# BigData

**An Introduction**

TLS

# Agenda



- The Data Voyage
- Understanding BigData
- BigData Enablers
- BigData use cases
- BigData Market

# The Data Voyage

TLS

# Understanding Data

What is data? A Customer's record, Subscriber's telephone # and address, a retail ticket, an insurance claim, an inventory list, a sales order, a flight ticket…etc.



This is how we perceive data today. But data is changing its form, even right now as we speak, data is changing.

# Data Amorphism

The manifestation of data has changed dramatically over the last few years. Data today is quite different from what data used to be a few years ago!

# The Data Challenge

- This new data is different from the regular data in a couple of respects, not only is it often unstructured, it has high volumes and in many cases changes rapidly

- As Gartner defines it, this new data has

- High Volume

- High Velocity

- High Variety

- High Complexity

# The Data Challenge… contd.

The industry has tried to tackle this "Data Challenge" from different technology and process perspectives. Some methods to improve the handling of data and others to directly improve the subsequent usage (Analytics) of it.

- Grid Computing (Volume)

- Cloud Computing (Velocity & Volume)

- Using Open Source Software (overall cost)

- Virtualization (Time to Test/Deploy, Processing Capacity)

- Innovations in BI – Realtime & Operational (Velocity)

- Dw Appliances (Volume)

- Industry specific models and solutions (Time to Deploy, Complexity)

**TLS**

# The Paradigm Shift

- Clouds are slow, Grids are expensive, Open Source less robust, Virtual Servers tend to slow things down, Dw Appliances bring in a lot of hardware…



Hadoop is an <u>Open Source</u> data storage and processing Architecture.

# The Paradigm Shift…contd.

By the early part of the last decade, Internet based Cos. Like Google, Yahoo, Linkedin, Amazon and others adopted Bigdata and invented programming models, frameworks to suit their needs and fill the knowledge gap.

With Hadoop & its associated technologies they could

- Handle terabytes & petabytes of data

- which was unstructured – text, audio, video & various formats

- on commodity hardware

- assuring fail-safe (read and write) operations

-  with high performance

- and  economy (remember its Open Source!)

More on Hadoop in the coming slides.

# Some facts - Wikibon

## Big Data Market Forecast, 2011-2026 ($US B)

# Some facts - eWeek

- Allied Market Research forecast that the global Hadoop market value will reach **$50.2 billion by 2020**
- Hadoop Will Be Used for Over 10 Percent of Data Processing and Storage
- Hadoop Will Lead in Infrastructure Spending
- Hadoop Will Not Be the Greatest Expense of the IT Budget
- Hadoop Will Democratize Data
- Hadoop Will Accelerate Big Data Adoption
- Hadoop Will Power the Startup Economy
- Hadoop Will Be Used for Critical Day-to-Day
- Hadoop Will Advance the Internet of Things
- Hadoop Will Be Used for Processing and Storing Highly Sensitive Data
- Hadoop Will Power the Connected World

http://www.eweek.com/database/slideshows/hadoop-2020-the-future-of-big-data-in-the-enterprise.html#sthash.FJ5FW6Kt.dpuf

# Understanding Bigdata

# BigData



Big Data = Transactions + Interactions + Observations

Source: Contents of above graphic created in partnership with Teradata, Inc.

# Bigdata Technology Landscape

# Conventional BI Landscape

| Access | Analytics | Data Repositories | Data Integration | Data Sources |
|---|---|---|---|---|
| Web Browser | **Business Applications** / Collaboration | Operational Data Stores | Extraction | Enterprise |
| Portals | Query and Reporting | Data Warehouses | Initial Staging | Unstructured |
| Devices | Data Mining | Data Marts | Data Quality | Informational |
| Web Services | Modeling | Staging Areas | Calculation and Splits | External |
| | Scorecard | Metadata | Clean Staging | |
| | Visualization | | Processing and Enrichment | |
| | Embedded Analytics | | Target Filtering | |
| | | | Load-Ready Publish | |
| | | | Load | |

**Data Flow and Workflow**

Data Governance

Metadata

Data Quality

Network Connectivity, Protocols and Access Middleware

Hardware and Software Platforms

**TLS**

# What happens to the old data order?

# Is this the end of the road for Conventional BI?

**TLS**

# Hadoop Clusters–Storage blocks of Bigdata



Cluster of machines running Hadoop at Yahoo (Source: Yahoo)

# The Hadoop Ecosystem

**Data Analytics**
**(Business Intelligence and Analytic tools )**

| Oozie (Workflow) | Chukwa (Monitoring) | Flume (Monitoring) | Zookeeper (Mgmt) |
|---|---|---|---|

| Hive (Sql) | Pig (Data Flow) | Avro (Serialization) | Mahout (Machine Learning) | Sqoop Sql to Hadoop |
|---|---|---|---|---|

