



Powering the Future of Data



Agenda

- **Solution Overview - 35 mins**
- **Sentiment Analysis Lab Demo - 15 mins**
- **QA – 10 mins**
- **Lab – 120 mins**



About Hortonworks

About Hortonworks

Customer Momentum

- ◆ +1 000 customers
- ◆ Doubled customer based in 2015
- ◆ Publicly traded on NASDAQ: HDP

The Leader in Connected Data Platforms

- ◆ Hortonworks DataFlow for data in motion
- ◆ Hortonworks Data Platform for data at rest
- ◆ Powering new modern data applications

Partner for Customer Success

- ◆ Leader in open-source community, focused on innovation to meet enterprise needs
- ◆ Unrivaled support subscriptions

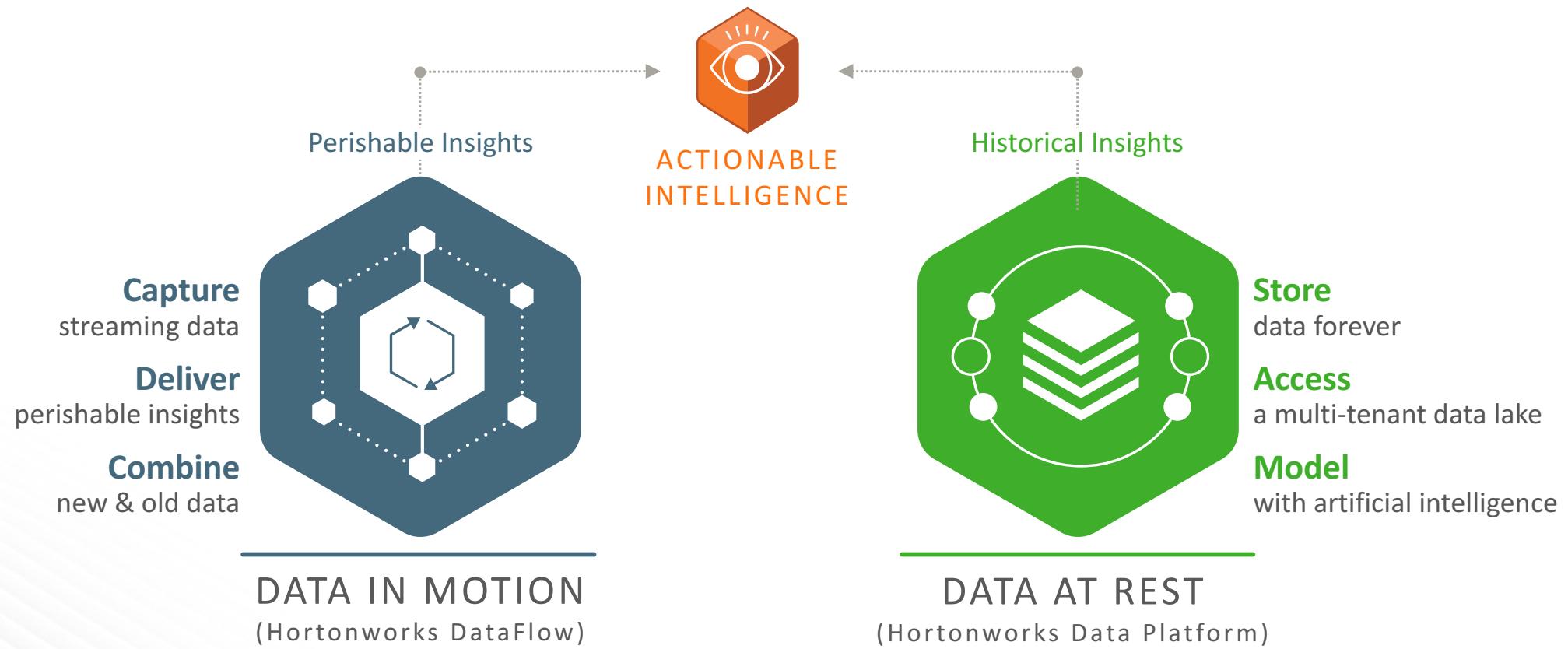


1000+
EMPLOYEES

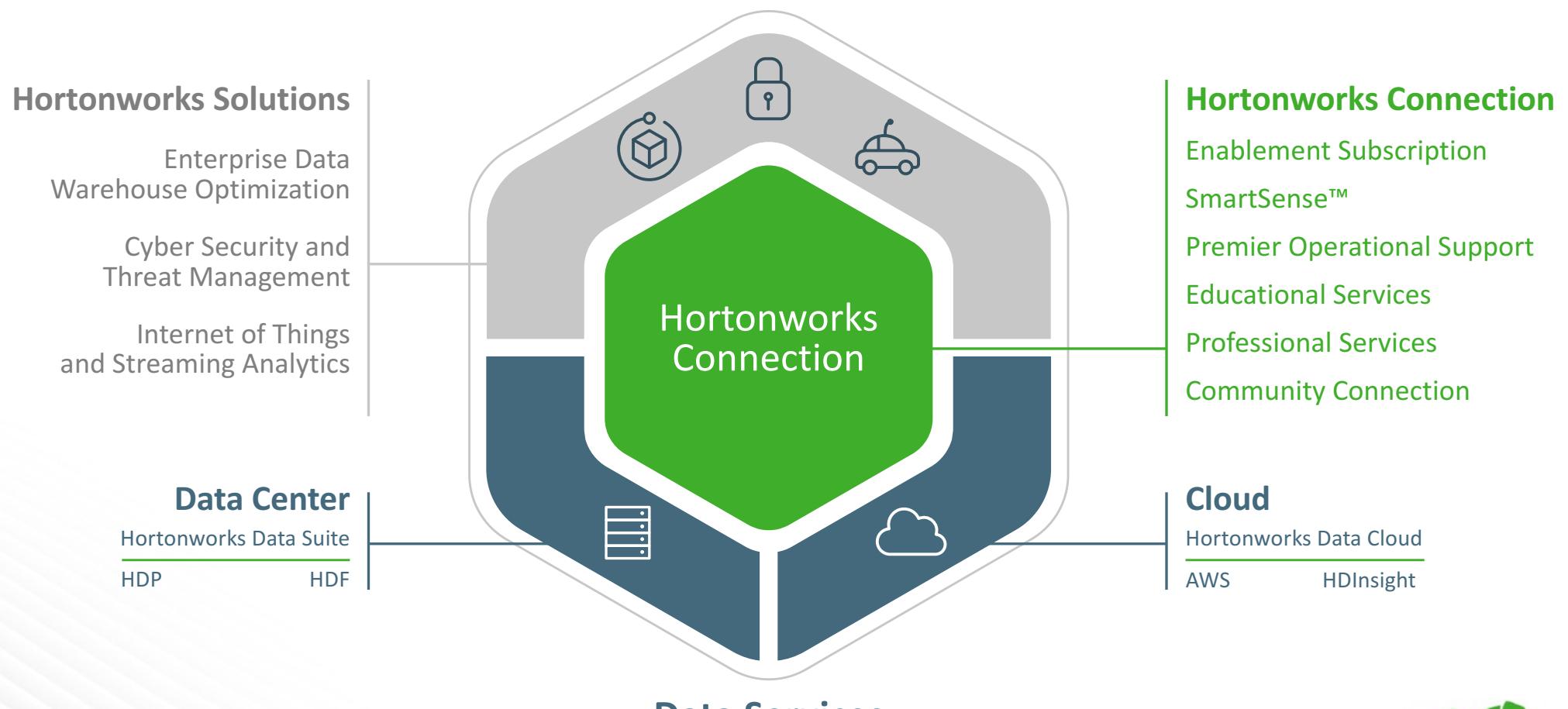
2600+
ECOSYSTEM
PARTNERS



A Connected Data Strategy Solves for All Data



Hortonworks Connection: Services and Solutions for Your Success



GA Cloud Offerings: Microsoft Azure and Amazon Web Services



Azure HDInsight

- ◆ Market-leading, feature-rich
- ◆ Microsoft branded & supported
- ◆ Hourly pricing (with annual option)
- ◆ Billed through customer's Azure account
- ◆ 24x7 Microsoft support with Azure pricing
- ◆ Hortonworks Community Connection (HCC)

