

Apache Hadoop Present & Future

Hadoop in China 2012

Eric Baldeschwieler

CTO Hortonworks

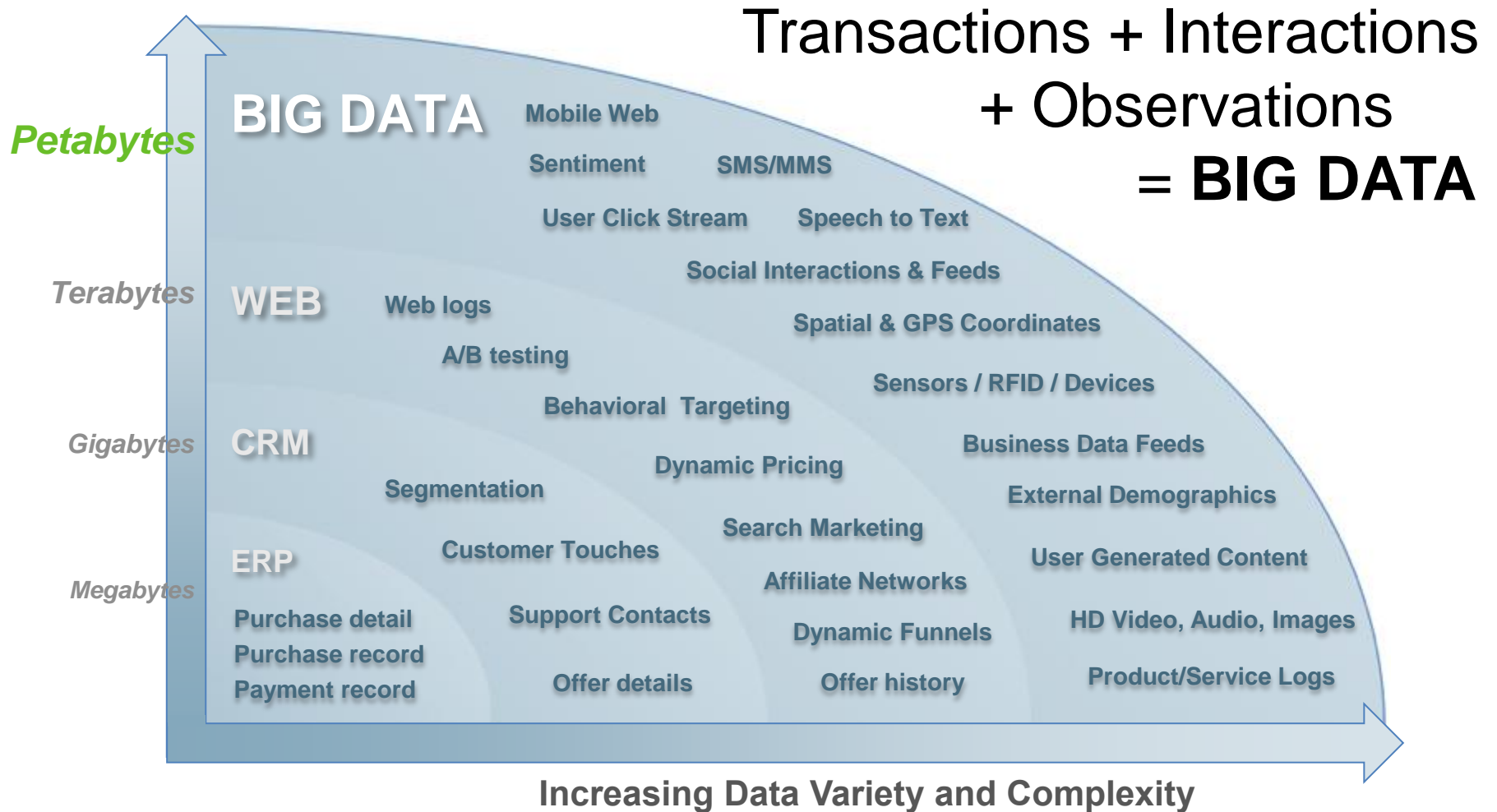
Twitter: @eric14 , @Hortonworks



What is Apache Hadoop?



