# Agenda

Súčasný stav

Hadoop = HDInsight

Moderný DWH = APS

Analytika = AzureML, PowerBI

Scenáre pre BigData

# Tradičný prístup

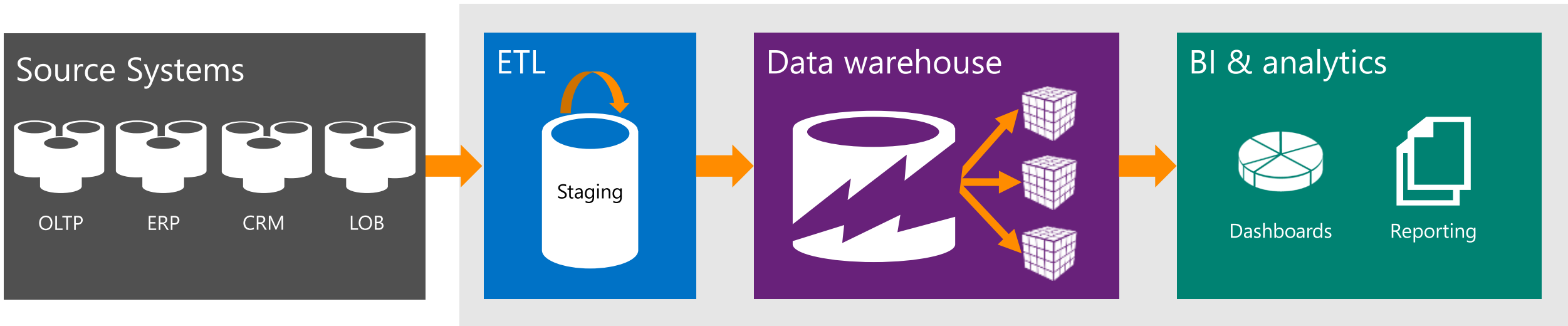① 1110 1110 1110 1010 1010 1010 1010 1010 1010
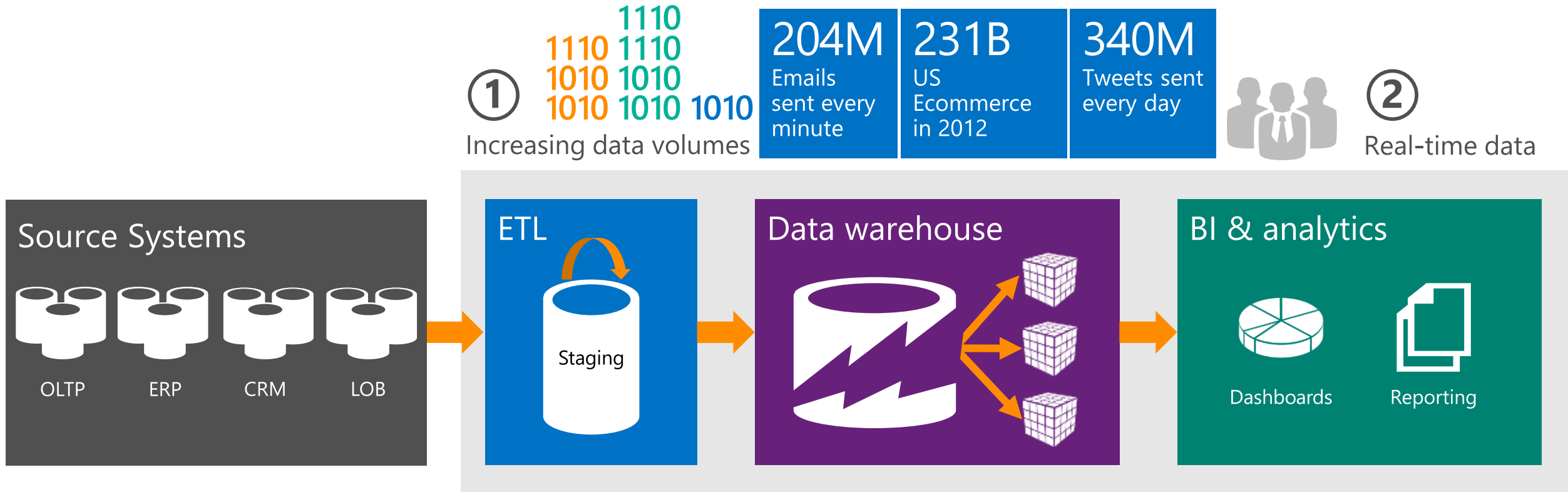Increasing data volumes

| 50x | 1Trillion | 40ZB |
|---|---|---|
| Data growth 2010-2020 | Web pages | Digital Universe 2020 |