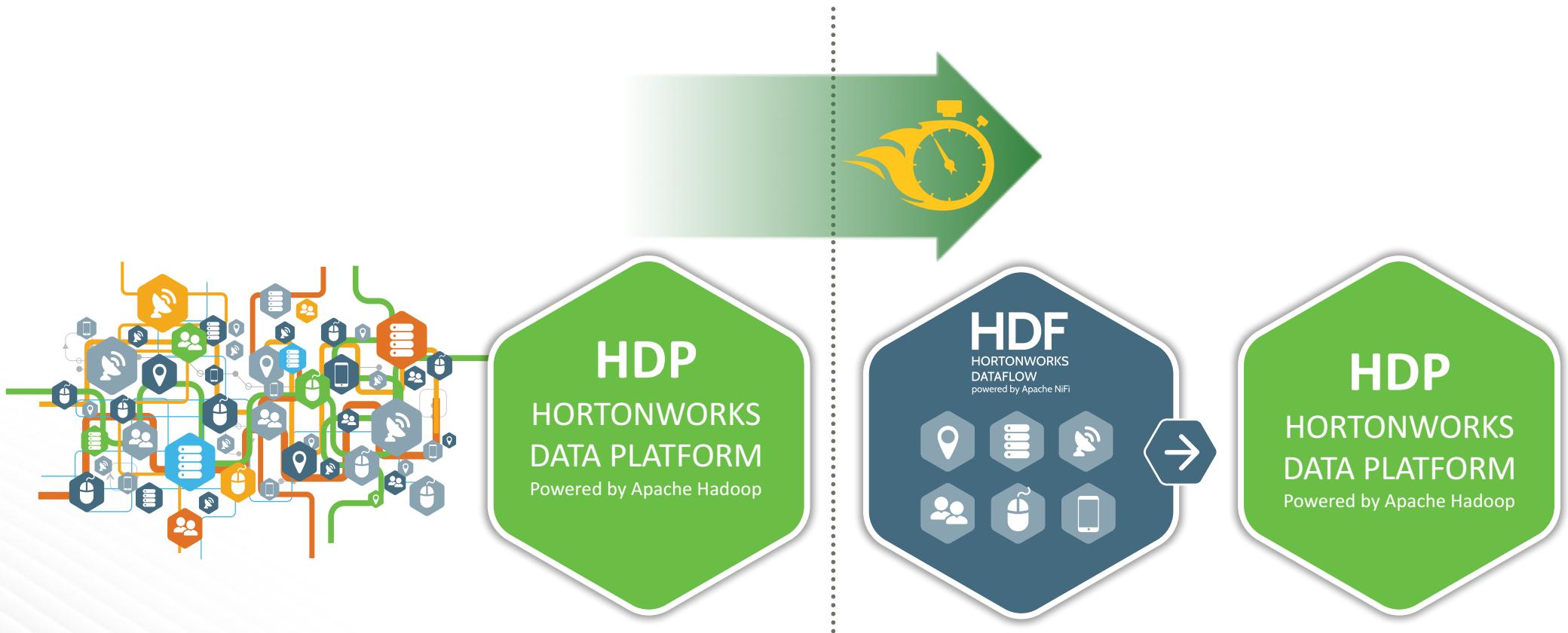


Hortonworks Data Cloud for AWS

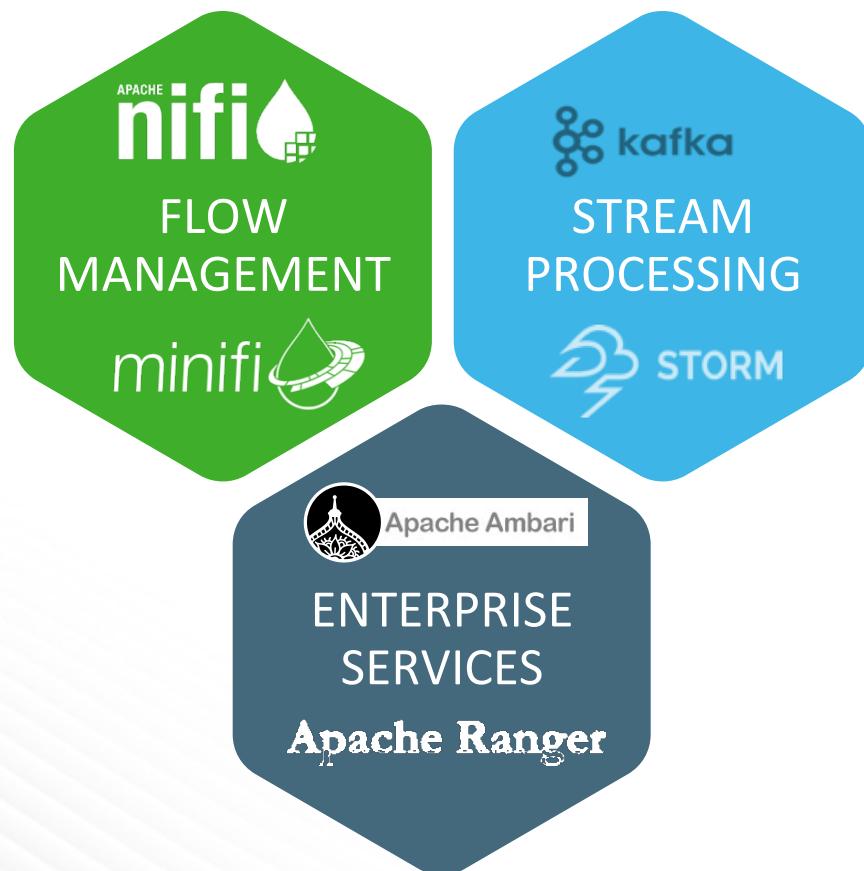
- ◆ Enterprise security, governance and operations
- ◆ Hortonworks branded & supported
- ◆ Hourly pricing (with annual option)
- ◆ Billed thru customer's AWS account
- ◆ 24x7 Hortonworks enablement sold separately
- ◆ Hortonworks Community Connection (HCC)

Hortonworks Data Flow Overview

HDF Makes Big Data Ingest Easy



Hortonworks DataFlow Powered by Apache NiFi, Storm & Kafka



Streaming Analytics

Capture perishable insights from data-in-motion

Visual Control Over Data Flows

to manage who can see and touch data in transit

End-to-End Security

to encrypt, decrypt, and filter data on its journey

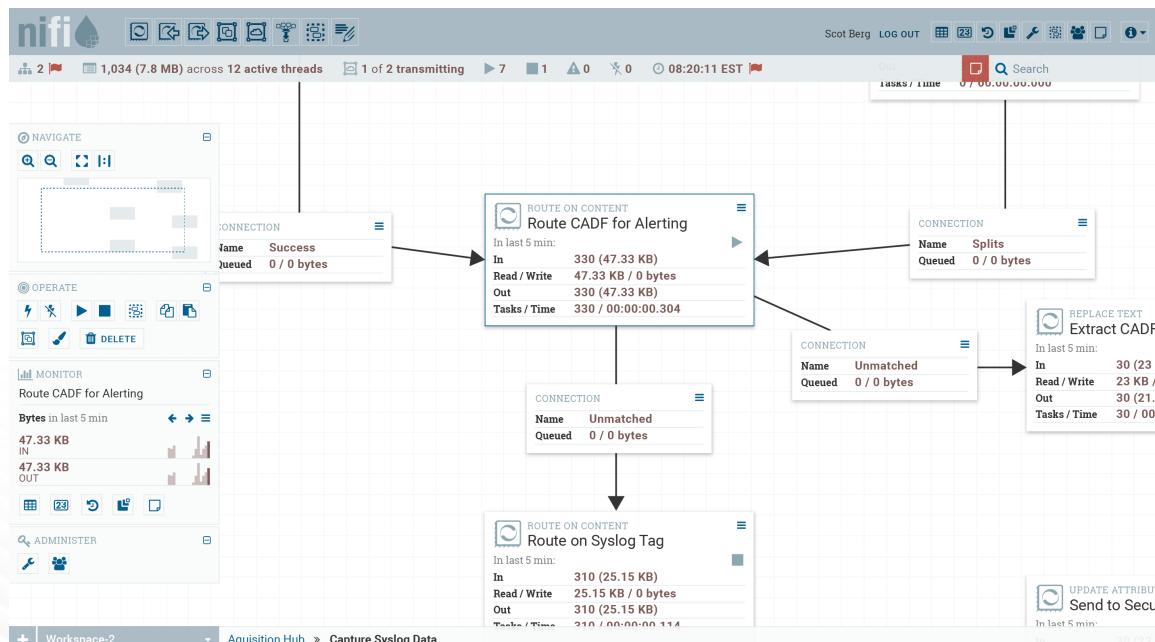
Real-Time Traceability

Rich metadata and contextual detail helps troubleshoot security issues



TM

Apache NiFi Key Features



◆ Visual User Interface

- Drag and drop for efficient, agile operations

◆ Immediate Feedback

- Start, stop, tune, replay data flows in real-time

◆ Adaptive to Volume and Bandwidth

- Any data, big or small

◆ Event Level Data Provenance

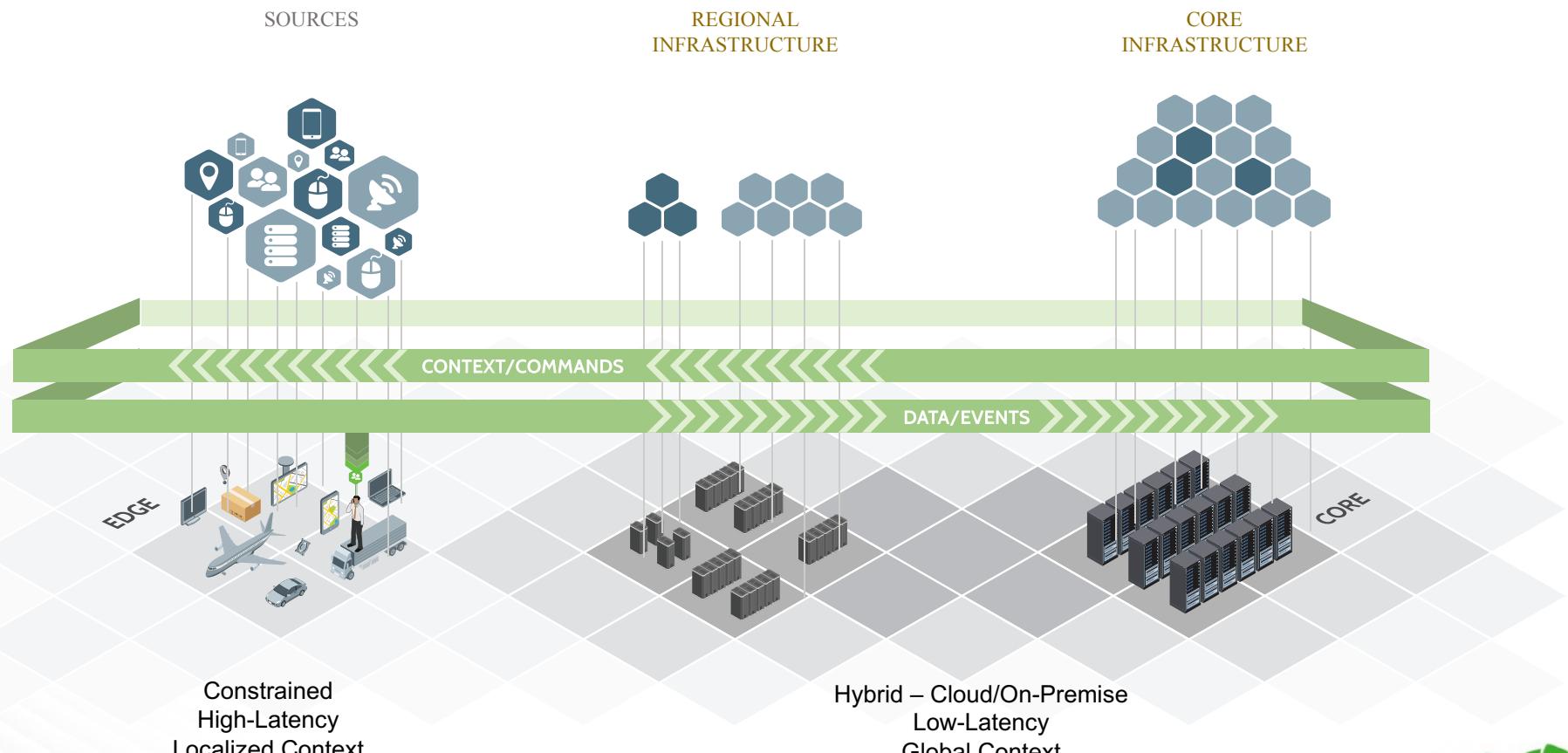
- Governance, compliance & data evaluation

