# *Insights to HDInsight*
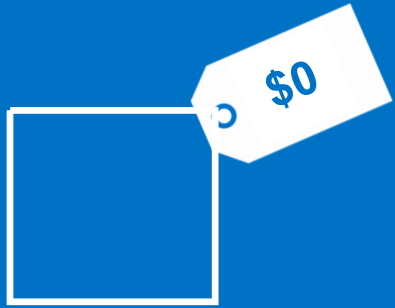
Ioannis Stavrinides, Microsoft

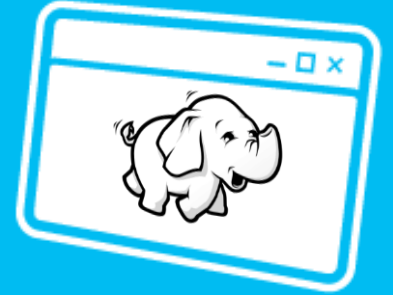# Why Hadoop in the Cloud?

No HW costs

Unlimited scale

Pay what you need

Deployed in minutes

# No hardware costs

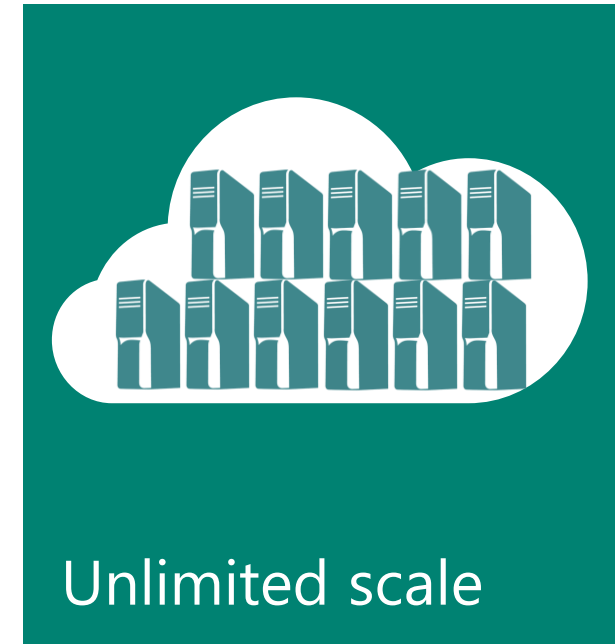## Hadoop in the Cloud bypasses hardware costs

Hardware acquisition
Hardware maintenance
Performance tuning



No HW costs

# Unlimited Scale

## Hadoop in the Cloud bypasses capacity planning

Spin up any number of Hadoop nodes on-demand

Go from tens of nodes to thousands of nodes



Unlimited scale

# Pay for What You Need

## Hadoop is billed by usage

Billed for usage

Clusters can be deleted when no longer used

Pay what you need

# Deployed in minutes
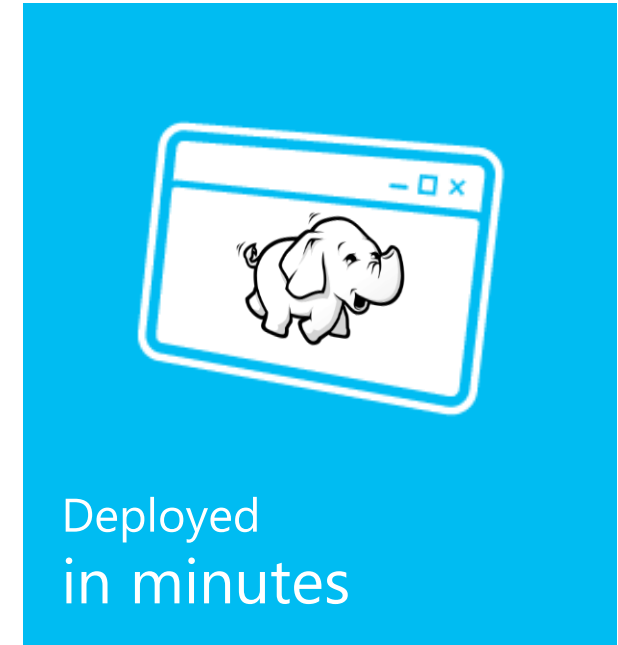
## Hadoop in the Cloud Bypasses deployment expertise

Hadoop is non-trivial to install and get up and running on multi-nodes

Education gap in IT community regarding Hadoop

## Hadoop is deployed in minutes
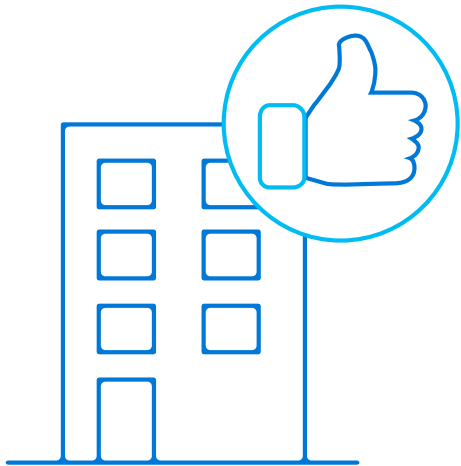
Spin up any number of Hadoop nodes on-demand

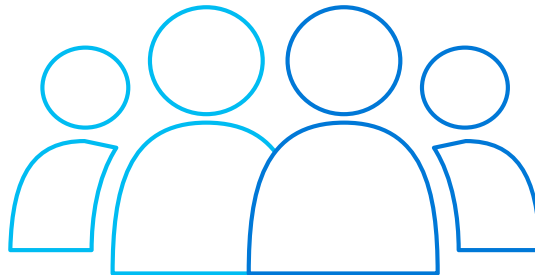Up and running in a few clicks (and within minutes)

