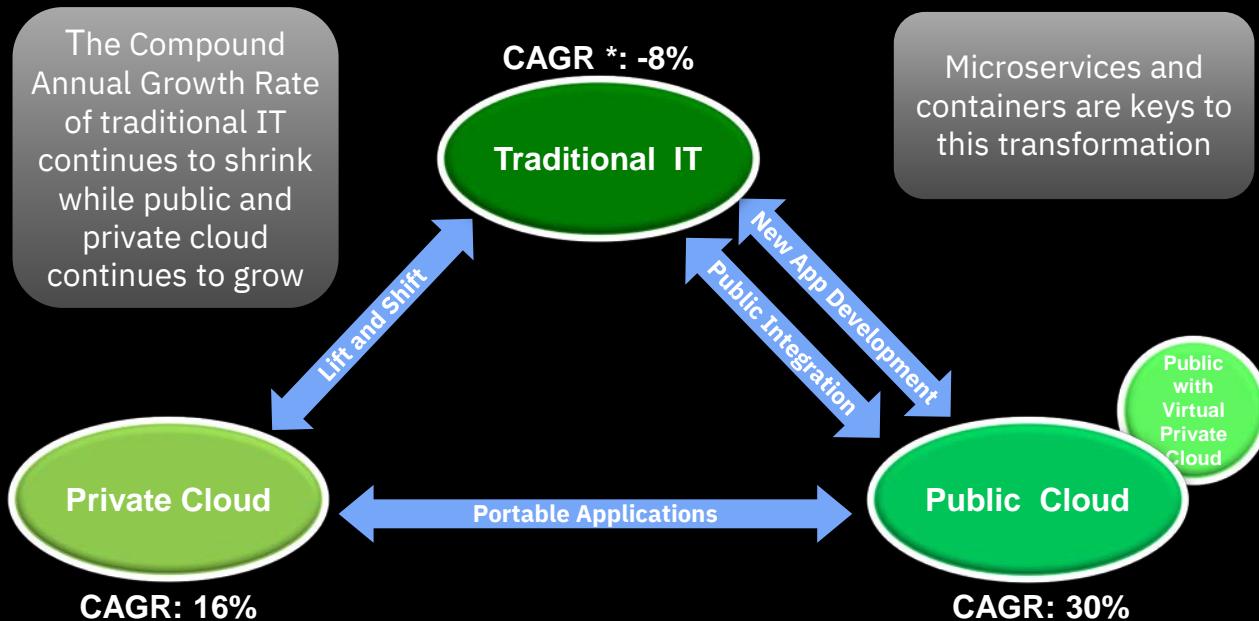
Cloud Pak for Data!

- Comprehensive, integrated, and extensible platform to make data ready for AI
- Built on enterprise grade **private cloud** that is reliable, efficient, scalable and portable
- **The right technology for the right job!**

Cloud Pak for Data: built on cloud-native architecture

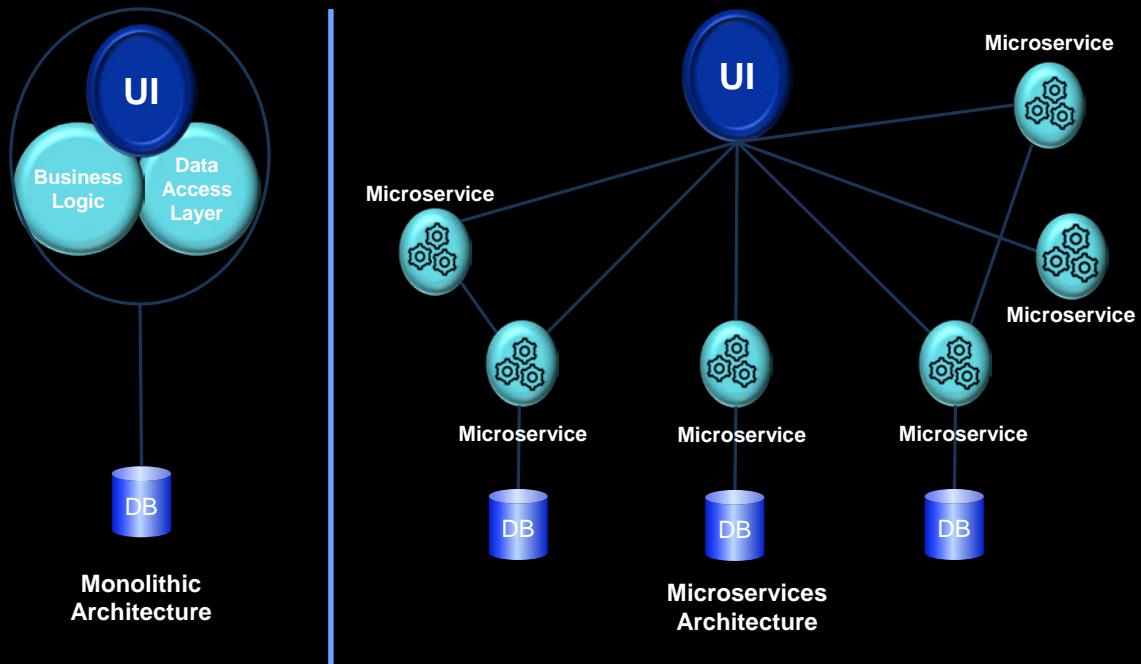


**Multi-cloud is being driven by cloud-native architectures
Microservices and containers are changing IT**



Microservices – the first key to cloud native applications

Making development & deployment more efficient



Microservices benefits *

- ***Improved fault isolation:***
Larger applications can remain largely unaffected by the failure of a single module
- ***Technological flexibility:***
Try out a new technology stack on an individual service and roll it back if required
- ***Easier development:***
A new developer can more easily understand the functionality of a service
- ***Optimized deployment:***
Auto provision, auto scale and provide auto-redundancy

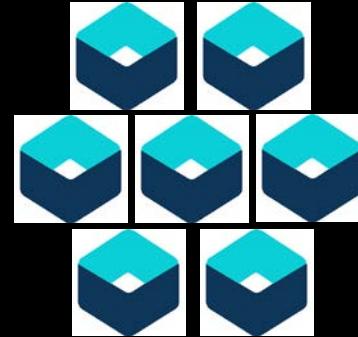
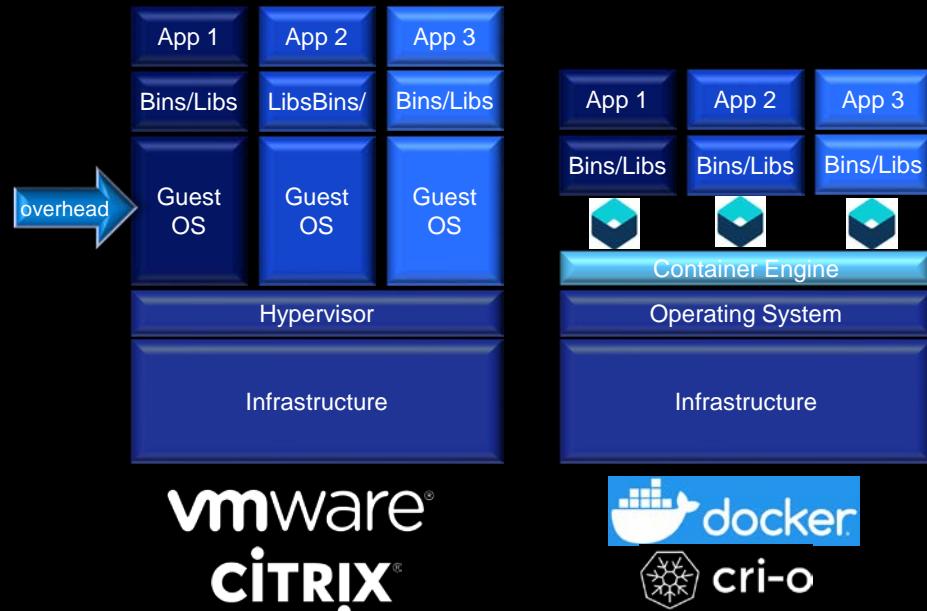
* This is not a claim that a microservice-based application approach is always better for every use case scenario

Containers – the second key to cloud native applications

Reducing operational and development costs



Virtual machines vs. containers *



Containers can be 2 – 3 times more resource efficient than virtual machines

On average Docker developers ship software 7x more frequently

* Containers virtual software in the way that virtual machines have virtualized hardware

Container automation and orchestration is essential

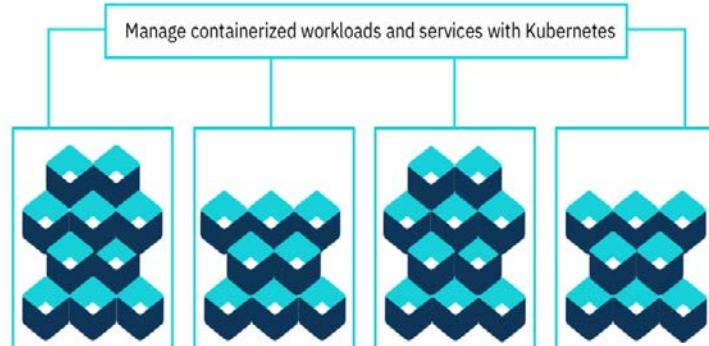
Enter: Kubernetes



**Containers are revolutionizing IT
But they require orchestration**

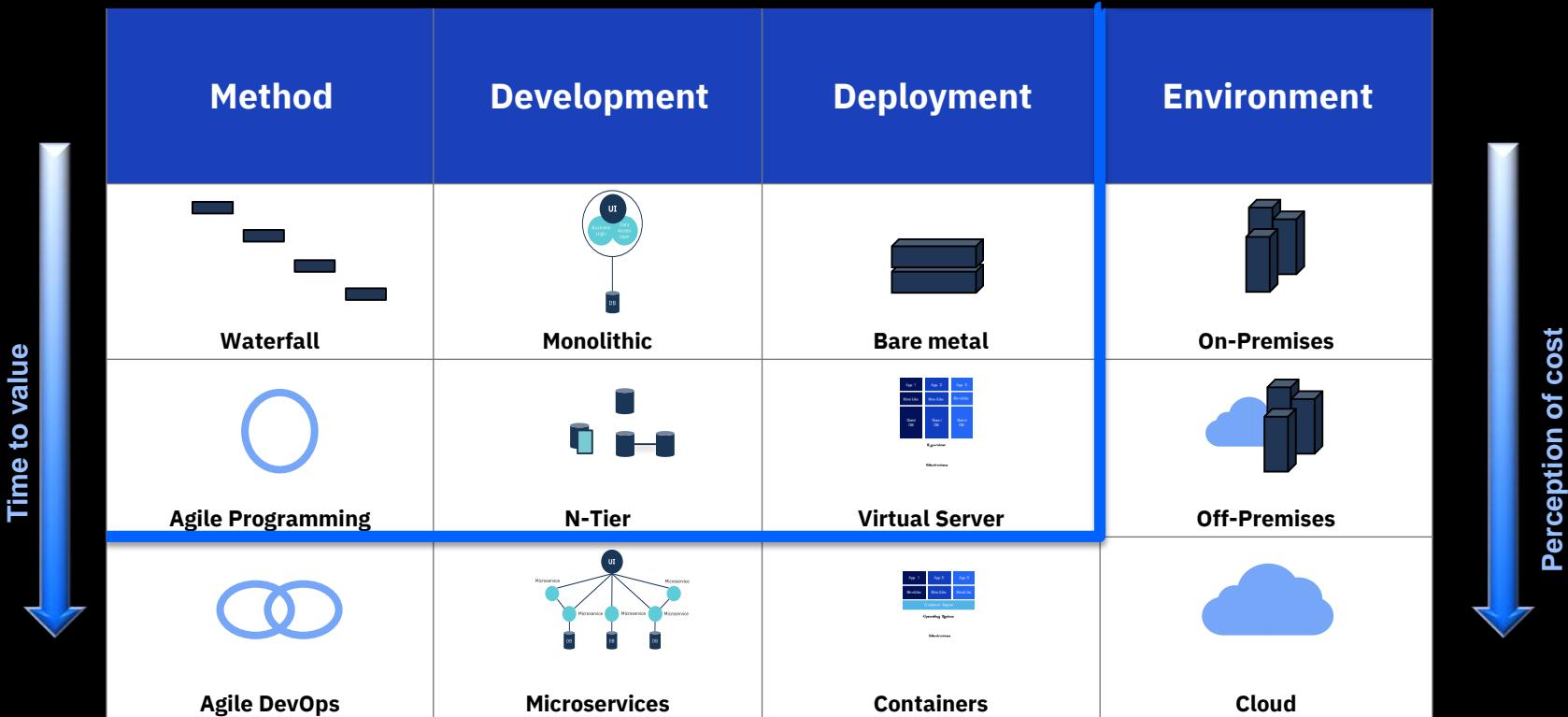


**Kubernetes - κυβερνήτης
Means “helmsman” or “pilot”**



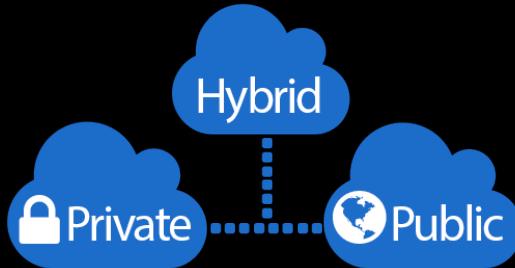
Private Clouds address the new IT reality

Created by digital transformation



Public Cloud + Private Cloud = Hybrid Cloud *

Different cloud options



	Public Cloud	On-Premises Private Cloud	Hosted Private Cloud	Hybrid Cloud
Hardware Deployment and Management	Vendor	Customer	Vendor	Shared between vendor and customer
Hardware Sharing Model	Shared	Dedicated	Dedicated	Partially shared and partially dedicated
Scalability	High	Medium	High	High
Low Cost	Yes	Sometimes	Sometimes	Sometimes
Predictable Cost	No	Yes	Yes	No
Utility Billing	Yes	No	Depends on vendor	Partial
Flexibility	Yes	Limited	Limited	Yes
Customization Capabilities	No	Yes	Depends on vendor	Partial
Enhanced Security and Compliance	No	Yes	Yes	Yes
Instant Provisioning	Yes	Yes	Yes	Yes

* A "Hybrid Cloud" is a highly orchestrated environment, where all sources act as one

A "Multi-cloud" environment simply refers to the use of multiple cloud sources of any kind, without necessarily being orchestrated

Why care about Private Clouds?

Adoption brings agility and efficiency



Cost Efficient & Scalable Infrastructure

50% Benefit

Data Center

System Utilization & Server Reduction

75% Benefit

Manage Performance

Elasticity, Bursting, High Availability

35% Benefit

DevOps

Faster Deployments

30% Benefit

Deployment Efficiency

Containers & Microservices

50% Benefit

Improved Security

Management & Risk Reduction



Accelerate Time to Market



Manage Data at Scale

The Private Cloud Platform Market Leader

Red Hat OpenShift Kubernetes Engine



OPENSIFT

Advanced Cluster Manager

OpenShift Container Platform

OpenShift Kubernetes Engine

Multi-cluster Management

Discovery : Policy : Compliance : Configuration : Workloads

Manage Workloads

Build Cloud-Native Apps

Developer Productivity

Platform Services

Service Mesh : Serverless Builds : CI/CD Pipelines Full Stack Logging Chargeback

Application Services

Databases : Languages Runtimes : Integration Business Automation 100+ ISV Services

Developer Services

Developer CLI : VS Code extensions : IDE Plugins Code Ready Workspaces CodeReady Containers

Cluster Services

Automated Ops : Over-The-Air Updates : Monitoring : Registry : Networking : Router : KubeVirt : OLM : Helm

Kubernetes

Red Hat Enterprise Linux & RHEL CoreOS



Physical



Virtual



Private cloud



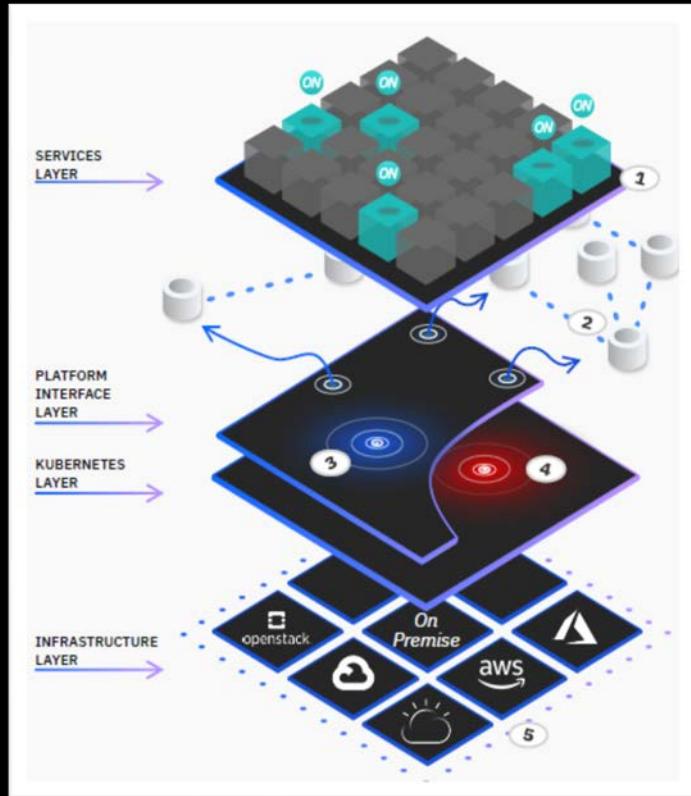
Public cloud



Managed cloud
(Azure, AWS, IBM, Red Hat)

Cloud Pak for Data (CPD)

Architectural overview



1. Services Ecosystem

At a click, access and deploy an ecosystem of 45+ services and templates from IBM and third parties.

2. Data Virtualization

Query across multiple data sources fast and easy without moving your data.

3. Control Plane (Platform Interface)

Speed time to value with a single platform that integrates data management, data governance and analysis for greater efficiency and improved use of resources

4. OpenShift (Kubernetes Engine)



Leverage the leading hybrid cloud, enterprise container platform for an innovative and fast deployment strategy.

5. Any Cloud (Infrastructure)

Avoid lock-in and leverage all cloud infrastructure with our “Any Cloud” mentality.

CPD Services Ecosystem



1. Cloud Pak for Data Base Services

Collect		Organize		Analyze		Deploy / Infuse	
Data Virtualization		Watson Knowledge Catalog (With IA, IGC, Refinery, InstaScan)		Watson Studio	Analytics Engine	Data Science: Model Design & Deployment	
Db2 Warehouse	PostgreSQL	Open Source Management		Dashboards	IBM Streams	Watson OpenScale	
Db2 Event Store	IBM Streams			Industry Accelerators (many)			
Db2 Big SQL	NPS			Watson Machine Learning			
OpenShift / Control Plane (Lite)							

2. Premium Services (purchase license or BYOL)

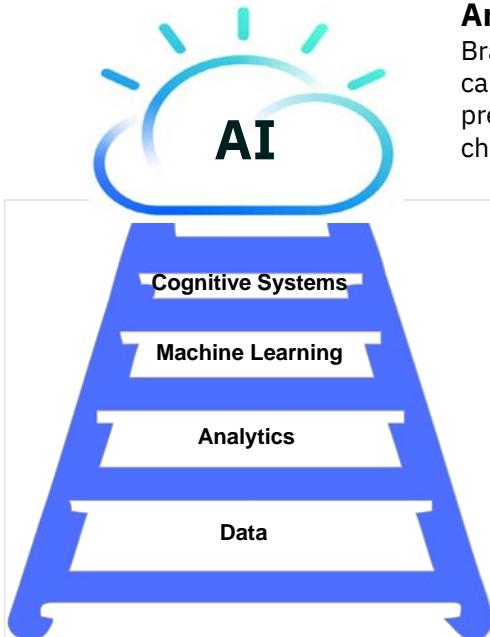
Collect	Organize	Analyze	Deploy / Infuse
Db2 AESE Virtual Data Pipeline	Infosphere DataStage Edition Infosphere Regulatory Accelerator Infosphere multi-cloud Data Mvmnt Infosphere Entity Resolution Master Data Management	Cognos Analytics SPSS Modeler Decision Optimization Watson Explorer Planning Analytics	Watson Assistant / Discovery Watson API Kit (Speech to Text, Text to Speech, Natural Language Understanding) Watson Financial Crimes Insights Planning Analytics

3. Third Party Extension Services



The Original* AI Ladder

Described: AI vs. Cognitive systems vs. Machine Learning vs. Analytics



Artificial Intelligence

Branch of computer science dealing with the simulation of intelligent behavior in machines. It is the capability of machines imitating intelligent (usually human) behavior. AI has been around for many years, predating machine learning and the term “cognitive systems.” e.g. IBM Deep Blue beat the world champ at chess in 1996 by using brute force computer programming.

Cognitive systems

Simulation of human thought processes in a computerized model with self-learning systems that use data mining, pattern recognition and natural language processing to mimic the way the human brain works. e.g. IBM Watson beats Jeopardy! Champs in 2011 using an ML based predictive analytics driven by NLP. Modern AI applications can also be described as “Cognitive.”

Machine Learning

Branch of AI that employs statistical models and advanced algorithms to enable computers to become "intelligent" by "learning" from data in lieu of being specifically programmed. Can include predictive analytics, NLP, facial recognition, search engines. e.g. Google's DeepMind beats Go champion by using ML neural networks and search.

Analytics

The discovery, interpretation, and communication of meaningful patterns in data.

Data

Includes structured, semi-structured and unstructured datatypes in on-premises, public or private cloud, and hybrid environments.