Big Data



The Big Data Problem

Platform for Big Data

Capture

- Collect data from all sources - structured and unstructured data
- all speeds batch, async, streaming, real-time

Process

- Transform, refine, aggregate, analyze, report

Exchange

- Deliver data with enterprise data systems
- Share data with analytic applications and processing

Operate

- Provision, monitor, diagnose, manage at scale
- Reliability, availability, affordability, scalability, interoperability

Across all deployment models

Operating
Systems

Virtual
Platforms

Cloud
Platforms

Big Data
Appliances

Big Data Use Cases

Vertical	Refine	Explore	Enrich
Social and Web	<ul style="list-style-type: none"> MDM CRM Ad service models 	<ul style="list-style-type: none"> Paths Feature usefulness 	<ul style="list-style-type: none"> Friends and associations Content recommendations Ad service
Retail	<ul style="list-style-type: none"> Loyalty programs Cross-channel customer MDM CRM 	<ul style="list-style-type: none"> Referrers Brand and Sentiment Analysis Paths Taxonomic relationships 	<ul style="list-style-type: none"> Dynamic Pricing/Targeted Offer
Intelligence	<ul style="list-style-type: none"> Threat Identification 	<ul style="list-style-type: none"> Person of Interest Discovery 	<ul style="list-style-type: none"> Cross Jurisdiction Queries
Finance	<ul style="list-style-type: none"> Risk Modeling & Fraud Identification Trade Performance Analytics 	<ul style="list-style-type: none"> Surveillance and Fraud Detection Customer Risk Analysis 	<ul style="list-style-type: none"> Real-time upsell, cross sales marketing offers
Energy & Utility	<ul style="list-style-type: none"> Production Optimization Smart Meters & Devices 	<ul style="list-style-type: none"> System wide analysis Fine grained reporting 	<ul style="list-style-type: none"> Automatic system recovery Interactive customer service
Manufacturing	<ul style="list-style-type: none"> Supply Chain Optimization 	<ul style="list-style-type: none"> Customer Churn Analysis 	<ul style="list-style-type: none"> Dynamic Delivery Replacement parts
Healthcare & Payer	<ul style="list-style-type: none"> Electronic Medical Records (EMPI) 	<ul style="list-style-type: none"> Clinical Trials Analysis 	<ul style="list-style-type: none"> Insurance Premium Determination

Apache Hadoop, Big Data Platform



*Open Source data management
with scale-out storage &
distributed processing*

Storage

HDFS



- Distributed across “nodes”
- Natively redundant
- Name node tracks locations

Processing

Map Reduce



- Splits a task across processors “near” the data & assembles results
- Self-Healing, High Bandwidth Clustered Storage