Deployed
in minutes

# Azure HDInsight
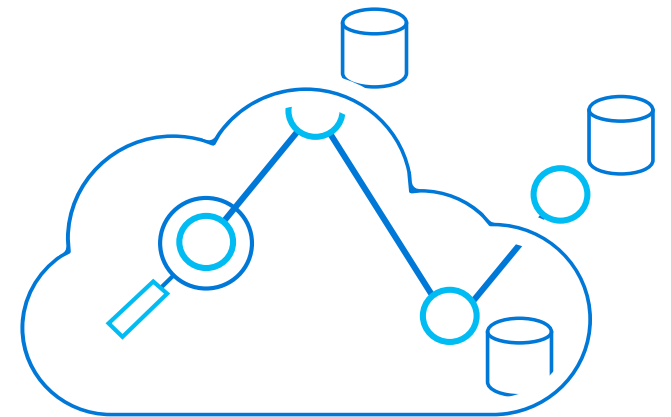## Big Data made easy

| Enterprise Ready | Easier and more productive for all users | Hybrid |
|---|---|---|

# Azure HDInsight
## Big Data made easy
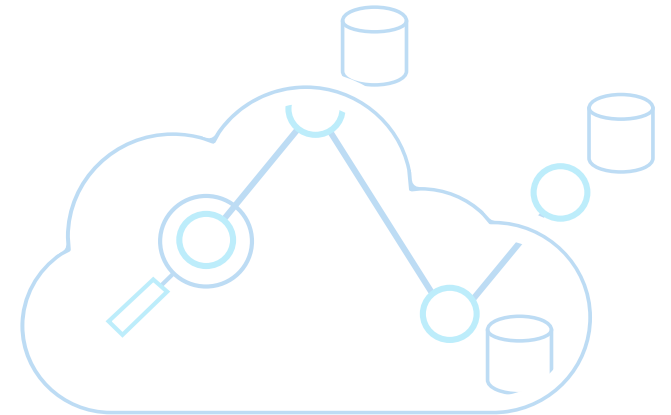
| Enterprise Ready | Easier and more productive for all users | Hybrid |
|---|---|---|

# Highly Available – Designed for the cloud ground up



- HDInsight provides primary and secondary headnodes allowing for better reliability

- Have invested in making entire stack including Resource Manager, HiverServer2 HA ready

- HDInsight stack includes Zookeeper nodes at no extra charge to customer

# Highest availability guarantee in the industry for peace of mind

99.9% SLA

- Managed, monitored and supported by Microsoft

- Enterprise-leading SLA—99.9% uptime for both VM connectivity and Hadoop running in VMs

- No IT resources needed for upgrades and patching

- Microsoft monitors your deployment so you don't have to

Microsoft

# Always encrypted, Role-based security & Auditing

- Always encrypted; in motion using SSL, and at rest using keys in Azure Key Vault

- Single sign-on, multi-factor authentication and integration of on-premises identities w/Active Directory integration

- Fine-grained ACLs for role-based access controls with Apache Ranger

- Auditing every access / configuration change with Apache Ranger

Microsoft

# Alerting, monitoring, and pre-emptive actions

- Enhanced workload protection through integration with Microsoft Operations Management Suite (OMS)

- Threat detection, monitoring, and management

# Petabyte size files and Trillions of objects



Store

EBs

TBs

- Store data in it's native format

- PB sized files, 200x larger than anyone else

- Scalable throughput for massively parallel analytics

- No need to redesign application or reparation data at higher scale

Microsoft

# Backed by Microsoft and Hortonworks



- Microsoft + Hortonworks has **37 committers** for Hadoop Core; more than all managed cloud Hadoop vendors combined

- Uniquely ready to support your deployment

- Can fix and commit code back to Hadoop

Microsoft

# Runs in the most datacenters worldwide

North Central US
*Illinois*

Central US
*Iowa*

West Europe
*Netherlands*

China North*
*Beijing*

China South*
*Shanghai*

Japan East
*Tokyo, Saitama*

North Europe
*Ireland*

West US
*California*

East US
*Virginia*

Japan West
*Osaka*

India Central
*Pune*

East US 2
*Virginia*

South Central US
*Texas*

East Asia
*Hong Kong*

SE Asia
*Singapore*

Azure doubling compute
and storage every 6 months

Australia East
*New South Wales*

Brazil South
*Sao Paulo State*

Australia South East
*Victoria*

Microsoft

# Lower total cost of ownership



- No hardware

- Hadoop support included with Azure support

- Pay only for what you use

- Independently scale storage and compute

- No need to hire specialized operations team

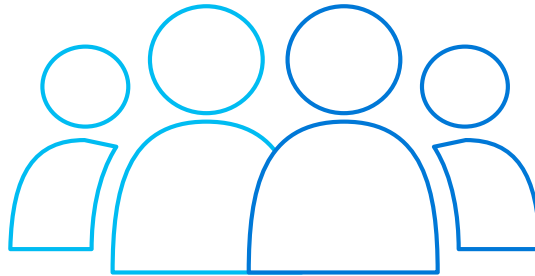- 63% lower total cost of ownership than on-premises*