* Presented in 2017 by Rob Thomas, IBM Senior Vice President, Cloud and Data Platform

Cloud Pak for Data (CPD)

Make your data ready for AI



There is
no **AI**
without **IA**



Infuse - Deploy trusted AI-driven business processes



Analyze - Scale insights with ML everywhere



Organize - Create a trusted analytics foundation



Collect - Make data simple & accessible

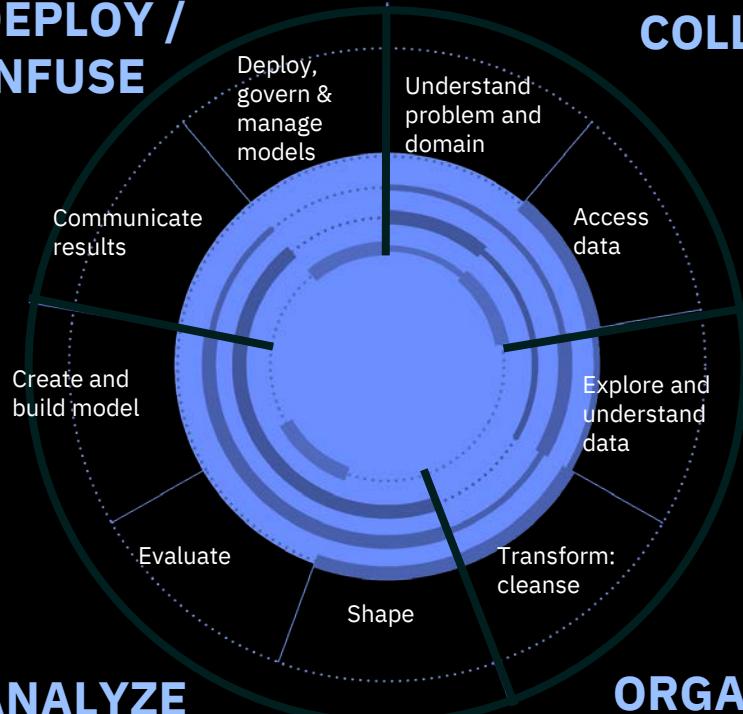
Strong Foundation – Built on “Cloud native architecture”

Cloud Pak for Data (CPD)

Increases workforce productivity across the analytics lifecycle



DEPLOY / INFUSE



COLLECT

ORGANIZE

Explore and understand data

Transform: cleanse

Shape

Evaluate

Create and build model

Communicate results

Deploy, govern & manage models

Understand problem and domain

Access data

Administrator / Architect

Ensures the usability of the compute, network, storage, etc.



Data Engineer

Architects data pipelines & ensures operability



Data Steward

Governs data and ensures regulatory compliance



Business Analyst

Works with data to apply insights to business strategy



Data Scientist

Dives deep into the data to draw insights for the business



Application Developer

Plugs into analysis and code to build applications



CPD Administration

Administer and manage the platform



The screenshot illustrates the CPD Administration interface, showing how to manage the platform through various sections:

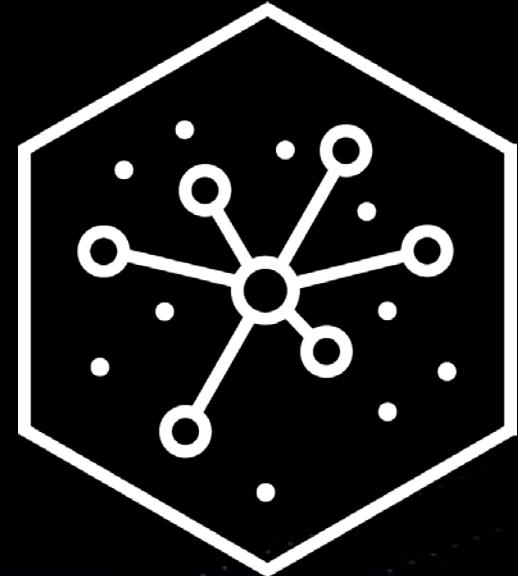
- Left Sidebar (Administer):** Contains links for "Manage platform" (highlighted with a green oval), "Configure platform", "Gather diagnostics", "Manage users" (highlighted with a green oval), and "Customize branding".
- Deployments View:** Shows a list of deployments with columns for Name, Type, Installed on, Service instances, vCPU, and Memory (GB). It includes a search bar and filters for All types and Clear all.
- Users View:** Shows a list of users with columns for Name, User ID, and Username.
- Pods View:** Shows a list of pods with columns for Name, Pods, vCPU, and Memory (GB).

Name	User ID	Username
admin	1000330999	admin
Business Analyst	1000331009	businessanalyst
CPD User	1000331002	cpduser
Data Engineer	1000331003	dataengineer

Name	Pods	vCPU	Memory (GB)
Total	8	0.14 of 1.65	1.50 of 2.80
db2wh-1590588027600-ibm-unified-console-api	5	0.10 of 1.00	1.30 of 2.00
db2wh-1590588027600-ibm-unified-console-influxdb	1	0.00 of 0.10	0.08 of 0.25
db2wh-1590588027600-ibm-unified-console-ucgoapi	1	0.00 of 0.25	0.07 of 0.10

Business Use Case

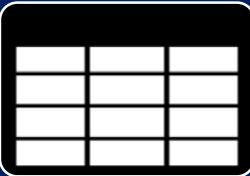
Lab 02 – Business Use Case: Customer Churn



Trade Co. Challenges

- Customer retention problem leading to declining revenue
- Underperforming rules-based system to identify separation (churn) risk
- Lack of centralized, vetted, and reliable data to ensure accuracy of analytics
- Disparate analytical tools for reporting and model development
- No simple way to infuse machine learning models into the customer facing Stock Trader Application

Separation (Churn) Risk: Current Rules Based System



Built Using Limited Data

Rules are developed using a single source of data that contains customer demographic information.



Manual Process to Develop Rules

Rules are manually developed based on the past experience of the marketing team. Rules are only updated once a year.



Low Overall Predictive Accuracy

Low overall predictive accuracy. We are both missing identifying customers who ultimately separate and incorrectly assigning high risk to customers who ultimately stay.

Separation (Churn) Risk: New Data Driven Approach



Incorporate Multiple Data Sources

Use vetted centralized transactional data along with customer demographics to understand separation behavior. Also, include the outcomes of the rules-based system for each customer where an accurate prediction was rendered.



Data Driven Process to Develop Machine Learning Models

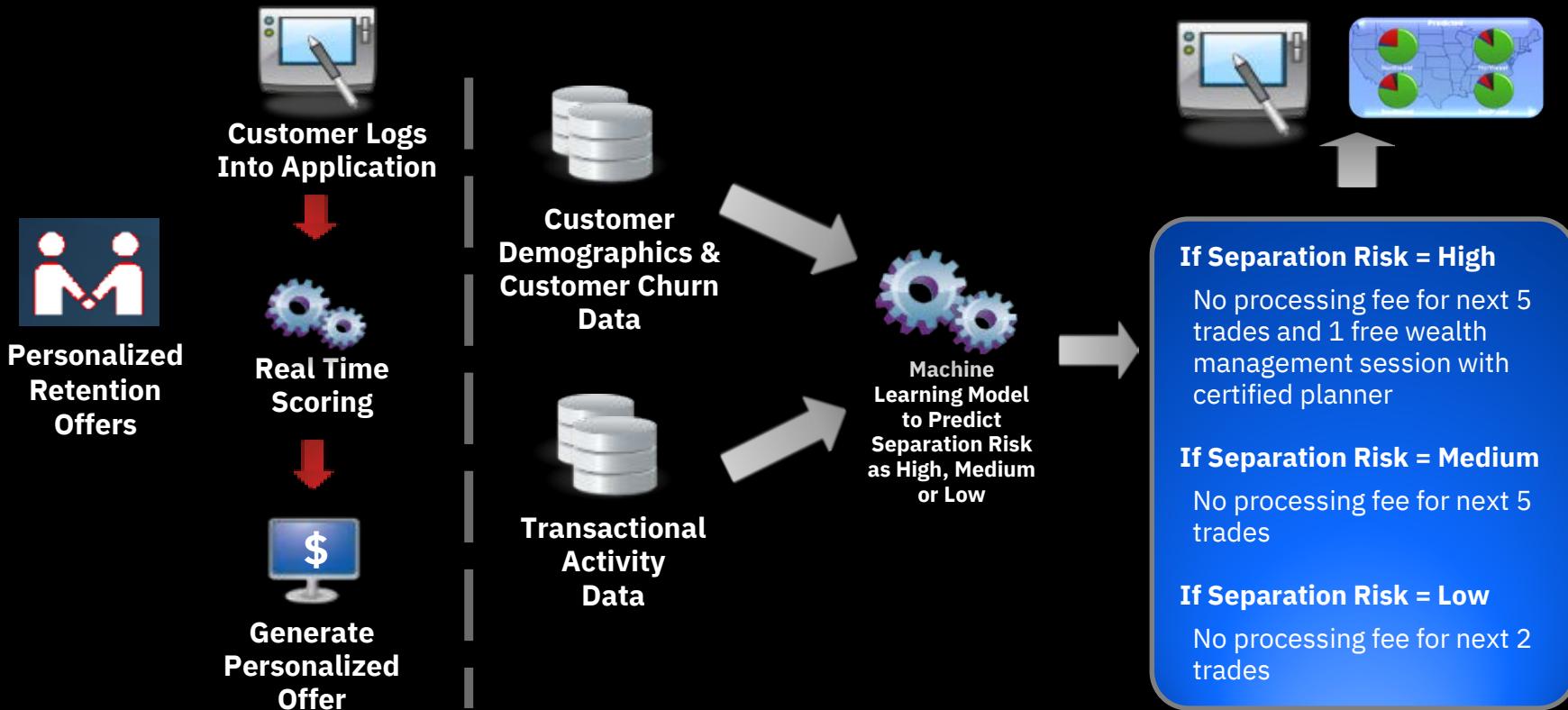
Develop predictive models for separation risk that automatically discover and incorporate all the patterns in the data including interactions and contingencies.



High Accuracy from Adaptive Machine Learning

Models will classify separation risk with a higher overall accuracy and will adapt to changing patterns in risk to maintain that accuracy. Machine Learning models will incorporate all the understanding from the rules-based system and build on that to develop highly complex set of predictive conditions.

Deployment: Stock Trader App. Integrated with AI





Stock Trader Application



Stock Trader Application with Infused ML



Flat File of Monthly Sales Performance



Cloud Pak for Data

1)



Dashboard of Sales Performance (Monthly Metrics)

4)



Dashboard of Churn Risk (Demographics Discovery)

2)



Organize Data: Discover, Govern and Catalog



Transform Data: Merge and Prepare data for Analysis

5)



Build Machine Learning Model for Churn Risk (AutoAI and Notebook)

7)



Integrate Model into Application (Stock Trader)



Dashboard of Business Impact (Monthly Metrics after AI)



Customer Demographics

Historic Churn Risk Results



Customer Activity



Combined Data:
• Demographics
• Churn Risk
• Activity



Post analysis customer data

ORGANIZE

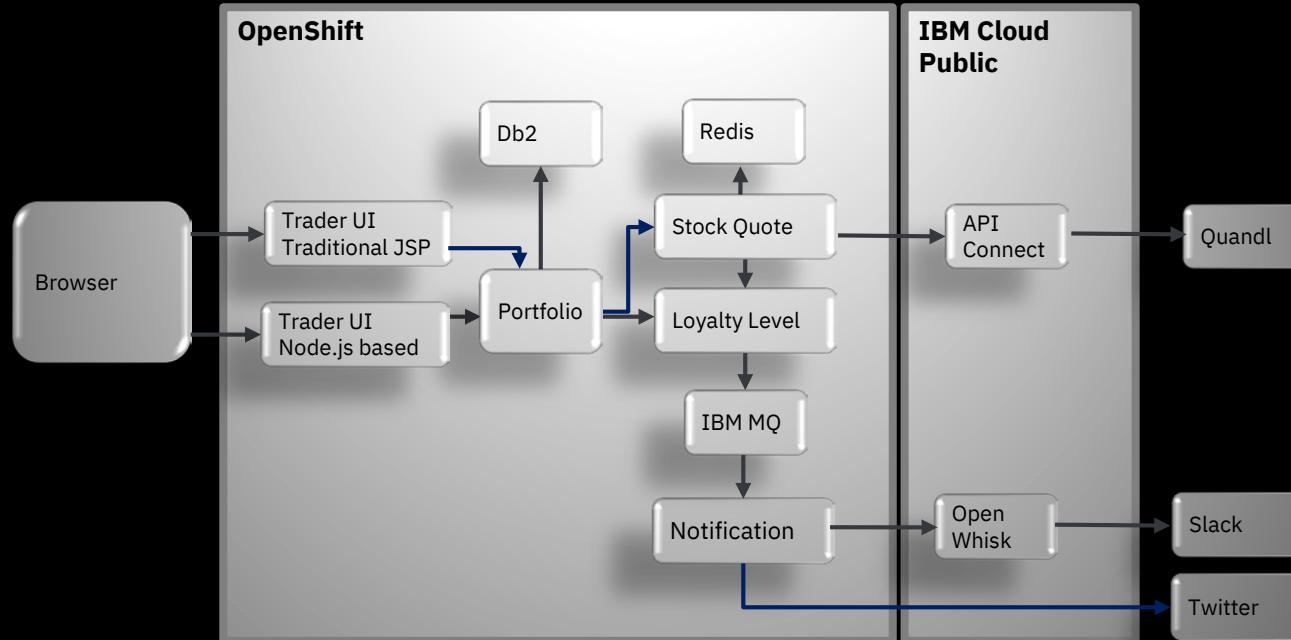
ANALYZE

DEPLOY / INFUSE

COLLECT

Trade Co. Application “Stock Trader” – Before Microservices Application

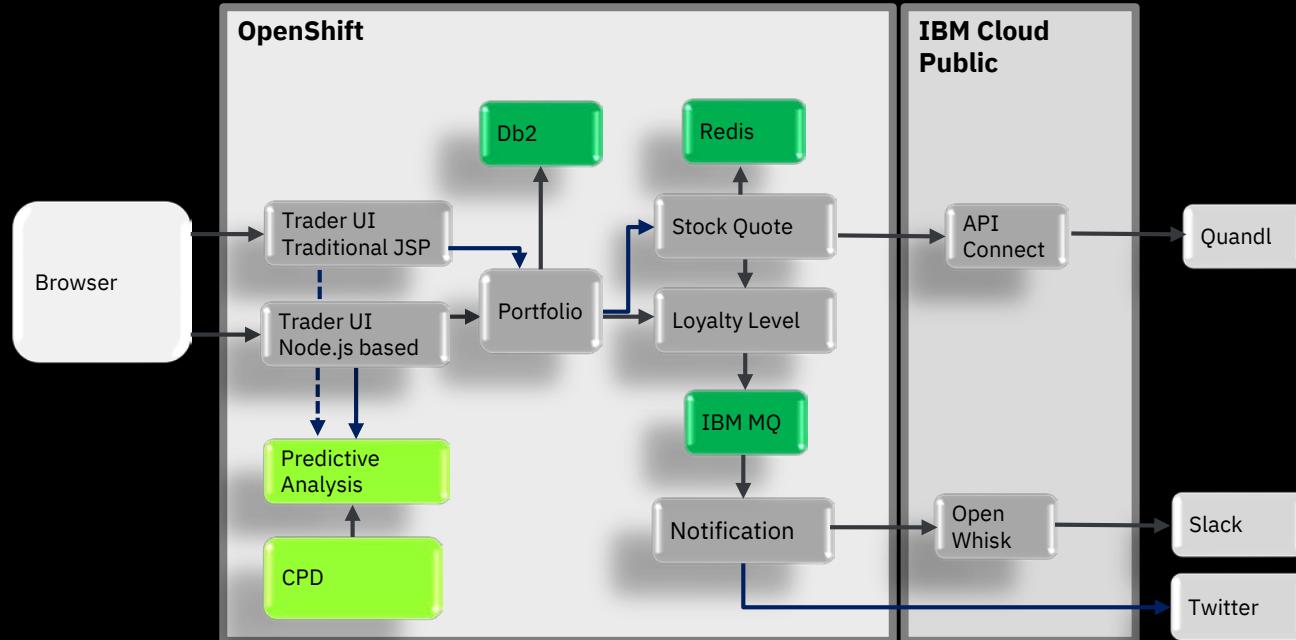
Stock Trader as a modern microservice application



Trade Co. Application “Stock Trader” – After: Infused with ML

Microservices Application

Stock Trader (enhanced) as a modern microservice application



Stock Trader – After Monetizing the ML model

IBM TRADER

Home Summary Add Portfolio Predictive Analysis Change User

Summary

Welcome to IBM Trader powered by ICP for Data

- Create a new portfolio
- Retrieve selected portfolio
- Update selected portfolio (add stock)
- Delete selected portfolio

Owner	Total	Loyalty Level
TechStocks	\$115,670	Gold

Submit Change User

Though looking simple - a lot has gone through to provide machine learning predictive model scoring service.

no processing fee for next 5 trades

Advertisement

IBM Cloud Pak for Data

- Cloud agile
- Lightning fast
- AI-ready

No assembly required

Trade Co. Dashboards

Before and After deploying the CPD developed ML model

Before AI



After AI



IBM Analytics Modernization Workshop

Part 1

<ul style="list-style-type: none">• Introduction• Business Use Case	<ul style="list-style-type: none">• Lab 01• Lab 02
<ul style="list-style-type: none">• Collect: Connect• Organize• Collect: Virtualize	<ul style="list-style-type: none">• Lab 03• Lab 04• Lab 05
<ul style="list-style-type: none">• Analyze• Deploy• Infuse – OpenScale• Infuse – Cognos Analytics• Wrap up	<ul style="list-style-type: none">• Lab 06• Lab 07• Lab 08• Lab 09• Lab 10

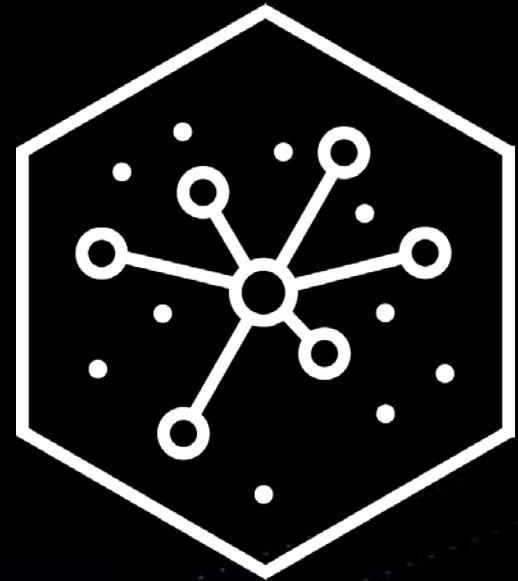
IBM Analytics Modernization Workshop

Part 2

- | | | |
|---------------|---|--|
| | <ul style="list-style-type: none">• Introduction• Business Use Case | <ul style="list-style-type: none">• Lab 01• Lab 02 |
| Part 2 | <ul style="list-style-type: none">• Collect: Connect• Organize• Collect: Virtualize | <ul style="list-style-type: none">• Lab 03• Lab 04• Lab 05 |
| | <ul style="list-style-type: none">• Analyze• Deploy• Infuse – OpenScale• Infuse – Cognos Analytics• Wrap up | <ul style="list-style-type: none">• Lab 06• Lab 07• Lab 08• Lab 09• Lab 10 |

Collect

*Lab 03 – Collect: Connections
Lab 05 – Collect: Virtualize*



CPD Collect

1. Provision in-cluster databases

Provision, host, and manage these data sources directly on the CPD cluster

 CockroachDB Partner Premium Ultra-resilient distributed SQL clusters designed for global business.	 Data Virtualization IBM Enabled ✓ Query many data sources as one.	 IBMDb2 IBM Premium Relational database that delivers advanced data management and analytics capabilities for transactional and warehousing workloads.	 Db2 Event Store IBM In-memory data store capable of extremely high speed ingest and deep, real-time analytics.	 Db2 Warehouse IBM Enabled ✓ Data warehouse designed for high-performance, in-database analytics.
 EDB Postgres Partner Premium Object-relational database designed for developers.	 IBM Db2 for z/OS IBM Create databases in Db2 for z/OS and work directly with the data from IBM Cloud Pak for Data	 MongoDB Enterprise Advanced Partner Premium Scalable, NoSQL database for enterprise deployments.	 Virtual Data Pipeline IBM Premium Access all the data you need for analytics and application testing without impacting production databases.	

CPD Collect

2. Connect to existing data sources: IBM *

* Note: This list is constantly updated and shows what exists as of August 27, 2020.

Connect directly to these data sources and perform the CPD component functionality shown

IBM Data Sources	Cognos Dashboards	DataStage Edition	Data Virtualization	WKC	Watson Studio
Analytics Engine HDFS				✓	✓
Classic Federation		✓			
Cloud object storage (IBM)	✓			✓	✓
Cloud object storage (infra)				✓	✓
Cloudant				✓	✓
Cognos Analytics				✓	✓
Compose for MySQL				✓	✓
Data Set		✓			
Data Virtualization	✓			✓	✓
Data Virtualization Mgr z/OS			✓		
PostgreSQL databases				✓	✓
Db2	✓	✓	✓	✓	✓
Db2 Big SQL			✓	✓	✓
Db2 Event Store			✓	✓	
Db2 for i			✓	✓	✓
Db2 for z/OS	✓	✓	✓	✓	✓
Db2 Hosted				✓	✓
Db2 on Cloud	✓	✓	✓	✓	✓
Db2 Warehouse	✓	✓	✓	✓	✓
Db2 Warehouse on Cloud	✓	✓	✓	✓	✓
Distributed Transactions		✓			
DRS		✓			

IBM Data Sources	Cognos Dashboards	DataStage Edition	Data Virtualization	WKC	Watson Studio
External Source				✓	
External Target				✓	
HDFS via Hadoop					✓
Hierarchical				✓	
Hive via Hadoop					✓
Impala via Engine for Hadoop					✓
Informix				✓	✓
Informix Enterprise / Load				✓	
ISD Input / Output				✓	
Java Integration				✓	
Lookup File Set				✓	
Netezza				✓	✓
Planning Analytics					✓
Obj. Strg. OpenStack Swift					✓
PureData for Analytics					✓
WebSphere MQ				✓	
Z/os DVM sources (VSAM, IMS, Adabas, etc.)				✓	

CPD Collect

2. Connect to existing data sources: Third-party *

* Note: This list is constantly updated and shows what exists as of August 27, 2020.