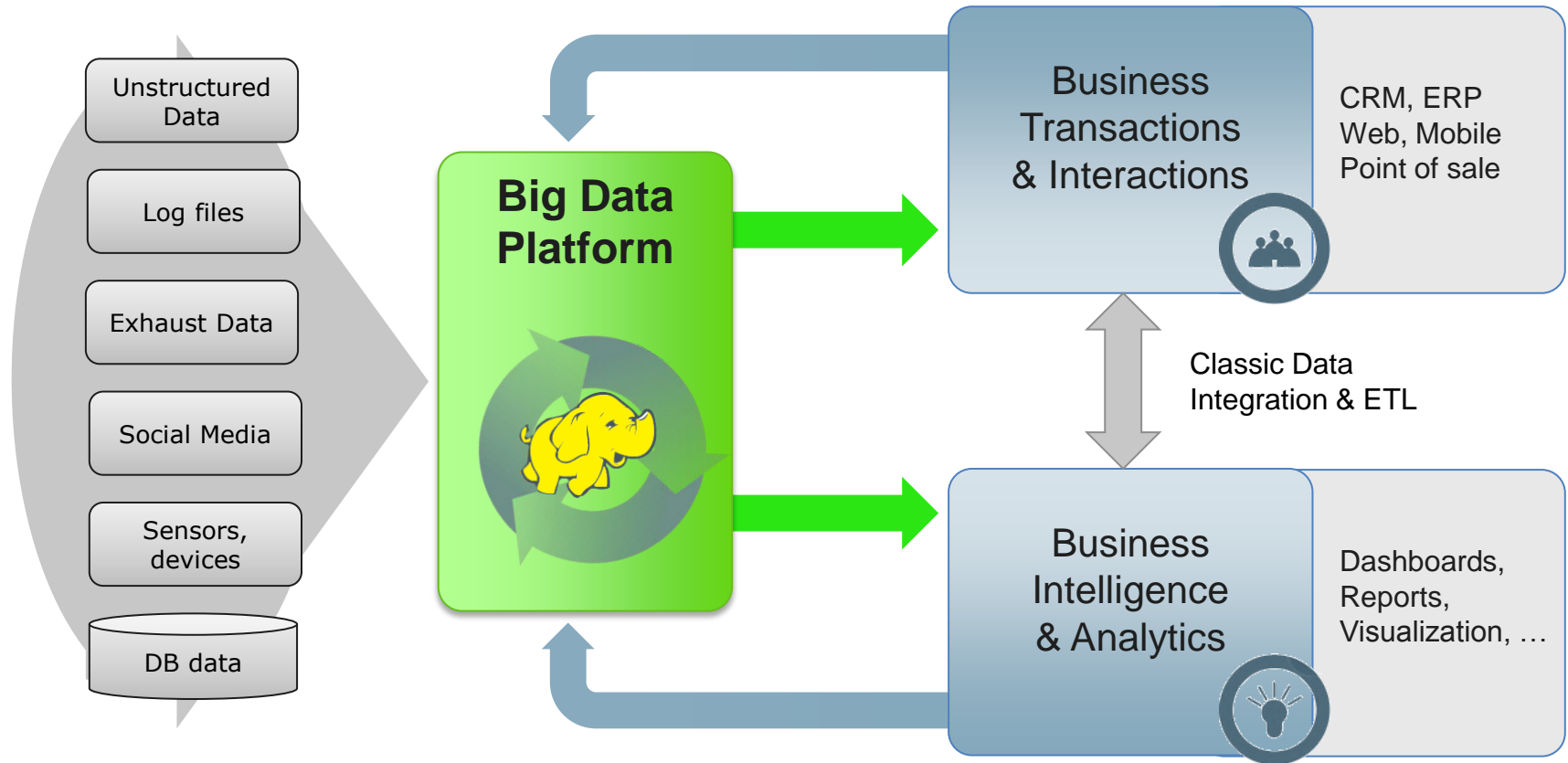
Key Characteristics

- **Scalable**
 - Efficiently store and process petabytes of data
 - Linear scale driven by additional processing and storage
- **Reliable**
 - Redundant storage
 - Failover across nodes and racks
- **Flexible**
 - Store all types of data in any format
 - Apply schema on analysis and sharing of the data
- **Economical**
 - Use commodity hardware
 - Open source software guards against vendor lock-in