Microsoft

# Azure HDInsight
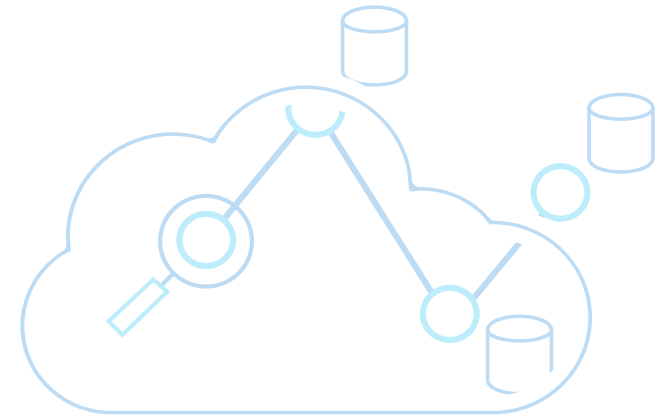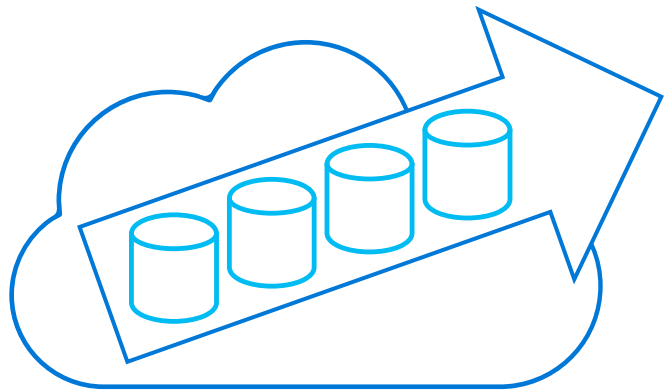## Big Data made easy

| Enterprise Ready | Easier and more productive for all users | Hybrid |

# Easy for administrators to spin up quickly

- Deploy big data projects in minutes

- No hardware to install, tune, configure or deploy

- No infrastructure or software to manage

- Scale to tens to thousands of machines instantly

# Debug and Optimize your Big Data programs with ease



- Deep integration with IDEs for developer productivity: Visual Studio, Eclipse, & IntelliJ
- Integrated with Hive, Pig, Storm, and Spark
- Visually see execution of Hive jobs ran by the Tez execution engine
- Full Intellisense

Microsoft

# Easy notebook experience for data engineers



- Most popular notebooks, Jupyter and Zeppelin out-of-the-box

- Combine code, statistical equations and visualizations

- Worked w/ Jupyter community to enhance kernel to allow Spark execution through REST endpoint

# Easy for data scientists with familiar R language



*Applies to HDInsight only

## R Server for HDInsight

- Largest portable R parallel analytics library

- Terabyte-scale machine learning—1,000x larger than in open source R

- Up to 100x faster performance using Spark and optimized vector/math libraries

- Enterprise-grade security and support

# Easy for business analysts with interactive reports over big data

- Interactive BI with big data

- Spark 2.0 integration

- Interactive Hive with LLAP-keeps data compressed running in-memory 25x faster

- ODBC driver to use Power BI or third party tools (Tableau, SAP, Qlik, etc.)

Microsoft

# Azure HDInsight
## Big Data made easy

| Enterprise Ready | Easier and more productive for all users | Hybrid |
|---|---|---|

# On-premises and cloud

- Uses Hortonworks Data Platform (HDP)

- Move projects from on-premises to cloud without code rewrite

- Hybrid scenarios supported like Dev/Test, burst, back up, disaster recovery

# Recognized by top analysts



## Forrester Wave for Big Data Hadoop Cloud

- Named industry leader by Forrester with the most comprehensive, scalable, and integrated platforms*

- Recognized for its cloud-first strategy that is paying off*

*The Forrester WaveTM: Big Data Hadoop Cloud Solutions, Q2 2016.

# Hadoop Workloads

Microsoft

# Hadoop is a platform with portfolio of projects

Governed by Apache Software Foundation (ASF)

Comprises core services of MapReduce, HDFS, and YARN

In addition to the core, includes functions across:

Governance and integration, Tools, Data Access, Security, and Operations

| Governance and integration | Tools | Security | Operations |
|---|---|---|---|

**Governance and integration**

**Data lifecycle and governance**

Falcon
Atlas

**Data workflow**

Sqoop
Flume
Kafka
NFS
WebHDFS

**Tools**

Zeppelin          Ambari User Views

**Data access**

| **Batch** | **Script** | **SQL** | **Nosql** | **Stream** | **Machine Learning** | **Others** |
|---|---|---|---|---|---|---|
| Map reduce | Pig | Hive Spark SQL | Hbase Accumulo Phoenix | Kafka Storm Spark | Sparkl Mlib Mahout | ISV engines |