Connect directly to these data sources and perform the CPD component functionality shown

Third-party Data Sources	Cognos Dashboards	DataStage Edition	Data Virtualization	WKC	Watson Studio	Third-party Data Sources	Cognos Dashboards	DataStage Edition	Data Virtualization	WKC	Watson Studio
Amazon Redshift			✓	✓	✓	Looker				✓	✓
Amazon S3		✓		✓	✓	MariaDB			✓		
Apache Cassandra	✓					MSFT Azure Blob & File		✓			
Apache Derby			✓			MSFT Azure Lake Store				✓	✓
Apache Hbase		✓				MSFT Azure SQL DB				✓	✓
Apache HDFS				✓	✓	MSFT SQL Server	✓	✓	✓	✓	✓
Apache Hive			✓	✓	✓	Minio				✓	✓
Apache Kafka	✓					Mongo			✓	✓	
Azure Storage	✓					MySQL			✓	✓	✓
BDFS	✓					ODBC		✓			
Cloudera Impala			✓	✓	✓	OData				✓	✓
Dropbox				✓	✓	Oracle		✓	✓	✓	✓
Filesystem	✓			✓	✓	Pivotal Greenplum		✓		✓	✓
FTP Enterprise	✓					PostgreSQL	✓		✓	✓	✓
FTP	✓			✓	✓	Salesforce.com		✓		✓	✓
Generic JDBC				✓	✓	SAP HANA			✓		
Google BigQuery	✓		✓	✓	✓	SAP Data Object		✓		✓	✓
Google Cloud Storage	✓			✓	✓	Snowflake		✓	✓	✓	✓
HDFS Generic web-HDFS				✓		Sybase Enterprise		✓	✓	✓	✓
HDFS HttpFS				✓		Sybase IQ / OC		✓		✓	✓
Hive JDBC	✓		✓			Tableau				✓	✓
Hive JDBC CDH	✓			✓		Teradata		✓	✓	✓	✓
Hive JDBC HDP	✓			✓							
Hortonworks HDFS				✓	✓						

CPD Collect

2. Connect to existing data sources: Data Files

Connect directly to these data sources and perform the CPD component functionality shown

Other Data Sources	Cognos Dashboards	DataStage Edition	Data Virtualization	WKC	Watson Studio
Comma Separated Value (CSV) files	✓	✓	✓		✓
Microsoft Excel Spreadsheets			✓		✓
Sequential File		✓			
Tab Separated Value (TSV) files			✓		

Connecting to any of the data sources by service:

Service	Comment
Cognos Dashboards	You can use the local and remote data sets that already exist in your analytics projects. You can create connections by selecting <i>Add data</i> source from the analytics dashboard menu.
DataStage Edition	You can transform data that is in a catalog by searching for the data that you want to use and selecting <i>Transform</i> .
Data Virtualization	You can create connections that can be used to virtualize data from the following locations: 1) <i>Connections</i> page, 2) <i>Data Sources</i> page in the Data Virtualization service
Watson Knowledge Center	You can create connections that can be used in the catalog and connections that can be used to curate data.
Watson Studio	Ideally, you should use data that is already in a catalog by searching for the data you want there and add it to an analytics project. Alternatively, you can create connections that can be used in analytics projects from the following locations: 1) <i>Connections</i> page, 2) <i>Assets</i> page of the analytics project. You can also <i>Add data</i> from files.

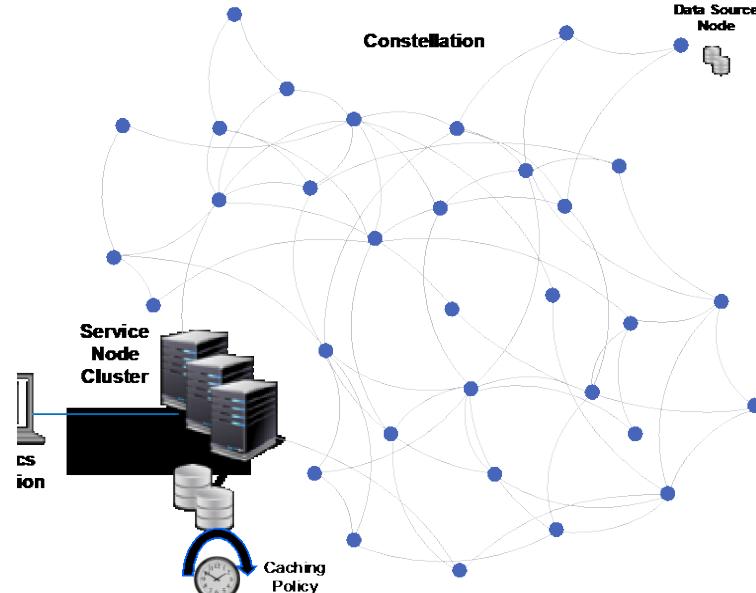
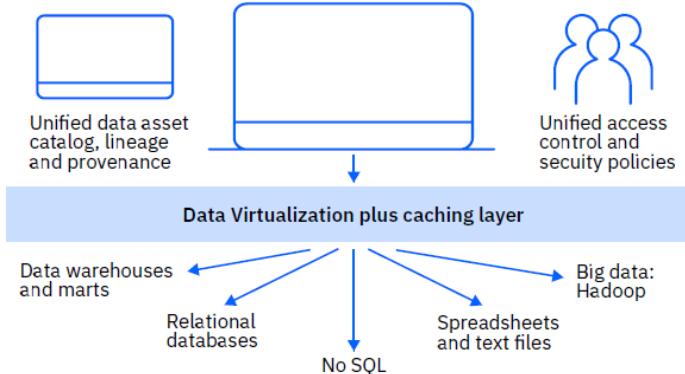
CPD Collect

3. Data Virtualization - Overview

Query across multiple databases and big data repositories which appear as one to an application

Data Virtualization in Cloud Pak for Data is a unique new technology that connects many data sources into a single self-balancing collection of data sources or databases referred to as a *constellation*

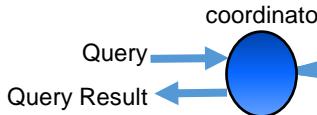
- Analytics are submitted against the data at the source (data is not copied)
- New algorithms are applied that support distributed SQL execution across heterogenous data sources



CPD Data Virtualization

Constellation “Computational Mesh” benefit

Classic Federation & Edge Computing



Query issued against the system

A coordinator receives the request and fans the work out to edge nodes

Edge nodes individually perform as much work as they can based on their own data. Individual results are sent back to the coordinator for final merging and remaining analytics.

Coordinator receives intermediary results from all edge nodes, merges results, and performs remaining analytics

To be clear: Federation is a form of Data Virtualization and has been used successfully for many years in IBM products like Db2

CPD Data Virtualization uses a new Computational Mesh * approach which meets the performance demands of today's modern data access requirements

New Computational Mesh



Query issued against the system

A coordinator receives the request and fans the work out to edge nodes

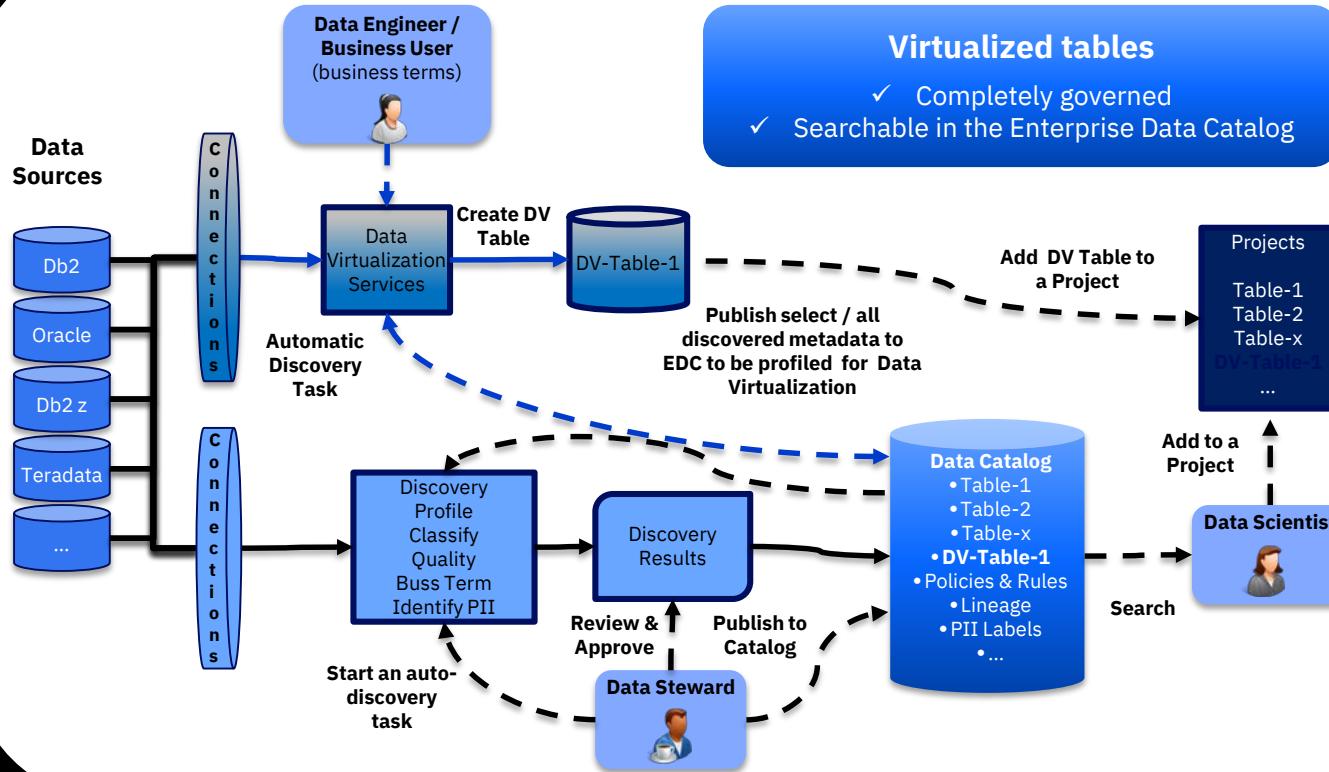
Edge nodes self organize into a constellation where they can communicate with a small number of peers. Nodes collaborate to perform almost all analytics, not only analytics on their own data.

Coordinator receives mostly finalized results from just a fraction of nodes. Completes the final work for the query result.

* Note: this is a work in progress. Remote Connectors with data source support is available today.

CPD Data Virtualization

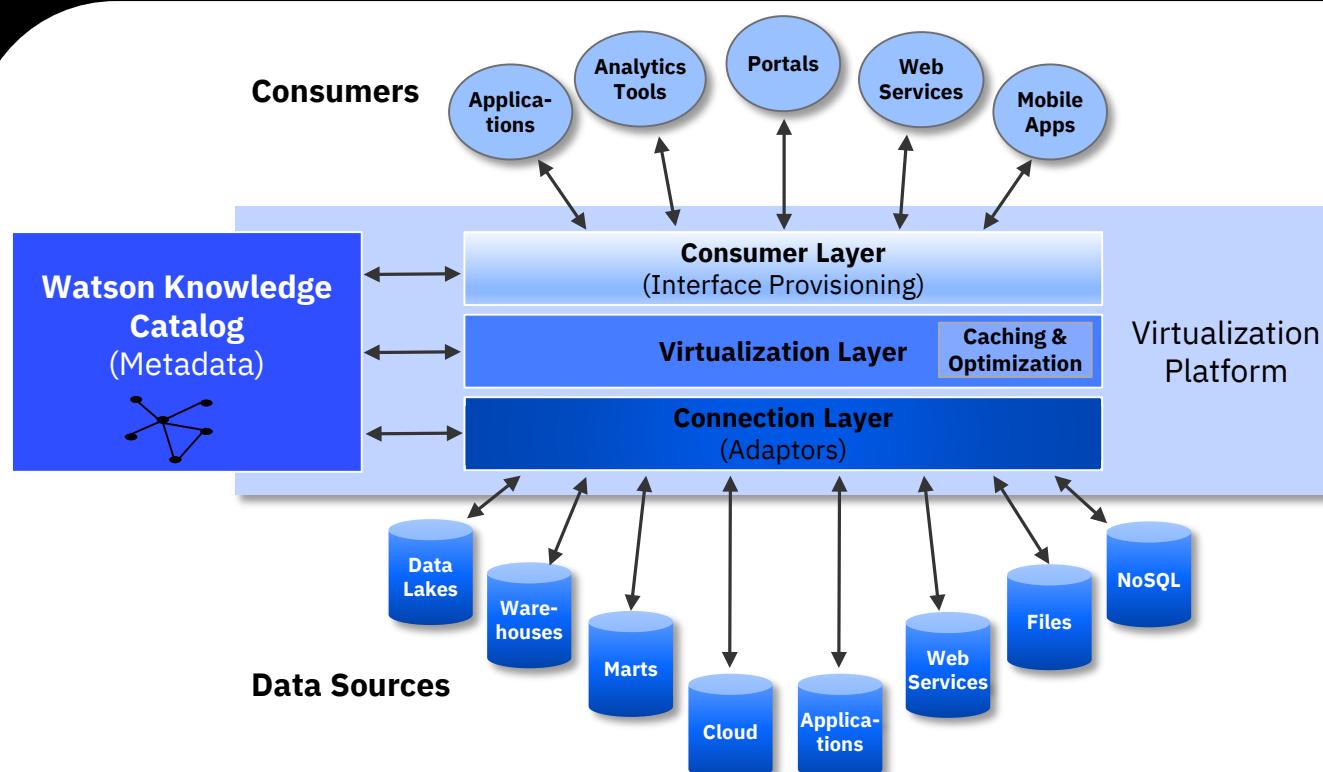
Data Governance is built in



- By using this flow, the DV Table will be published to catalog with fingerprint information collected (quality, profile, classification, assign business terms, Sensitive Information Detection)
- If any PII data is present, the system will automatically flag the table in the catalog with a PII Label
- Data Masking Services can de-identify the sensitive information when user attempts to view the data

CPD Data Virtualization

With Watson Knowledge Catalog (WKC) built in



- Provides the ability to search, view, access, manipulate, and analyze data
- No need to know or understand its physical format or location
- No need to move or copy it

CPD Data Virtualization

Benefits and use cases

Benefits

Simple:

- Self-discovering, self-organizing cluster
- Joins provide a one source input to analytics

Flexible:

- Once established, it is easy to add new sources to the constellation
- Integrates disparate data assets with simple automation, providing seamless access to data as one

Scalable:

- Can access thousands of sources, IOT and edge devices

Cost Effective:

- Leverages the compute resources of source systems to execute the SQL

Secure:

- Inherits privileges & masking policies of the data sources
- Built in governance, security, and access control

Use Cases

Data Scientists:

- Significant productivity increase getting access to sources discovery and assembly of data sets

Current State answer requirement:

- Current state required for up-to-date analytics
- One time access to data, then throw it away
 - e.g. “How much cash is ‘on hand’ across our branches worldwide?” “What is our current ‘claims’ liability?”

ETL and/or Data Governance saturation

- Self-service – In the event that Data Engineers cannot keep up with business demands for access to data



CPD Collect

4. External Data Sets

The Weather Company



Historical Weather Data

- 3 Years of Historical Weather Data
 - Current Weather Conditions
 - Weather Forecast Data
 - Location Look-up Services
 - Industry Accelerators (Retail & Manufacturing)
- 90-day trial

Equifax



Demographic Consumer Data

- Ability to Pay
- Economic Cohorts
- Credit Styles Pro
- Income 360
- Wealth Complete Data
- IXI-Data

Premium

People Data Labs



People Data

- Dataset of ~1.5B profiles consisting of both b2b and b2c data on each person
- Data is accessible via an Enrichment API or Data License

Premium

BCC



Real Time Stock Data

- Real Time Stock Market Data

Premium



CPD Collect

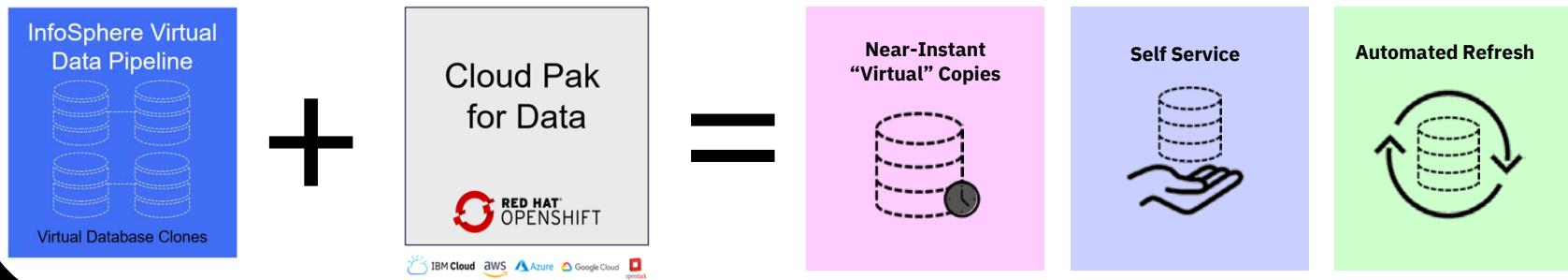
5. Premium Service: Virtual Data Pipeline – Overview

Provision and refresh analytics and test data in minutes

Virtual Data Pipeline in Cloud Pak for Data is a service that allows users to *instantly provision virtual database copies* to work with near real-time for data analytics and application testing, AI model training and testing, and data virtualization.

Accessing production data quickly and securely can be a challenge:

- *Risk:* Accessing Production data can impact operations and present a security and data privacy risk
- *Time Consuming:* Moving data to data warehouses and marts or creating duplicate copies can result in stale data
- *Cost:* Creating multiple copies and versions of data for each functional area uses a lot of storage and increases costs

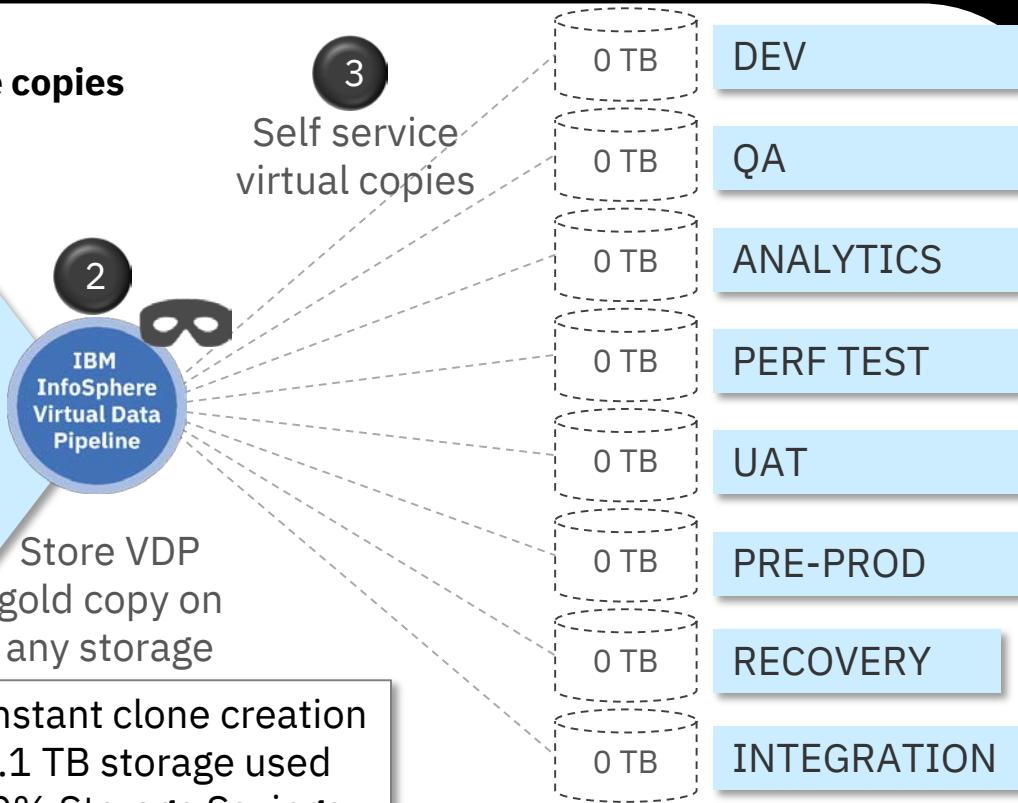
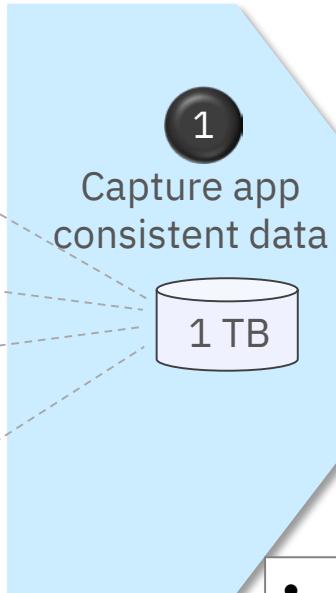




Virtual Data Pipeline

Cost savings example

Capture, clone and serve virtual database copies





Virtual Data Pipeline

Multi-cloud data management example

Manage Multi-cloud data sources from a single pane of glass



1

Capture hybrid
multi-cloud
data sources
for CPD

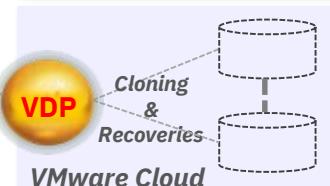
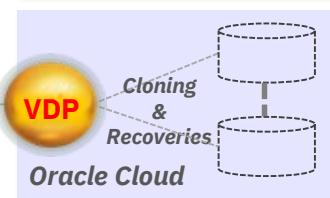
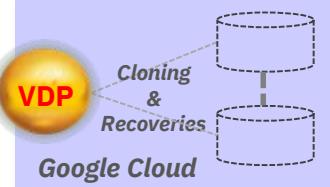
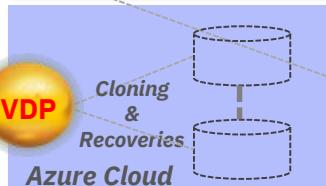
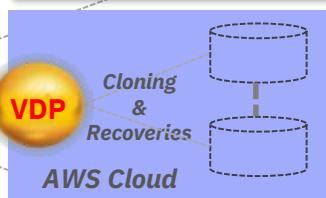
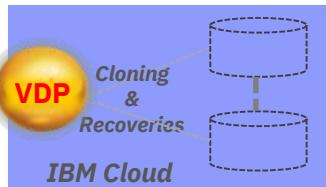
2

VDP

Store VDP gold
copy on any
storage

3

Replicate to public cloud
or another data center



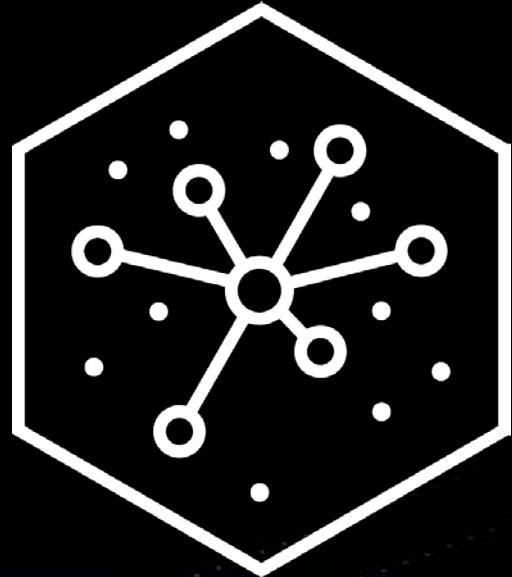
DATABASE AGNOSTIC

STORAGE AGNOSTIC

CLOUD AGNOSTIC

Organize

Lab 04 – Organize





CPD Organize

Watson Knowledge Catalog (WKC)



Data Scientists

Data Analysts

Business

Analysts

- Search and find relevant data
- Prepare data for consumption and analysis
- Consume and analyze the data
- Rate, comment on and share the data