Deploying Apache Hadoop

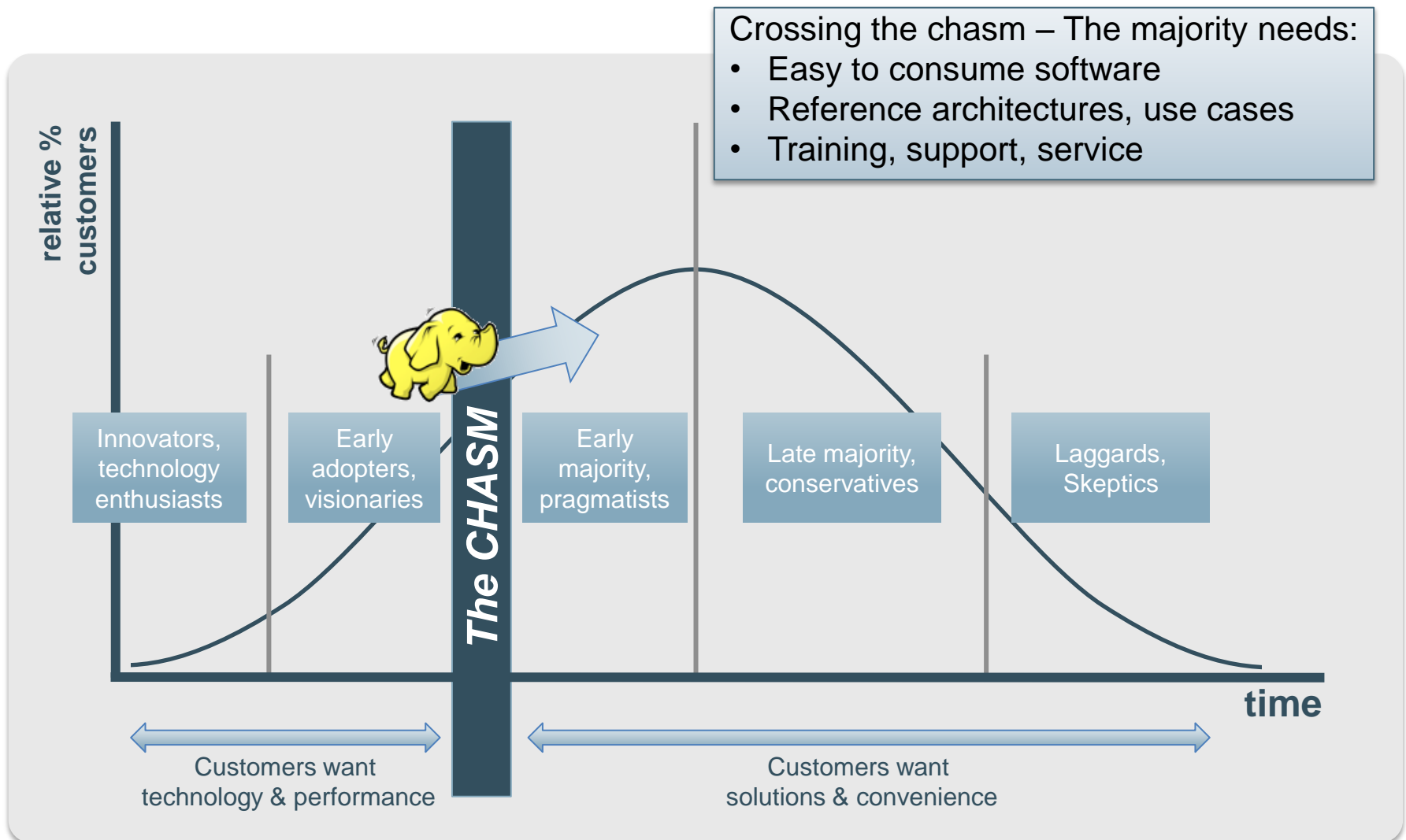


Hadoop in Enterprise Big Data Flows



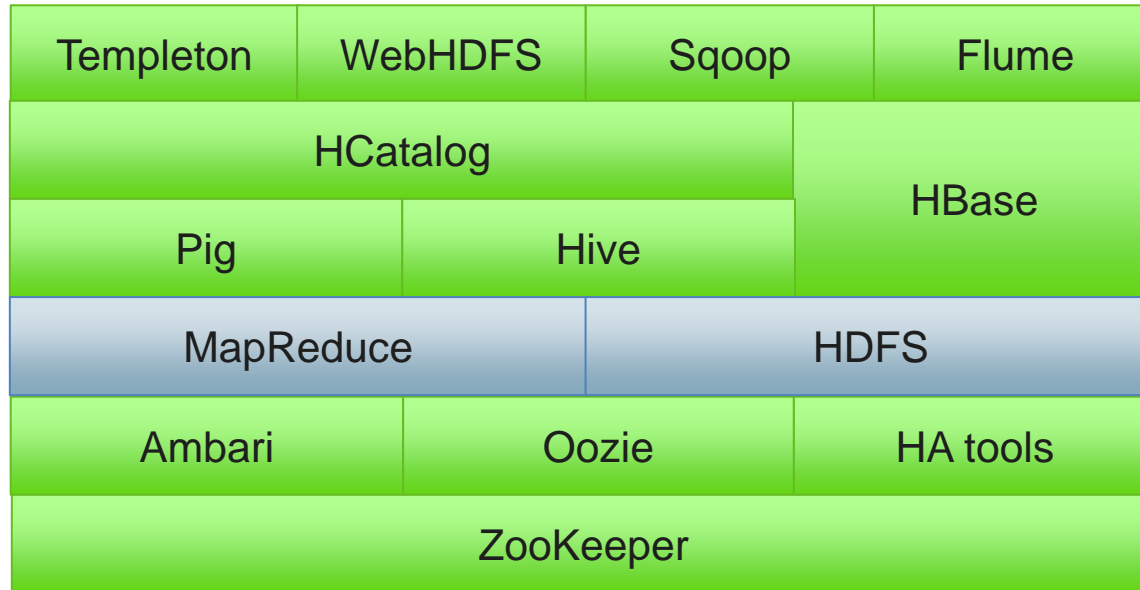
- 1 Capture Big Data**
Collect data from all sources structured & unstructured
- 2 Process**