## Source Systems

OLTP    ERP    CRM    LOB

## ETL

Staging

## Data warehouse

## BI & analytics

Dashboards    Reporting

# Breaking points of traditional approach

**①** Increasing data volumes

1110
1110 1110
1110 1110
1010 1010
1010 1010 1010

**204M** Emails sent every minute

**231B** US Ecommerce in 2012

**340M** Tweets sent every day

**②** Real-time data

## Source Systems

OLTP    ERP    CRM    LOB

## ETL

Staging

## Data warehouse

## BI & analytics

Dashboards    Reporting

# Breaking points of traditional approach

① Increasing data volumes

② Real-time data

**Source Systems**

OLTP  ERP  CRM  LOB

**ETL**

Staging

**Data warehouse**

**BI & analytics**

Dashboards  Reporting

**New Data**

Devices  Web  Sensors  Social

③ New data types

1110
1110 1110
1010 1010
1010 1010 1010

**15x** Machine generated data 2020

**2.4M** Facebook content per minute

**1.3M** Hours on Skype per hour

# Breaking points of traditional approach

## Source Systems

- OLTP
- ERP
- CRM
- LOB

## New Data

- Devices
- Web
- Sensors
- Social

① Increasing data volumes

② Real-time data

③ New data types

④ Cloud-born data

### ETL
Staging

### Data warehouse

### BI & analytics
- Dashboards
- Reporting

$100B spend on cloud

40% CRM sold are SaaS

50% large orgs have hybrid by 2017

1110
1110 1110
1010 1010
1010 1010 1010

# HDInsight v Azure (Haddop)

# HDInsight – krok 1/3

# HDInsight – krok 2/3

# HDInsight – krok 3/3

# 15' .... Hadoop cluster - running

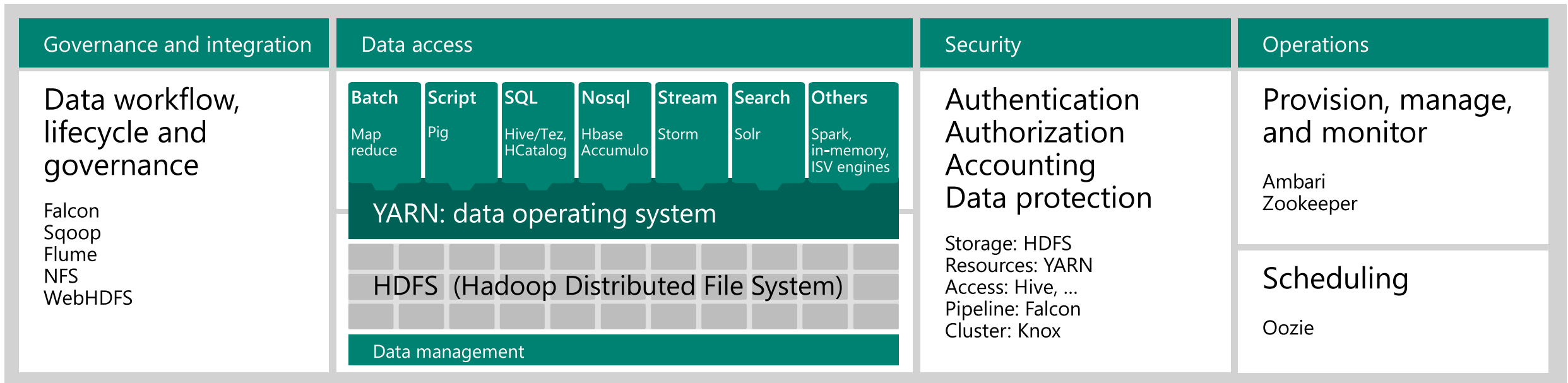# Hadoop is a platform with portfolio of projects

Governed by Apache Software Foundation (ASF)

Comprises core services of MapReduce, HDFS, and YARN

In addition to the core, includes functions across:

Data services which allow you to manipulate and move data (Hive, HBase, Pig, Flume, Sqoop)

Operational services which help manage the cluster (Ambari, Falcon, and Oozie)

| Governance and integration | Data access | | | | | | | Security | Operations |
|---|---|---|---|---|---|---|---|---|---|

**Governance and integration**

Data workflow, lifecycle and governance

Falcon
Sqoop
Flume
NFS
WebHDFS

**Data access**

| Batch | Script | SQL | Nosql | Stream | Search | Others |
|---|---|---|---|---|---|---|
| Map reduce | Pig | Hive/Tez, HCatalog | Hbase Accumulo | Storm | Solr | Spark, in-memory, ISV engines |

YARN: data operating system

HDFS  (Hadoop Distributed File System)

Data management

**Security**

Authentication
Authorization
Accounting
Data protection

Storage: HDFS
Resources: YARN
Access: Hive, …
Pipeline: Falcon
Cluster: Knox

**Operations**

Provision, manage, and monitor

Ambari
Zookeeper

Scheduling

Oozie

# A Hadoop distribution is a package of projects
## Tested for consistency across entire package

| | Hadoop and YARN | Tez | Pig | Hive and HCatalog | HBase | Phoenix | Accumulo | Storm | Mahout | Solr | Falcon | Sqoop | Flume | Ambari | Oozie | Zookeeper | Knox |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| HDP 2.1 April 2014 | 2.4.0 | 0.4.0 | 0.12.1 | 0.13.0 | 0.98.0 | 4.0.0 | 1.5.1 | 0.9.1 | 0.9.0 | 4.7.2 | 0.5.0 | 1.4.4 | 1.4.0 | 1.5.1 | 4.0.0 | 3.4.5 | .0.4.0 |
| HDP 2.0 October 2013 | 2.2.0 | | 0.12.0 | 0.12.0 | 0.96.1 | | | | 0.8.0 | | | 1.4.4 | 1.3.0 | 1.4.4 | 3.3.2 | 3.4.5 | .0.4.0 |
| HDP 1.3 May 2013 | 1.1.2 | | 011.0 | 0.11.0 | 0.94.6 | | | | 0.7.0 | | | 1.4.3 | 1.3.1 | 1.2.5 | 3.3.2 | 3.4.5 | .0.4.0 |

Data management | Data access | Governance and integration | Operations | Security