**YARN: data operating system**

**HDFS  (Hadoop Distributed File System)**

**Data management**

**Security**

Authentication
Authorization
Accounting
Data protection

Ranger
Knox
Atlas
HDFS Encryption

**Operations**

**Provision, manage, and monitor**

Ambari
Zookeeper
Cloudbreak

**Scheduling**

Oozie

# HDFS

## HDFS is a distributed file system

From a few nodes to thousands of nodes
Files can be spread out over multiple nodes
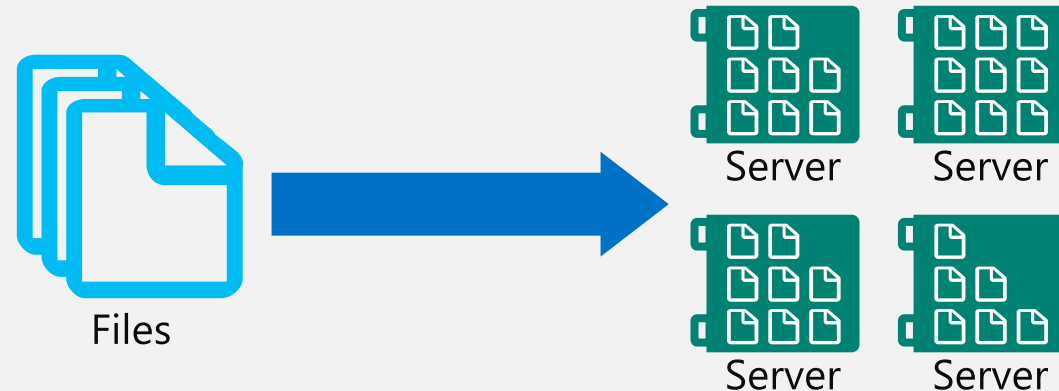
## HDFS stores large amounts of data

Very large files are supported including those larger than the capacity of a single node

## HDFS stores non-relational files

Files

Server    Server

Server    Server

# MapReduce

Batch
Map reduce
Script
ng
SQL
Hive
Spark SQL
Nosql
Hbase
Accumulo
Phoenix
Stream
Kafka
Storm
Spark
Machine
Learning
Sparkl Mlib
Mahout
Others
ISV engines

YARN: data operating system

HDFS (Hadoop Distributed File System)
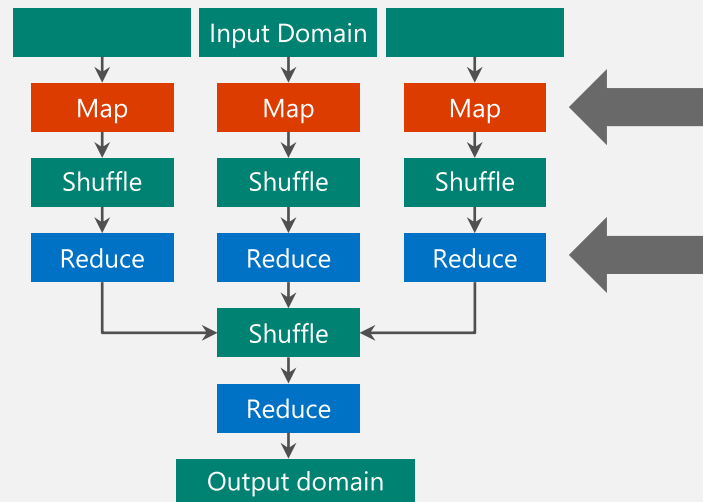
## Takes processing to where data is

Distributed processing: instead of serializing processing through one pipe, distributes computing locally where data is

Brings back only the resultant data

Scales linearly as you add nodes

## Three-step execution

Map: Developer writes map functions to the data

Shuffle / Distributes: Framework automatically shuffles for you (networking, synchronization, recovery, scheduling)

Reduce: Developer writes reduce functions to bring resultant data back



```javascript
// Map Reduce function in JavaScript

var map = function (key, value, context) {
var words = value.split(/[^a-zA-Z]/);
for (var i = 0; i < words.length; i++) {
            if (words[i] !== "")
context.write(words[i].toLowerCase(),
1);}
}};

var reduce = function (key, values, context) {
var sum = 0;
while (values.hasNext()) {
sum += parseInt(values.next());
    }
context.write(key, sum);
};
```

# Hive

## SQL-like queries on Hadoop data in HDFS

HiveQL is a SQL-like language (subset of SQL)