Transform, refine, aggregate, analyze, report
- 3 Exchange Results**
Interoperate and share data with applications/analytics

Hadoop: Poised for Rapid Growth



Source: Geoffrey Moore - Crossing the Chasm

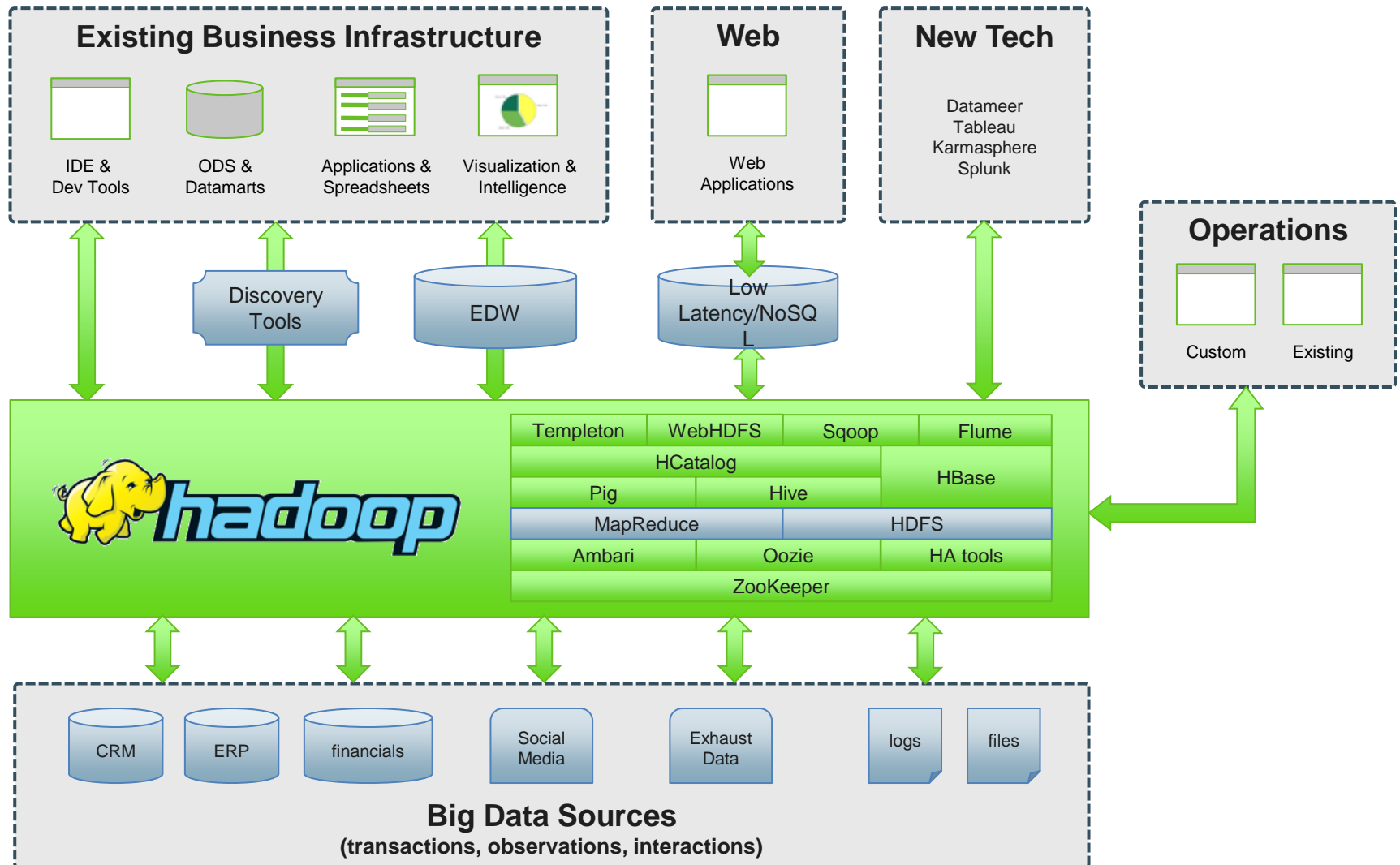
What is a Hadoop “Distribution”