**MapReduce**

**Hbase**

**Hadoop Distributed File System(HDFS)**

**TLS**

# Bigdata Enablers

TLS

# Oracle Bigdata Reference Landscape



Oracle Integrated Software Solution Stack

# Oracle Bigdata Reference Landscape

**Oracle Engineered Solutions**

Data Variety

↑

Unstructured

**Big Data Appliance**
- Hadoop
- NoSQL Database
- Oracle Loader for hadoop
- Oracle Data Integrator

**Oracle Exadata**
- Image
- XML
- Text
- Semantic  Graphs
- Spatial

**Exalytics**
- Speed of Thought Analytics

Schema

↓

Information Density

Acquire          Organize          Analyze

ORACLE

# Oracle Bigdata Reference Landscape

**Another perspective of Oracle BigData Solution**

# Oracle Bigdata Appliance Architecture

- The Oracle Bigdata Appliance software includes
- Full distribution of Cloudera's version of Hadoop including Apache Hadoop (CDH4)
- Oracle Bigdata Appliance Plug-In for Enterprise Manager
- Cloudera Manager to administer all aspects of Cloudera CDH
- Oracle distribution of the statistical package R
- Oracle NoSQL Database Community Edition2
- Oracle Enterprise Linux operating system and Oracle Java VM
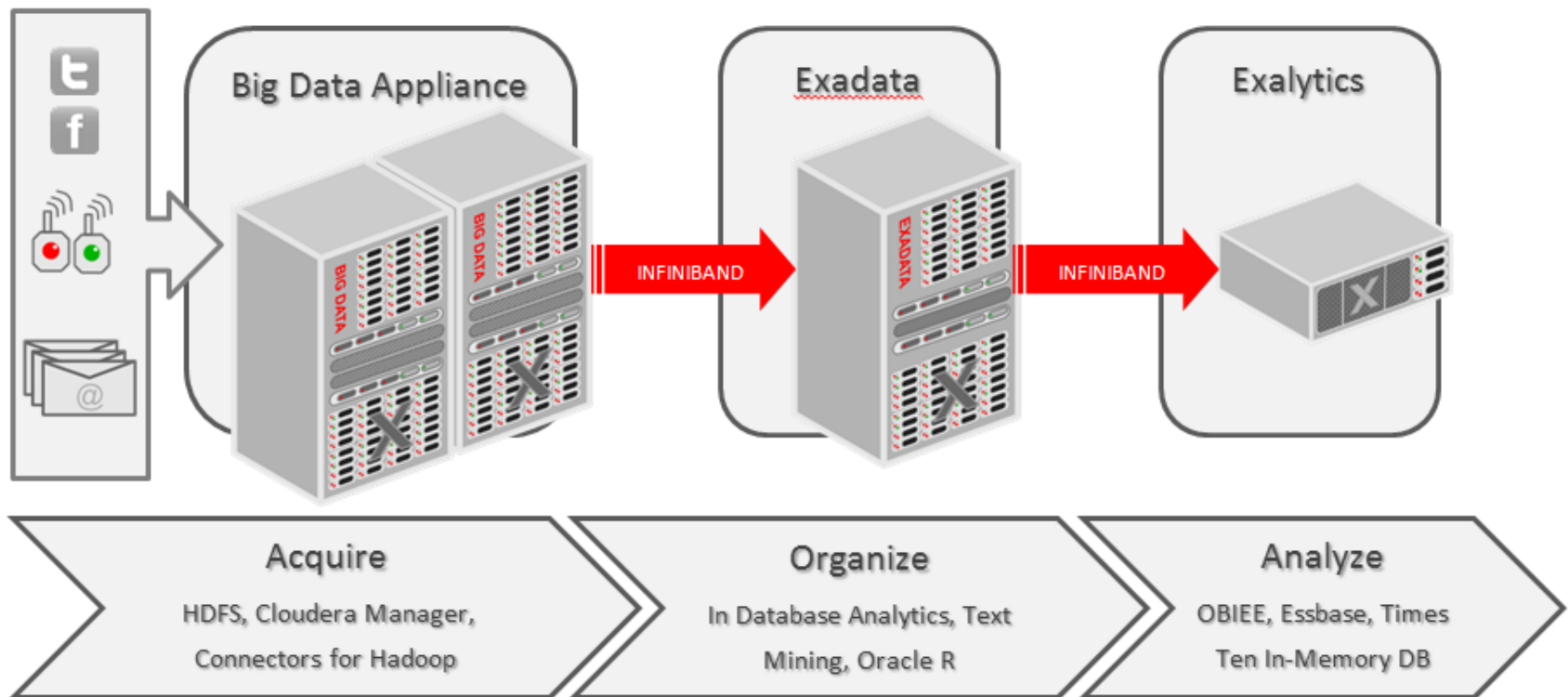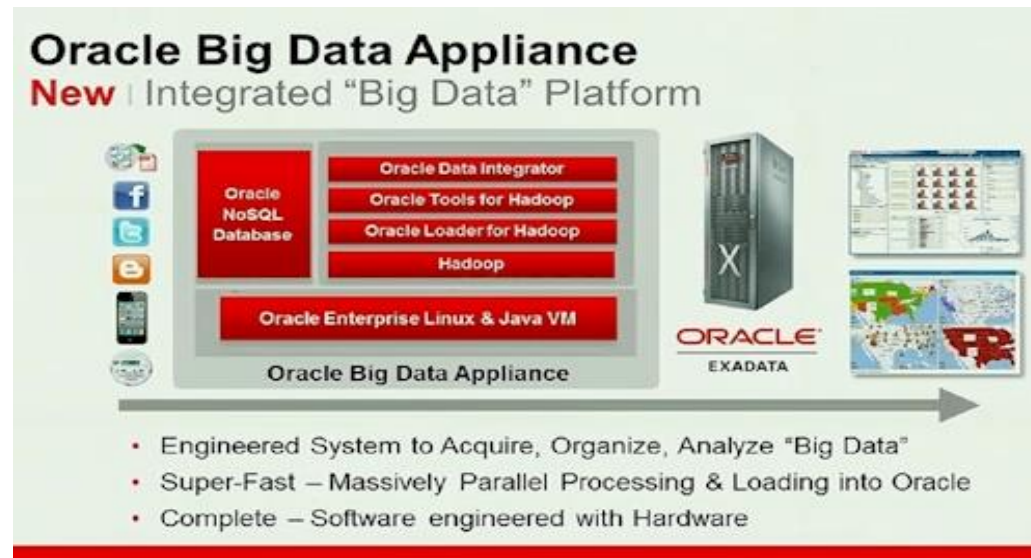


**Oracle Big Data Appliance**
**New** | Integrated "Big Data" Platform

Oracle NoSQL Database

- Oracle Data Integrator
- Oracle Tools for Hadoop
- Oracle Loader for Hadoop
- Hadoop
- Oracle Enterprise Linux & Java VM

Oracle Big Data Appliance

ORACLE EXADATA

- Engineered System to Acquire, Organize, Analyze "Big Data"
- Super-Fast – Massively Parallel Processing & Loading into Oracle
- Complete – Software engineered with Hardware

# Teradata Aster Big Analytics Appliance

The Teradata Aster Big Analytics appliance offers a Hybrid architecture that includes Aster Database, Aster SQL-mapreduce and Apache Hadoop



SQL-H interfaces with the Apache HCatalog to provide a mechanism for users to directly access the data in Hadoop from Aster Database

# Aster Database Architecture

Aster Database not only stores large volumes of data but also processes data and analytic applications in-database to deliver faster, deeper insights