# Microsoft Analytics Platform System
Appliance pre moderný Datawarehouse

Analytics Platform System

ETL/ELT with SSIS, DQS, MDS

ERP  CRM  LOB  APPS

ETL/ELT with DWLoader

Hadoop / Big Data

Ad hoc queries

PDW

PolyBase
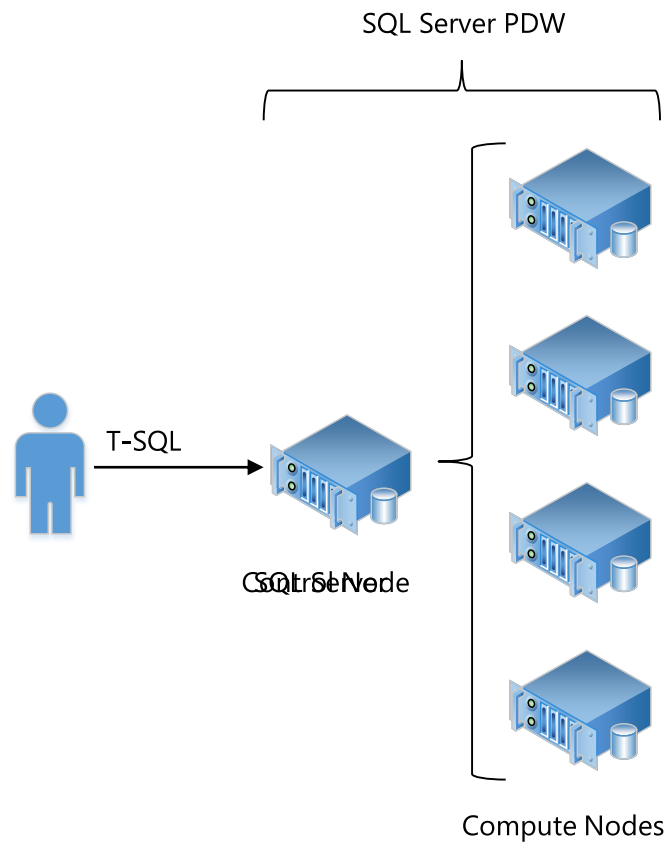
HDInsight

SQL Server SMP

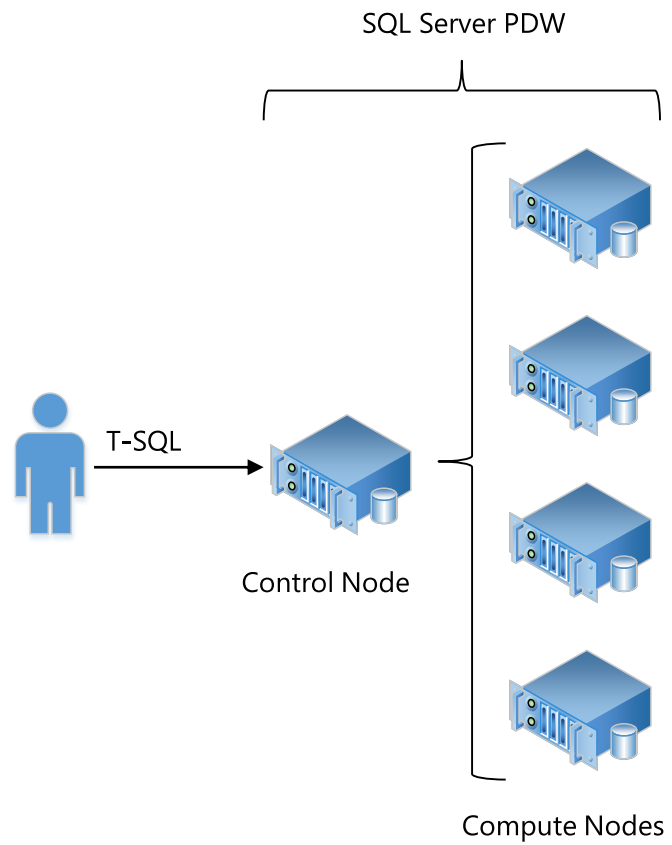SSRS / SSAS

BI Tools

# Čo je Parallel Data Warehouse?

- ## Shared-nothing parallel database system
  - » Massively parallel processing (MPP)
  - » A "Control" server that accepts user queries, generates a plan, and distributes operations in parallel to compute nodes
  - » Multiple "Compute" servers running SQL Server
  - » A "Management" server for administering the system
  - » A "Data Movement Service" that facilitates parallel SQL operations