Data Stewards

CDO

LOB Risk

LOB Product

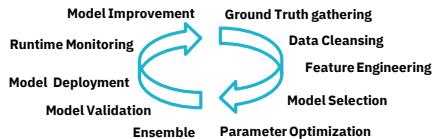
- Manage metadata repository
- Manage Reference Data
- Data governance workflow
- Discover metadata assets
- Classify data assets
- Data stewardship
- Build data glossary
- Data ownership
- Data lineage



Quality Analysts

- Profile data
- Understand, monitor and remediate data quality
- Apply validation rules

AI Lifecycle



Enterprise Data Consumption

Enterprise Data Governance

Enterprise Data Quality

WKC is installed when you install any of the following services:

Watson Studio, Watson Machine Learning, Watson OpenScale

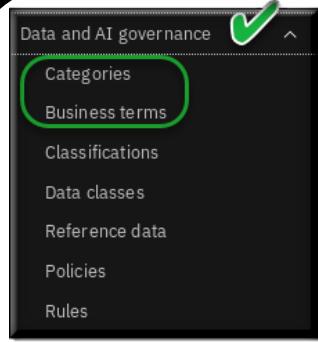
IBM Watson Knowledge Catalog on Cloud Pak for Data

End-to-End Platform for Business-Ready Data



CPD Organize

Data and AI governance: Categories and Business terms



Categories provide logical structure to a business glossary

Business terms standardize definitions of business concepts

1. Manually create Categories and Business Terms
2. Import Categories and Business Terms from CSV or XML files
3. Import a Glossary from an **industry accelerator**





CPD Organize

Data and AI governance: Policies and Rules

Data and AI governance ✓

- Categories
- Business terms
- Classifications
- Data classes
- Reference data
- Policies
- Rules

Policies describe how to control data and consist of one or more rules

Rules describe the criteria for compliance with business objectives

Policies

Published	Draft
<input type="text"/> Find policies	
Data Privacy Company-wide data privacy policy for se Data Privacy Last modified: May 28, 2020	
Net Gains and Net Losses ar A customer can only have a value in Net Customer Churn Category Last modified: May 28, 2020	

Rules

Published	Draft
<input type="text"/> Find rules	
Sort by: Name Show:	
All Credit Card Information Must be Protected All constructs of a credit card must be protected to ensure that those who should not be able to view the information are not allowed to. The information can be redacted since it typically is not used as a unique identifier. This includes the Credit... Data Privacy Governance rule Last modified: May 28, 2020	
All Email Addresses Must be Protected All Email Addresses must be protected to ensure that those who should not be able to view the information are not allowed to. The information must be masked (obfuscated) where the original format and validity of the email address is preserved... Data Privacy	



CPD Organize

Data and AI governance: Classifications and Data classes

The sidebar menu includes:

- Categories
- Business terms
- Classifications** (highlighted with a green oval)
- Data classes** (highlighted with a green oval)
- Reference data
- Policies
- Rules

Classifications describes the sensitivity level of data

Data classes describe the contents of data in a column in a structured data set

Classifications interface:

- Published tab is selected.
- Search bar: Find classifications.
- Sort by: Name.
- Items listed:
 - Confidential**:
Confidential data is data that if compromised in some form, is likely to result in significant and/or long-term harm to individuals whose data it is. Access to confidential information is restricted to those who need it.
[uncategorized]
Last modified: May 27, 2020
 - Personally Identifiable Information**:
Personally identifiable information (PII) is defined as any data that could potentially identify a specific individual or organization, and which can be used to distinguish one person from another. This type of information can be considered PII.
[uncategorized]
Last modified: May 27, 2020

Data classes interface:

- Published tab is selected.
- Search bar: Find data classes.
- Items listed:
 - Account Number**:
A value representing an Account Number.
[uncategorized]
Last modified: May 27, 2020
 - Address Line 3**:
Address Line 3 of a multi-line address.
[uncategorized]
Last modified: May 27, 2020



CPD Organize

Data and AI governance: Reference Data

The sidebar menu includes: Categories, Business terms, Classifications, Data classes, Reference data (highlighted with a green border), Policies, and Rules.

Reference Data Sets define list of permissible values that are allowed for use within a data field.

May be referenced by Business Terms, Policies, Rules and Data Classes

The interface shows: Reference data, Published, Draft, Find reference data search bar, State and Province Codes entry (highlighted with a green border), and Customer Churn Category link.

The details page shows: State and Province Codes, Published, Overview, Related content, and a search bar asking "What are you looking for today?". The table lists codes and their values:

	Code	Value
▼	AA	Armed Forces (the) Americas
▼	AB	Alberta
▼	AE	Armed Forces Europe
▼	AK	Alaska



CPD Organize

Auto-discover assets

Data Discovery

CUSTOMER_DEMOGRAPHICS			
DOB	100%	1	Date of Birth 100% ▾
ESTINCOME	100%	2	NoClassDetected 100% ▾
GENDER	100%	3	Gender 100% ▾
HOMEOwner	100%		Indicator 100% ▾
ID	100%		NoClassDetected 100% ▾
LATITUDE	100%		Latitude 100% ▾
LONGITUDE	100%		Longitude 100% ▾
STATE	92%		US State Code 92% ▾
STATUS	100%		Code 100% ▾
TAXID	93%		US Social Security Num... 93% ▾
ZIP	92%		US Zip Code 92% ▾

Use machine learning based auto-discovery to:

- ① Analyze data quality
- ② Analyze columns (Classify data)
- ③ Assign Business terms

You can perform discovery with data sampling to allow for self-service data access with a search.



CPD Organize

Publish to a catalog

Catalog and govern your assets

Catalogs / CPD Workshop Catalog

CPD Workshop Catalog

Browse Assets Access Control Settings

What assets are you looking for?

Any type Any source Any tag

Showing 8 of 8 items

<input type="checkbox"/> Name	Owner	Tags
<input type="checkbox"/> Customer Activity	CPD User	
<input type="checkbox"/> Customer Churn	CPD User	
<input type="checkbox"/> Customer Churn	CPD User	
<input type="checkbox"/> Customer Demographics	CPD User	
<input type="checkbox"/> Db2Warehouse	CPD User	global...

Watson Recommend: Highly Rated Recently Added

Warehouse Data asset Customer Demographics Data asset Customer Activity Data asset

CPD User CPD User CPD User CPD User CPD User CPD User

May 28, 2020 10:38 AM May 28, 2020 10:41 AM May 28, 2020 10:42 AM

0 reviews 1 review 0 reviews

★ ★ ★ ★ ★ ★ ★ ★ ★ ★ ★ ★ ★ ★ ★

Showing 8 of 8 items

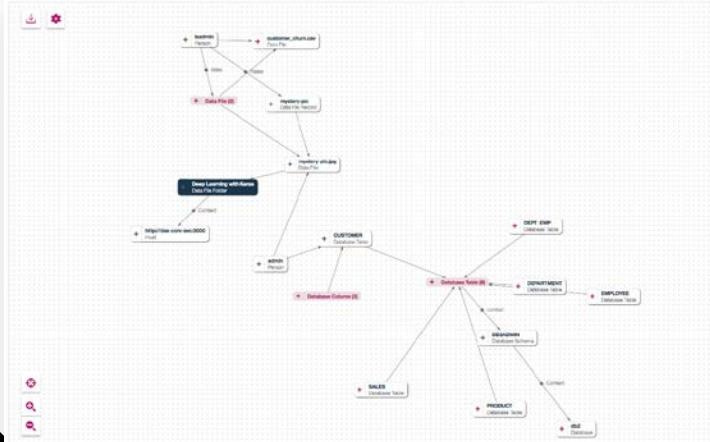
Checklist



CPD Organize

Relationship graph with explorer

- Explore relationships between data assets, terms, analytic assets, users, etc.
- Gain in-depth understanding of metadata through crowdsourcing (e.g., ratings, comments) and machine learning



This screenshot shows a simplified view of the 'mystery-pic.jpg' data asset within the CPD Relationship Graph Explorer. On the left, there are two cards for 'mystery-pic.jpg': one for the Data File and one for the Data File Record. Both cards show a 5-star rating (1 rating), a description field ('Deep Learning with Keras'), and a comment section. On the right, a smaller graph visualization shows the 'mystery-pic.jpg' node connected to other nodes like 'Deep Learning with Keras' (Data File Folder) and 'badadmin Person' (User), with relationship types like 'Context' and 'multiple relationships'.

Explore deeper to understand context and usage patterns



CPD Organize

Refine data with visualizations

Refine can cleanse and shape tabular data with a graphical flow editor using functions and logical operators.

Use it to remove data that is incorrect, incomplete, improperly formatted, etc.

Shape the data by filtering, sorting, combining or removing columns. You can create a Data Refinery flow as a set of ordered operations on the data to run repeatedly any time.



My Projects / CPD Workshop Analytics Project / Customer Demographics

Preview Profile Lineage

Schema: 18 Columns | 2066 rows
Preview: 1000 rows Last refresh: 1 day ago

ID Smallint GENDER String STATUS String CHILDREN Smallint ESTINCOME Decimal HOMEOWNER String AGE Smallint TAXID String

ID Smallint	GENDER String	STATUS String	CHILDREN Smallint	ESTINCOME Decimal	HOMEOWNER String	AGE Smallint	TAXID String
Identif... ▾	Gender ▾	Code ▾	Code ▾	Not clas... ▾	Indicator ▾	Code ▾	US So... ▾
481	F	M	2	28267	N	30	386283240
482	F	M	2	36725.1	N	56	162447113
483	M	S	1	94188.3	N	58	673845765
484	F	M	2	91861	Y	42	209619292



Data Refinery also includes a graphical interface to profile data to validate it with 20+ customizable charts that give perspective and insights into the data.



CPD Organize

Profile data

The **Profile** of a data asset includes generated metadata and statistics about the textual content of the data.

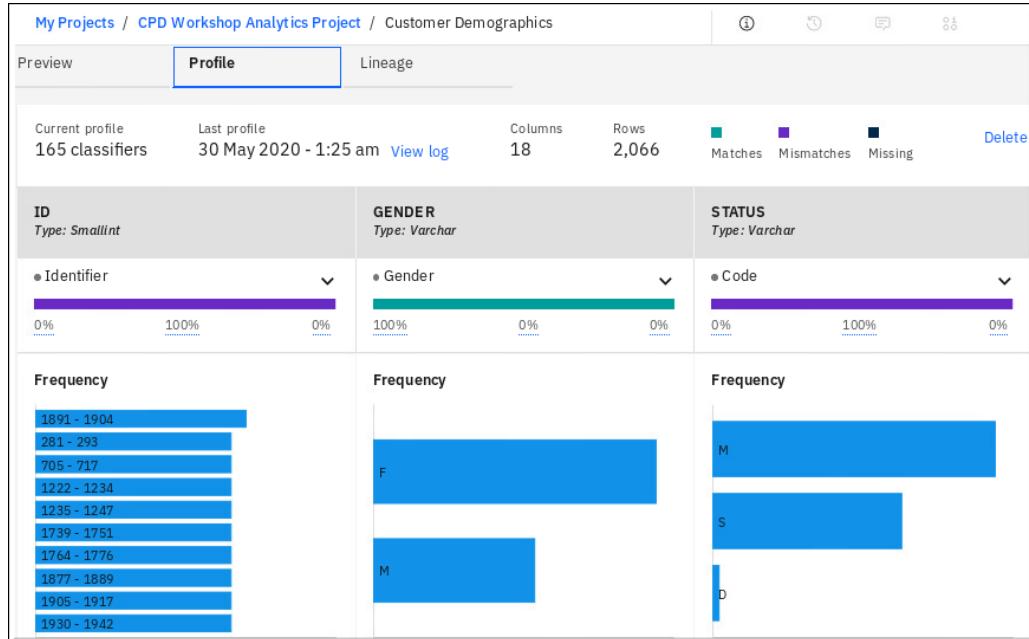
All catalog or project members can see data asset profiles.

Profiles are automatically created:

- In catalogs, profiles for unstructured data assets are created automatically, regardless of whether policies are enforced
- In governed catalogs, profiles for structured data assets are created automatically

Profiles can be manually created:

- In ungoverned catalogs for structured data assets
- In projects for both structured and unstructured data assets



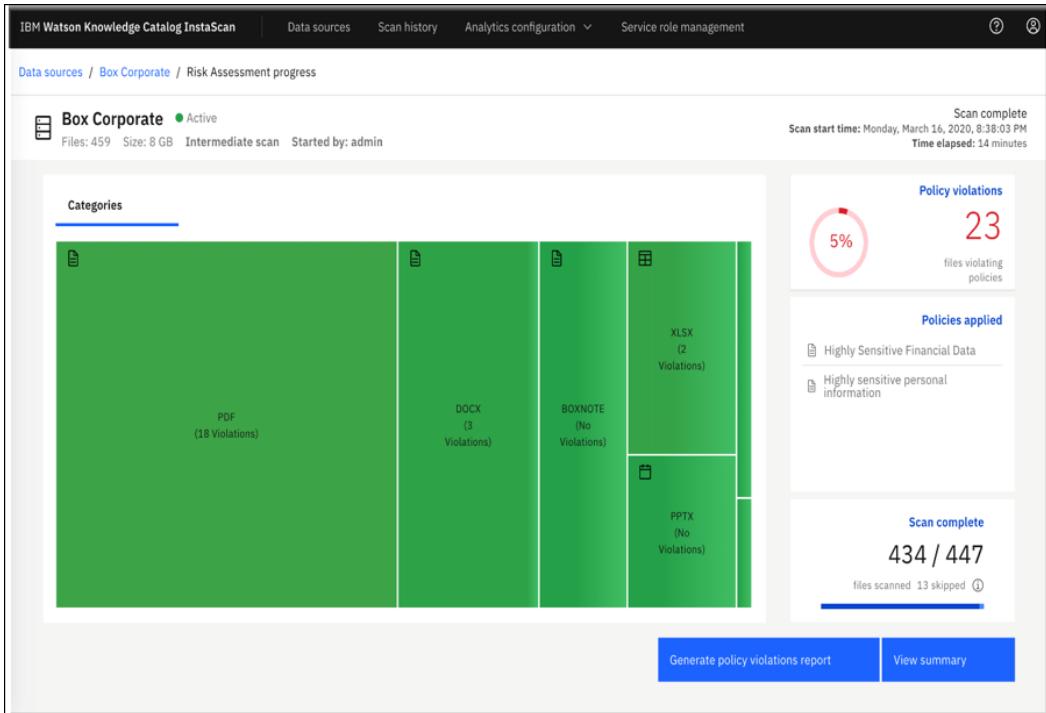


CPD Organize

InstaScan

InstaScan can perform risk assessments and compliance checks of unstructured data

- Scan email, PDFs, word processor documents and images
- Quickly determine which areas have high concentration of sensitive information and prioritize hot spots
- Automatically apply classification labels to data that violates corporate policies
- Create and share ongoing compliance reports with CISO or regulators
- Integrates with Box Shield to provides a comprehensive data privacy solution for unstructured data in Box





CPD Organize

Search for data

The screenshot shows the CPD Organize interface with a search bar at the top containing the term "churn". Below the search bar, there are filters for "Any type", "Any tag", and "Steward/Owner". The results section displays "Showing 19 of 19 items" and lists three assets:

Name	Type
Gender Categories > Customer Churn Category Customer Churn Category	Business term
Customer Churn All catalogs > CPD Workshop Catalog	Data asset
Customer Churn All projects > CPD Workshop Analytics Project	Data asset

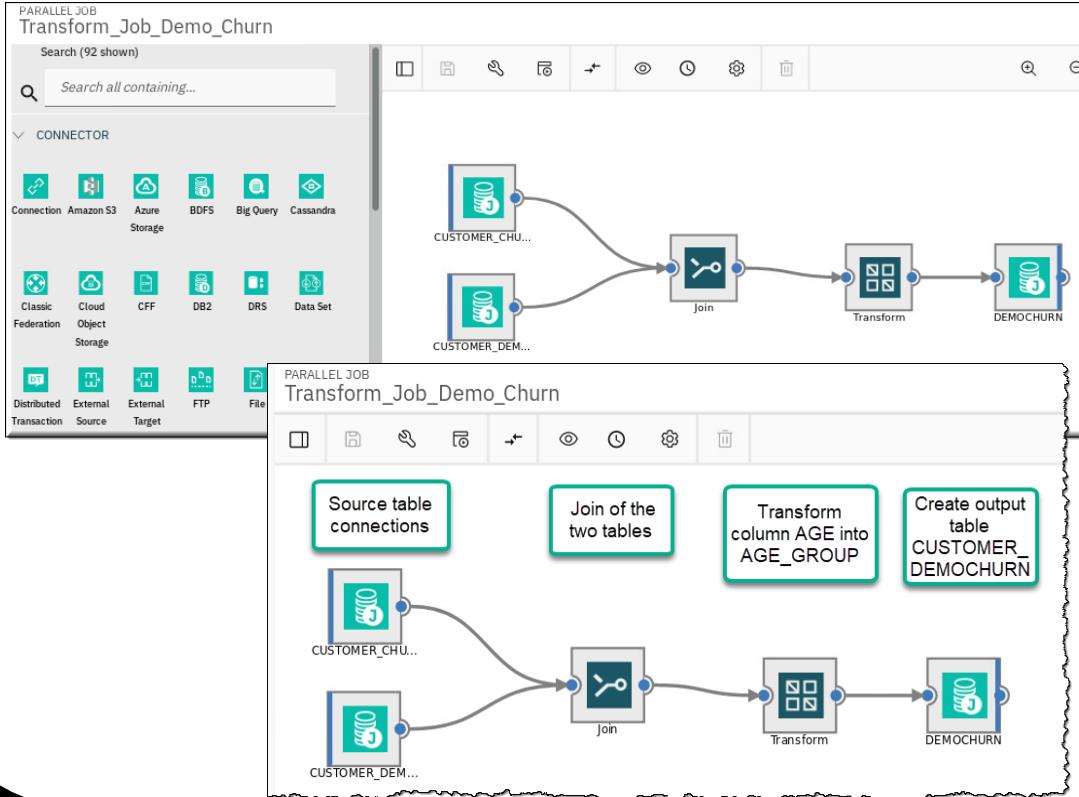
Relevancy of search results factors:

<i>Text match</i>	The provided text is searched in the asset name & description, where name contributes more to the higher place in the result list.
<i>Asset rating</i>	The higher the average rating the asset has, the higher it is on the results list.
<i>Comments</i>	The higher the number of comments, the higher the asset is on the results list.
<i>Context match</i>	The search results list might contain the closest neighbors of assets that are returned based on the text match.
<i>Modification date</i>	The assets that were modified recently are more likely to be returned in the search results.
<i>Quality score</i>	The higher the score, the higher the asset is on the results list. Quality score applies to database tables, views & columns, design tables, views & columns, data file records & fields.
<i>Usage</i>	The more relationships of type uses an asset has, the higher it is on the results list.



CPD Organize

DataStage - Transform and migrate data, build and execute ETL jobs at scale



- Use powerful data transformation capabilities
- ML infused Smart Job Clustering, Smart Job Assist and automated job sequencing
- Design once, run on any-cloud using Kafka
- Remotely execute job for co-located access to data, satisfy geopolitical requirements and save costs
- In-line data quality and governance to build trusted data when the data is being delivered to a target environment such as a data lake
- Create CI/CD pipelines with GitHub

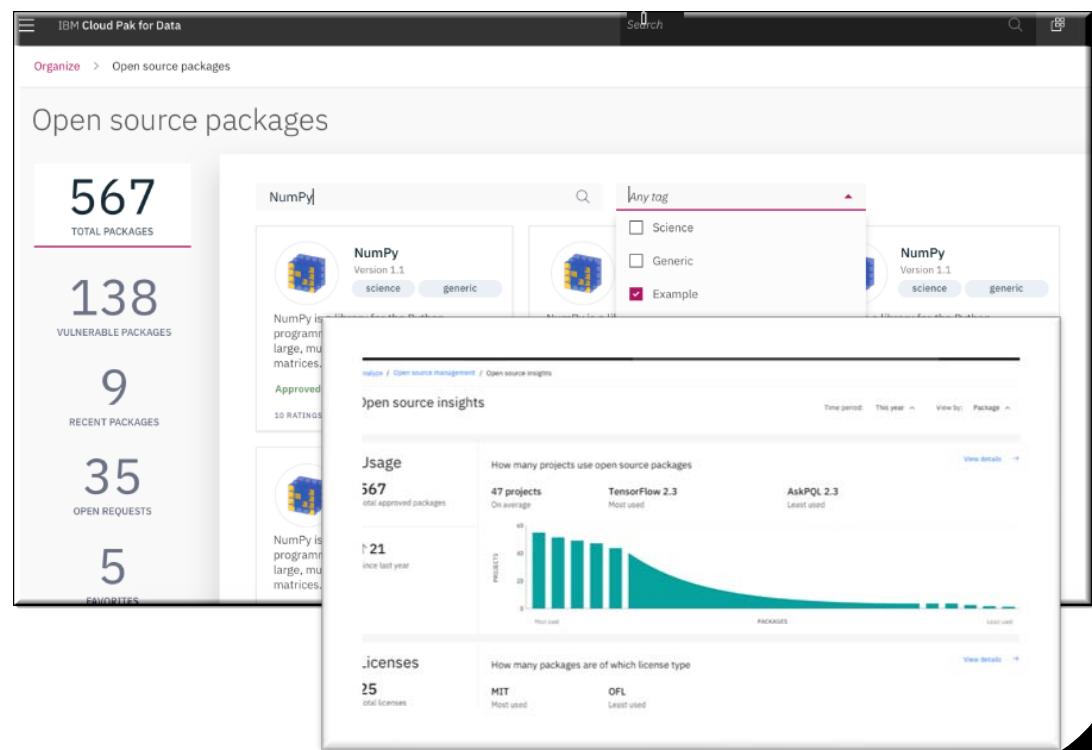


CPD Organize

Open Source Software Management

Open Source Software Management is:

- A *centralized inventory* of approved open source packages to ensure code quality & security
- A vehicle for *self-service open source consumption* and workflow requests for developers
- A *correlation of vulnerabilities* across open source portfolio: highlight security risks
- An *infusion of collaboration* by allowing developers to rate and comment on opensource packages
- A *set of dashboards* for CIOs to manage risks and accelerate innovation with OSS





CPD Organize

Master Data Connect

Master Data Connect is an extension of:

- ✓ IBM Master Data Management Advanced Edition
- ✓ IBM Master Data Management Standard Edition

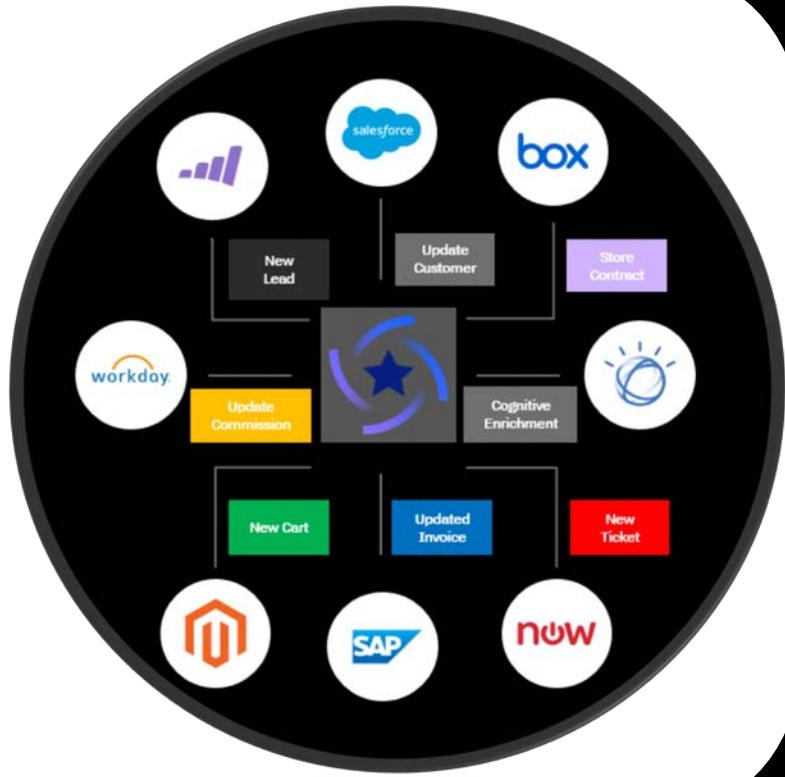
Services the majority of the MDM system workload - specifically reading MDM data.