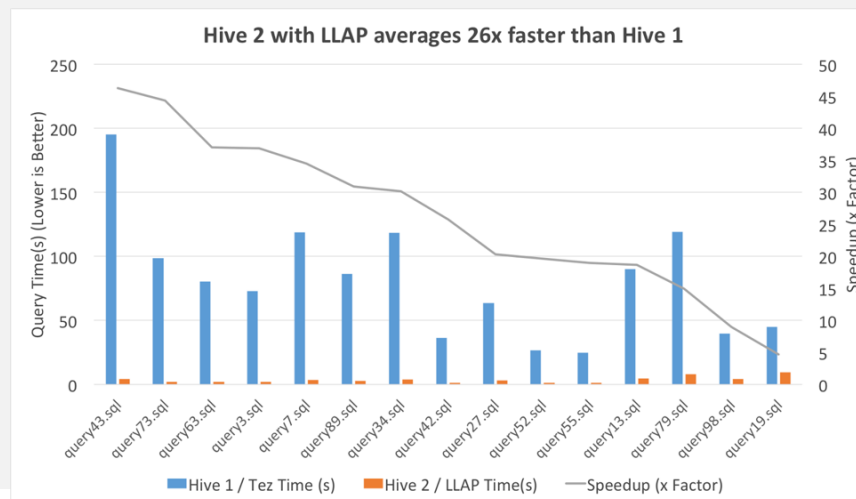Hive structures include well-understood database concepts such as tables, rows, columns, partitions

Compiled into MapReduce jobs that are executed on Hadoop

## Dramatic performance gains with Hive w/LLAP
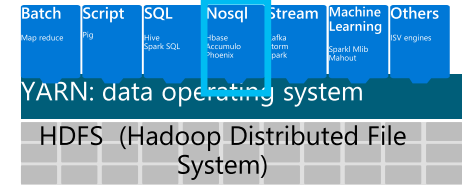
Performance gains up to 25x

ODBC drivers to integrate with Power BI, Tableau, Qlik, etc.

Opens up scenarios to do interactive BI and reporting on big data

| Batch | Script | SQL | Nosql | Stream | Machine Learning | Others |
|-------|--------|-----|-------|--------|------------------|--------|
| Map reduce | Pig | Hive Spark SQL | Hbase Accumulo Phoenix | Kafka Storm Spark | Sparkl Mllib Mahout | ISV engines |

YARN: data operating system

HDFS (Hadoop Distributed File System)

**Hive 2 with LLAP averages 26x faster than Hive 1**

Query Time(s) (Lower is Better) / Speedup (x Factor)

query43.sql query73.sql query63.sql query3.sql query7.sql query89.sql query34.sql query42.sql query27.sql query52.sql query55.sql query13.sql query79.sql query98.sql query19.sql

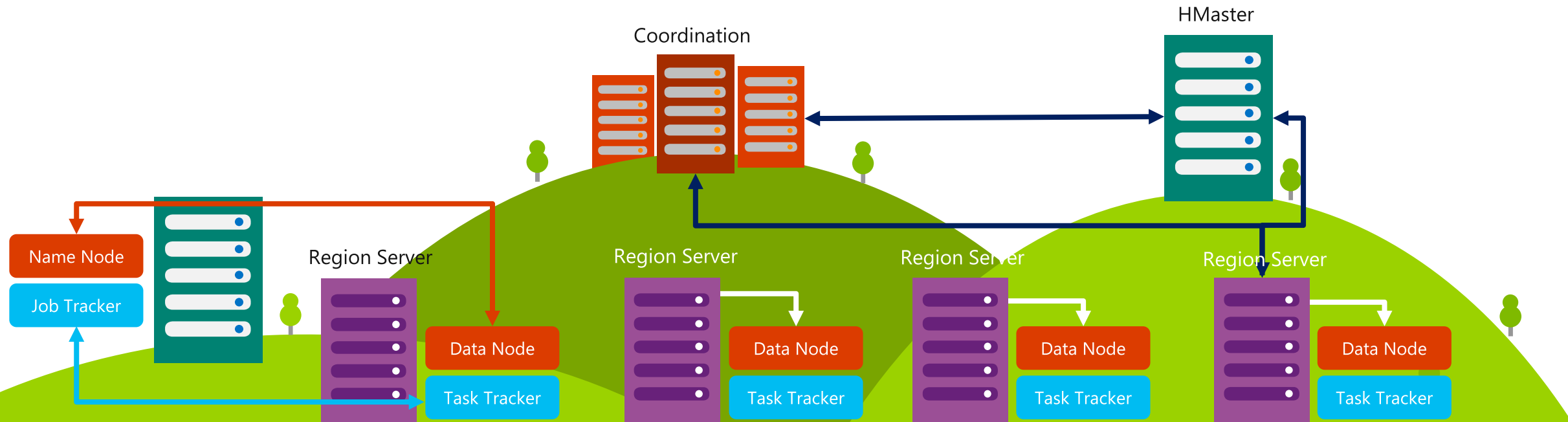■ Hive 1 / Tez Time (s)  ■ Hive 2 / LLAP Time(s)  — Speedup (x Factor)

# HBase

## NoSQL database on data in HDFS

Columnar, NoSQL database

Runs on top of the Hadoop Distributed File System (HDFS)