◆ Secure Data Acquisition & Transport

- Fine grained encryption for controlled data sharing and selective data democratization



HDF Manages Dataflow



Hortonworks Data Platform Overview

What is Apache Hadoop?

Hadoop was designed for Big Data

HDFS – Hadoop Distributed File System

- Data broken into blocks and replicated 3x
- Automatically replaces lost data / computers
- Very high bandwidth, not IOPs optimized

YARN – Distributed Computation Layer

- Distributed execution on HDFS
- Many programming models
 - MapReduce, SQL, Streaming, ML...
- Multi-users, with queues, priorities, etc...

◆ Cost Effective

- Commodity hardware
- Open source software

◆ Scalable

- Efficiently store and process petabytes of data
- Grows linearly by adding commodity computers

◆ Reliable

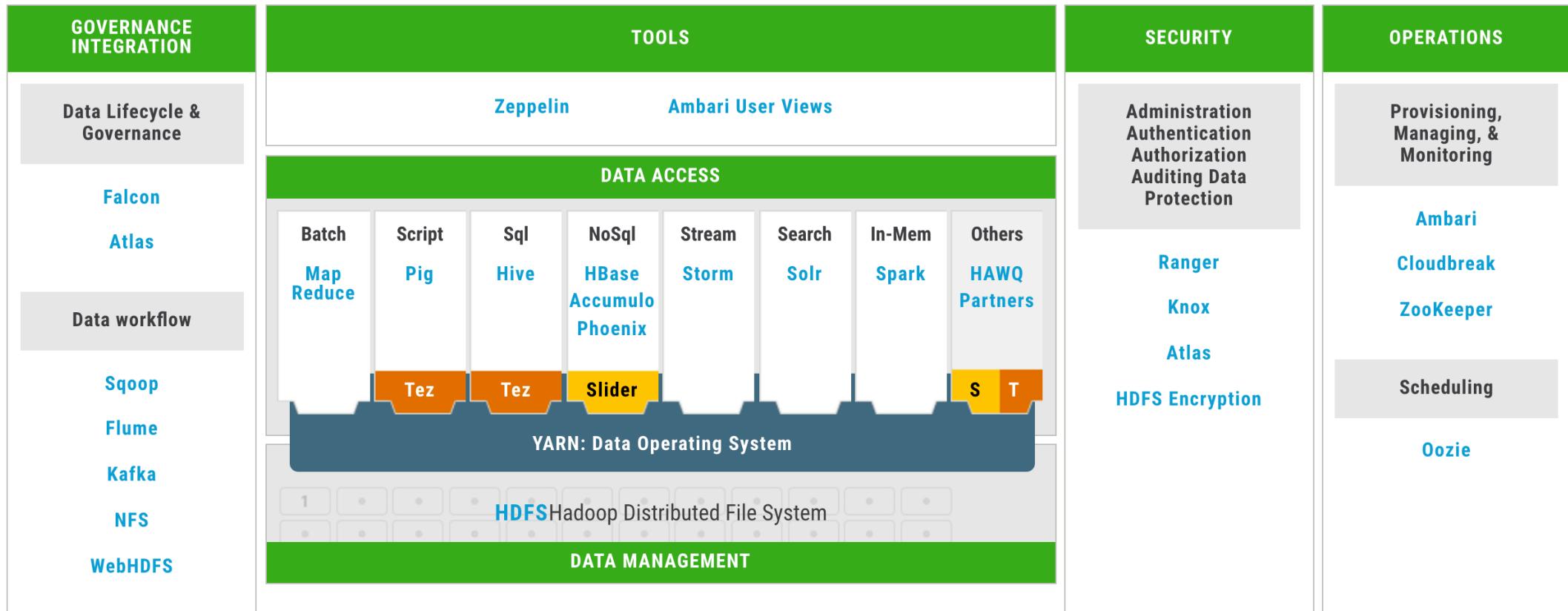
- Self healing as hardware fails or is added

◆ Flexible

- Store all types of data in many formats
- Schema-on-read



Hortonworks Data Platform



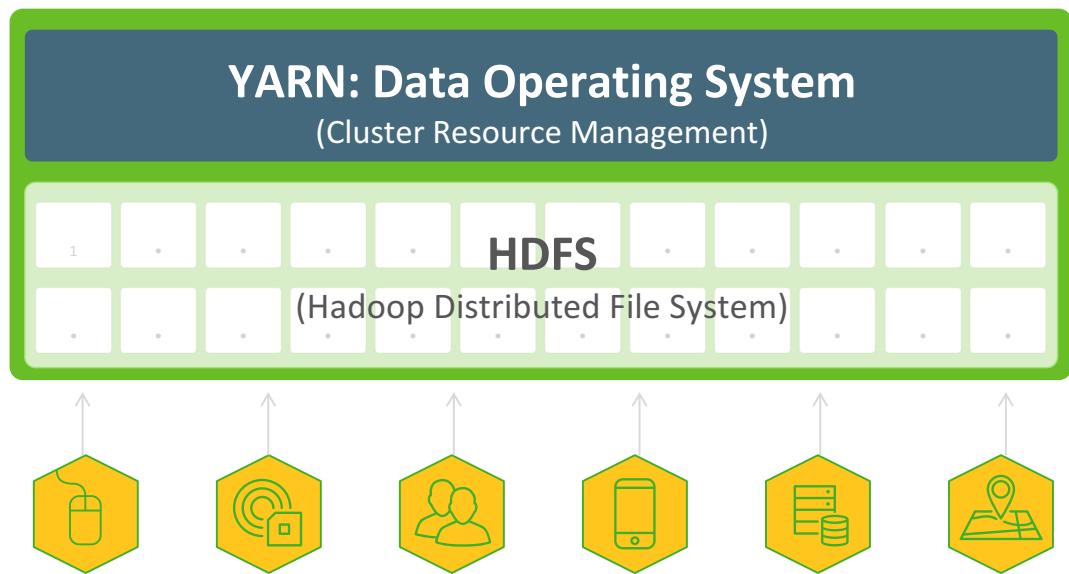
Core Hadoop: HDFS + YARN

Distributed Computing of Storage, CPU and Memory

Centralized Architecture

Highly Scalable

Cost Effective



HDFS: Name Nodes & Data Nodes

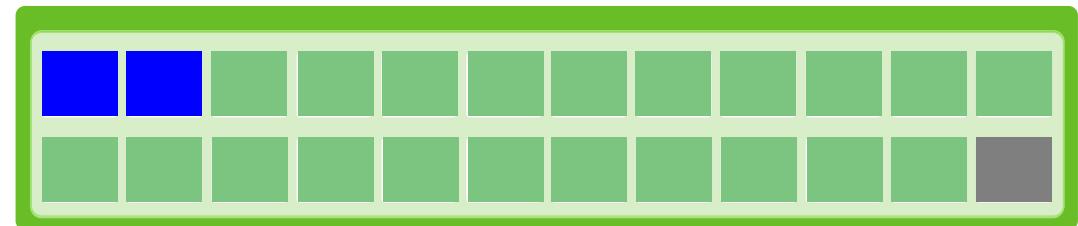
Data and its Namespace

Name Nodes

- Data Catalog, a “Namespace”
- Primary & Secondary Nodes for HA
- Manual or Automated Failover
- Quorum Journaling (States & Edits)

Data Nodes

- “Where the data lives”
- Data stored in replicas of 3 (default)
- Rack aware
- Heterogeneous Storage Options
- Large block sizes (64MB default)

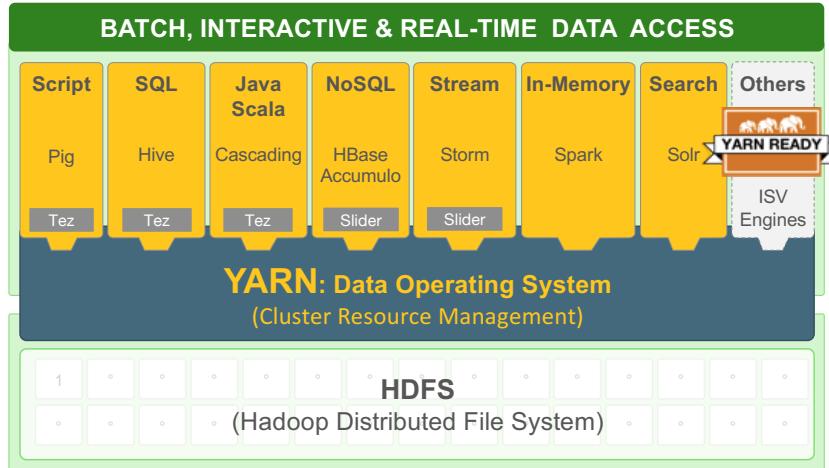


Edge & Utility Nodes

- Cluster Administration
- 3rd Party Tools

YARN: A Data Operating System

Yet Another Resource Negotiator



YARN Ready Applications

Facilitates ongoing innovation and enterprise adoption via ecosystem of new and existing “YARN Ready” solutions

YARN The Architectural Center of Hadoop

- Common data platform, many applications
- Support multi-tenant access & processing
- Batch, interactive & real-time use cases
- Supports 3rd-party ISV tools
(ex. SAS, Syncsort, Actian, etc.)