A collection of Apache Hadoop and related technologies to make it easy to install and operate a Hadoop Big Data Platform

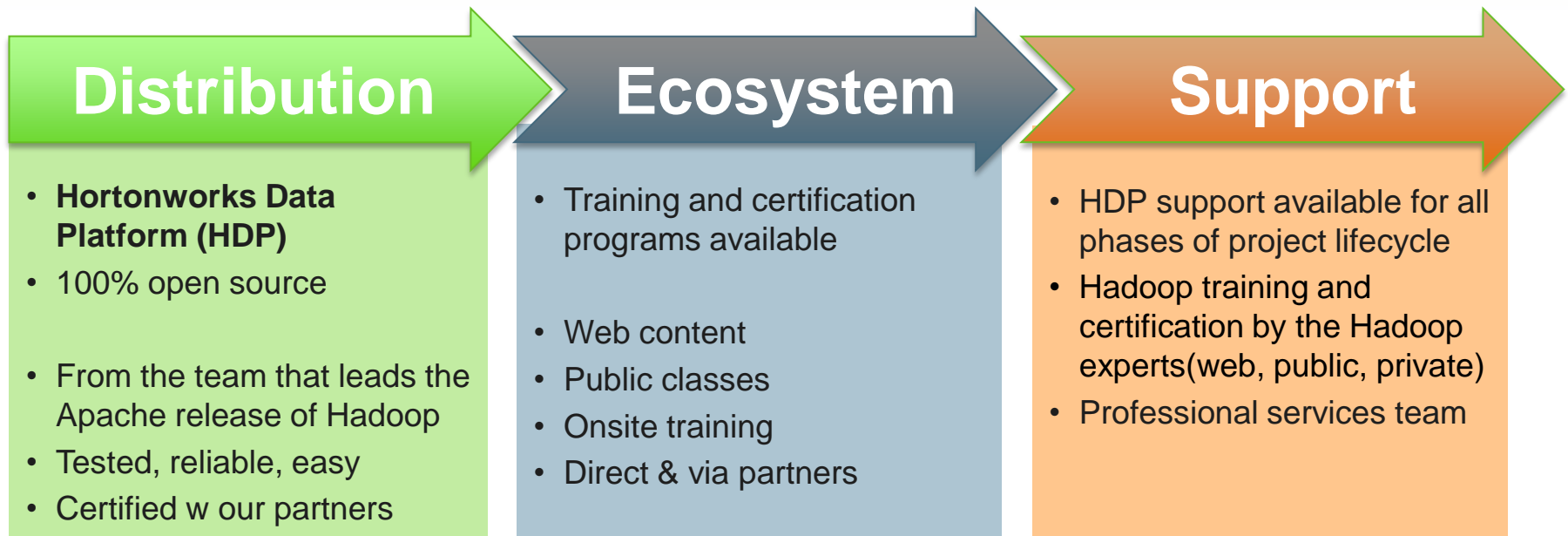
- Solves the tricky problem of choosing the right versions of many components that all have different release cycles
- Tested and pre-packaged to ease installation and usage