- ## Delivered as an appliance
  - » Balanced and pre-configured software and industry standard hardware from *HP*
  - » Single Call Support
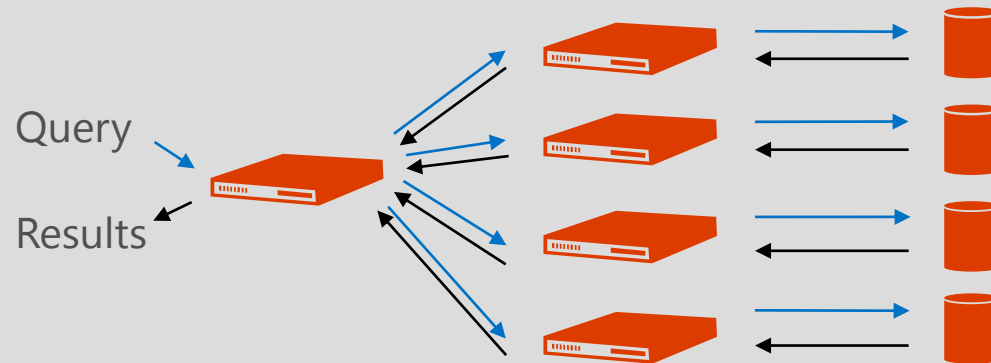  - » Fastest Time to Market
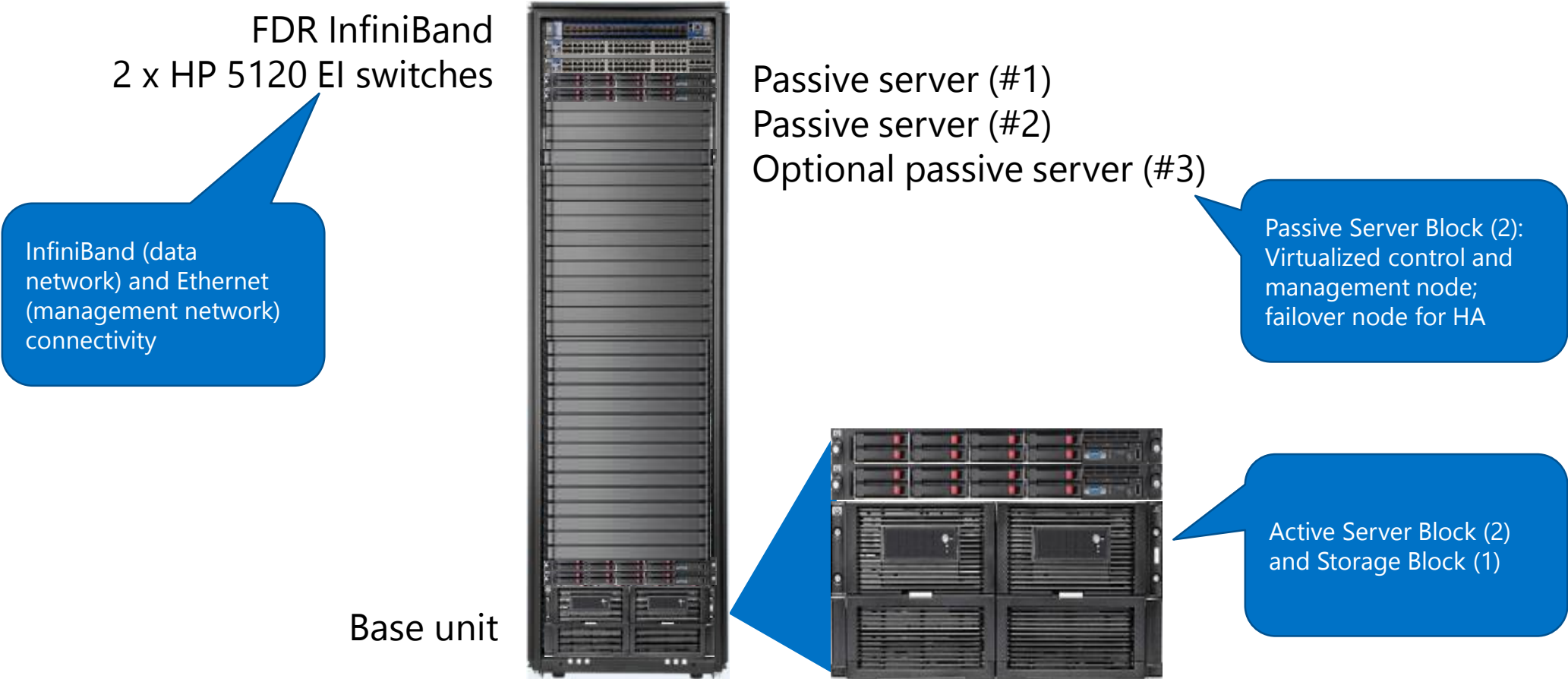  - » Scales from 2 to 56 Nodes

# SQL Server Parallel Data Warehouse

SQL Server PDW

T-SQL

Control Node
SQL Server

Compute Nodes

# SQL Server Parallel Data Warehouse

# HP ConvergedSystem 300 for Microsoft Analytics Platform Base unit

FDR InfiniBand
2 x HP 5120 EI switches

InfiniBand (data network) and Ethernet (management network) connectivity

Passive server (#1)
Passive server (#2)
Optional passive server (#3)

Passive Server Block (2): Virtualized control and management node; failover node for HA

Active Server Block (2) and Storage Block (1)

Base unit

# Microsoft Analytics Platform System

Appliance pre moderný datawarehouse

Analytics Platform System

ETL/ELT with SSIS, DQS, MDS

ERP · CRM · LOB · APPS

ETL/ELT with DWLoader

Hadoop / Big Data

Ad hoc queries

PDW

PolyBase

HDInsight

SQL Server SMP

SSRS / SSAS

BI Tools

# SQL Server Parallel Data Warehouse – Appliance Update 1
## Extending the distributed Data Warehouse further

PDW 2008 R2   PDW 2012     PDW 2012 AU1
PDW 2012 AU1

T-SQL

Control Node
SQL Server

Compute Nodes

PolyBase

HDFS

Windows Azure Blob Storage

Hadoop on-premise

MapReduce

# Query Hadoop data with T-SQL using PolyBase

Bringing the worlds or big data and the data warehouse together for users and IT

Select...

Result set

SQL Server Parallel Data Warehouse

PolyBase

Microsoft HDInsight

Windows Azure HDInsight

Cloudera

Hortonworks (Windows, Linux)

Single T-SQL query model for PDW and Hadoop with rich features of T-SQL including joins without ETL

Leverages the power of MPP to enhance query execution performance

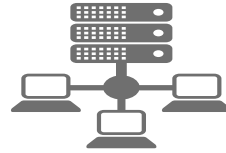Supports Windows Azure HDInsight to enable new hybrid cloud scenarios

Query non-Microsoft Hadoop distributions such as Hortonworks and Cloudera

# Access Hadoop on different cluster (cloud or on premise)
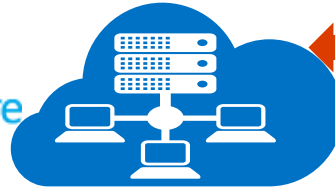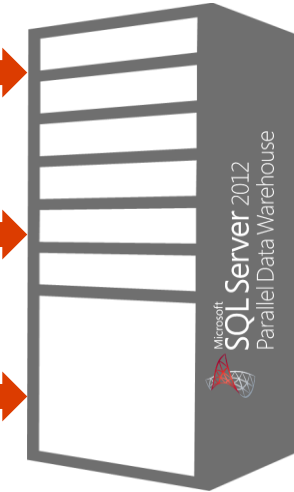