YARN in Production

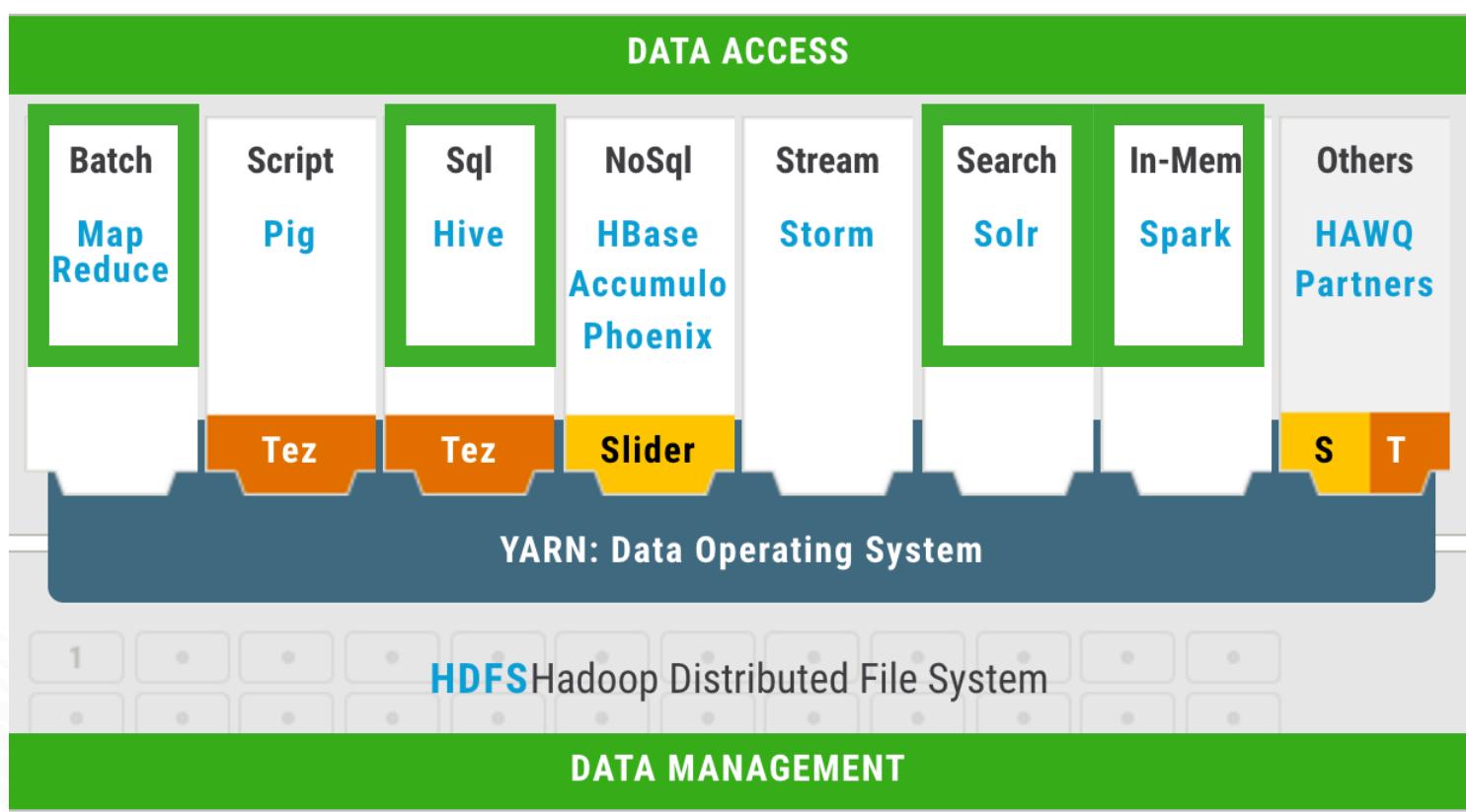
- Yahoo: ~40,000 nodes, multiple clusters running YARN across over 365PB of data
- Spotify, Progressive, Kohls, UHG, Sprint, JPMC, Target, AIG, Samsung

Data Access on Hadoop

Hortonworks Data Platform

MapReduce, Hive, Spark, Solr





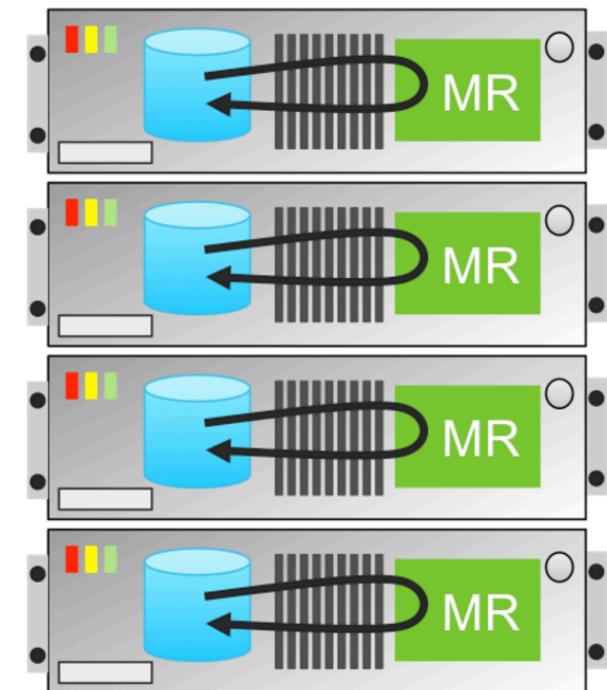
MapReduce

Hadoop architecture is based on a cluster of machines:

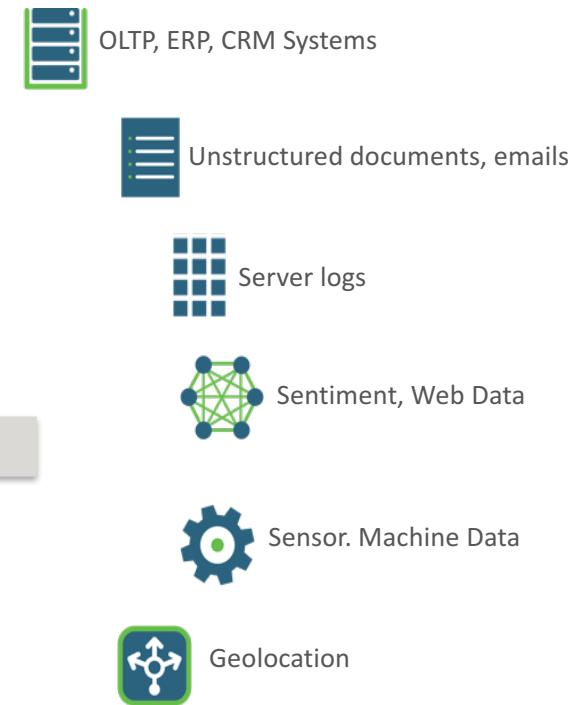
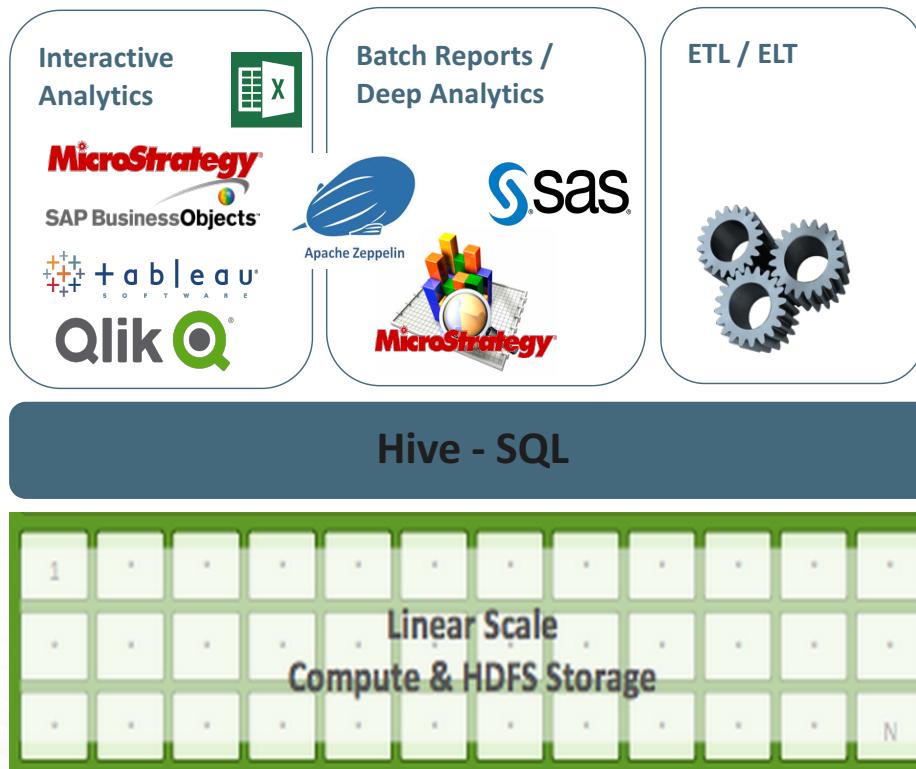
- Scale-out, shared nothing architecture, distributed processing.
- Requires distributed processing programs.