# Advantage Teradata (+Aster database)
## A combined view of the Teradata and Aster database suite

# Paraccel

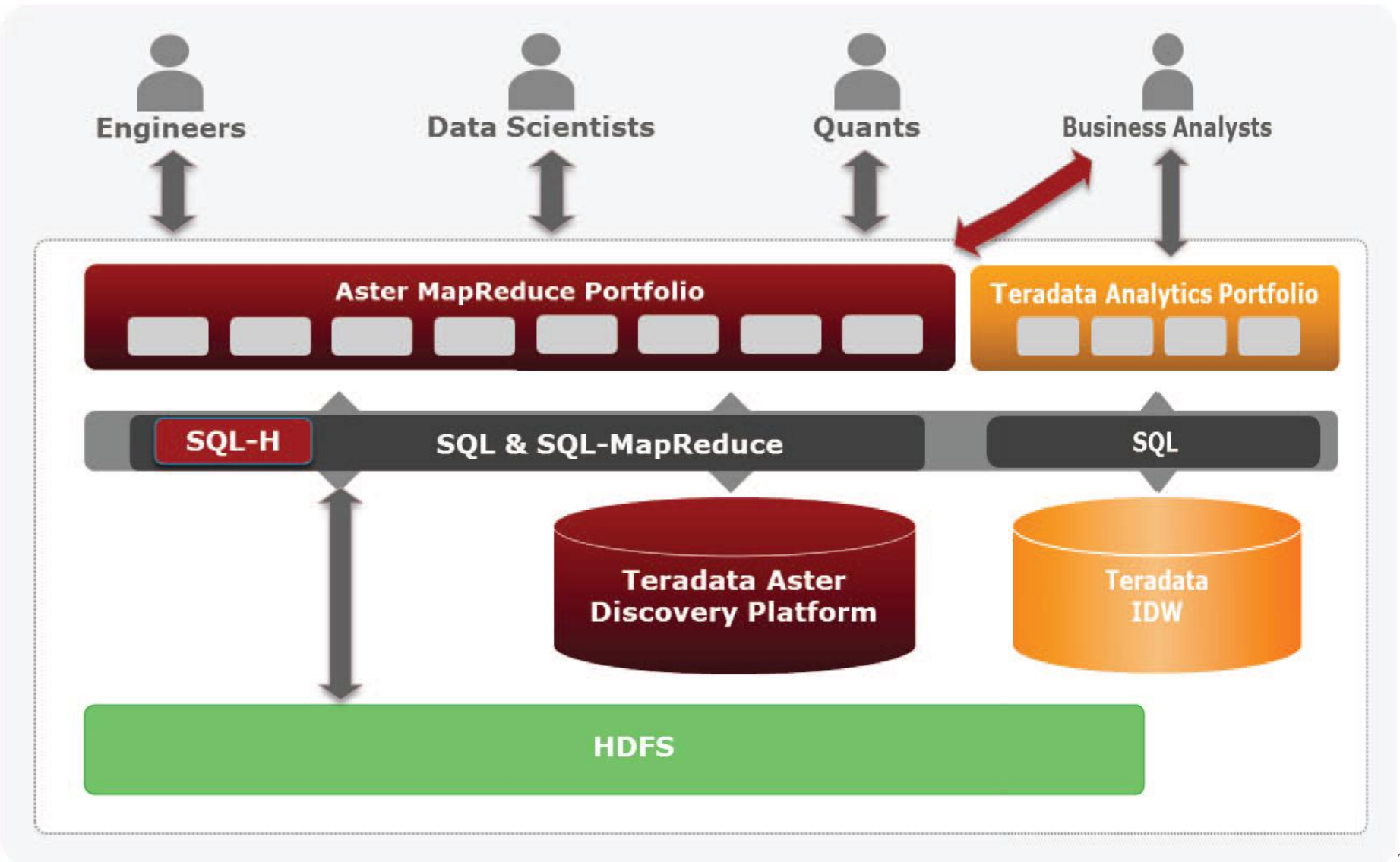## Paraccel Analytic Platform

# Paraccel Hadoop Analytics

ParAccel Hadoop Analytics is a visual end-to-end data preparation and analytics model development and execution environment



- Paraccel Hadoop Analytics automatically detects & utilizes all cores and nodes available at runtime up to a settable limit. It runs on any hardware, laptop, desktop, super-server or cluster

# Paraccel Hadoop Analytics

- Display data graphically in interactive charts that can be modified on-the-fly. Scatter plots, line graphs, bar graphs, dashboards & more allow visualizing data to instantly spot trends. R provides a wide variety of open source options.
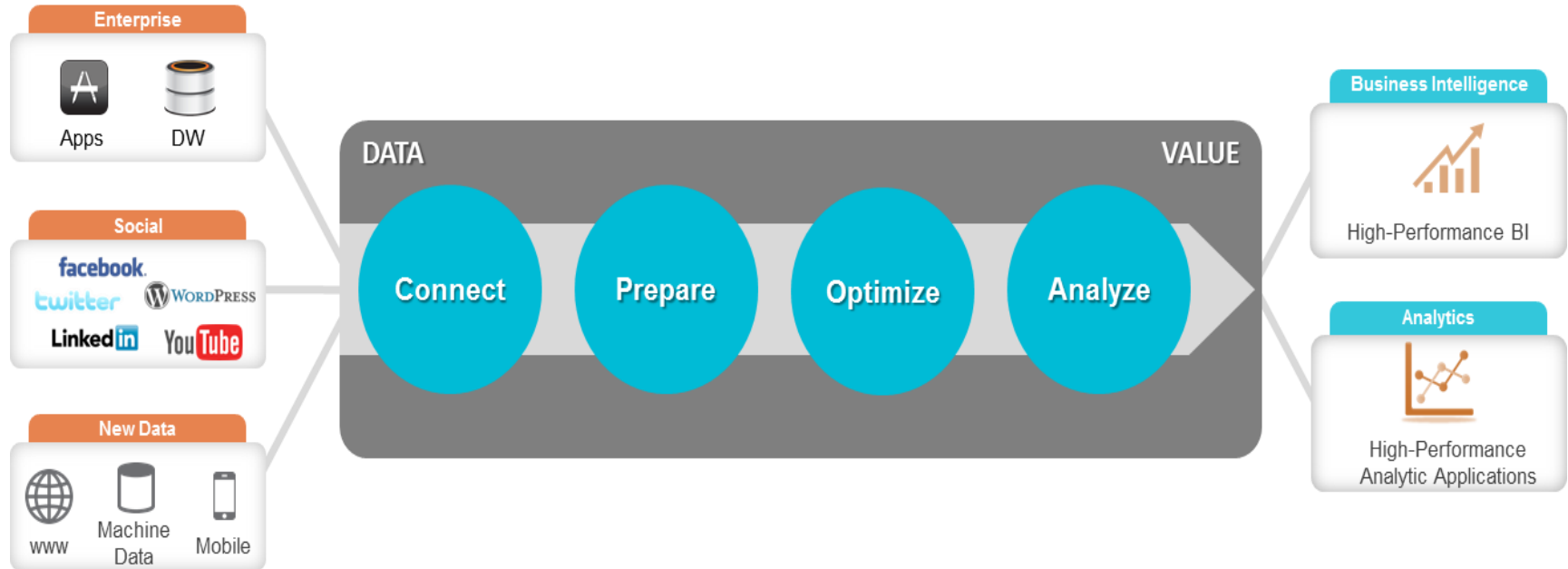