Hortonworks Data Platform (HDP) on Windows or Linux
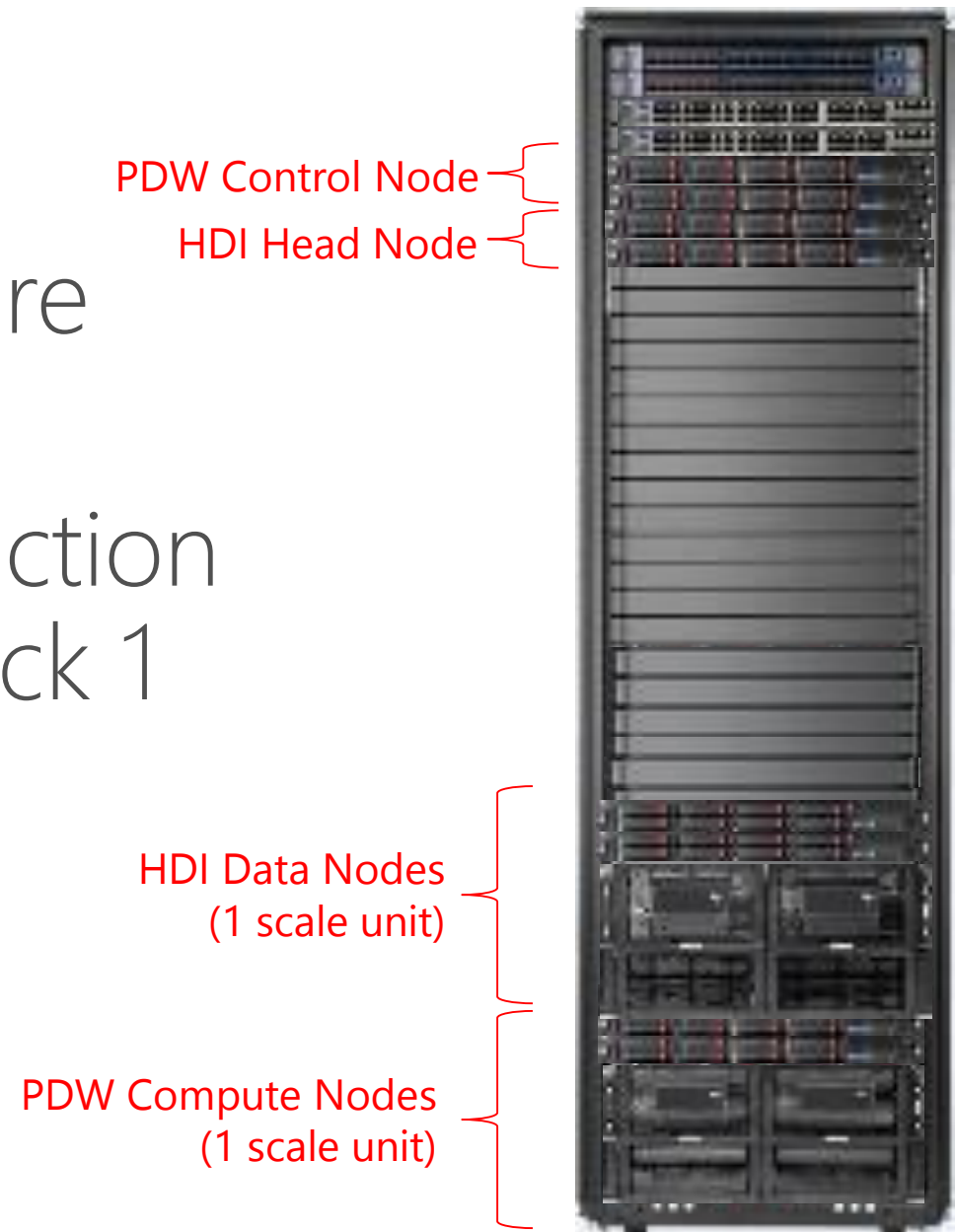
HDInsight on Azure (rebranded HDP)