MapReduce:

- MapReduce Framework offers APIs to write distributed processing programs:
 - 3 Tasks/APIs: *Map, Shuffle, Reduce*
- Highly Scalable: Principle of Data Locality
- Best suited for batch processing.



Hive – Single tool for all SQL use cases



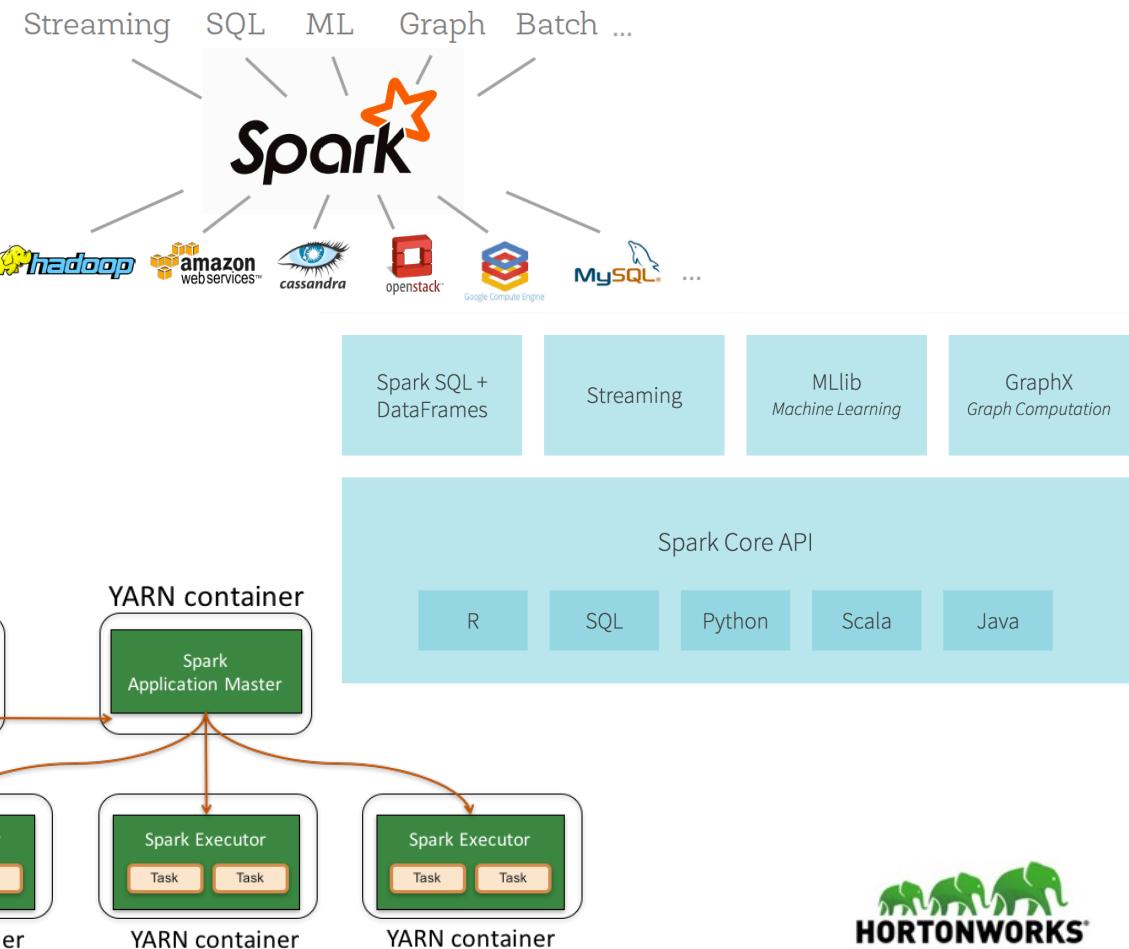
Apache Spark

What is Apache Spark ?

- Apache open source project originally developed at AMPLab (University of California Berkeley)
- Unified data processing engine that operates across varied data workloads and platforms

Why Spark ?

- Elegant Developer APIs
 - Single environment for data munging and Machine Learning (ML)
- In-memory computation model – Fast!
 - Effective for iterative computations and ML
- Machine Learning
 - Implementation of distributed ML algorithms
- YARN Ready



Apache Solr

What is Lucene?

- Apache Lucene is a high-performance, **full-featured text search engine library**
- **Provides API** to add search and indexing to your applications
- Provides scalable, High-Performance Indexing
- Efficient **Search Algorithms** (scoring, phrase, wildcard, sorting, geospatial, faceting....)
- When using Lucene you need to write code to do all of this

What is Solr?

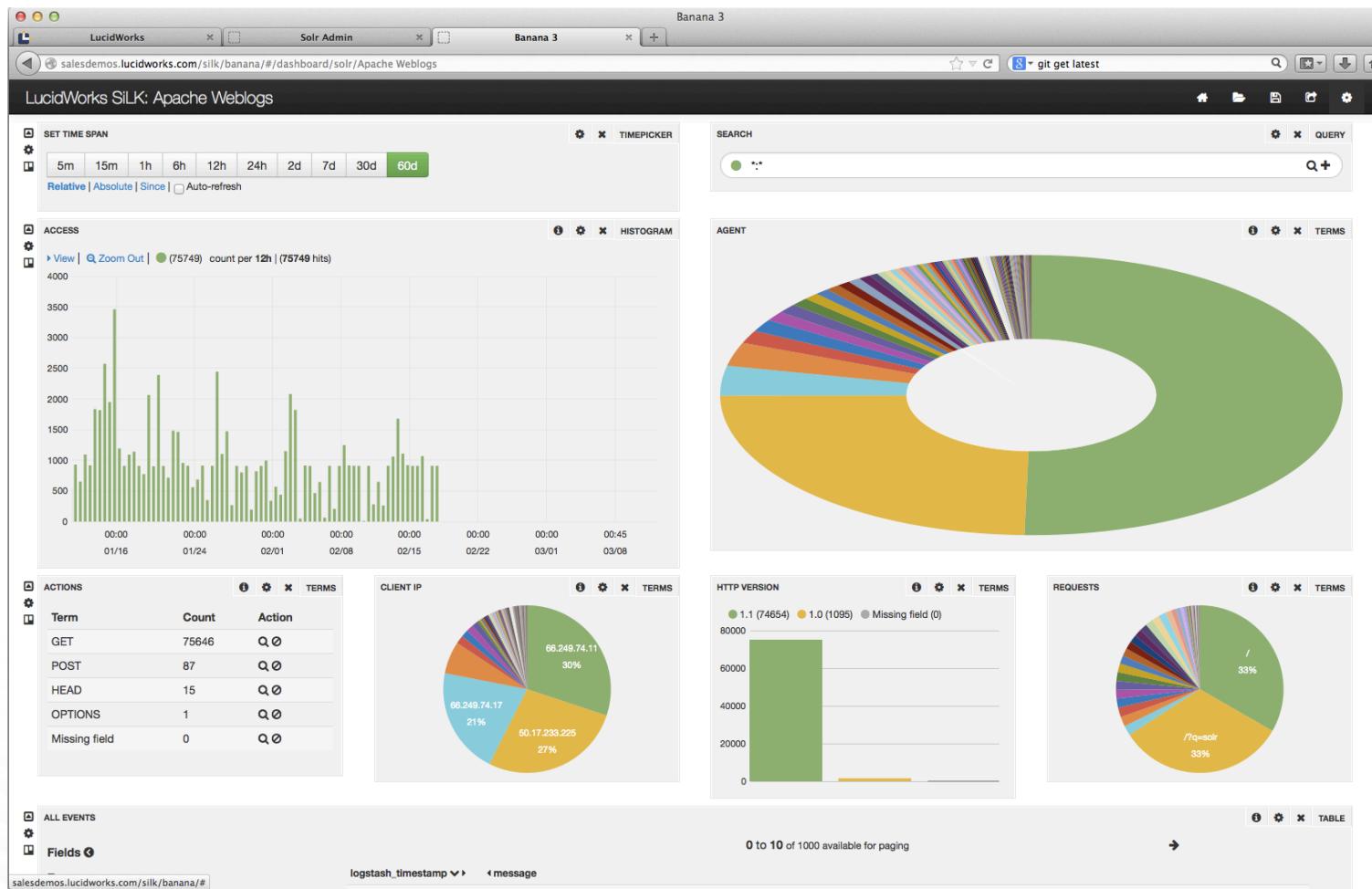
- **Search server built on top of Apache Lucene**
 - Provides API to access Lucene over HTTP
 - Ability to add more features on top of Lucene
- Majority of programming tasks in Lucene are XML configurations in Solr
- Provides **SolrCloud** (Scalability, Replication, Load Balancing)
- Provides **features** like Faceting, Clustering, Data Import Handlers, Multiple Language support, Rich document support



Key Features

- **Advanced Full-Text Search Capabilities**
- **Optimized for High Volume Traffic**
- **Standards based open interfaces: XML, JSON, HTTP**
- **Comprehensive Administration interface**
- **Near Real-Time, Batch and Real-Time Indexing capabilities**
- **Linearly scalable, highly available and extensible plugin architecture**
- **Deep Hadoop integration**

Banana: Visualization and Reporting



Apache Zeppelin GA: The Data Science Notebook



Features

- Ad-hoc experimentation
- Deeply integrated with Spark + Hadoop
- Interactive data ingestion, data exploration, visualization, sharing and collaboration

Data engineers, data analysts and data scientists



Use Cases

- Data exploration and discovery
- Visualization
- Interactive snippet-at-a-time experience
- “Modern Data Science Studio”



Security, Governance & Operations

Hortonworks Data Platform



HDP Security: Comprehensive, Complete, Extensible

Administration

Central management and consistent security

Single administrative console to set policy across the entire cluster: [Apache Ranger](#)