Provides flexibility in that new columns can be added to column families at any time
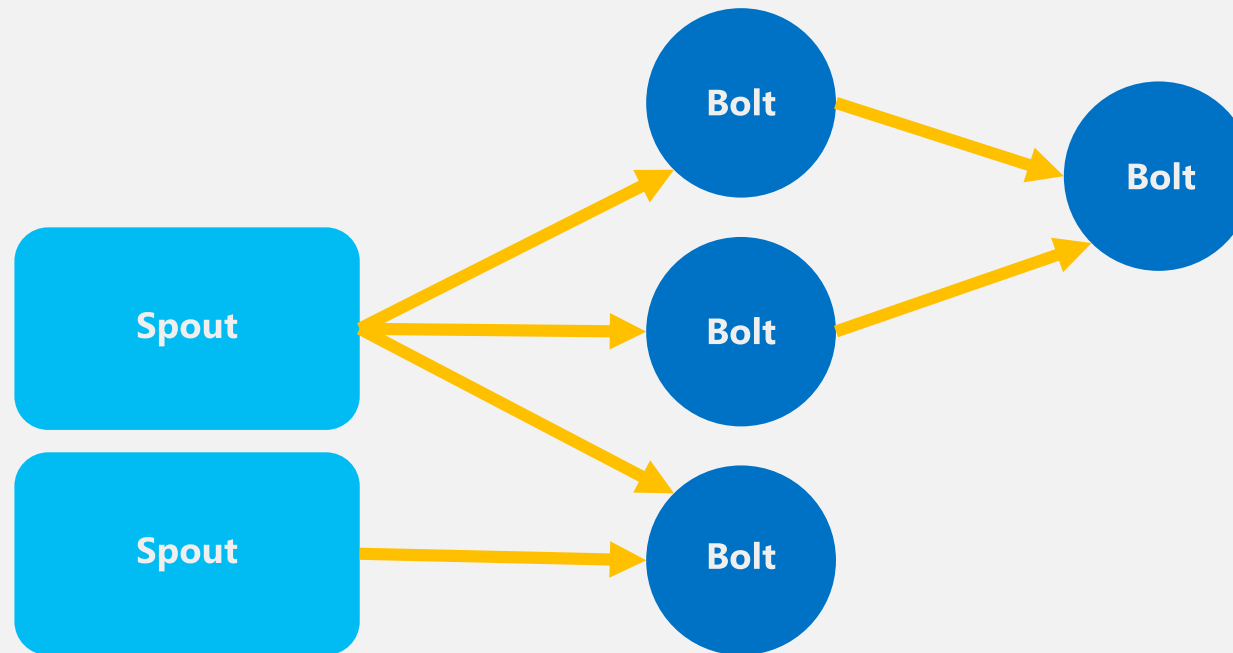
# Storm
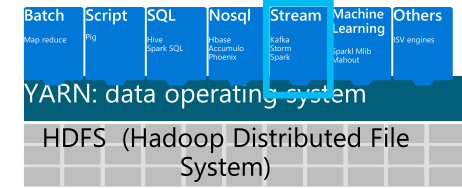


## Stream analytics for Near-Real Time processing

Consumes millions of real-time events from a scalable event broker (i.e.; Apache Kafka, Azure Event Hub)

Performs time-sensitive computation

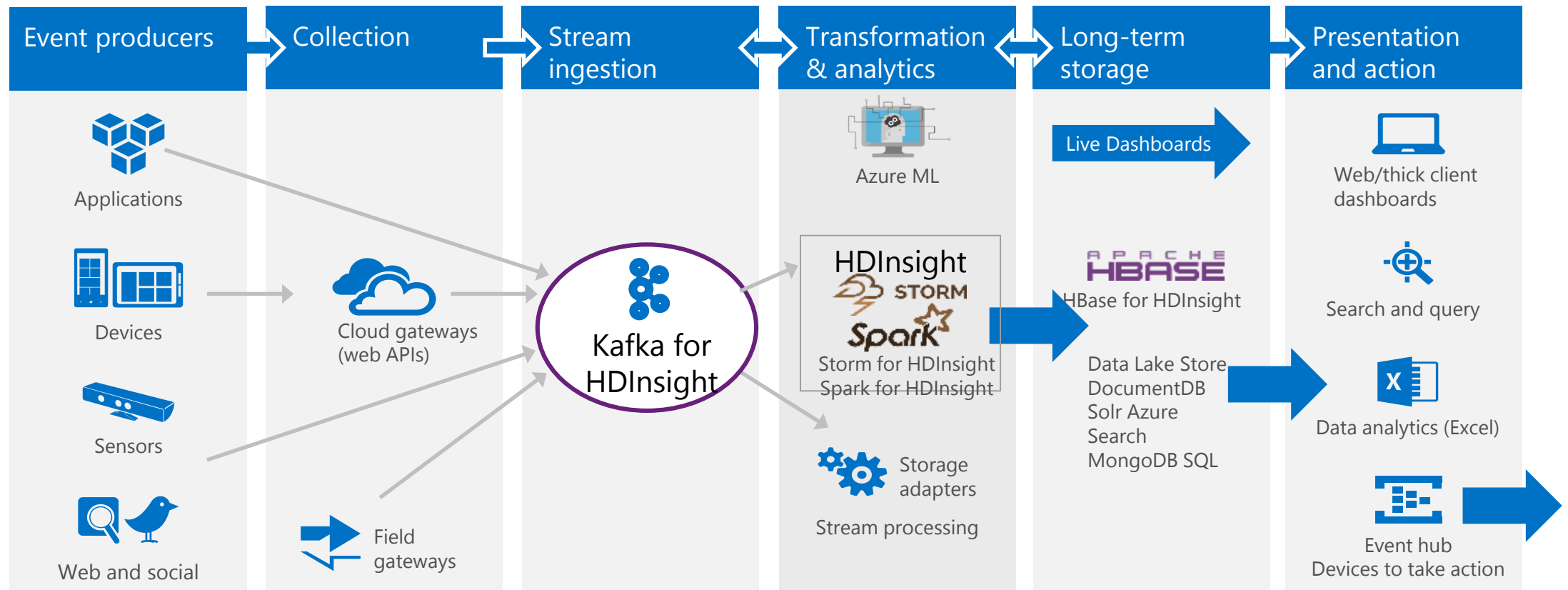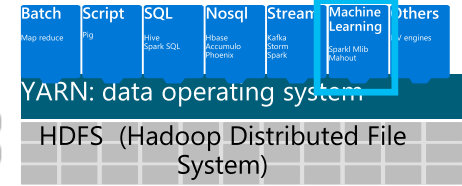Output to persistent stores, dashboards or devices