Provides clients the opportunity to deploy a read-only instance of their master data entities on Cloud Pak for Data.

Majority (typically 80%-95%) of MDM workload is read/inquiry.

Mastered entities represent the highest quality data enterprises have about key domains like person, patient or organization.

Connect master data to your applications via RESTful APIs or IBM App Connect.



IBM Analytics Modernization Workshop

Part 2

- | | |
|---|--|
| <ul style="list-style-type: none">• Introduction• Business Use Case | <ul style="list-style-type: none">• Lab 01• Lab 02 |
| <ul style="list-style-type: none">• Collect: Connect• Organize• Collect: Virtualize | <ul style="list-style-type: none">• Lab 03• Lab 04• Lab 05 |
| <ul style="list-style-type: none">• Analyze• Deploy• Infuse – OpenScale• Infuse – Cognos Analytics• Wrap up | <ul style="list-style-type: none">• Lab 06• Lab 07• Lab 08• Lab 09• Lab 10 |

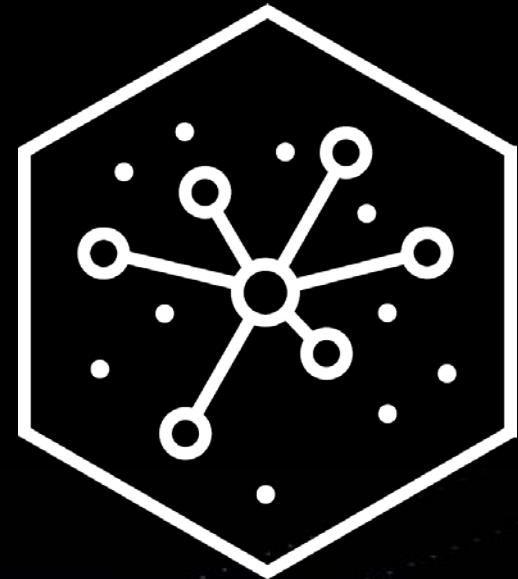
IBM Analytics Modernization Workshop

Part 3

- | | | |
|--|---|--|
| | <ul style="list-style-type: none">• Introduction• Business Use Case | <ul style="list-style-type: none">• Lab 01• Lab 02 |
| | <ul style="list-style-type: none">• Collect: Connect• Organize• Collect: Virtualize | <ul style="list-style-type: none">• Lab 03• Lab 04• Lab 05 |
| | <ul style="list-style-type: none">• Analyze• Deploy• Infuse – OpenScale• Infuse – Cognos Analytics• Wrap up | <ul style="list-style-type: none">• Lab 06• Lab 07• Lab 08• Lab 09• Lab 10 |

Analyze

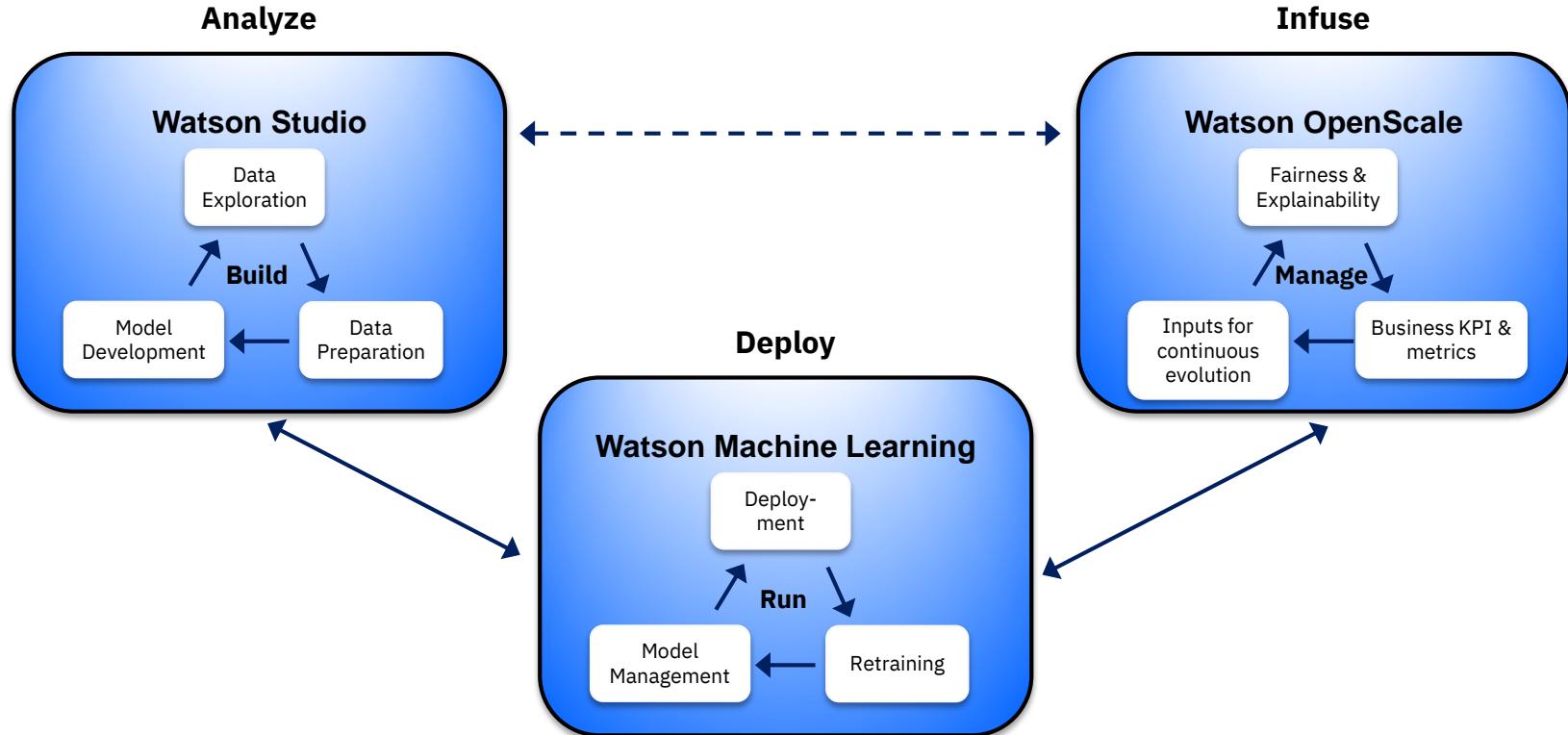
Lab 06 – Analyze: AutoAI & Notebooks





Analyze

The Data Science Lifecycle: Overview



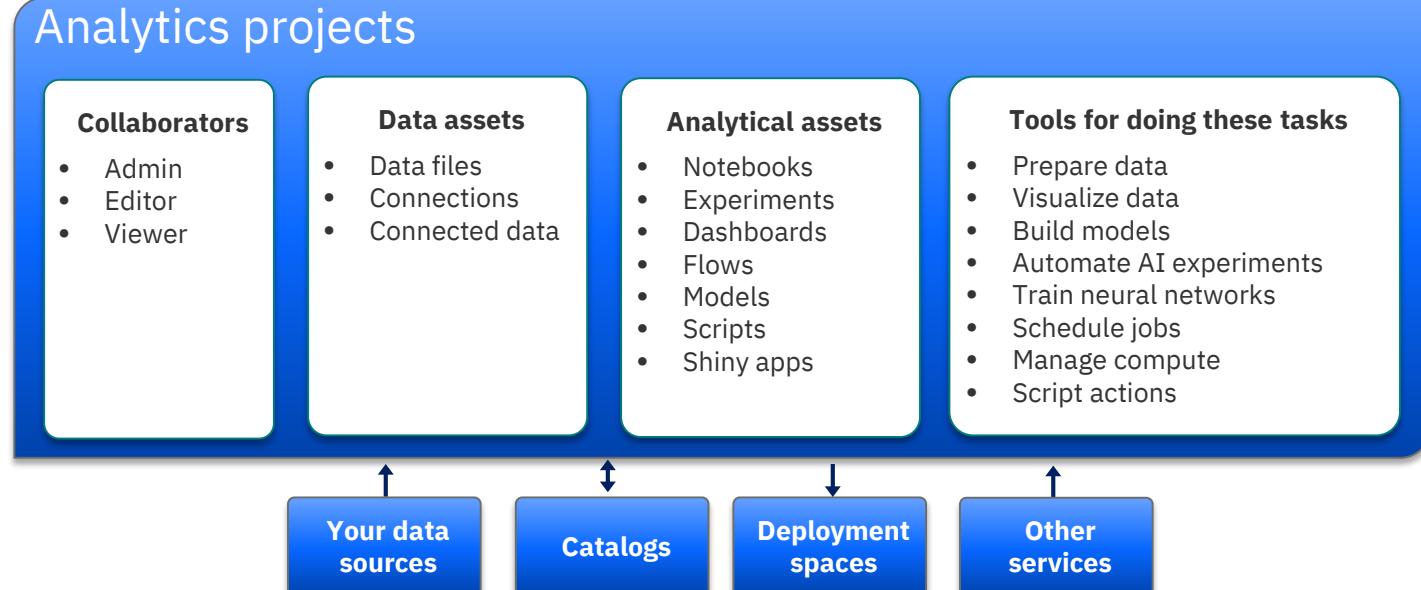


Analyze

Watson Studio: Collaborating with Analytics Projects

Watson Studio provides the environment and tools to collaborate on business problems.

Watson Studio is centered around the *Analytics Project*. Data scientists and business analysts use analytics projects to organize resources and analyze data with various tools.



Analyze

Watson Studio Integration



Watson Studio integrates with these services:

Watson Knowledge Catalog

- Easily move assets between projects and catalogs
- Catalogs and projects support the same types of data assets
- Data protection rules are enforced on catalog assets that you add to projects

Watson Machine Learning:

- You can easily move assets between analytics projects and deployment spaces

Watson Studio service includes these tools:

- Data Refinery
- Jupyter notebook editor
- JupyterLab IDE

Watson Studio projects can manage these separately installed service assets:

- Watson Machine Learning AutoAI experiments
- Watson Machine Learning Accelerator DL experiments
- Cognos Dashboards Embedded
- IBM Streams flows
- SPSS Modeler flows
- Decision Optimization models
- RStudio R Shiny apps

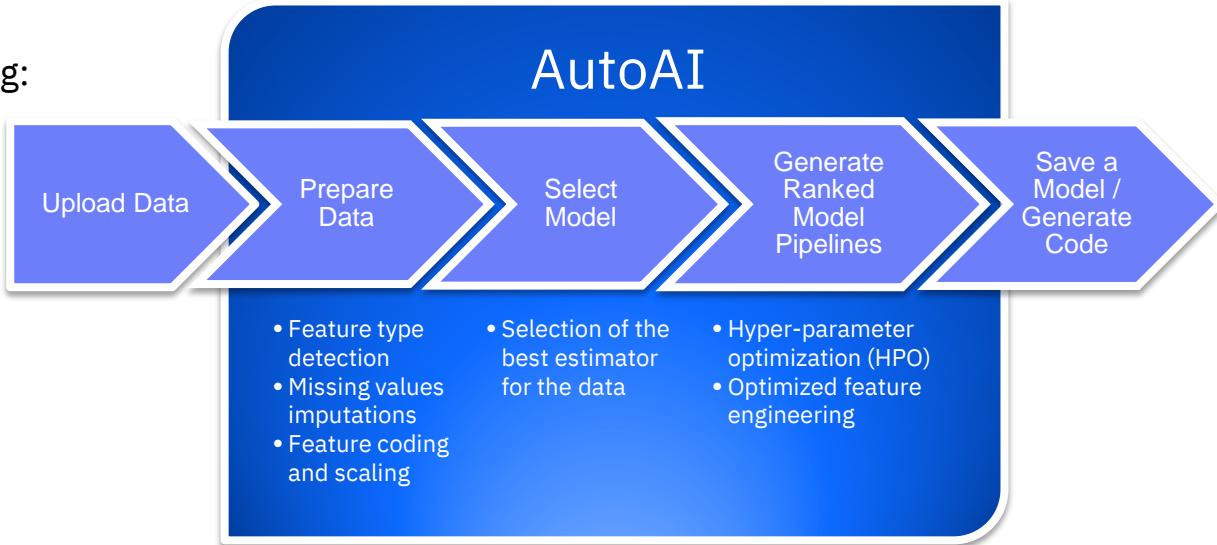


Analyze

AutoAI * – Overview

AutoAI is an award-winning technology that simplifies the Machine Learning model creation and AI lifecycle by automating the following:

- **Data preparation**
- **Model development**
- **Feature engineering**
- **Hyper-parameter optimization**



AutoAI delivers training feedback visualizations for real-time model performance results with:

- **Binary, Multiclass, and Regression support**
- **One-click model deployment**

* AutoAI is enabled with the Watson Machine Learning service install, but it is driven through a Watson Studio Analytics Project

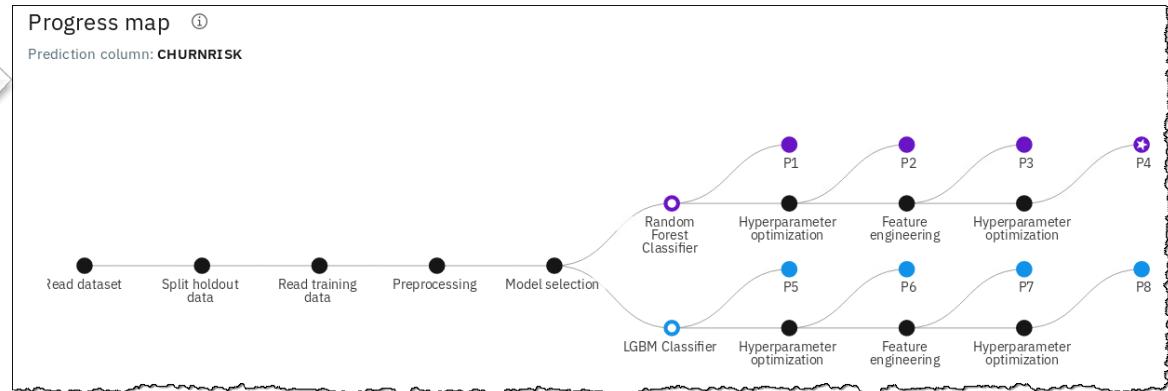


Analyze

AutoAI – Infographics

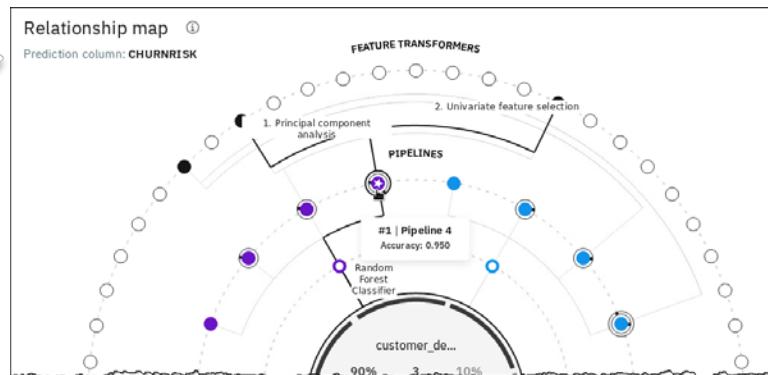
AutoAI Progress map

Displays a progress each step as it creates the best model for your data.



AutoAI Relationship map

Interactive infographic that shows the relationship of the pipelines, algorithms and the feature transformers.





Analyze

AutoAI – Pipelines

AutoAI pipeline leaderboard

Shows the ranking of the pipelines for each potential model, the higher the better.

Pipeline leaderboard

Rank ↑	Name	Algorithm	Accuracy (Optimiz...)	Enhancements
★ 1	Pipeline 4	Random Forest Classifier	0.950	HPO-1 FE HPO-2
2	Pipeline 8	LGBM Classifier	0.949	HPO-1 FE HPO-2
3	Pipeline 7	LGBM Classifier	0.946	HPO-1 FE

After AutoAI completes its model creation steps, you can drill into the pipeline(s) to understand how it came to its conclusion.

Save the pipeline in your project as a:

- **model**
- **notebook**





Analyze

AutoAI – Benefits

1 Speed Model Selection

Shortlist top performing models in minutes instead of days/weeks

Drastically reduce neural network search time

2 Jump the Skills Gap

Go live with better models using the skill sets you have

Increase repeatability and minimize human intervention

3 Drive Productivity

Get started with AI experiments without knowing how to code

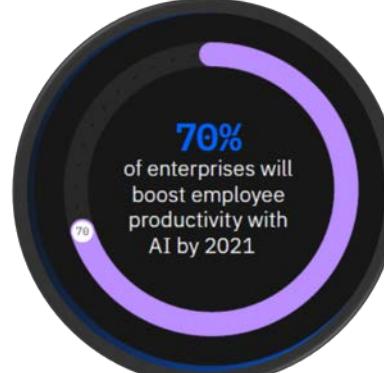
Do more innovative work instead of mundane tasks (e.g. lengthy feature selection process)



80%
of data scientists' time is spent on tedious tasks that could be automated



63%
of companies see availability of technical skills as a challenge to implementation



70%
of enterprises will boost employee productivity with AI by 2021



Analyze

Notebooks, RStudio and other tools

The default notebook environment:
Jupyter Notebook with Python 3.6

TradingCustomerChurn > Notebooks > 01TradingCustomerChurnClassifierSparkML

File Edit View Insert Cell Kernel Widgets Help

Trading Platform Customer Attrition Risk Prediction using SparkML

There are many users of online trading platforms and these companies would like to run analytics on and predict churn based on user activity elsewhere is key to maintaining profitability.

In this notebook, we will leverage Watson Studio Local (build into IBM Cloud Private for Data) to do the following:

1. Ingest merged customer demographics and trading activity data
2. Visualize merged dataset and get better understanding of data to build hypotheses for prediction
3. Leverage SparkML library to build classification model that predicts whether customer has propensity to churn

Developer tool services available:

- Jupyter Notebooks with Python 3.6 for GPU
- Jupyter Notebooks with R 3.6
- RStudio Server with R3.6
- Lightbend Platform
- OpenSource Management



 Jupyter Notebooks with Python 3.6 for GPU
Open Source

Optional development environment to create Jupyter Notebooks that use GPU-accelerated Python 3.6 libraries.

 Jupyter Notebooks with R 3.6
Open Source

Optional development environment to create Jupyter Notebooks that use R 3.6 libraries.

 Lightbend Platform
Partner Premium

Lightbend Platform makes it easy to deploy Reactive Microservices, real-time streaming and Machine Learning (ML).

 Open Source Management
IBM

Make it easy for developers and data scientists to find and access approved open source packages.

 RStudio Server with R3.6
Partner Enabled

Optional development environment for working with R.