Authentication

Authenticate users and systems

Authentication for perimeter and cluster; integrates with existing Active Directory and LDAP solutions: [Kerberos](#) | [Apache Knox](#)

Authorization

Provision access to data

Consistent authorization controls across all Apache components within HDP: [Apache Ranger](#)

Audit

Maintain a record of data access

Record of data access events across all components that is consistent and accessible: [Apache Ranger](#)

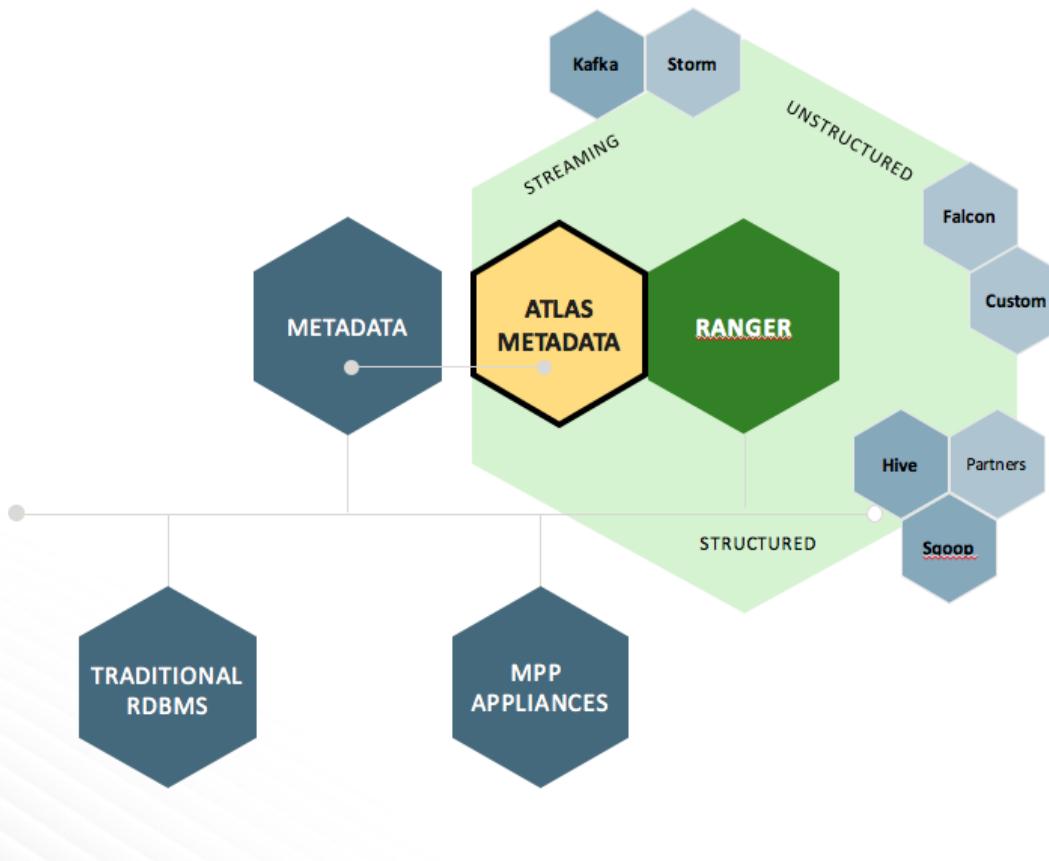
Data Protection

Protect data at rest and in motion

Secure data in motion and data at rest: [HDFS TDE w/ Ranger KMS + HSM](#), [Ranger Data Masking + Row Filtering](#), [Wire encryption + Partner Solutions](#)

Apache Atlas

Atlas is a scalable and extensible set of core foundational governance services – enabling enterprises to effectively and efficiently meet their compliance requirements within Hadoop and allows integration with the whole enterprise data ecosystem.

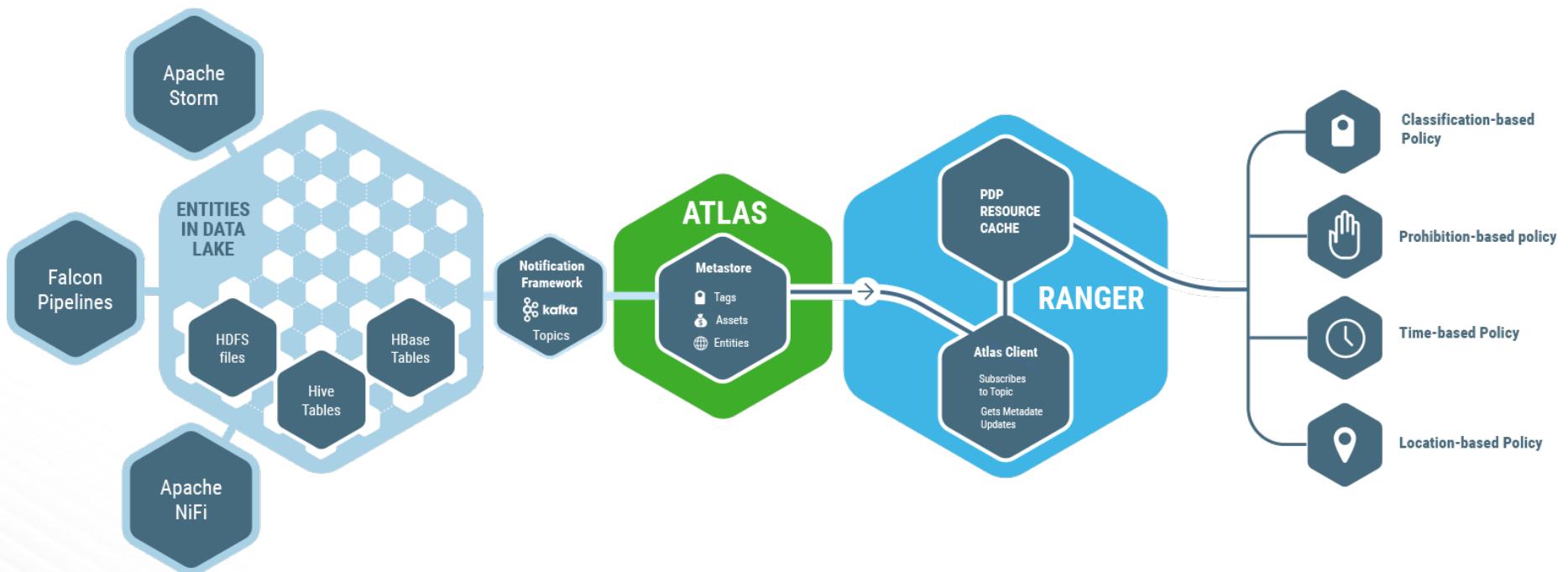


Key features

- Data Classification (taxonomy)**
- Security and policies engine**
- Data Lineage and Search over metadata**
- Centralized Auditing**
- Metadata Exchange**