Hadoop in Enterprise Data Architectures



Hortonworks and Hadoop in the Enterprise

- ***“We believe that by the end of 2015, more than half the world's data will be processed by Apache Hadoop.”***
- Hortonworks mission:
 - Provide a complete and 100% open source Apache Hadoop Distribution
 - Invest in Apache Hadoop to make it ***“The enterprise big data platform”***





The Road Ahead

Coming improvements
to Hadoop

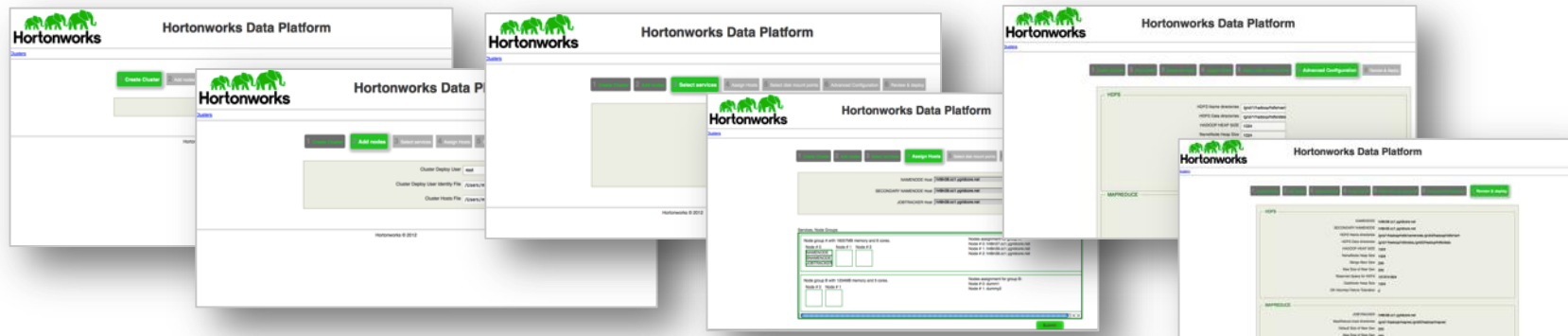
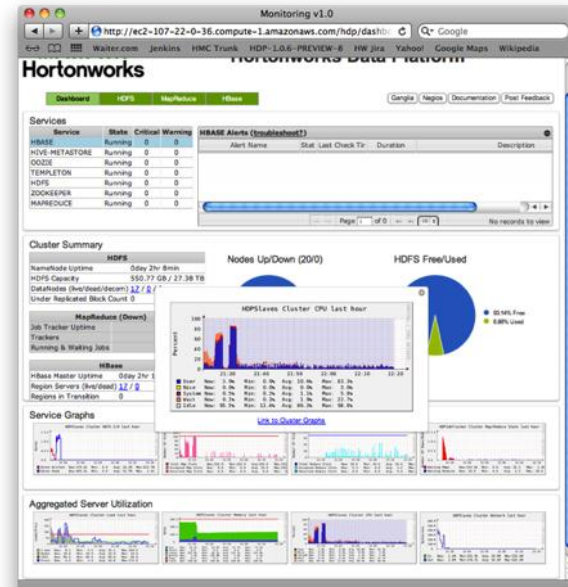
Ambari: Management & Monitoring Services