# Kafka



## High-throughput, low-latency for real-time data

Stream millions of events per second

Enterprise-grade management and control

| Event producers | Collection | Stream ingestion | Transformation & analytics | Long-term storage | Presentation and action |
|---|---|---|---|---|---|



Applications

Devices

Sensors

Web and social

Cloud gateways (web APIs)

Field gateways

Kafka for HDInsight

Azure ML

HDInsight
STORM
Spark
Storm for HDInsight
Spark for HDInsight

Storage adapters
Stream processing

Live Dashboards

APACHE HBASE
HBase for HDInsight

Data Lake Store
DocumentDB
Solr Azure
Search
MongoDB SQL

Web/thick client dashboards

Search and query

Data analytics (Excel)

Event hub
Devices to take action

# Mahout

Batch
Map reduce

Script
Pig

SQL
Hive
Spark SQL

Nosql
Hbase
Accumulo
Phoenix

Stream
Kafka
Storm
Spark

Machine Learning
Sparkl Mllb
Mahout

Others
V engines

YARN: data operating system

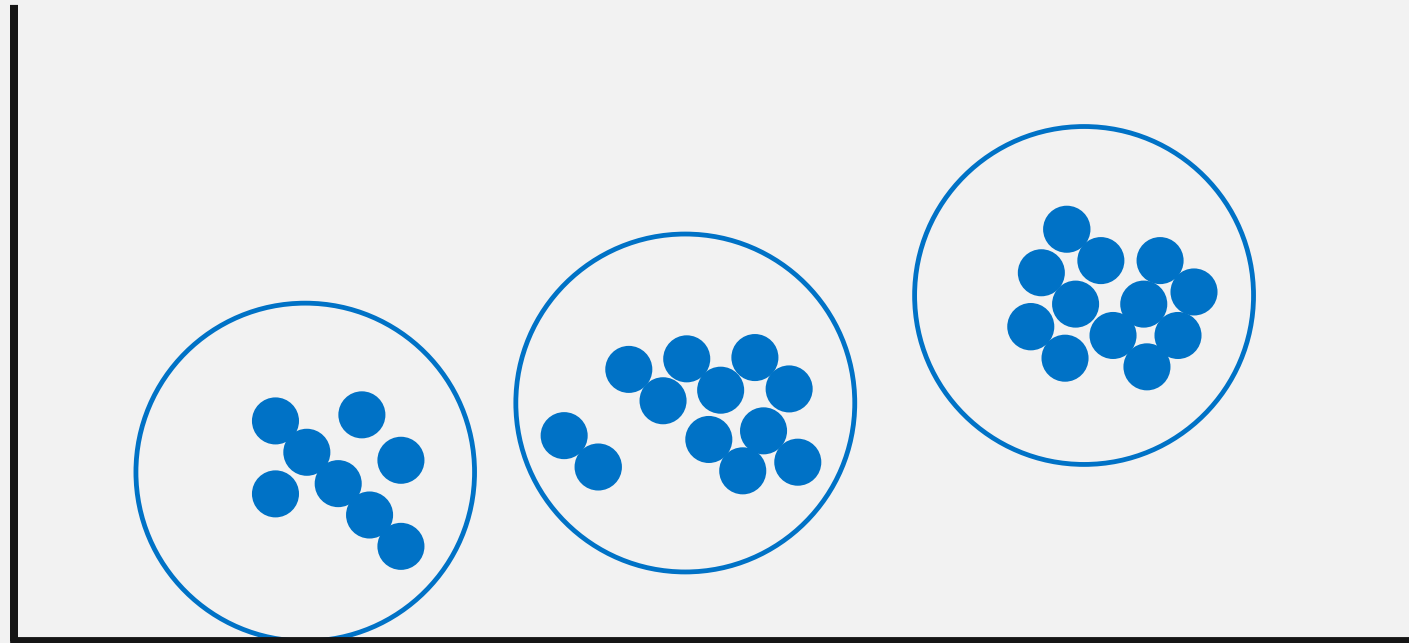HDFS (Hadoop Distributed File System)
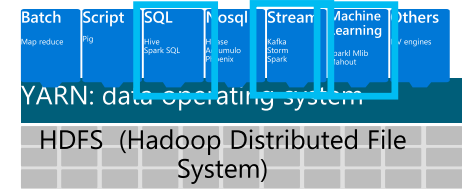
## Machine learning library