Analyze

Watson Studio notebooks: Build Data Science & Machine Learning models

We split original dataset into train and test datasets. We fit the pipeline to training data and apply the trained model to transform test data and generate churn risk class prediction

```
In [67]: # instantiate a random forest classifier, take the default settings
rf=RandomForestClassifier(labelCol="label", featuresCol="features")

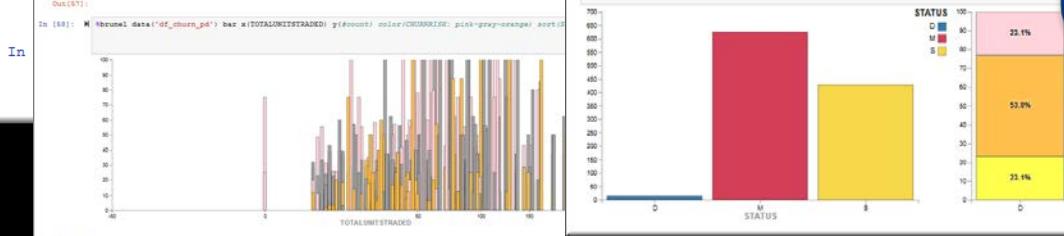
# Convert indexed labels back to original labels.
labelConverter = IndexToString(inputCol="prediction", outputCol="predictedLabel", labels=labelIndexer.labels)

stages += [labelIndexer, assembler, rf, labelConverter]

pipeline = Pipeline(stages = stages)

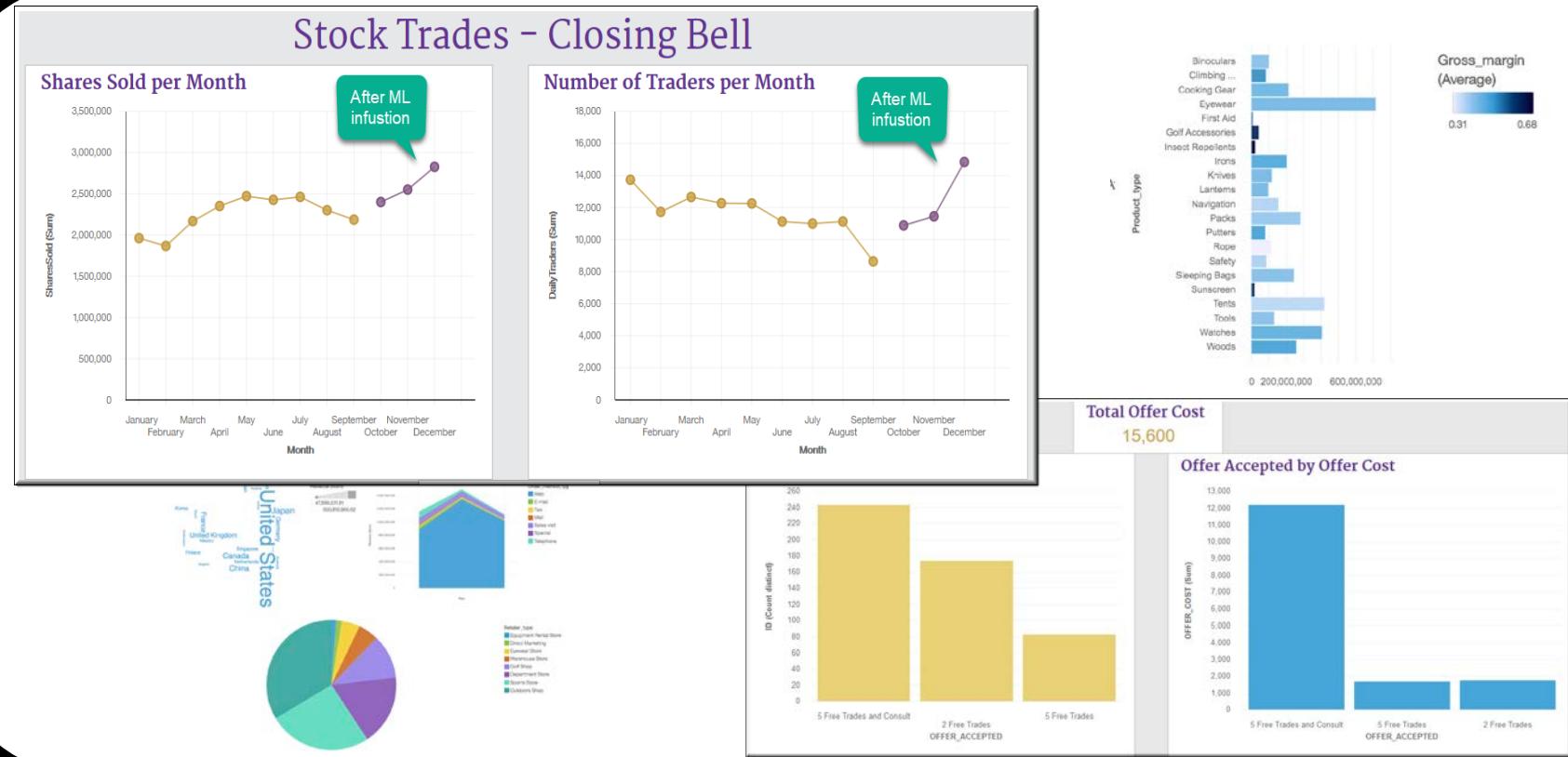
In [68]: # Split data into train and test datasets
train, test = df_churn.randomSplit([0.7,0.3], seed=100)
train.cache()
test.cache()
```

```
Out[68]: DataFrame[AGE: int, AGE_GROUP: string, CHILDREN: int, CHURNRISK: string, ESTINCOME: int, GENDER: string, INCOME: float, TOTALDODS: int, TOTALUNITSTRADED: int, STATUS: string, SMALLESTTRANSACTION: int, TOTALDOL: float]
```



- Data Scientists and Data Engineers collaborate with each other in CPD platform – while still maintaining data governance
- Collaboration using GitHub or BitBucket is integrated into the platform, which brings a cohesiveness to the work culture and helps to automate CI/CD pipe line
- Exploit GPUs for deep learning predictive ML models
- Programmatically build data visualizations and data wrangling
- Real-time or batch model scoring
- Evaluate model accuracy

Analyze Cognos Dashboards Embedded





Analyze IBM Streams

IBM Streams has built-in streaming analytics

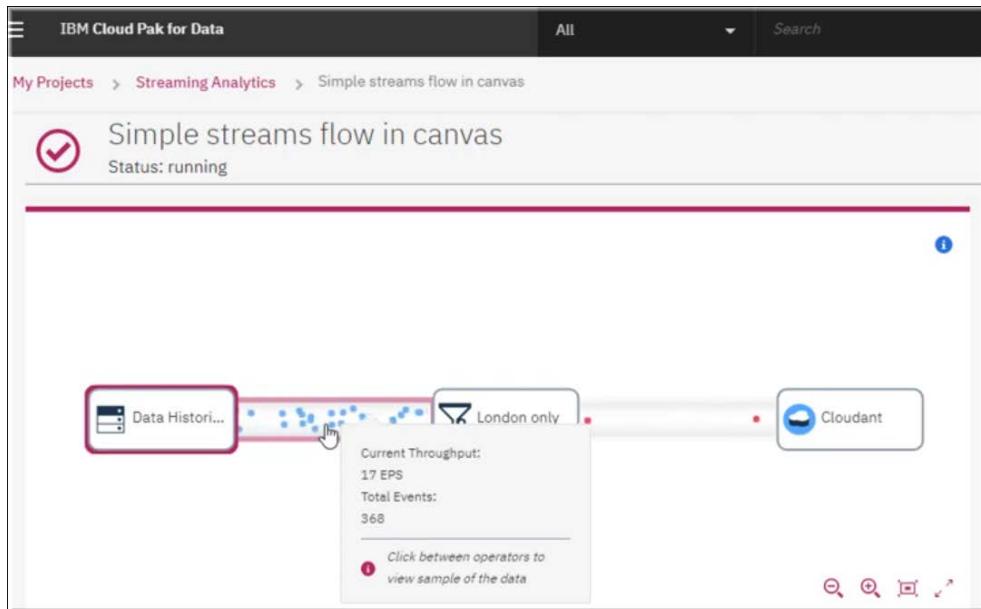
- No business disruption—run, score & update models continuously
- Machine learning, natural language, spatial-temporal, acoustic time series, etc.

Open architecture built for speed

- Millions of events per second for massive amounts of data analytics support
- Ultra-low latency clustered runtime
- Integrate via Kafka, JSON SQL/NoSQL & more

Rapid development

- Wizards, drag/drop development, performance dashboards, debugger
- Python, Java, Scala, PMML, R, C/C++ support
- VS Code and Atom plug-ins
- Export flows to a Python notebook



A streams flow consists of *operators*.

Every node on the streams flow canvas is an operator.

Operator types include: sources, targets, data processing, alerts and real-time analytics



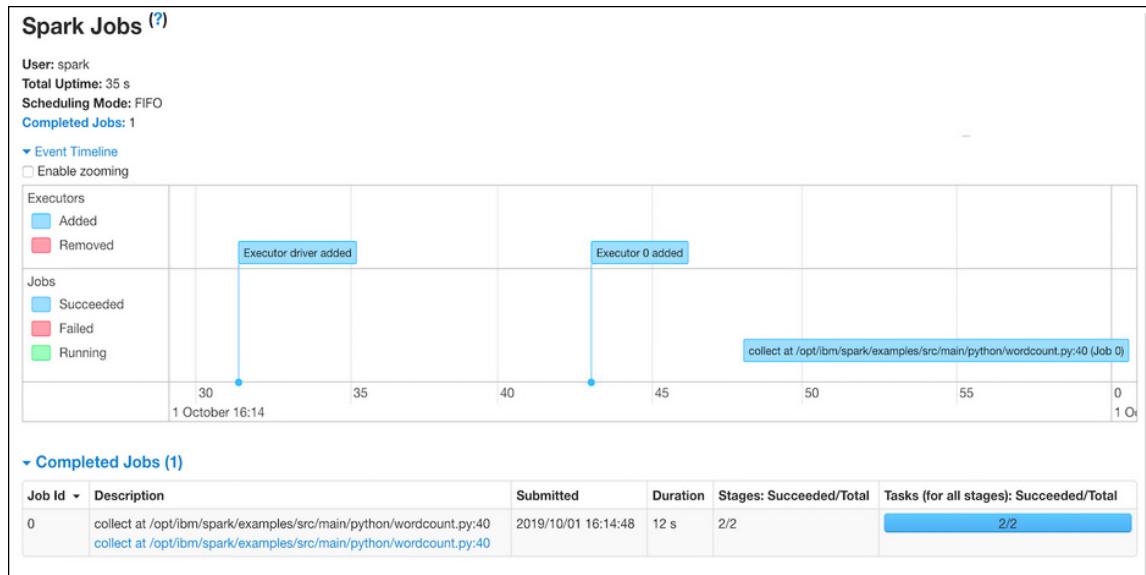
Analyze

Analytics Engine - powered by Apache Spark

Analytics Engine is a serverless, performant, customizable and dedicated Spark engine that is available in seconds.

100% open source, Analytics Engine can run a variety of workloads on the CPD cluster:

- Watson Studio notebooks that call Apache Spark APIs
- Spark application that run Spark SQL
- Data transformation jobs
- Data Science jobs
- Machine Learning jobs





Analyze

Premium Service: SPSS Modeler

The screenshot displays the IBM Watson Data Studio interface. On the left, a sidebar lists various modeling techniques: Association Rules, Auto Classifier, Auto Numeric, C5.0, C&R Tree, CHAID, GLE, Linear, Linear-AS, and Linear SVM. The main workspace shows a flowchart of data mining steps: UCI ML R... → Filter → Select → Data Audit → Partition → Target D... → Decision ... → Decision ... → Table. A central chart panel shows a histogram with a normal distribution curve overlaid. Below it is a spreadsheet view of a dataset with columns: Age, Sex, BP, Cholesterol, Na, K. The data rows are:

	Age	Sex	BP	Cholesterol	Na	K
1	23	F	HIGH	HIGH	0.79355	0.031258
2	47	M	LOW	HIGH	0.739309	0.056946
3	47	M	LOW	HIGH	0.697269	0.068944
4	28	F	NORMAL	HIGH	0.563682	0.072289
5	61	F	LOW	HIGH	0.659294	0.030969
6	22	F	NORMAL	HIGH	0.678901	0.078641
7	49	F	NORMAL	HIGH	0.789637	0.048588
8	41	M	LOW	HIGH	0.766335	0.069461
9	60	M	NORMAL	HIGH	0.777205	0.05123
10	43	M	LOW	NORMAL	0.526102	0.027164

A network diagram panel on the right shows relationships between variables like Age, BP, and Cholesterol.

SPSS Modeler

- A leading visual data science and machine-learning and predictive analytics solution
- Helps enterprises accelerate time to value and achieve desired outcomes by speeding up operational tasks for data scientists and business analysts
- Tap into data assets and modern applications, with complete algorithms and models that are ready for immediate use

Analyze

Premium Service: Decision Optimization



Decision Optimization (DO) enables data science teams to capitalize on the power of *prescriptive analytics* and build solutions using a combination of techniques like optimization and machine learning.

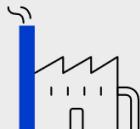
Integrated with Watson Studio, Decision Optimization can combine optimization techniques with coding and non-coding tools, model management and deployment – as well as other data science capabilities.

Decision Optimization evaluates millions of possibilities – balancing trade-offs and business constraints to find the best possible solution.

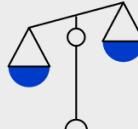
Insights that drive optimal decisions to complex problems



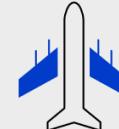
Determine location
and capacity
of warehouses



Determine which plant
should manufacture
which product



Build financial
portfolios by balancing
risks and rewards



Allocate aircraft
and crew to flights



Analyze

Watson Machine Learning : WML Accelerator for Deep Learning

The **Experiment Builder** GUI interface is available from a WML project when you install the **WML Accelerator**. It is the simplest method to perform Deep Learning experiments.

The Watson Machine Learning Accelerator has the following components:

- **IBM Spectrum Conductor Deep Learning Impact version 2.1.0:**
Provides robust, end-to-end workflow support for deep learning application logic for the complete lifecycle management from data ingest and preparation to building, optimizing, training and testing the model.
- **IBM Spectrum Conductor version 2.4.0:**
A *highly available* and resilient *multitenant distributed* framework, providing deep learning application lifecycle support, centralized management and monitoring, and end-to-end security.
- **Deep Learning Frameworks:**
TensorFlow, PyTorch, Keras, and Caffe
- **Deep Learning Rest API**

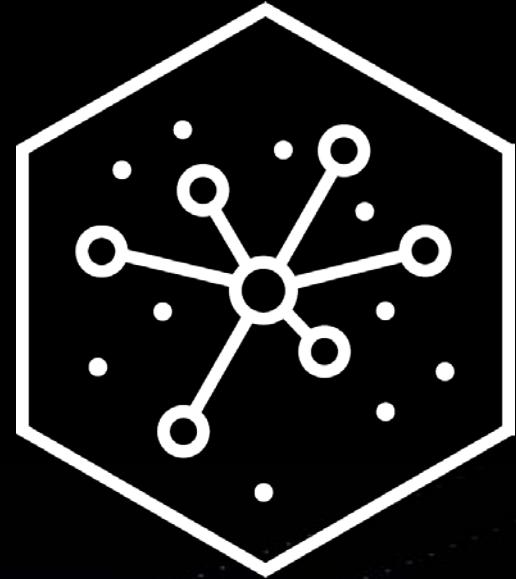


Key Features

- Supported AI frameworks with performance optimizations for GPU acceleration
- Transparent GPU Topology for Distributed Training
- Auto Hyper Parameter Optimization (HPO)
- Multi-Tenancy for Training and Inference
- Elastic Distributed Training (EDT)
- Resource Utilization, Monitoring, and Reporting
- Elastic Distributed Inference (EDI)
- Secure Deployment Model
- Role Based Access Control / Kerberos Authentication

Deploy

Lab 07 – Deploy





Deploy

Watson Machine Learning: Deployment Spaces

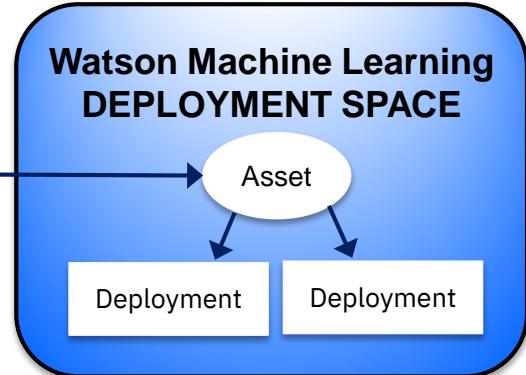
A **Deployment Space** is where you can:

- Promote and save models
- Create the deployments from the models
- Find the information you need to score the model and get a prediction back
- Embed the deployment in an app so you can interact with it programmatically

The screenshot shows two views of the Watson Studio interface. The top view is a navigation bar with tabs: Assets (highlighted with a green checkmark), Deployments, Access control, and Settings. Below this, under the 'Models' section, there is a table with one row: Name (churn_risk_model) and Type (scikit-learn). The bottom view is a detailed 'Deployments' page with tabs: Assets, Deployments (highlighted with a green checkmark), Access control, and Settings. Under the 'Deployments' tab, there is a table with one row: Name (churn_risk_model-deployment).



Prepare data and train



Configure, test, deploy and monitor



Deploy

Watson Machine Learning: Deployments

A **Deployment** is the last stage of the model development work. It means you put the model into production so that you can pass data to the model and return a score (or prediction).

After deploying a model, you can access the model *endpoint*, which you will need to make the model available for wider use in applications.

There are three type of WML deployments:

- **Online** – Provides an API endpoint needed to access the deployment programmatically to use in an application. Code snippets are provided in a variety of programming languages that illustrate how to access the deployment.
- **Batch** – Processes input data from a file and writes the output to a file.
- **Virtual** – For models downloaded into WML from different frameworks other than those built on the CPD platform.



Deploy

Watson Machine Learning: Supported Frameworks for deployments

Framework	Version	Framework	Version
Spark	2.3, 2.4	Tensorflow	1.15
PMML	3.0 to 4.3	Keras	2.2.5
Hybrid/AutoML	0.1	Caffe	1.0
SPSS	17.1, 18.1, 18.2	PyTorch	1.0, 1.1, 1.2
Scikit-learn	0.20, 0.22	Decision Optimization	12.9, 12.10
XGBoost 0.82	xgboost_0.82	Python Functions	0.1
XGBoost 0.90	xgboost_0.90	Python Scripts	1.0



Deploy

Watson Machine Learning: Online deployment testing

Overview Implementation **Test ✓** Lineage

Enter input data

age
41

job
accountant

marital
married

education

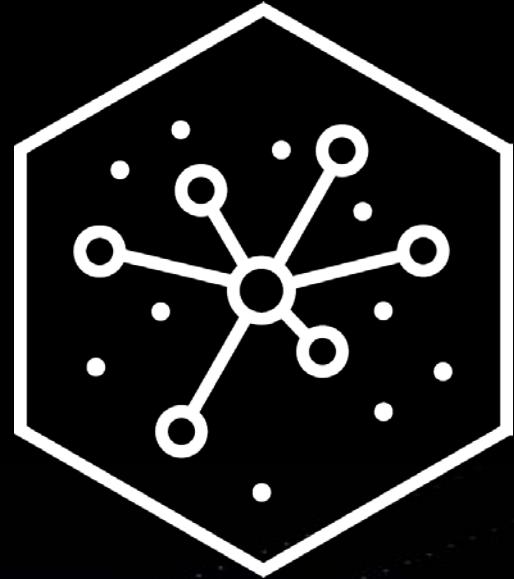
Predict ✓

{
 "predictions": [
 {
 "fields": [
 "prediction",
 "probability"
],
 "values": [
 [
 "no",
 [
 0.9984213709831238,
 0.0015786453150212765
]
]
]
]
]
}

- You can use the *form* in the *Deployment Space* to easily test an online deployment.
- On the *Test* tab of the Deployment details page, enter test data, and click *Predict* to see the result.

Infuse

*Lab 08 – Infuse: Watson OpenScale
Lab 09 – Infuse: Cognos Analytics*





Infuse

Watson OpenScale: Overview

Watson OpenScale:

- Automates and operates AI at scale across its entire lifecycle
- Delivers transparent, explainable outcomes freed from bias and drift
- Provides confidence in AI outcomes and spans the gap between the teams that operate AI and the business units that use these applications
- Monitors models developed in a 3rd party IDE, open source framework and hosted in a 3rd party or private model serve engine

Manage AI at Scale



Model build / train frameworks



Model serving environments





Infuse

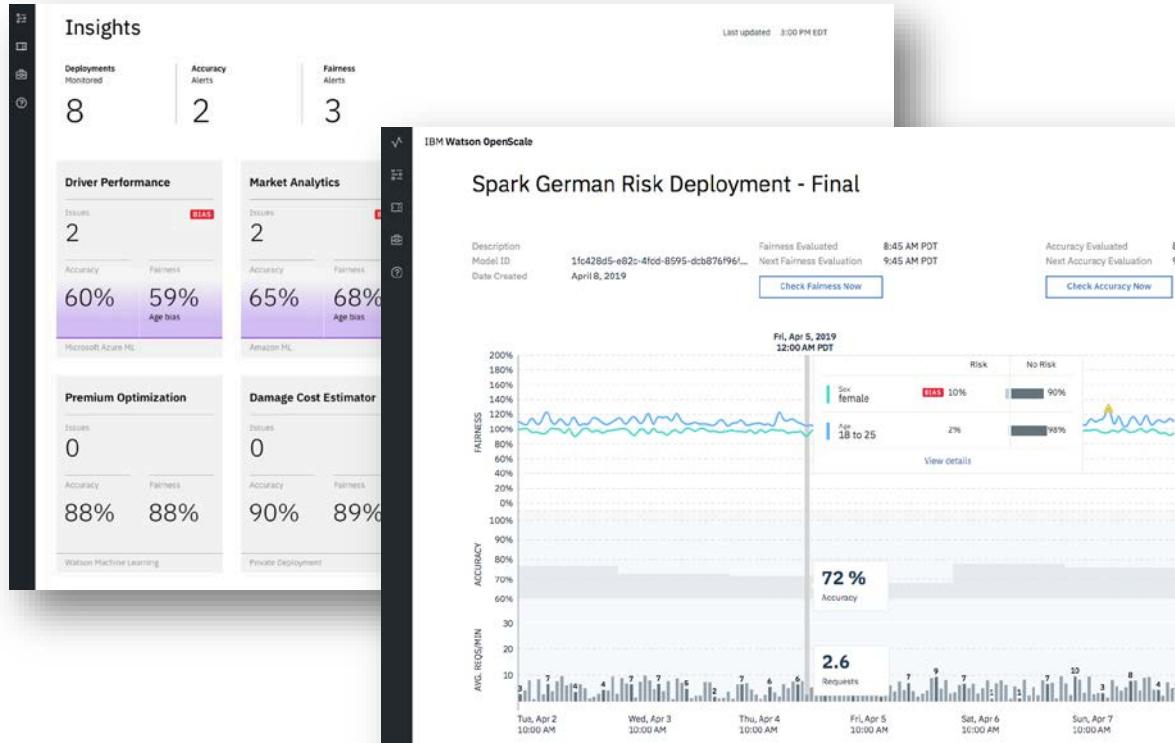
Watson OpenScale: Operations dashboard

Description:

Monitor deployed models in a single dashboard that can be filtered by deployment making it easy to manage AI in apps

Value:

- Configure alerts or actions to be triggered when KPIs exceed threshold, ensuring model quality for improve business outcomes
- Measure model accuracy as it pertains to its ability to deliver outcomes more accurate than knowledge workers
- Provides “continuous evolution” for your models





Infuse

Watson OpenScale: Model Fairness

Description:

Production Models need to make fair decisions and *can not be biased* in their recommendations

How it works:

- Outcomes are selected as “favorable or unfavorable”
- “Favored Populations” and “protected populations” are selected where majority and minority groups are found
- A score is calculated based on the probability of favorable outcome for minority vs. probability of favorable outcome for majority

The screenshot shows two sequential steps in the Watson OpenScale interface:

Step 1: Select the features to monitor

This step displays a grid of features for monitoring. The "Sex" feature is highlighted with a blue border. Other visible features include CheckingStatus, LoanDuration, CreditHistory, ExistingSavings, EmploymentDuration, InstallmentPercent, CurrentResidencDuration, OwnsProperty, Age, and InstallmentAmount.

Step 2: Specify the favorable outcomes

This step allows users to define favorable and unfavorable outcomes. Under "Favorable values", "No Risk" is listed. Under "Unfavorable values", "Risk" is listed.