- Output through JDBC to visualization tools like Tableau, Actuate, YellowFin, etc. to expand  visualization options without limits.
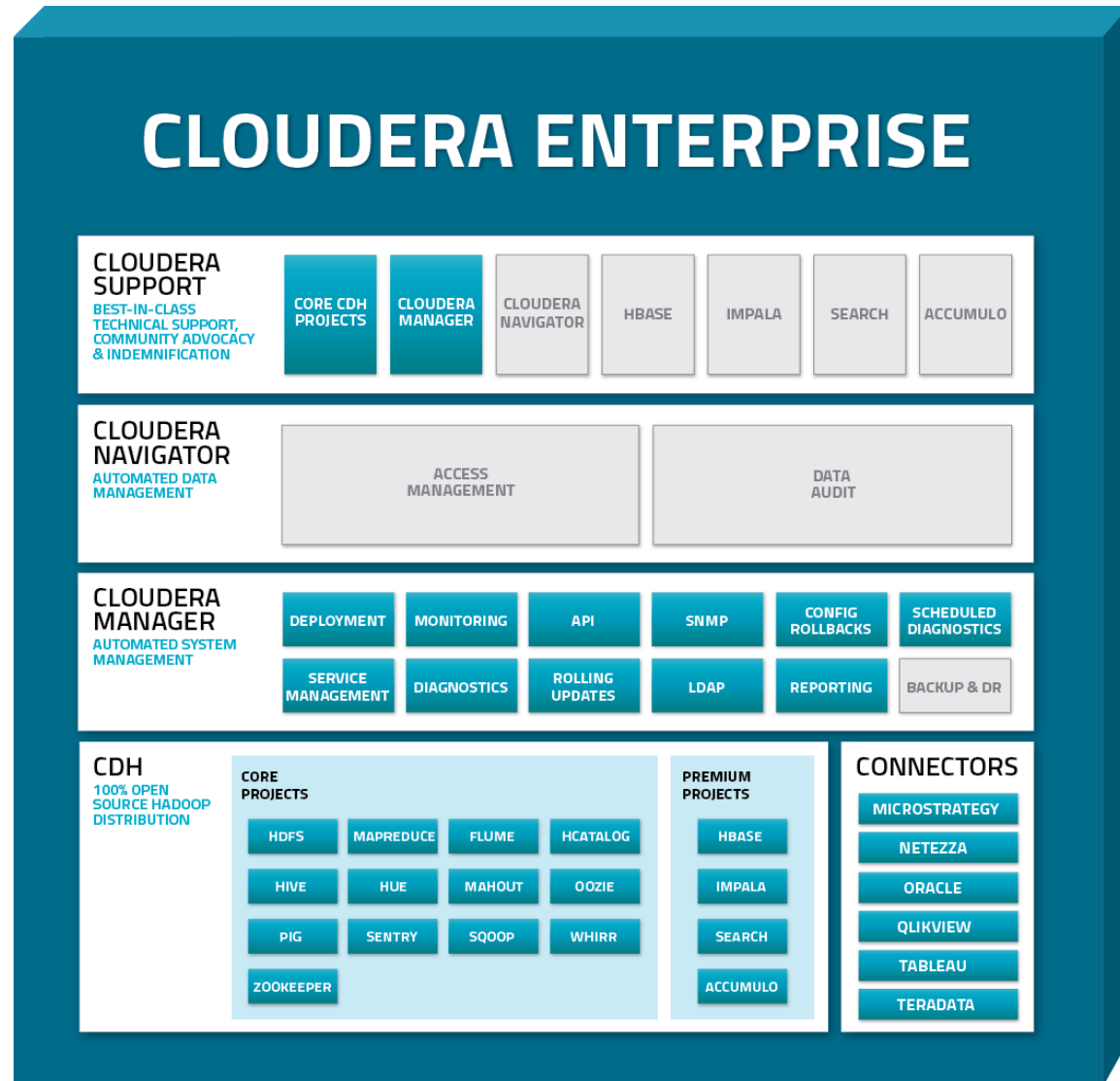
TLS

# The Paraccel Edge

- Accelerated Data Preparation - Processing 250 million health insurance claims to find fraud and claims mismanagement took 48 days. With ParAccel scaling across all cores on all nodes, it took less than a day.

- Accelerated Analytics Design - Six MapReduce coders took a week to build a machine learning predictive analytics workflow on Hadoop. One analyst with ParAccel Hadoop Analytics reproduced the workflow in ten minutes, with more accurate predictive results, and not one line of code.

- Accelerated Execution Speed per Dollar - A risk management solution a global bank was using had a processing time of 15+ hours. ParAccel executed the same solution in20 minutes, on half the hardware.