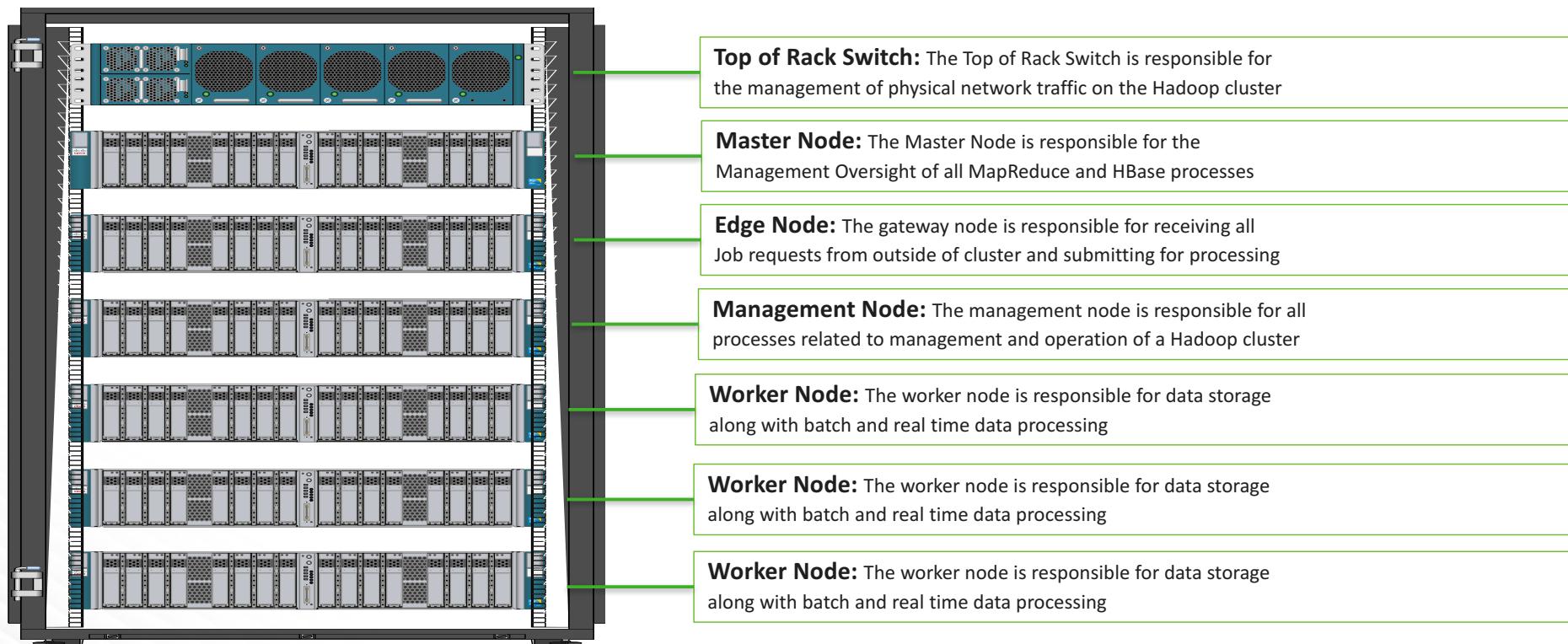


Atlas + Ranger: Metadata driven security model

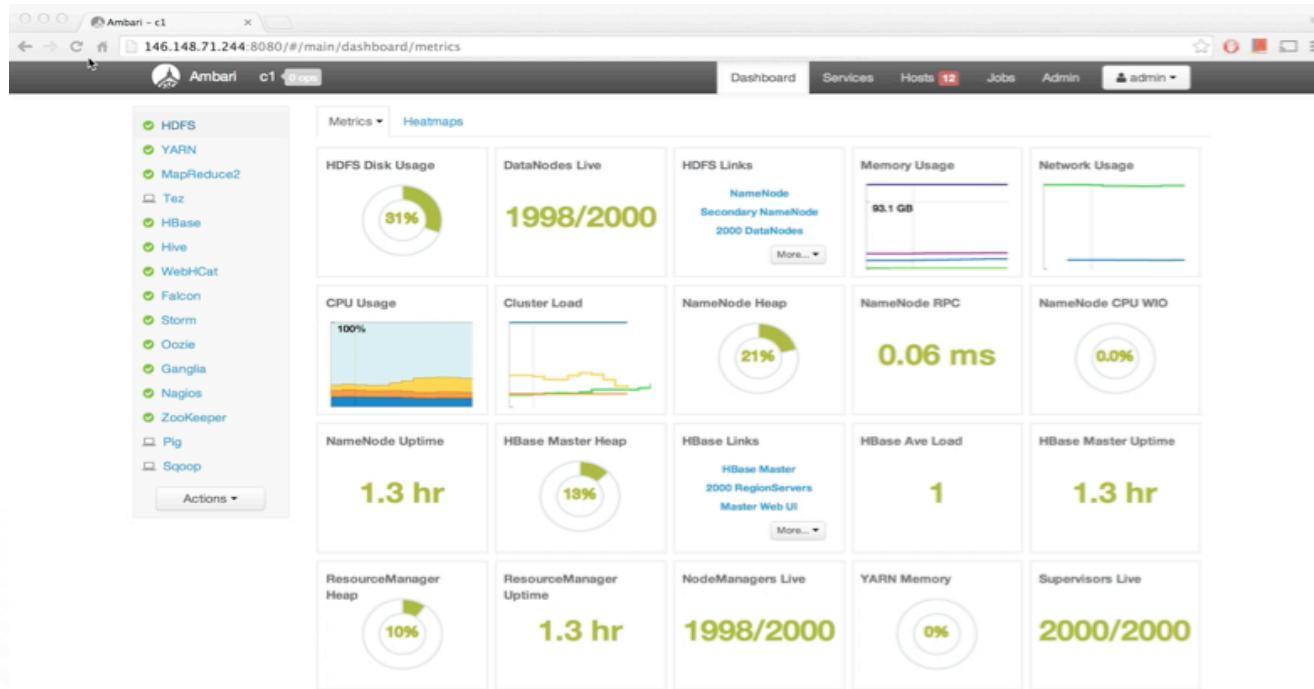


Typical Hadoop Rack Topology

Hadoop is deployed as a “Rack Aware” platform so that the data is securely replicated across multiple racks and the Hadoop services are also deployed across multiple racks, preventing cluster outages and job failures due to a single or multiple rack failures



Apache Ambari: Cluster Operations



A completely open source management platform for provisioning, managing, monitoring and securing Apache Hadoop clusters. Apache Ambari takes the guesswork out of operating Hadoop.

Simplified Installation,
Configuration and
Management

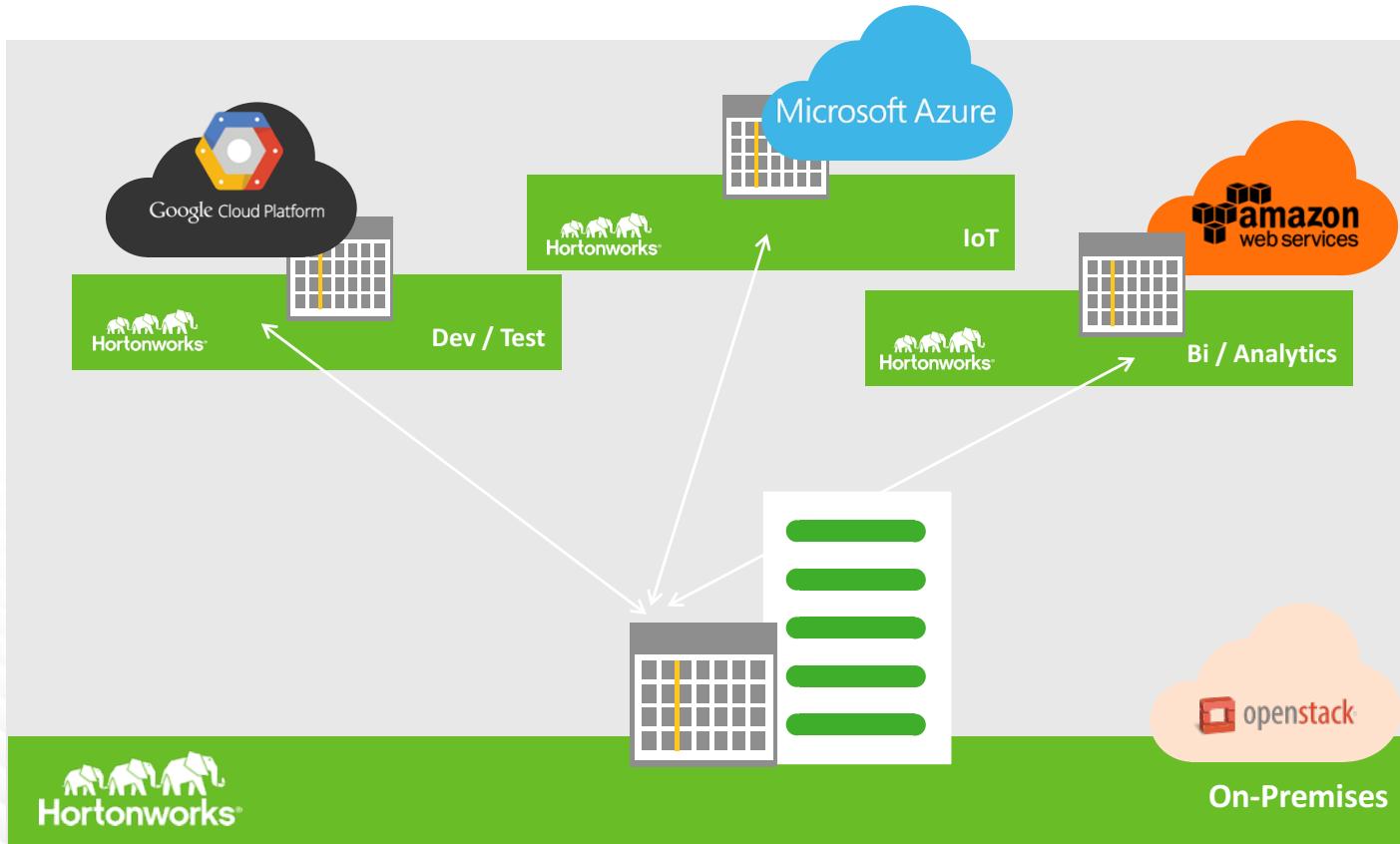
*Highly Extensible and
Customizable*

*Full Visibility into
Cluster Health*

*Centralized Security
Setup*

Cloudbreak: Cloud operations

Easily Launch HDP on Any Cloud with Cloudbreak



Cloudbreak is a tool for provisioning HDP clusters on cloud infrastructure. Cloudbreak allows enterprises to simplify the provisioning of clusters in the cloud and optimize their use of cloud resources as workloads change.

Cloudbreak: Cloud operations

The screenshot shows the Cloudbreak dashboard interface. At the top, there are three cluster status cards: AWS (11 nodes, 0 uptime), Azure (8 nodes, 0 uptime), and GCP (9 nodes, 0 23 uptime). A message at the top left says "20:21:19 Terminating the cluster and its infrastructure...". Below the cards is a navigation menu with the following items:

- manage security groups
- manage networks
- manage templates
- manage blueprints
- manage credentials
- manage recipes
- manage security configurations
- manage platforms

At the bottom of the dashboard is a footer bar with the text "SEQUENCEIQ".