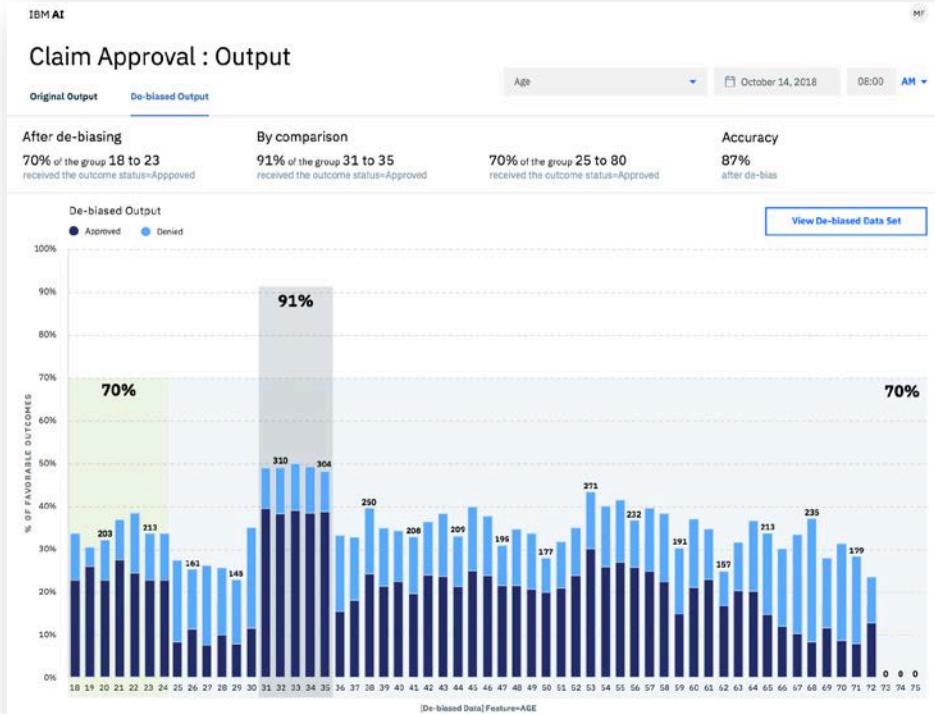
Watson OpenScale: Bias Mitigation

Description:

Fairness is enforced with automatic bias mitigation.

How it works:

- Calculated on an *hourly basis* (over a sliding window defined by the user)
- Optimizations identify the *right subset of data to perturb* (rather than perturbing all the data)
- Perturbed data is sent to the deployed model* to determine effect of perturbations
- An internal bias detection model (logistic regression) is built using perturbed data that *classifies whether new prediction will be biased or not*
- Users receive both the *original prediction* plus the *internal model's classification* of whether the monitored model's prediction is biased or not





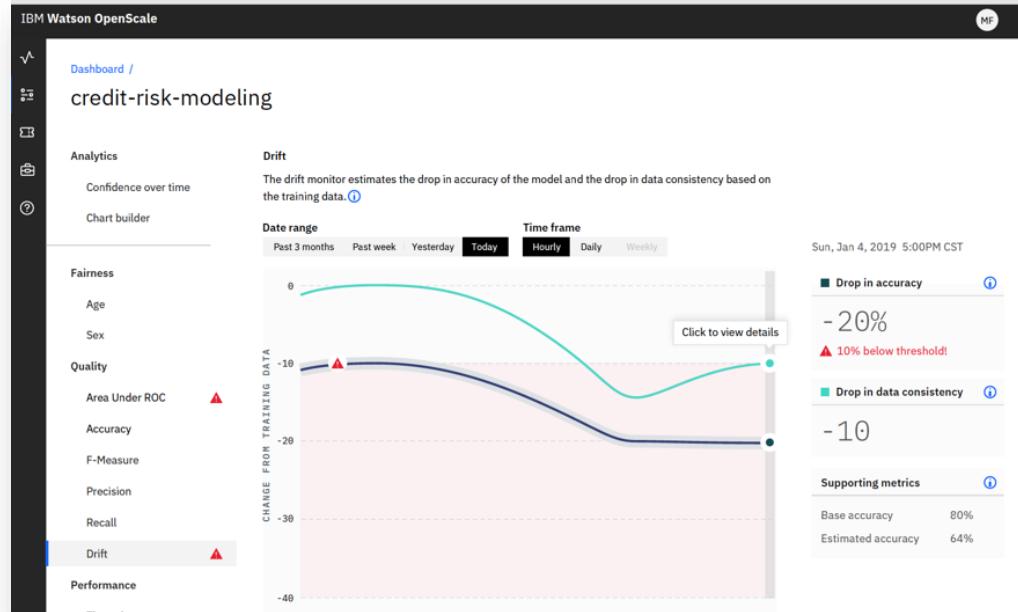
Infuse

Watson OpenScale: Drift detection

Description:

OpenScale monitors for two types of drift:

- **Drop in accuracy:** It estimates the drop in accuracy of the model at runtime. Accuracy could drop if there is an increase in transactions similar to those which the model was unable to evaluate correctly with the training data.
- **Drop in data consistency:** It estimates the drop in consistency of the data at runtime as compared to the characteristics of the data at training time.



OpenScale does drift detection on the entire payload data.

OpenScale measures the drift without requiring labeled data. Accuracy computation using labeled data can be expensive and might not be comprehensive



Infuse

Watson OpenScale: Explainability

Description:

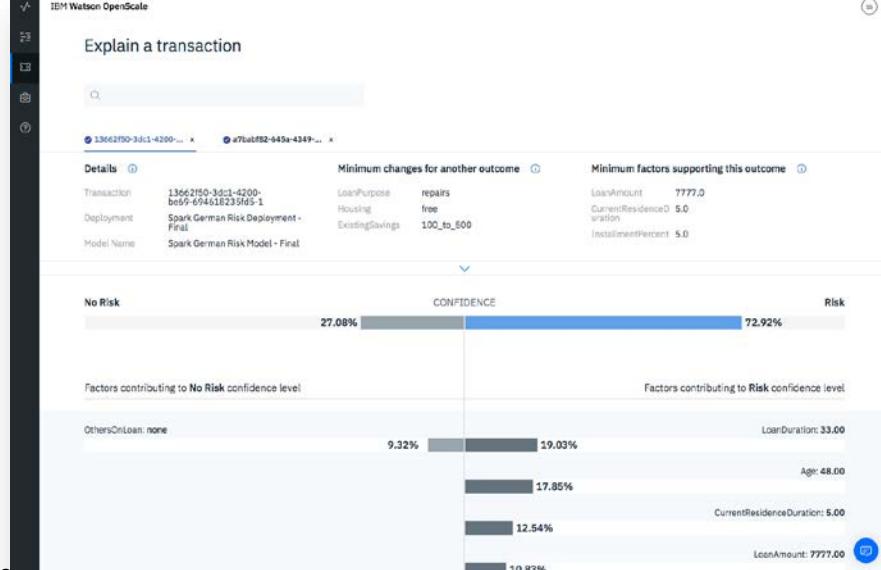
Allows you to understand which feature values of a model that are most influencing a prediction for a specific transaction

Example:

A loan is not approved by a model prediction - explainability will tell you why

How it works:

- Perturbation analysis on thousands of variations
- Risk model is created for two variations:
 - **LIME (local) Explanation:** set of features which played a positive or negative role in the prediction - also identifies the feature weights which helps to identify the most or least important features
 - **Contrastive Explanation:** Explains the behavior of the model in the vicinity of the data point whose explanation is being generated – assumption: the most common value is the least interesting from an explanation point of view





Infuse

Watson OpenScale: Business (Application) KPIs

Description:

The OpenScale UI dashboard contains an Application centric dashboard in addition to the Model centric dashboard:

- Bring in business event data into OpenScale and use it to compute Key Performance Indicators (KPIs)
- Monitor correlation between model monitors and KPIs - get alerts and recommendation from OpenScale
- Visualize the correlation through a plot

The screenshot illustrates the Watson OpenScale UI with four numbered callouts:

1. A browser window showing the URL aiopenscale.cloud.ibm.com/aione.
2. The "Credit Risk Application" dashboard, which includes sections for "Associated Models", "Event Details", "KPIs", and "Logging Endpoint". It also features a "Add Associated Models" button and a section titled "Associated Models Deployments".
3. The "Insights Dashboard" showing "Application Monitors beta" and "Model 1". Below this, a summary states: "German credit risk compliant deployment" and "Drift : Drop in accuracy".
4. A plot titled "Program Impact & Credit Risk Model Drift" showing the relationship between Accepted Credits (Y-axis) and CREDIT RISK MODEL DRIFT (X-axis). The plot includes a legend for correlation levels: "Strong correlation", "Some correlation", and "No correlation". A note indicates: "When drift magnitude rises by 0.88%, accepted credits falls by 5.27" and "Large correlation | V".



Infuse

Watson OpenScale: Payload logging

Description:

Payload logs capture (in Postgres) the request sent to the model or python function along with statistics about its health which when combined with feedback data provides insights into AI and application behavior

Value:

- Enable logging of payloads for *traceability* of business outcomes to AI recommendations
- Payload *data powers visualizations* for the OpenScale dashboard, making it easy to monitor health of deployed models
- Payload *fuels the Open Datamart* for custom reporting and business KPI integration

CARS4U - Postgres Connection	public
Schemas (1)	Tables (4)
▶ public	> cars4u_action_recommendati...
▶	cars4u_business_and_action_...
	cars4u_business_area_predict...
	cars4u_satisfaction_predictio...



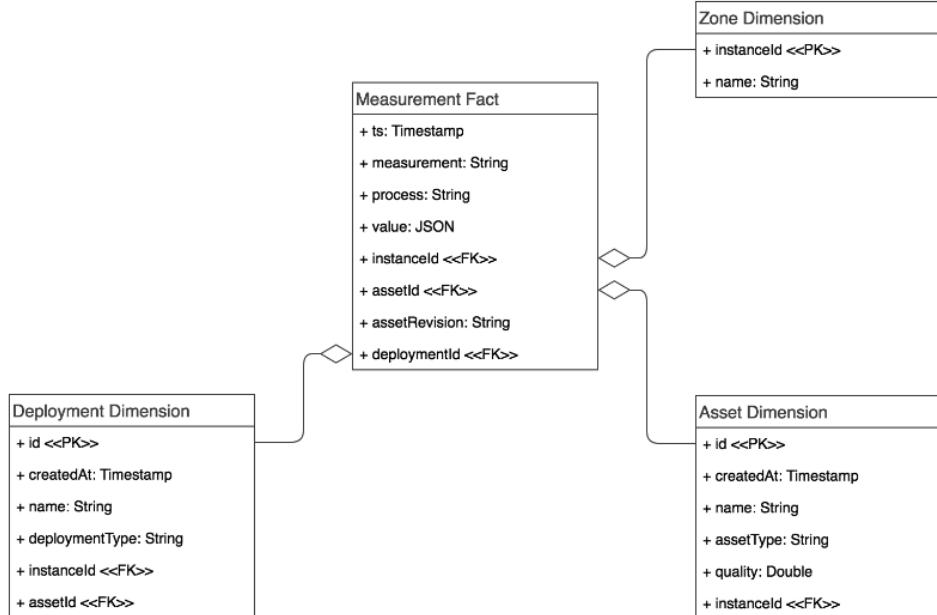
Infuse

Watson OpenScale: Data mart

Description:

The Data Mart is a data system to support the following use cases:

- Enable *OpenScale UI dashboards* for operations staff
- Enable *3rd party reporting tools* so operations staff can customize their own dashboards using the underlying OpenScale data with their own tools
- Enable *data engineers to integrate OpenScale data* with existing enterprise data marts, warehouses, and data lakes by exporting data into those systems
- Provide *data scientists with the actual runtime payload, scoring and feedback data* that can be utilized in continuous learning





Infuse

Premium Service: Cognos Analytics

Cognos Analytics is self-service analytics, infused with AI and machine learning.

- Enables you to create stunning visualizations to share your findings through *dashboards* and *reports*
- These can be embedded (infused) into your applications
- The Cognos Analytics service makes it easier for you to extract meaning from your data with features such as:
 - **Automated data preparation**
 - **Automated modeling**
 - **Automated creation of visualizations and dashboards**
 - **Data exploration**



Infuse

Premium Service: Planning Analytics



Use Planning Analytics to:

- Automatically create plans, budgets and forecasts
- Steer business performance by bridging operations and finance for any department allowing you to adapt to changing business conditions
- See impact before executing – explore what-if scenarios and assess impact to determine the best course of action
- Make changes in real-time – pivot plans, budgets, and forecasts quickly to meet changing demands and priorities

Differentiators:

- Adjust financial plans in real time across departments
- Protect your investment in Microsoft Excel while transcending limitations of spreadsheets
- Uncover deep insights through AI-infused planning, without the need for help from a data scientist

A screenshot of the IBM Cloud Pak for Data interface. At the top, there's a navigation bar with icons for dashboard, search, and user profile. Below it, the main area shows the "Services catalog / Planning Analytics". On the left, there's a sidebar with a blue icon of three horizontal lines with circles at the ends, followed by the text "Planning Analytics", "IBM", "Enabled ✓", and "Premium". To the right of the sidebar, there's a large image showing various charts and data tables related to budgeting and forecasting. Below the image, there's a "Description" section with text about creating good plans based on data from across the business using IBM Planning Analytics powered by TM1®. There's also a detailed description of Planning Analytics as an AI-infused solution that integrates with Microsoft Excel to build sophisticated multidimensional models.

Infuse

Premium Service: Watson Assistant



Watson Assistant enables you to build conversational interfaces (chat bots) into any application, device, or channel

- Intuitive tooling for dialog building
- Intent recommendations from chat logs for continuous improvements
- Digression handling when user changes topics mid-conversation
- Out-of-the-box integration with search capabilities when coupled with Watson Discovery
- Common misspellings are automatically handled by Watson Assistant.
- Contextual entities: Support for additional languages not yet in Cloud Pak
- Fuzzy matching: Support for additional languages not yet in Cloud Pak

Benefits:

- Get up and running quickly
- Easily improve customer experience
- Reduce costs and speed of resolution
- Increase value by integrating apps and channels
- Ensure security, resiliency and data privacy - anywhere



Infuse

Premium Service: Watson Discovery

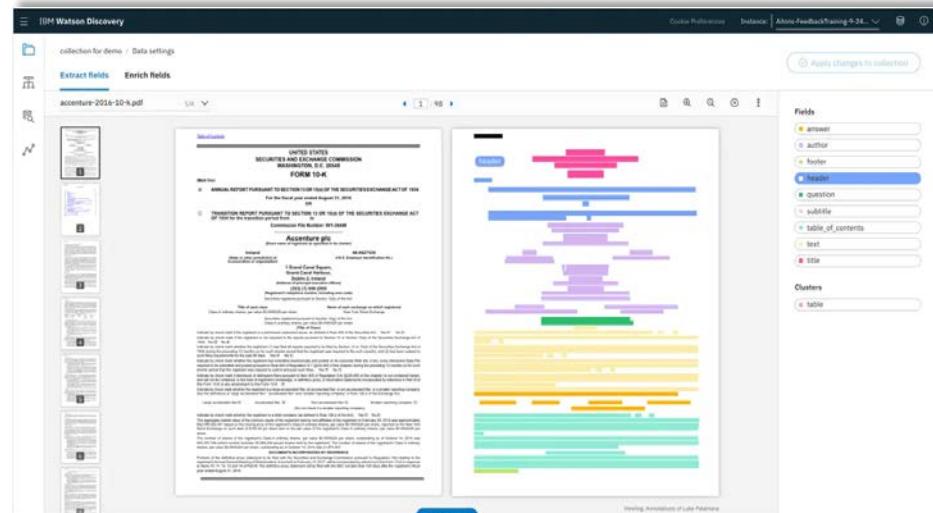


Watson Discovery Surface answers and rich insights from your enterprise data.

- Award-winning enterprise search and AI search technology that breaks open data silos and retrieves specific answers to your questions while analyzing trends and relationships buried in enterprise data.
- Applies the latest breakthroughs in machine learning, including natural language processing capabilities, and is easily trained on the language of your domain.

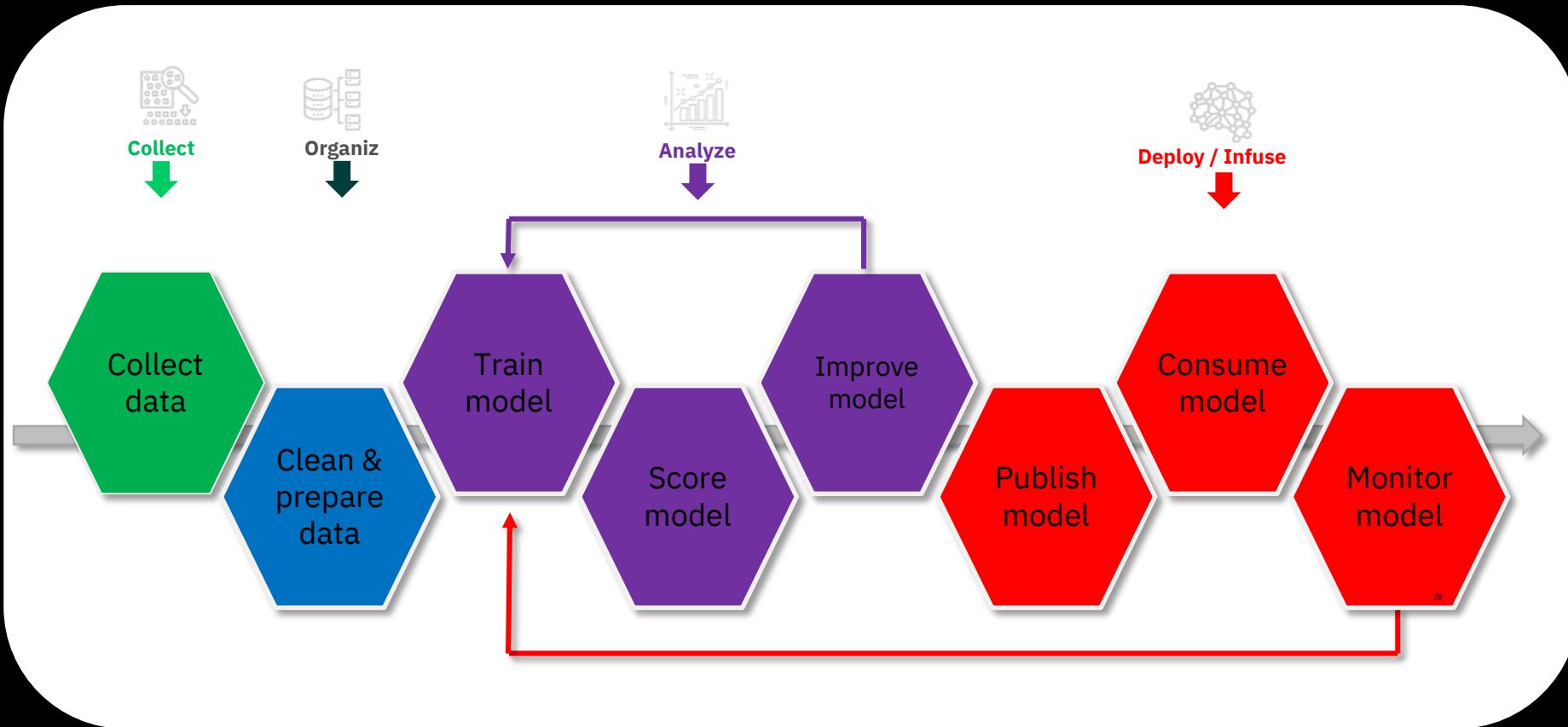
Benefits:

- Traditional enterprise search engines don't provide exact answers because they can't understand the nuances of phrases and acronyms in your industry and accurately search through your complex documents in a timely manner.
- Delivers specific answers to your queries while also serving up the entire document and supporting links, allowing your employees and customers to make informed decisions with confidence.



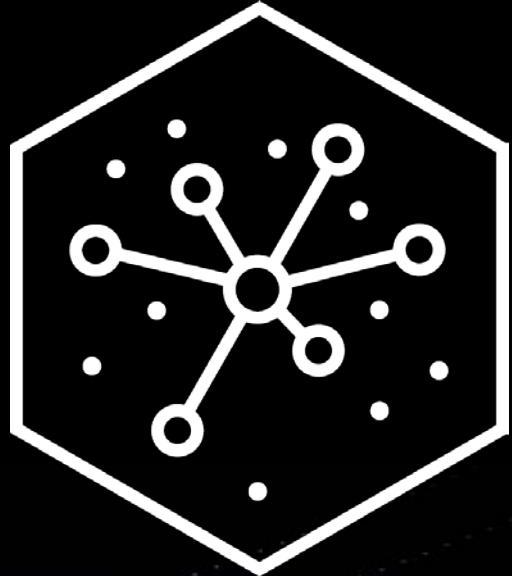
Machine Learning Model Lifecycle

CPD simplifies the entire process



Wrap up

Lab 10 – Wrap up



Cloud Pak for Data

Unified Experience

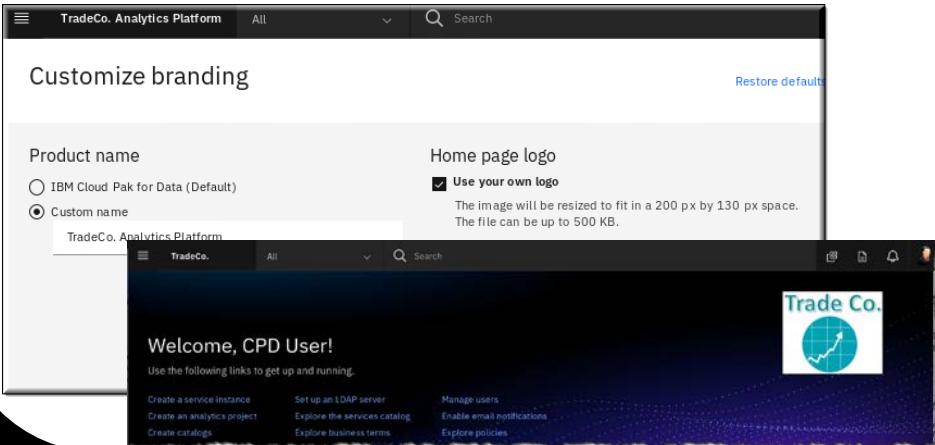


Customized Logo and Branding

Example: Customize the web client interface per tenant.

Customizable components in the Home Page:

- Product name
- Home page logo



Group 1 language support

All services in Cloud Pak for Data are translated for the following languages:

- Simplified Chinese
- Traditional Chinese
- Japanese
- French
- German
- Italian
- Spanish
- Brazilian Portuguese

Carbon 10

Modernized, modular and flexible open source design framework

- Consistent look and feel
- Reusable components

Cloud Pak for Data

Version 3.0.1 notable updates



Update	Comment
Support for POWER systems	Build an AI foundation for speed and scale with the premier, built-in GPU acceleration platform for faster time to AI.
OpenShift Container Storage	Software-defined Storage automated for quicker and efficient hybrid multi-cloud deployments and optimized for Red Hat OpenShift Platform.
Red Hat OpenShift 4.x support	Or continue to run on RHOC 3.11.
IAM Integration	IBM Cloud Platform Common Services (CPCS) IAM for all Cloud Paks provides authentication support via the OpenID Connect (OIDC) specification for SSO capability.
Fine-grained permissions	4 New fine-grained permissions: <ul style="list-style-type: none">Configure authenticationConfigure platformManage usersMonitor platform
Operations Management	New utilities for upgrade, backup and restore, import and export

Cloud Pak for Data Security Considerations



Security Features

- ✓ Security Architecture and Design
- ✓ Access Control, Authentication and Authorization (e.g. integrates with leading LDAPs)
- ✓ Data Protection
- ✓ Security Logging

Security Engineering

- ✓ Development trained in Secure Coding Practices
- ✓ Secure Engineering Development Practices: threat modeling, risk assessment, static and dynamic code analysis, penetration testing, container scanning, etc.