- Accelerated Analytics Deployment - A global data science consultancy took a week or more for each customer when deploying custom risk management analytics applications. Using ParAccel, they regularly deploy those applications in a matter of hours.
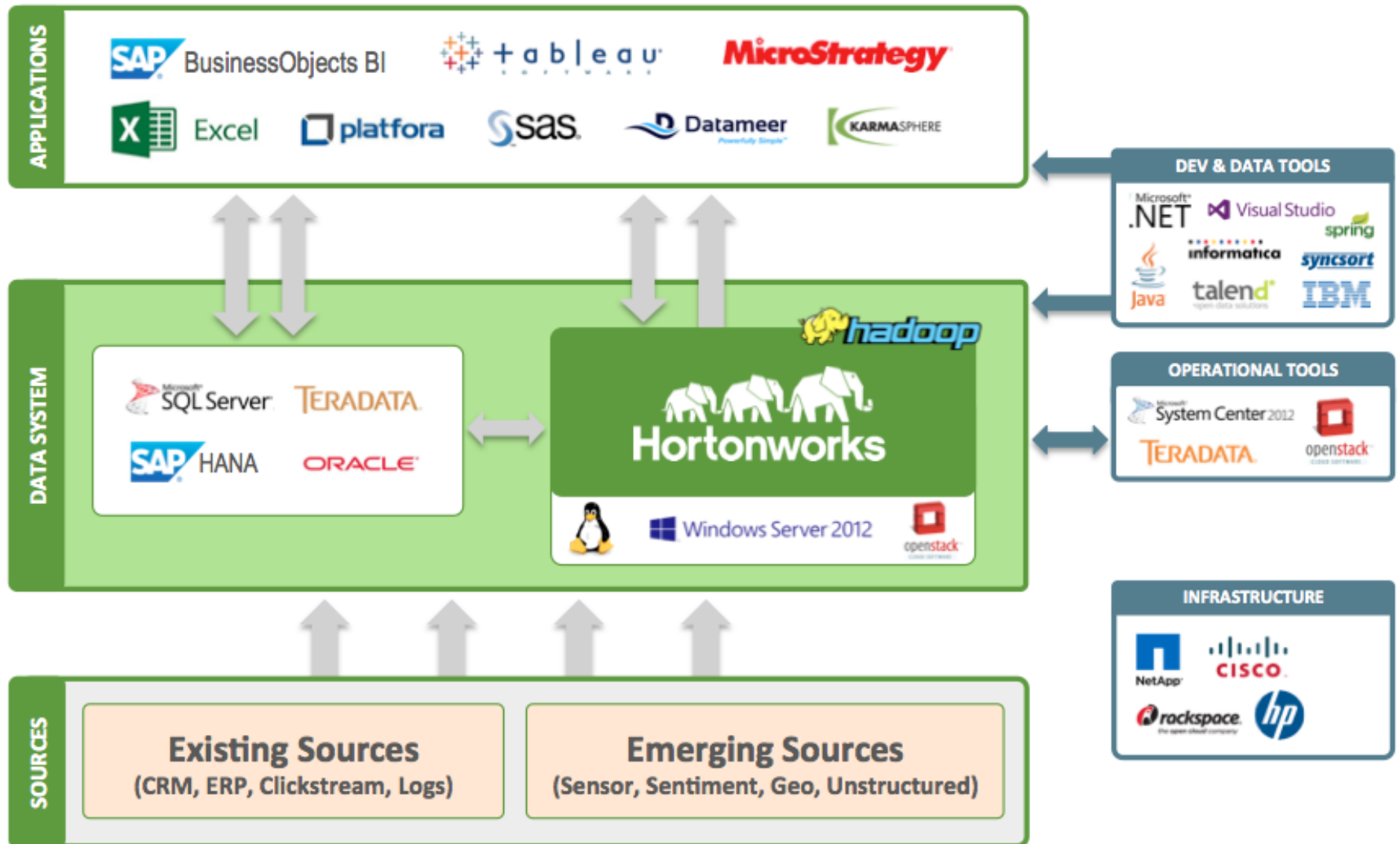
**TLS**

# Cloudera Enterprise
## The leader in the Hadoop market with 100+ customers

- Adds Cloudera Manager – Console for administering & managing Hadoop deployments and Enterprise Support

- Cloudera Manager offers wizard-based installation and configuration menus for deploying Hadoop. Also offers tools to help system managers monitor the health of the platform, diagnose problems, optimize performance, and make required configuration and security changes
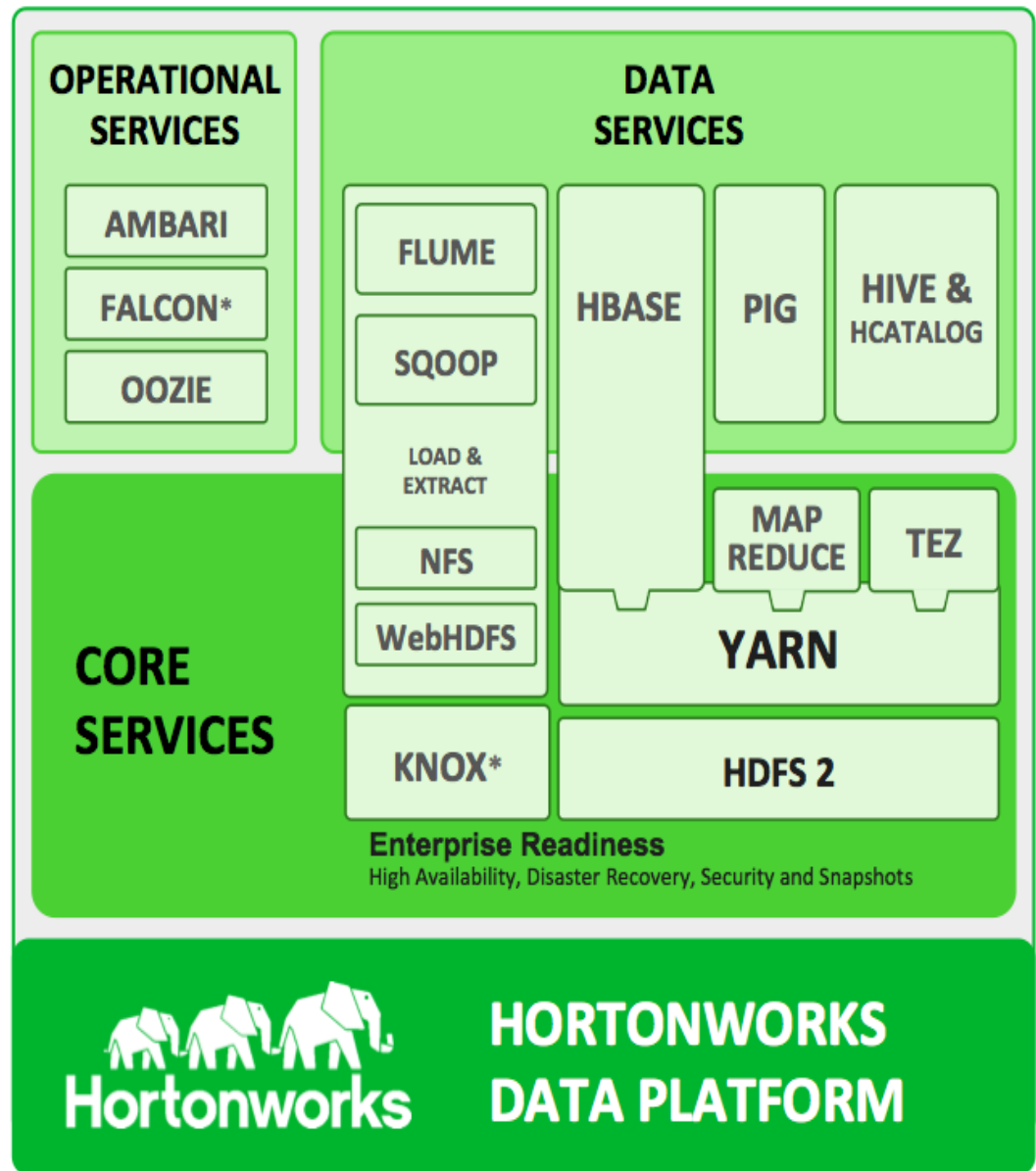
**CLOUDERA ENTERPRISE**