**HDInsight is the Microsoft branded Hortonworks Data Platform**
- **We made it work on Windows**
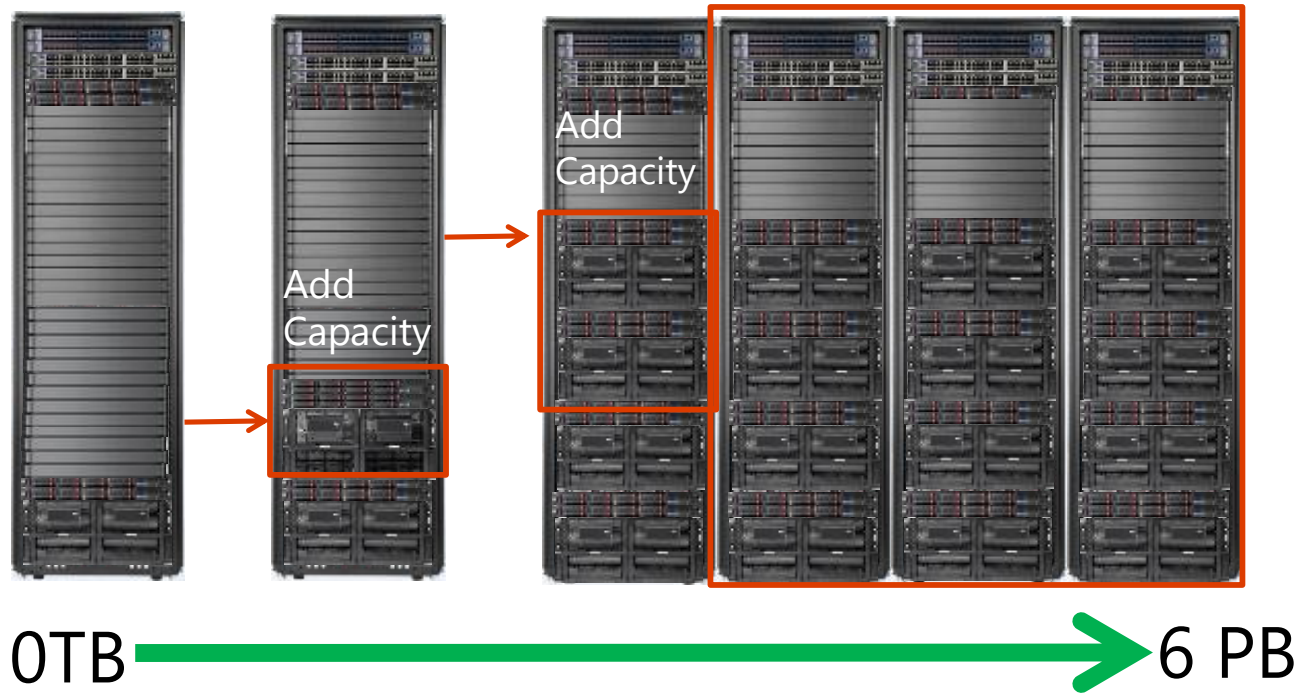- **We brougth SystemCenter support to it**

# Hardware topology overview



- Uses the same PDW hardware and topology
- The key difference is introduction of 2 additional servers on rack 1 for the HDI Head Node
    - 1 active server and 1 fail over server.

PDW Control Node

HDI Head Node

HDI Data Nodes
(1 scale unit)

PDW Compute Nodes
(1 scale unit)

# Seamlessly add capacity

Scale from a Quarter Rack with 2 Compute Nodes up to 56 Compute Nodes!

Add Capacity

Add Capacity

Add Capacity

0TB

6 PB

## Smallest (0TB) To Largest (6PB)

- Start small with a few Terabyte warehouse
- From 2 compute nodes to 56 compute nodes
- 1 quarter rack up to 7 full racks
- Add capacity up to 6 Petabytes

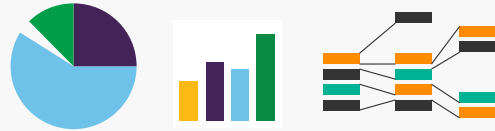Start Small And Grow

Largest Warehouse

**PB**

Minimal Downtime
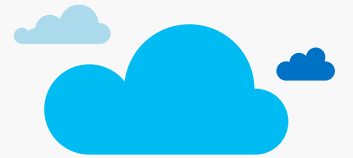
# Bringing Hadoop to a billion people