Security Operations

- ✓ Audit Log consolidation and analysis
- ✓ User access management
- ✓ Security Incident Management

Governance & Compliance

- ✓ Compliance Controls defined by Outside Agencies
- ✓ System Security Plans for maintaining compliance security postures

Compliance Best Practices:

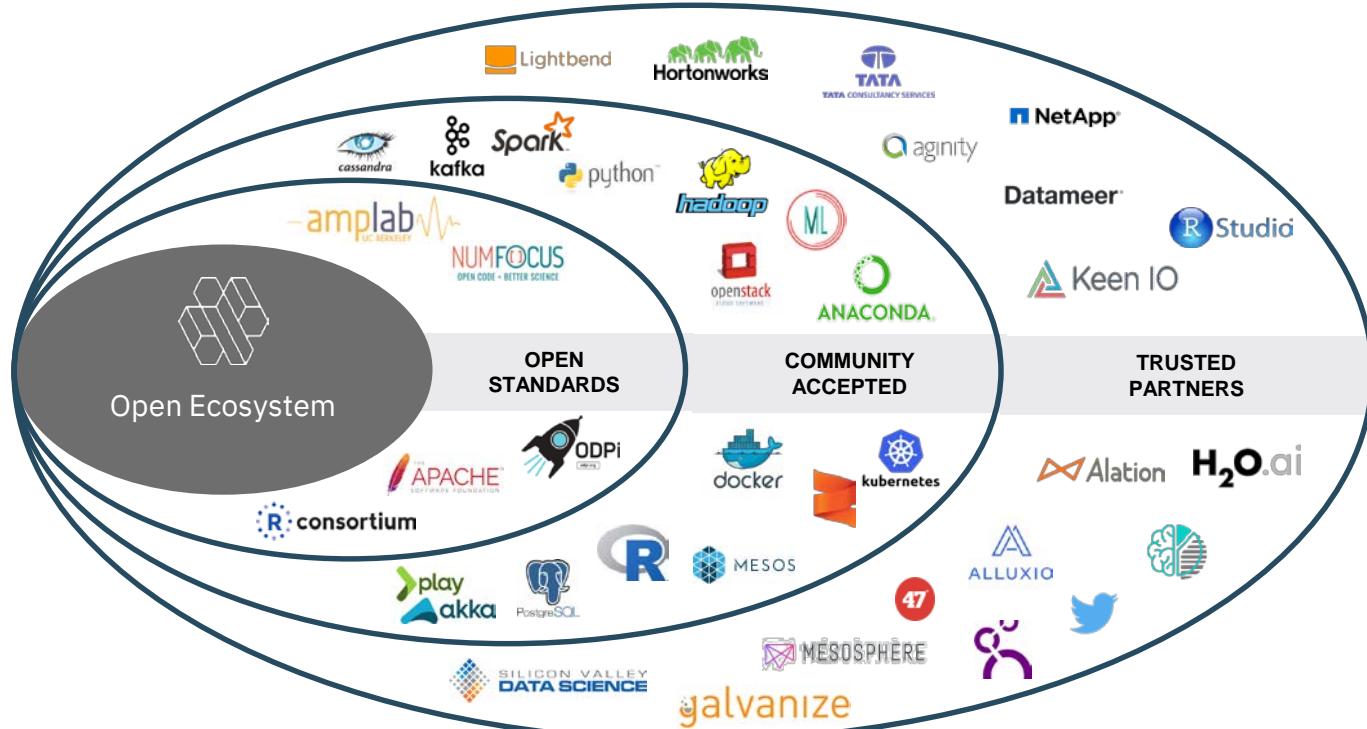
FISMA High “ready” with System Security Plans, spanning 350 controls:

- Risk Assessment
- Certification, Accreditation and Security Assessments
- System Services and Acquisition
- Security Planning
- Configuration Management
- System and Communications Protection
- Personnel Security
- Awareness and Training
- Physical and Environmental Protection
- Media Protection
- Contingency Planning
- System and Information Integrity
- Incident Response
- Identification and Authentication
- Access Control, Accountability and Audit

GDPR “readiness” considerations

CPD built on an open Ecosystem

Where IBM leads, partners and co-creates



IBM's approach to Open technology: <https://developer.ibm.com/articles/cl-open-architecture-update/>



Build Once - Run anywhere
In your own data center
Or the cloud infrastructure of your choice



Your Data Center



IBM Cloud



OPENSHIFT



openstack™



Helps avoid vendor lock in

CPD Use Cases

Industry agnostic



Create a Customer Focused Enterprise

- Rich profile for every customer kept up to date in real time as new customer behaviour is collected (360 view)
 - Spending Patterns
 - Behaviour
- Deliver tailored offerings based on segmentation
- Provide “next best action” in real-time

These use cases apply to most industry verticals

CPD Use Cases

Industry specific

Over 24 Data Science & AI use-cases across 15 industry verticals (*applicability varies by customer*)

Cloud Pak for Data differentiation :

- Operationalize models in a matter of minutes (Deploy, scale & manage models with minimum effort)
- Model governance (Lineage and provenance – Who created, when, what data was used, comments, ratings etc.)

Use Case(s)	Aerospace & Defense	Automotive	Banking	Chemicals & Petroleum	Consumer	Education	Electronics	Energy, Environmental & Utilities	Financial Markets	Government	Healthcare & Life Sciences	Insurance	Industrial Products	Telco, Media & Entertainment	Travel & Transportation
Predictive Maintenance	X							X							X
Real time analytics (IOT)	X	X		X											X
Customer Churn / Retention		X	X												X
Anomaly Detection	X	X						X						X	
Regulatory Compliance			X						X			X			
Anti-Money Laundering (AML)			X												
Cross Sell / Up-Sell			X		X										
Demand Forecasting				X			X	X							
Inventory Optimization					X										
Retention & Time to Degree								X							
Application modernization								X							
Student Safety							X								
Predictive Customer Insights				X				X							
Counter Fraud & Payments									X			X			
Counter Party Credit-risk										X					
Client Insights for Wealth Management										X					
Threat Prediction & Prevention											X				
Patient Diagnosis												X			
Data Privacy												X			
Client Risk Scoring			X										X		
Targeted Ads														X	
Intrusion Detection	X													X	
Route Optimization															X

Cloud Pak for Data

Ranked #1 by Forrester for “Enterprise Insight Platforms”



Enterprise Insight Platforms - Definition

- Enterprise insight platforms pre-integrate most — or all — of the technology required to build systems of insight and thus help business move faster. The need to move faster and change more easily is the driving force behind customer demand for these platforms.
- Vendors that can better support all the personas of an insight team with unified experiences that feature governance and can creatively enable hybrid cloud and multi-cloud delivery will win.

Forrester on “ICP for Data”

IBM has an impressive portfolio of individual data management and analytics capabilities that have consistently scored well on individual component Forrester Waves. With ICP for Data, IBM has pre-integrated capabilities that allow clients to be productive in a week or less. We were also impressed with its ML-assisted data cataloging and governance tools. IBM's platform uses Kubernetes to deploy on-premises or into the public cloud. Lastly, IBM's support for different insight team personas through tailored but unified experiences is commendable. Firms looking to unify the work of insight teams will do well on this platform.

Forrester on Microsoft’s Perceived Weakness – Azure Cloud Platform

While Microsoft offers AI services, its multimodal predictive analytics and machine learning (PAML) tools scored poorly in previous Forrester Waves. Finally, we found this offering to be too light on data governance capabilities and self-service data preparation tooling, both of which are critical insight team capabilities.

Report Preview : <https://ibm.box.com/s/bry68nm9alduszvrffo105cvjhmd7pn>

Cloud Pak for Data

Ranked #1 by Forrester for “Multimodal Predictive Analytics And Machine Learning”



Why IBM

- IBM is packed with AI lifecycle services everywhere — public, private, and on-prem. What do you get when you combine a full stack of data analytics capabilities — from data management to PAML to business intelligence — in a microservices framework that can run seamlessly on premises, in the private cloud, and in multiple public clouds?
- IBM Cloud Pak for Data, an offering that makes capabilities across the PAML lifecycle available when and where your users need them.
- The crown jewel of Cloud Pak for Data is Watson Studio, a PAML offering that combines easy-to-use, SPSS-inspired workflow capabilities with open source ML libraries and notebook-based interfaces.
- IBM continues to add innovations from IBM Research like fairness monitoring, bias mitigation, AutoML, and federated learning.
- IBM offers a compelling, scalable, increasingly integrated, and harmonized platform with differentiated capabilities that spans the entire PAML lifecycle and can be deployed anywhere across any cloud. Users will have to navigate some technical previews and quickly evolving features, but that will save them from otherwise having to stitch together a mismatch of proprietary solutions and open source software.

Report Preview : <https://ibm.box.com/v/Forrester2020>



Cloud Pak for Data Editions

Make your data ready for AI – Cloud Agility, Lightning Fast & AI-ready

Standard Edition

- Beginning with a minimum of 24 VPCs
- Expands in increments of 1 VPC
- Up to a *maximum* of 64 VPCs
- You cannot use separately priced premium services with this edition
- 50% Enterprise Edition list price

Enterprise Edition

- Deploy in any form: on premises, on public cloud or CPD System
- Beginning w/ recommended minimum of 48 VPCs
- Select 24 VPC configurations supported
- Expands in increments of 1 VPC
- No maximum

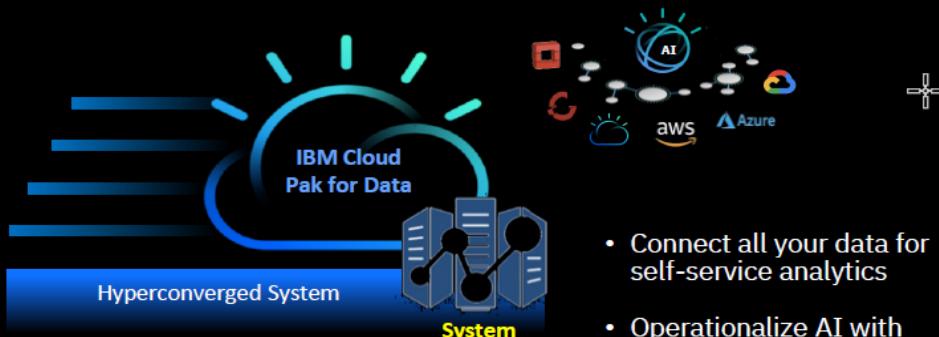
Non-Prod Edition

- Can only be used for non-production scenarios
- No restriction in terms capabilities or VPC quantities
- Needs to be in separate namespace from production licenses
- Need parallel and appropriately sized Standard or Enterprise Edition licenses for production use
- 50% of Enterprise Edition list price



IBM Cloud Pak for Data System

True plug-and-play enterprise data & AI in hours right out of the box



- Brings the elasticity & scalability of public clouds securely behind the firewall

- Connect all your data for self-service analytics
- Operationalize AI with trust & transparency
- Deploy dynamic cloud native data workloads

An **all-in-one** data & AI system with all the necessary systems and software components pre-integrated

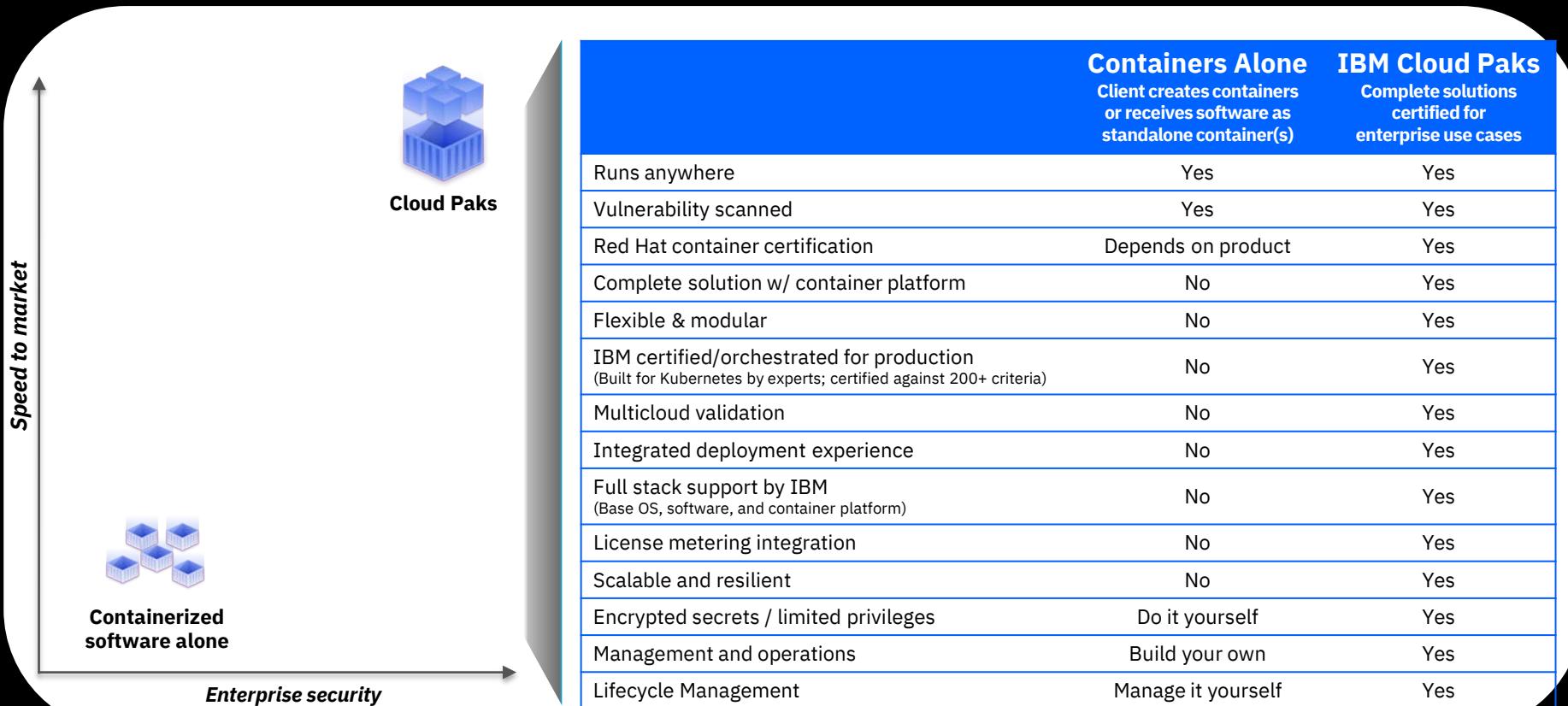
Deploy a complete private cloud in under 4 hours, with no assembly required

Dynamically scale compute, storage and networking resources with plug and play of new hardware nodes

Simplify management and optimization with a unified and intuitive dashboard

Cloud Pak for Data vs. Containers Alone

IBM Certified and production ready



Cloud Pak for Data

IBM Kubernetes Certified

<http://ibm.biz/cp-certify>



Production Grade	Security	Quality Assurance	Lifecycle Management
			

Consistency and Standards

	<ul style="list-style-type: none">• Consistent Packaging / Publishing• Supporting Operators and Helm• Consistent Entitlement management• Common management of OSS elements	<ul style="list-style-type: none">• UBI and Red Hat Certified• Consistent use of OCP and IBM Services• ~200 Code Standards enforced• Governed Best Practice / Anti Practices
---	---	---

 Red Hat OpenShift	 OPERATOR FRAMEWORK		<ul style="list-style-type: none">• Managed Image CVEs• Packaging• Publishing	<ul style="list-style-type: none">• Trusted Source• E2E Support
---	--	--	---	--

Cloud Pak for Data

Key differentiators (vs. Microsoft, AWS and Google)



Data virtualization

- Query all of your data sources as one
- Governance, security, and scalability by design
- 400X faster than federation
- *Unmatched by Microsoft, AWS, or Google (may only have simple federation at best)*

Data governance

- Data privacy & governance by design: data discovery & curation, with policy & rules management
- Metadata management and shopping for data
- Smarter compliance: Regulatory ML, industry accelerators, FISMA HIGH certification, etc.
- *Unmatched by Microsoft, AWS, or Google*

Model bias and drift detections w/ explainability

- Detect and mitigate model bias and drift
- Explainability of a model prediction for a single transaction
- *Unmatched by Microsoft, AWS, or Google (none have all three)*

Governing and operationalizing AI

- Governed AI lifecycle management
- CI/CD style pipelines for AI DevOps
- AI model trust and transparency
- *Unmatched by Microsoft, AWS, or Google*

Open source based hyperconverged system

- Query all of your data sources as one
- Governance, security, and scalability by design
- 40X faster than federation
- *Unmatched by Microsoft, AWS, or Google*



Turbocharged digital transformation

Read the story in any of these magazines:

Business Chief US

(front cover, story pages 12-23,
Cloud Pak for Data pages 16-17)

Business Chief Canada

(story pages 144-153, Cloud Pak for Data pages 146-147)

Gigabit Magazine

(front cover, story pages 12-23, Cloud Pak for Data pages 16-17)



Facing its own path toward digital transformation, Sprint started preparing its data for artificial intelligence (AI) with the goal of using machine learning algorithms to gain quicker insights and increase responsiveness to customers.

Sprint chose [Cloud Pak for Data](#) because it enables AI projects in weeks rather than months through unifying and simplifying three critical stages in the journey to AI: the collection, organization and analysis of data.



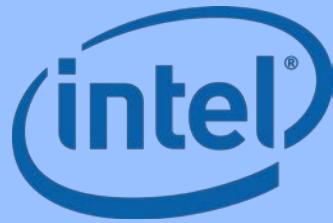
Cloud Pak for Data

Industry: Telecommunications
Geography: North America

Simplify and automate how your organization turns data into insights within a unified, all-in-one design.

"Cloud Pak for Data enabled Sprint to digest high volumes of data for near real-time ML/AI analysis, and the trial results have shown potential to take Sprint to the next phase of digital transformation."

Michelle Gehl
VP Networks OSS Applications and Operations, Sprint



Is AI your priority? Start with a data strategy

[Read the blog and](#)
[watch the video](#)

Intel and IBM are great partners and closely aligned in becoming more data-centric.

Intel's participation and contribution is meaningful because customers can run [Cloud Pak for Data](#) at speed on their Intel-based infrastructure. The union between IBM and Intel is supercharging the ability of data scientists to drive better insight and better business outcomes in a way that has never been seen before.

Cloud Pak for Data

Industry: Technology

Geography: North America

Simplify and automate how your organization turns data into insights within a unified, all-in-one design.

"Cloud Pak for Data is really important because it helps to do a couple of things that are mind blowing for data scientists — auto discovery of data and rapid integration of hyper-relevant data."

Melvin Greer

Senior Principal Engineer and Chief Data Scientist
- Americas, Intel Corporation

Government | MEA (United Arab Emirates)

Accelerating the Customs Process

With Cloud Pak for Data's end-to-end platform, Dubai Customs is building trusted AI models to help identify risky goods entering the country and reduce the number of false positives.

- **When:** Current | **Duration:** 3 months
- **Product:** Cloud Pak for Data, Watson OpenScale
- **Pillars:** DSAI



Cloud Pak for Data

Industry: Government
Geography: MEA

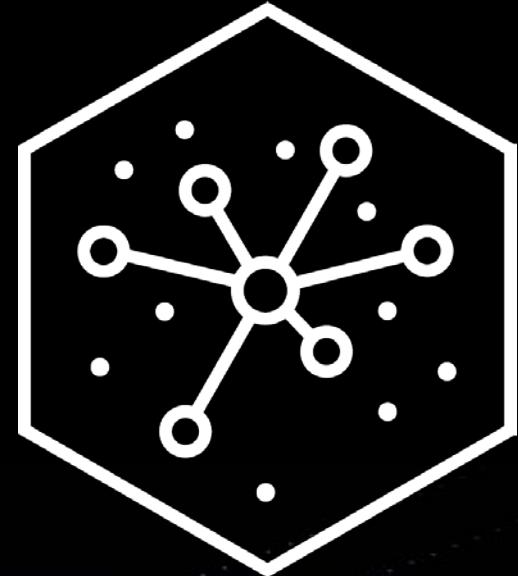
IBM Analytics Modernization Workshop

Part 3

- | | |
|---|--|
| <ul style="list-style-type: none">• Introduction• Business Use Case | <ul style="list-style-type: none">• Lab 01• Lab 02 |
| <ul style="list-style-type: none">• Collect: Connect• Organize• Collect: Virtualize | <ul style="list-style-type: none">• Lab 03• Lab 04• Lab 05 |
| <ul style="list-style-type: none">• Analyze• Deploy• Infuse – OpenScale• Infuse – Cognos Analytics• Wrap up | <ul style="list-style-type: none">• Lab 06• Lab 07• Lab 08• Lab 09• Lab 10 |

Thank you for your time!

**Begin your journey now on the
IBM Platform built for AI...**



We appreciate your feedback.

Copyright and trademarks

© Copyright IBM Corporation 2020

IBM Corporation
Route 100
Somers, NY 10589

Produced in the United States of America
July 2020

IBM, the IBM logo, ibm.com, API Connect, Db2, Elastic Storage, FlashCore, POWER, Spectrum Scale, UrbanCode, WebSphere and IBM Z are 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” at www.ibm.com/legal/copytrade.shtml

Intel, Intel logo, Intel Inside, Intel Inside logo, Intel Centrino, Intel Centrino logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium, and Pentium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both. Microsoft, Windows, Windows NT, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

VMware, the VMware logo, VMware Cloud Foundation, VMware Cloud Foundation Service, VMware vCenter Server, and VMware vSphere are registered trademarks or trademarks of VMware, Inc. or its subsidiaries in the United States and/or other jurisdictions.

This document is current as of the initial date of publication and may be changed by IBM at any time. Not all offerings are available in every country in which IBM operates.

THE INFORMATION IN THIS DOCUMENT IS PROVIDED “AS IS” WITHOUT ANY WARRANTY, EXPRESS OR IMPLIED, INCLUDING WITHOUT ANY WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND ANY WARRANTY OR CONDITION OF NON-INFRINGEMENT. IBM products are warranted according to the terms and conditions of the agreements under which they are provided.