**CLOUDERA SUPPORT**
BEST-IN-CLASS TECHNICAL SUPPORT, COMMUNITY ADVOCACY & INDEMNIFICATION

| CORE CDH PROJECTS | CLOUDERA MANAGER | CLOUDERA NAVIGATOR | HBASE | IMPALA | SEARCH | ACCUMULO |
|---|---|---|---|---|---|---|

**CLOUDERA NAVIGATOR**
AUTOMATED DATA MANAGEMENT

| ACCESS MANAGEMENT | DATA AUDIT |
|---|---|

**CLOUDERA MANAGER**
AUTOMATED SYSTEM MANAGEMENT

| DEPLOYMENT | MONITORING | API | SNMP | CONFIG ROLLBACKS | SCHEDULED DIAGNOSTICS |
|---|---|---|---|---|---|
| SERVICE MANAGEMENT | DIAGNOSTICS | ROLLING UPDATES | LDAP | REPORTING | BACKUP & DR |

**CDH**
100% OPEN SOURCE HADOOP DISTRIBUTION

CORE PROJECTS

| HDFS | MAPREDUCE | FLUME | HCATALOG |
|---|---|---|---|
| HIVE | HUE | MAHOUT | OOZIE |
| PIG | SENTRY | SQOOP | WHIRR |
| ZOOKEEPER | | | |