Primary Use Cases

Deploy on Public or Private Clouds

Dynamically configure and manage clusters on public or private clouds (Amazon Web Services, Microsoft Azure, Google Cloud Platform and OpenStack)

Automated Scaling

Seamlessly manage elasticity requirements as cluster workloads change

Secured Cluster Access

Supports configuration for Kerberos, defining network boundaries and configuring security groups



Wrap-up

Hortonworks Data Platform



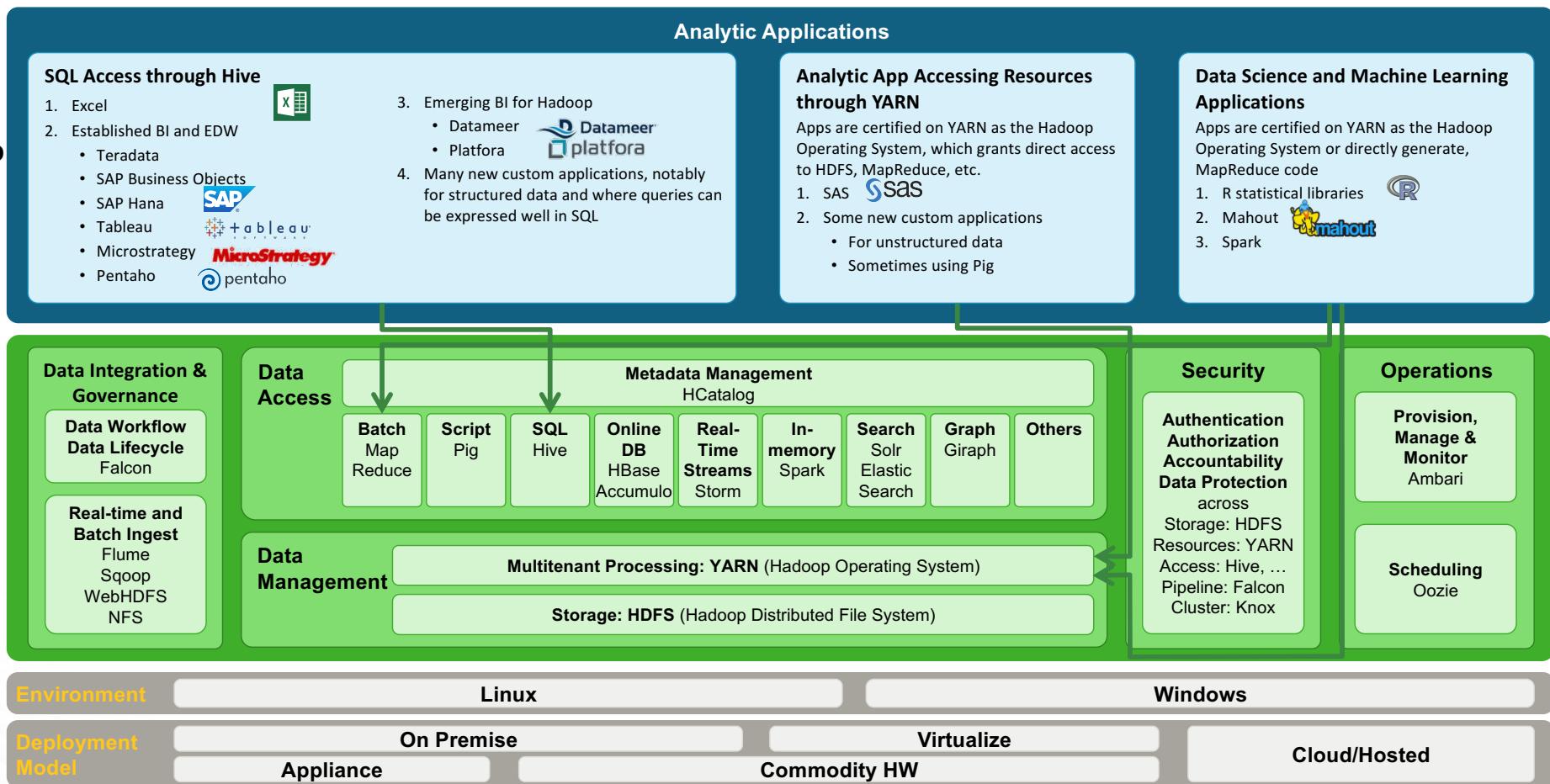
The big picture: Ways of Consuming Hadoop

= Hadoop

How is Hadoop consumed?

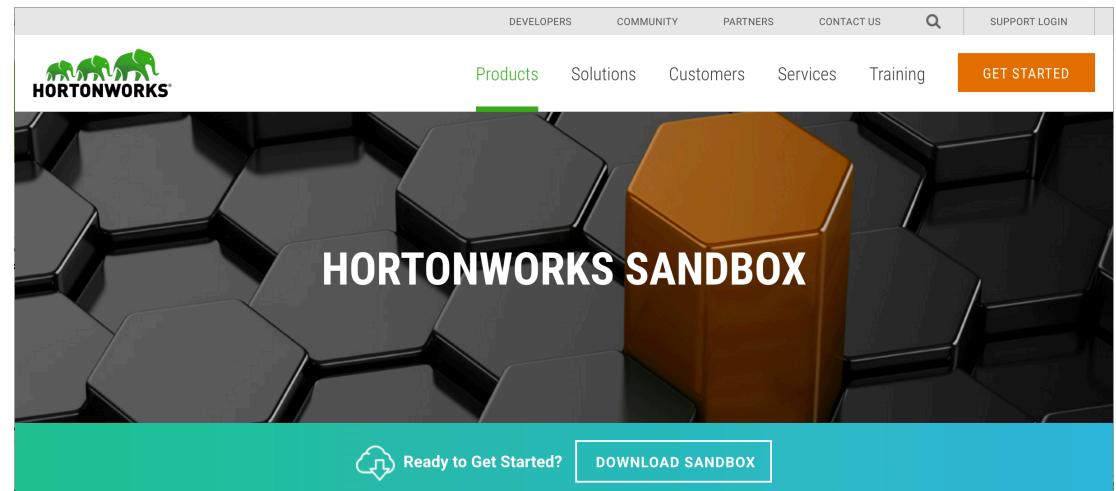
What is Hadoop?

Where does Hadoop run?



Next Steps: Try Hortonworks Data Platform

- Download HDP Sandbox:
<http://hortonworks.com/products/sandbox/>
- Deploy the HDP Sandbox on Azure:
<http://hortonworks.com/hadoop-tutorial/deploying-hortonworks-sandbox-on-microsoft-azure/>
- Learning the ropes:
<http://hortonworks.com/hadoop-tutorial/learning-the-ropes -of-the-hortonworks-sandbox>



A screenshot of the Microsoft Azure website. The top navigation bar includes links for Sales (1-800-867-1389), My Account, Portal, and Search. The main content area features a large blue banner with the text "Free one-month trial" and "Sign up for free and get \$200 to spend on all Azure services". It includes a "Try it now >" button, a "Or buy now >" link, and links for Frequently Asked Questions and More questions? Call us: 1-800-867-1389. To the right of the banner is a list of services with checkmarks: Virtual Machines, SQL Databases, Websites, Hadoop, Mobile Push, Media Streaming, Active Directory, and Everything else... .



**DATA
WORKS
SUMMIT**



**HADOOP
SUMMIT**

Presented by Hortonworks and Yahoo!

DATAWORKS SUMMIT AND HADOOP SUMMIT | SAN JOSE

JUNE 13–15 2017

**CALL FOR ABSTRACTS
IS NOW **OPEN****

<https://dataworkssummit.com/san-jose-2017/abstracts/submit-abstract>

**CALL FOR ABSTRACTS
CLOSES FEBRUARY 10**

[https://dataworkssummit.com/
san-jose-2017/agenda/#tracks](https://dataworkssummit.com/san-jose-2017/agenda/#tracks)

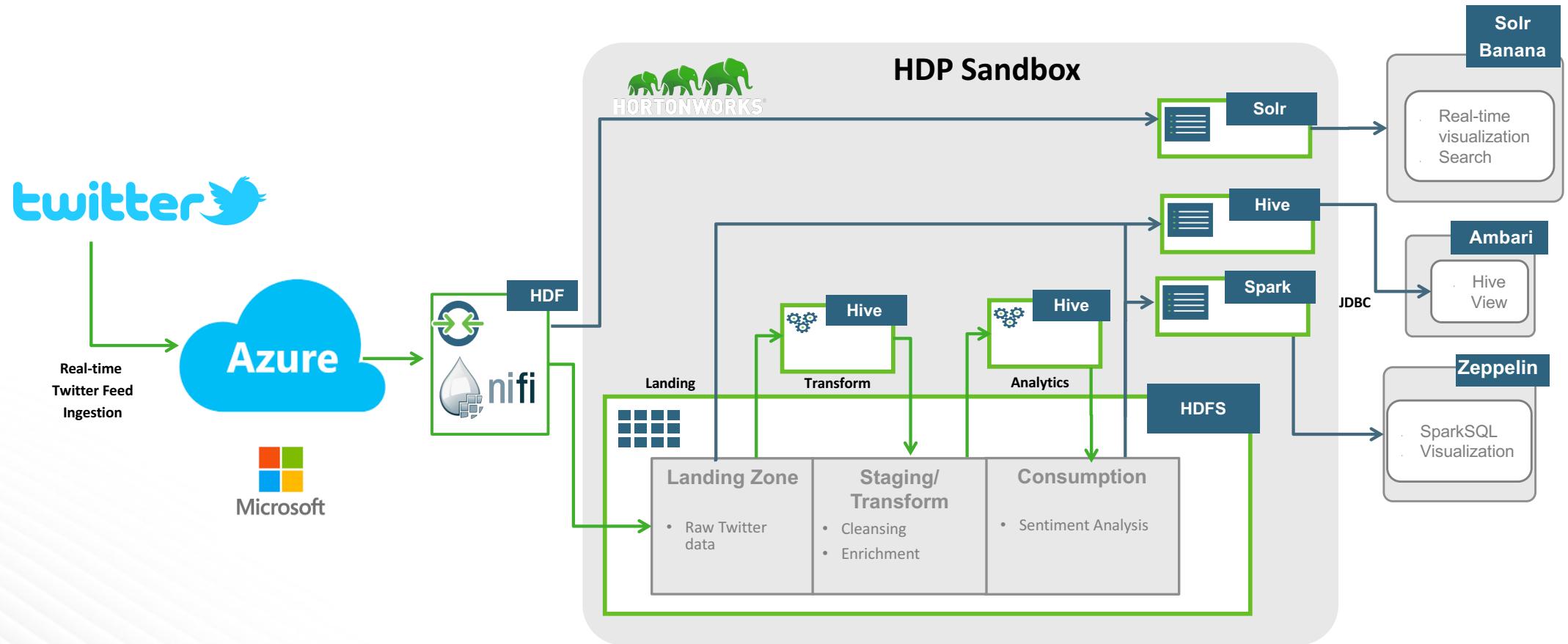


Demo of HDP and HDF

Sentiment Analysis on Twitter



Sentiment Analysis Demo Overview





Hortonworks

Thank You !