A library of machine learning algorithms to execute on data in HDFS

Algorithms are not dependent on size of data and can scale with large datasets

Library includes: Collaborative Filtering, Classification, Clustering, Dimensionality Reduction, Topic Models

# Spark



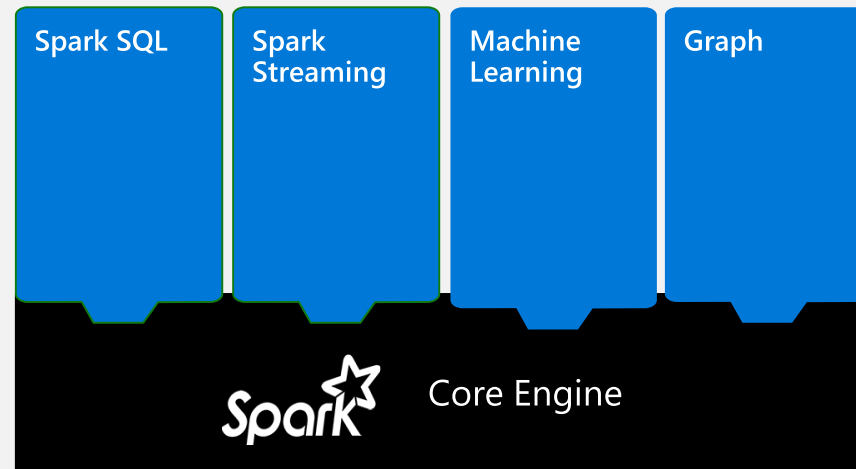# Massive data processing framework built on in-memory

Single execution model for multiple tasks

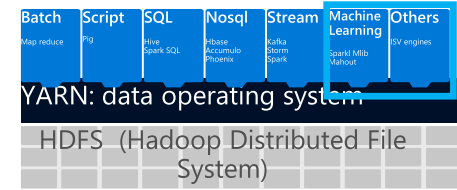Processing up to 100x faster performance

Developer friendly (Java, Python, Scala)
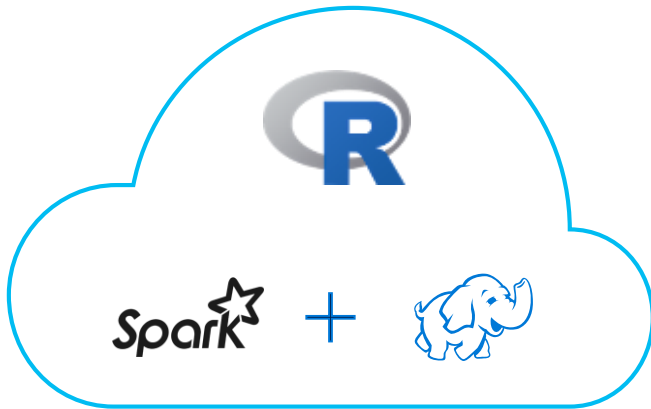
BI tool of choice (Power BI, Tabelau, Qlik, SAP)

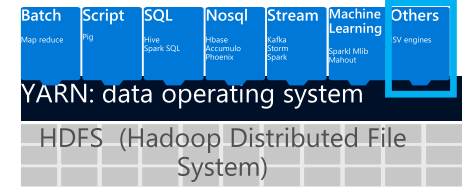Notebook experience (Jupyter & Zeppelin)

# R Server

## Predictive analytics, machine learning, and statistical modeling for big data

- Largest portable R parallel analytics library

- Terabyte-scale machine learning—1,000x larger than in open source R

- Up to 100x faster performance using Spark and optimized vector/math libraries

- Enterprise-grade security and support

# ISV Integration

Batch | Script | SQL | Nosql | Stream | Machine Learning | Others
Map reduce | Pig | Hive Spark SQL | Hbase Accumulo Phoenix | Kafka Storm Spark | Sparkl Mlib Mahout | SV engines

YARN: data operating system

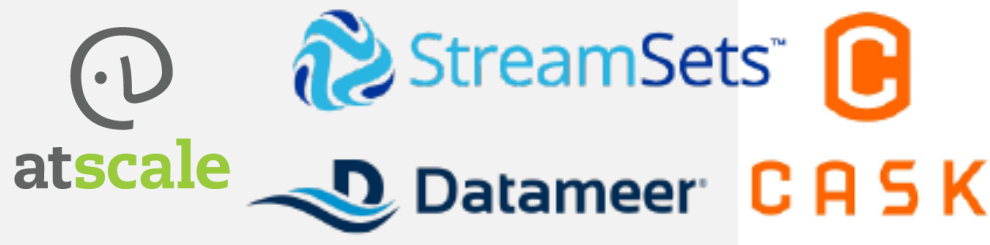HDFS (Hadoop Distributed File System)

## Integration with leading productivity applications

Spin up Hadoop and Spark clusters pre-integrated and pre-tuned with ISV applications out-of-the-box

Runs on the HDInsight clusters; does not require separate VMs

Fast and easy way to spin up applications

# Demo

Spark on HDInsight

# Business in action with Cloud solutions

Digital Business Conference, Malta 2017