PREMIUM PROJECTS

| HBASE |
|---|
| IMPALA |
| SEARCH |
| ACCUMULO |

**CONNECTORS**

| MICROSTRATEGY |
|---|
| NETEZZA |
| ORACLE |
| QLIKVIEW |
| TABLEAU |
| TERADATA |

**TLS**
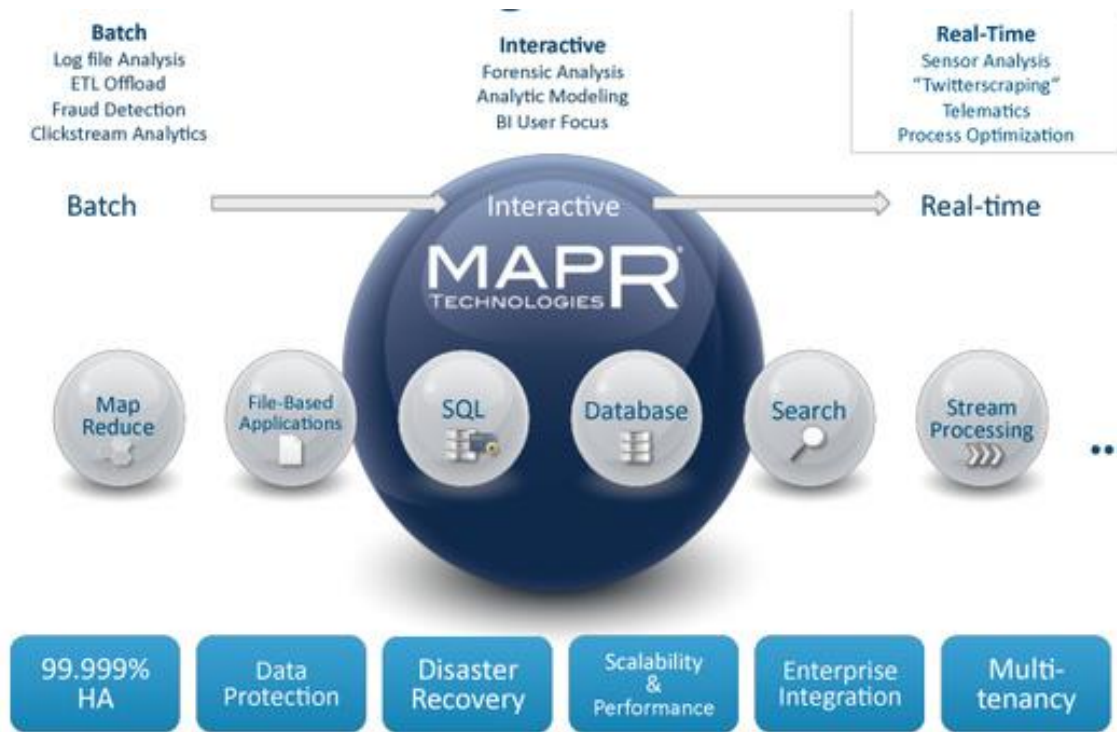
**CLOUDERA ENTERPRISE**

# HortonWorks

# HortonWorks

- Comes in 3 editions –
  - HDP (2 or 1.3)
  - HDP for Windows
  - HDP Sandbox
- Include Apache Ambari (similar to Cloudera Manager)
- BCP improvements via snapshots for recovery
- Better connectors to Oracle, Netezza
- Improved functionality and performance in Sqoop
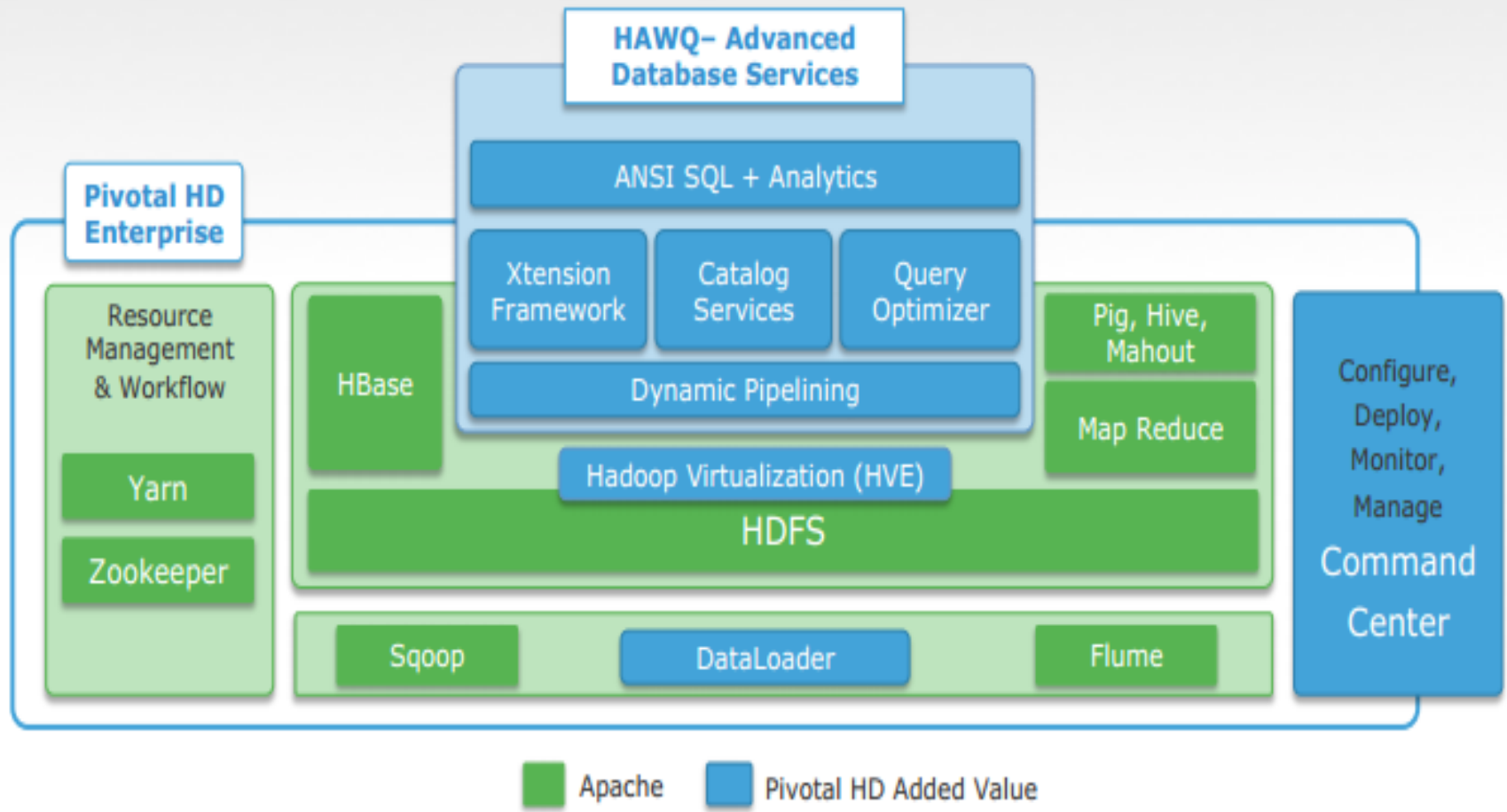- Improve SQL support using Hive



**OPERATIONAL SERVICES**
- AMBARI
- FALCON*
- OOZIE

**DATA SERVICES**
- FLUME
- SQOOP
- HBASE
- PIG
- HIVE & HCATALOG
- LOAD & EXTRACT
- NFS
- WebHDFS

**CORE SERVICES**
- MAP REDUCE
- TEZ
- YARN
- KNOX*
- HDFS 2

**Enterprise Readiness**
High Availability, Disaster Recovery, Security and Snapshots

**Hortonworks**
**HORTONWORKS DATA PLATFORM**

## TLS

# MapR

## Arguably the strongest challenger to Cloudera



- Comes in 3 editions – M3, M5 and M7

- M3 is Free Hadoop, M5 is Enterprise Hadoop and M7 Enterprise NoSQL (Hbase) and Hadoop