**Excel** as the BI tool for everyone

**Power** BI for collaboration & new experiences

1 Billion Microsoft Office users
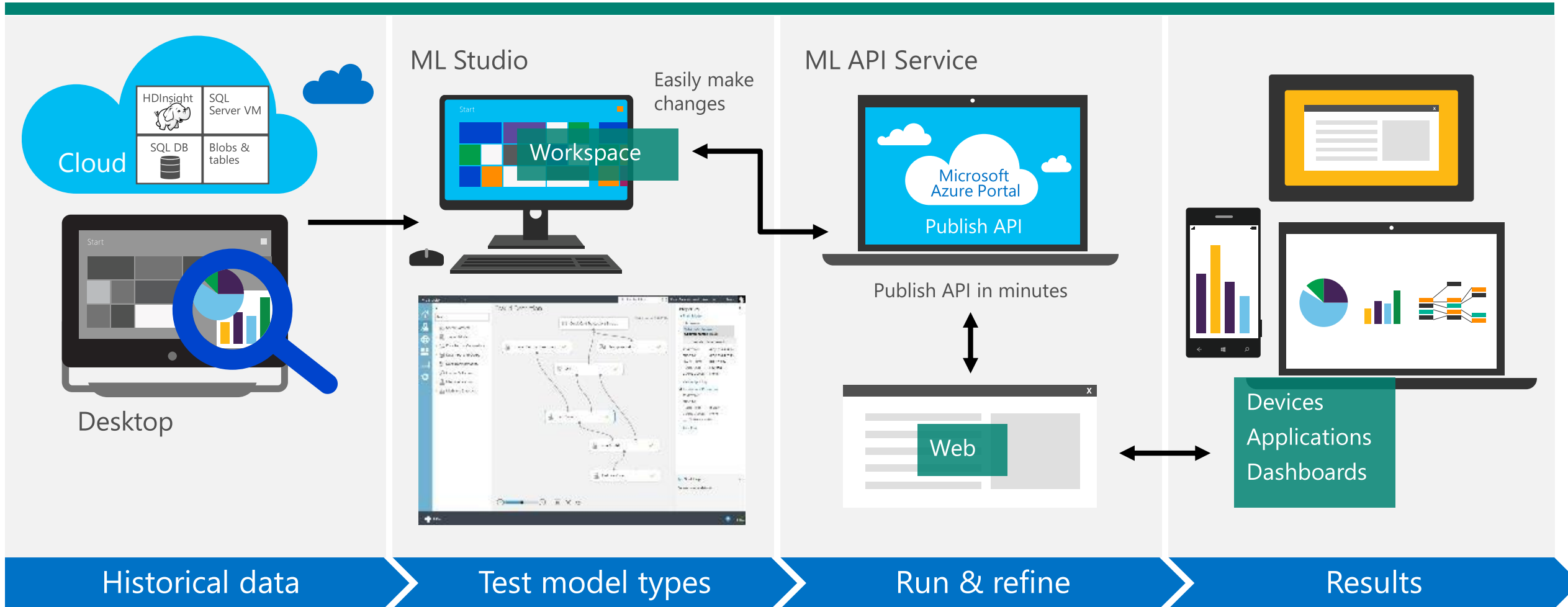- Connect to HDInsight
- Analyze
- Visualize

Office 365 is our fastest-growing commercial product ever
- Share
- Ask
- Access

Scalable, manageable, trusted

# Making advanced analytics accessible to Hadoop
## Microsoft Azure Machine Learning



Cloud

HDInsight | SQL Server VM

SQL DB | Blobs & tables

Desktop

ML Studio

Start

Workspace

Easily make changes

ML API Service

Microsoft Azure Portal

Publish API

Publish API in minutes

Web

Devices
Applications
Dashboards

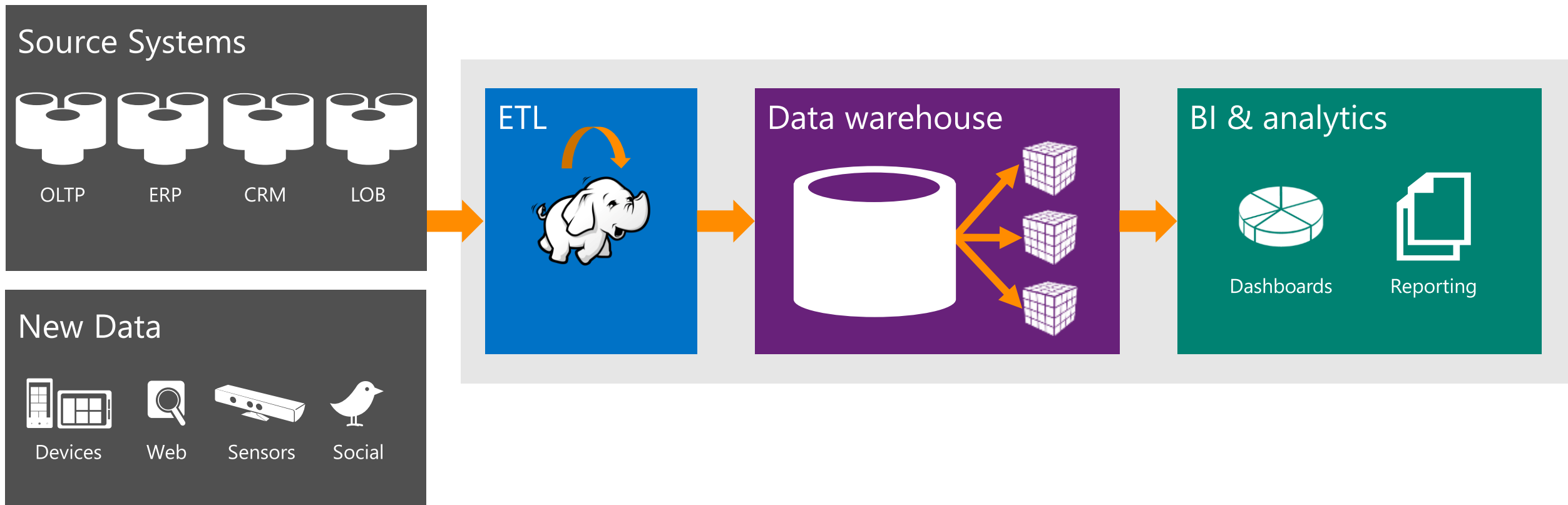| Historical data | Test model types | Run & refine | Results |

# PowerBI, Excel PowerMap, PowerQuery, ...

# Hadoop Scenario 1: pre-process ETL

Shift the pre-processing of ETL in staging data warehouse to Hadoop
Shifts high cost data warehousing to lower cost Hadoop clusters

**Source Systems**

OLTP    ERP    CRM    LOB

**New Data**

Devices    Web    Sensors    Social

**ETL**

**Data warehouse**
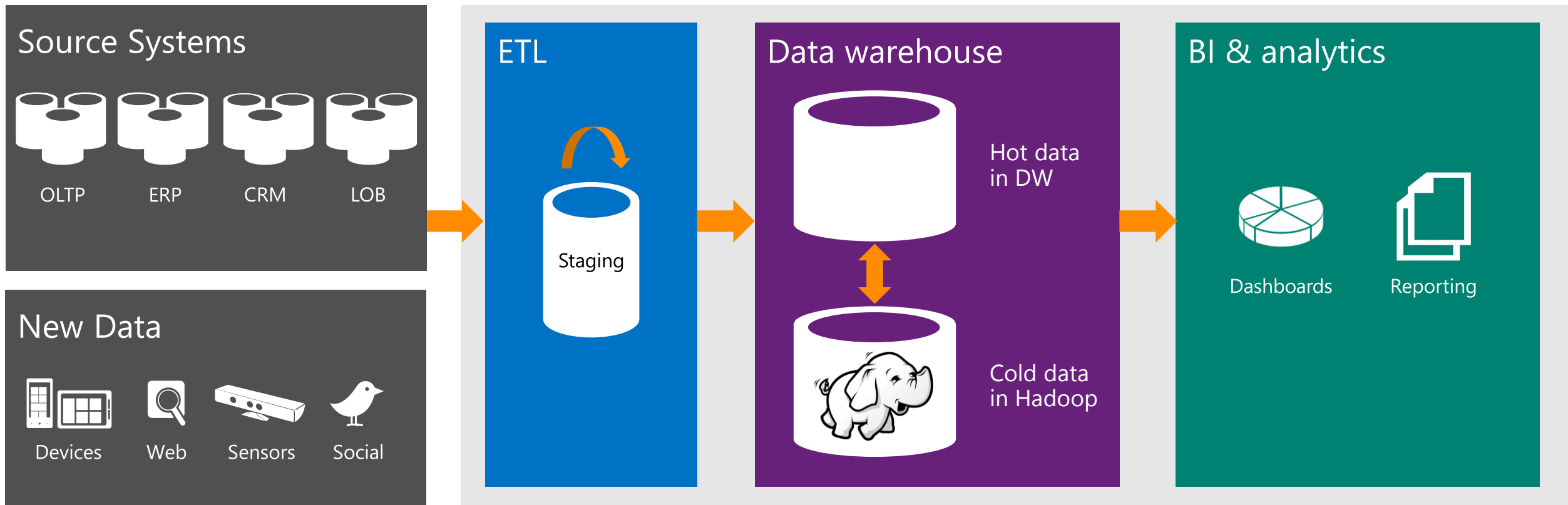
**BI & analytics**

Dashboards    Reporting

# Hadoop Scenario 2: hot and cold storage

Offloading large volume of historical data into cold storage with Hadoop

Keep data warehouse for hot data to allow BI and analytics

When data from cold storage is needed, it can be moved back into the warehouse



**Source Systems**

OLTP   ERP   CRM   LOB

**New Data**

Devices   Web   Sensors   Social

**ETL**

Staging

**Data warehouse**

Hot data in DW

Cold data in Hadoop

**BI & analytics**
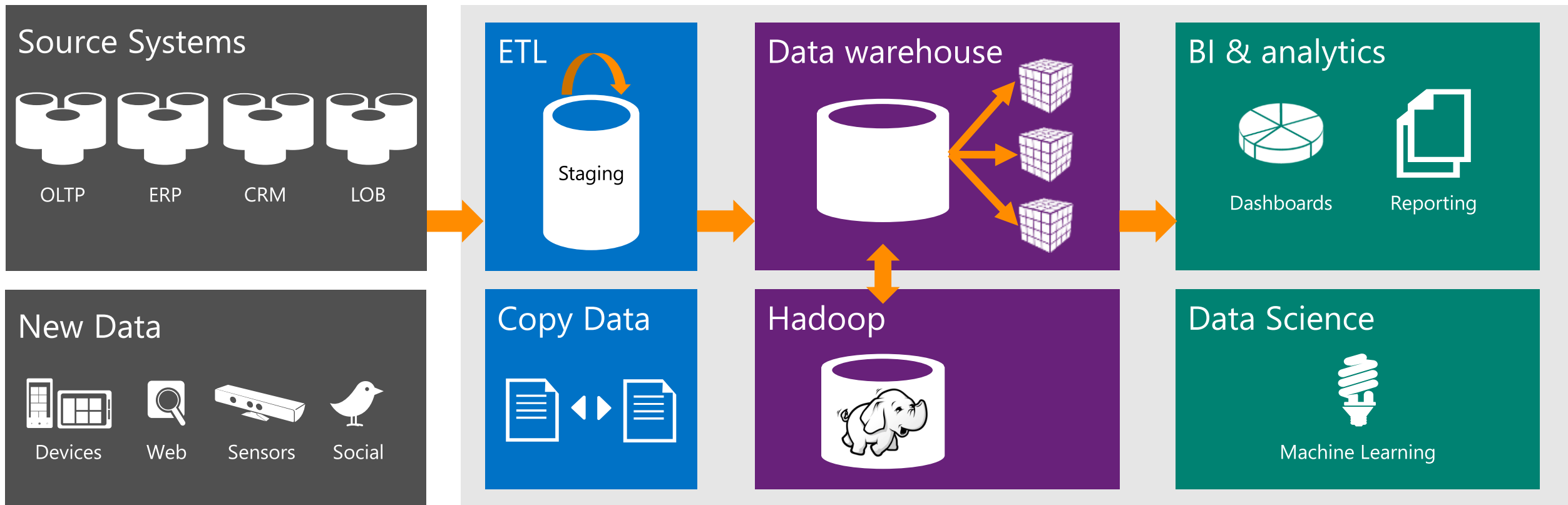
Dashboards   Reporting

# Hadoop Scenario 3: true data discovery

Keep data warehouse for operational BI and analytics
Allow data scientists to gain new discoveries on raw data (no format or structure)
Operationalize discoveries back into the warehouse

# Industry Use Cases of Hadoop

## Financial services
New account risk screens
Fraud prevention
Trading risk
Maximize deposit spread
Insurance underwriting
Accelerate loan processing

## Retail
360˚ view of the customer
Analyze brand sentiment
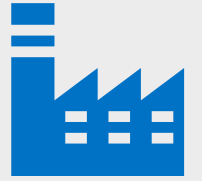Localized, personalized promotions
Website optimization
Optimal store layout

## Telecom
Call detail records (CDRs)
Infrastructure investment
Next product to buy (NPTB)
Real-time bandwidth allocation
New product development

## Manufacturing
Supplier consolidation
Supply chain and logistics
Assembly line quality assurance
Proactive maintenance
Crowd source quality assurance

## Healthcare
Genomic data for medical trials
Monitor patient vitals
Reduce re-admittance rates
Store medical research data
Recruit cohorts for pharmaceutical trials

## Utilities, oil and gas
Smart meter stream analysis
Slow oil well decline curves
Optimize lease bidding
Compliance reporting
Proactive equipment repair
Seismic image processing

## Public sector
Analyze public sentiment
Protect critical networks
Prevent fraud and waste
Crowd source reporting for repairs to infrastructure
Fulfill open records requests

**Microsoft**

# Get started today!

- For more information visit: http://azure.microsoft.com/en-us/services/hdinsight/