- Proprietary Mountable file system

- High availability name nodes & job trackers

- Amazon's Elastic MapReduce now includes MapR M7

- Was the underlying Hadoop distro used in EMC Greenplum (called Greenplum HD)
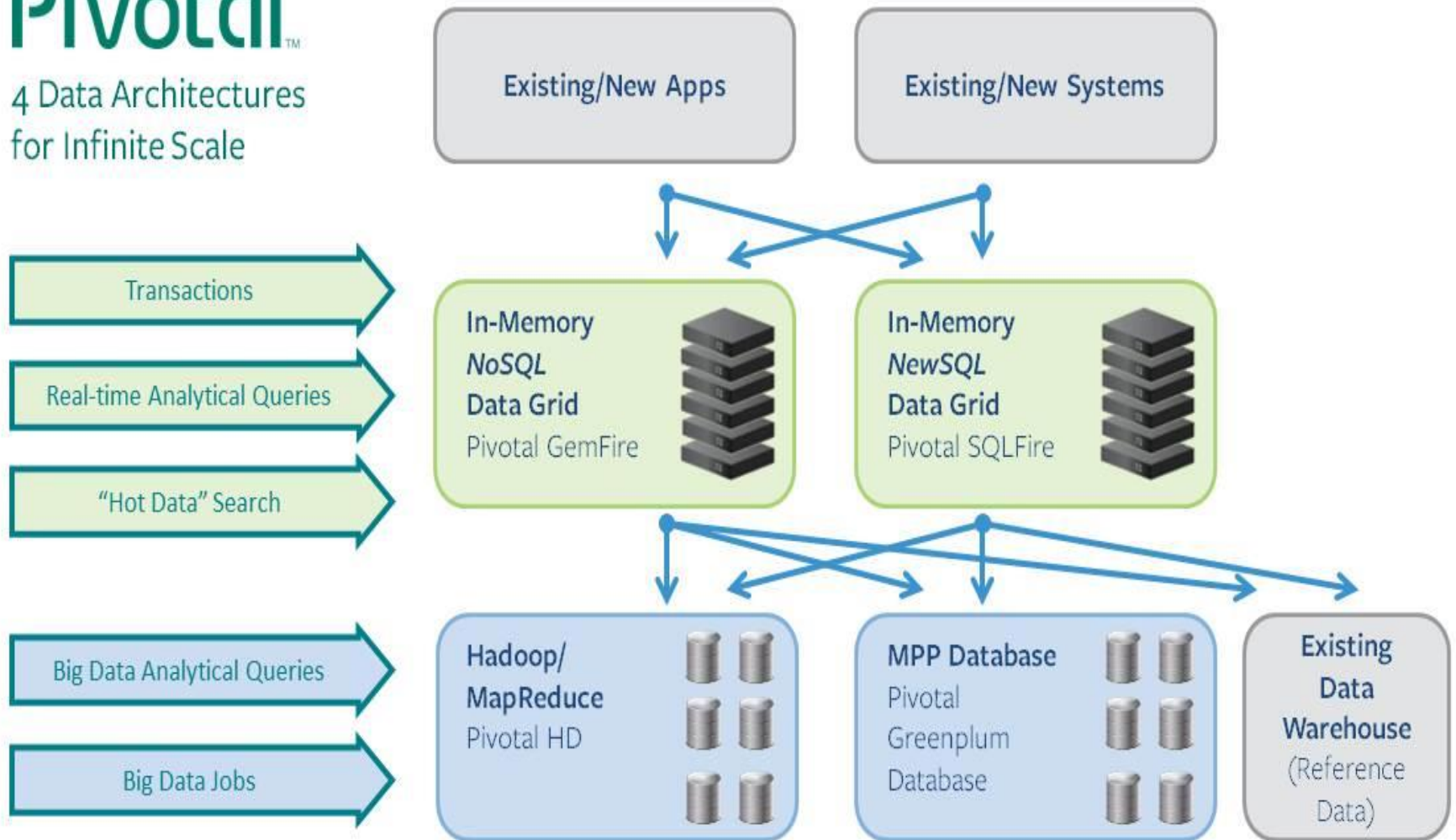
**TLS**

# Pivotal Analytics

Pivotal analytics is a unique end-to-end analytics platform that fills the need for an elastically scalable, real-time analytics solution that delivers big answers.

# Greenplum HD Pivotal – HAWQ

## HAWQ: The Crown Jewels of Greenplum

**HAWQ**

- High-Performance Query Processing
  - Multi-petabyte scalability
  - Interactive and true ANSI SQL support
  - Programmable analytics

- Enterprise-Class Database Services
  - Column storage and indexes
  - Workload Management

- Comprehensive Data Management
  - Scatter-Gather Data Loading
  - Multi-level Partitioning
  - 3rd Party Tool & Open Client Interfaces

## HAWQ Benchmarks

| | HAWQ | HIVE | |
|---|---|---|---|
| User intelligence | 4.2 | 198 | 47X |
| Sales analysis | 8.7 | 161 | 19X |
| Click analysis | 2.0 | 415 | 208X |
| Data exploration | 2.7 | 1,285 | 476X |
| BI drill down | 2.8 | 1,815 | 648X |

**TLS**

# Bigdata Use Cases

TLS

# Bigdata Use Cases – HealthCare

**HealthCare**

- Clinical Decision Support System

- Remote Patient Monitoring

- Patient Profiling

- Comparative Effective Research

TLS

# Bigdata Use Cases – Retail

**Marketing**

- Cross-Selling

- Location-based Marketing

- Analyzing in – store behavior

- Sentiment Analysis

TLS

# Bigdata Use Cases – Manufacturing

**Manufacturing**

- Sensors Driven Operations

- Supply Chain & Inventory Management

- Shorten Design to Value Cycle

TLS

# Bigdata Use Cases – Public Services

**Public Sector Administration**

- Traffic Decongestion

- Civic Compliance – Noise, Air & Water Pollution

- Locality and Area Development

**TLS**

# Bigdata Use Cases – Telecom

**Telecom**

- Geo – Targeted Advertising

- Electronic Toll Collection

- Emergency Response

- Remote Personal Car safety & monitoring

- Urban Planning

**TLS**

# Bigdata Use Cases – Oil & Gas

## Oil & Gas Sector

- Digital Oil Field

- Remote Drilling Operations

**TLS**

# The Bigdata Market

TLS

# Gartner Hype Cycle on Bigdata

# Some Unanswered Questions

- Too much emphasis on the 3Vs & 1C, what about Value?

- Conventional BI has mature and elaborate processes and definitions for data validation, data integrity, data governance, data extraction, metadata management and ETL, what about these disciplines in Bigdata?

- Convergence of Structured and Unstructured data (single query) how will it be handled?

- Challenges with data update and data erasure

- Writing MapReduce in Java is cumbersome and inefficient, need quick evolution of IDEs and GUIs

- Convergence of BI capability & Java programming – Cross skilling, multi-skilling challenges

TLS

# Thank You

**Disclaimer**