- **Powerful monitoring and alerting dashboards**

- View topology, health & utilization of cluster
- Detailed view of cluster operations, server & storage utilization, job status, and performance levels
- Get alerts to critical events

- **Simple installation & provisioning**

- Easy configuration process
- One-click deployment for clusters of all sizes
- Analyzes/recommends optimal services configuration
- Automatically configures mount points in the cluster



Hive + HCatalog

- **ODBC / JDBC support**

- Support for BI tools such as Tableau & Excel

- **HCatalog**

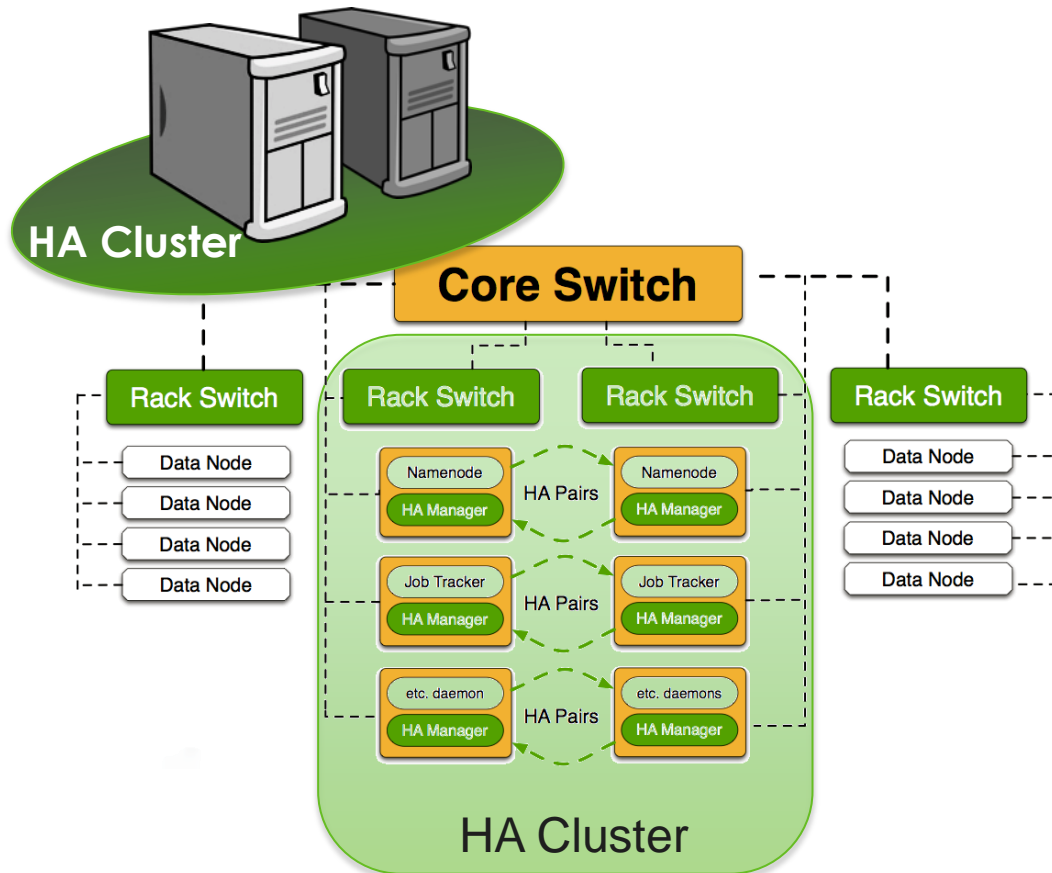
- Allow the use of Hive tables with all Hadoop tools (Pig, MapReduce...)
- Allow use of Hadoop data with MPP DBs (Aster SQLH)
- Providing a new higher level of abstraction for Hadoop data

- **Performance**

- New column oriented file formats
- Lots of inner loop performance improvements
- Improved query planning
- Yarn integration

Full Stack High Availability

Proven HA solutions with proven Hadoop 1.0 & 2.0



Failover and restart for

- NameNode
- JobTracker
- Other services to come...

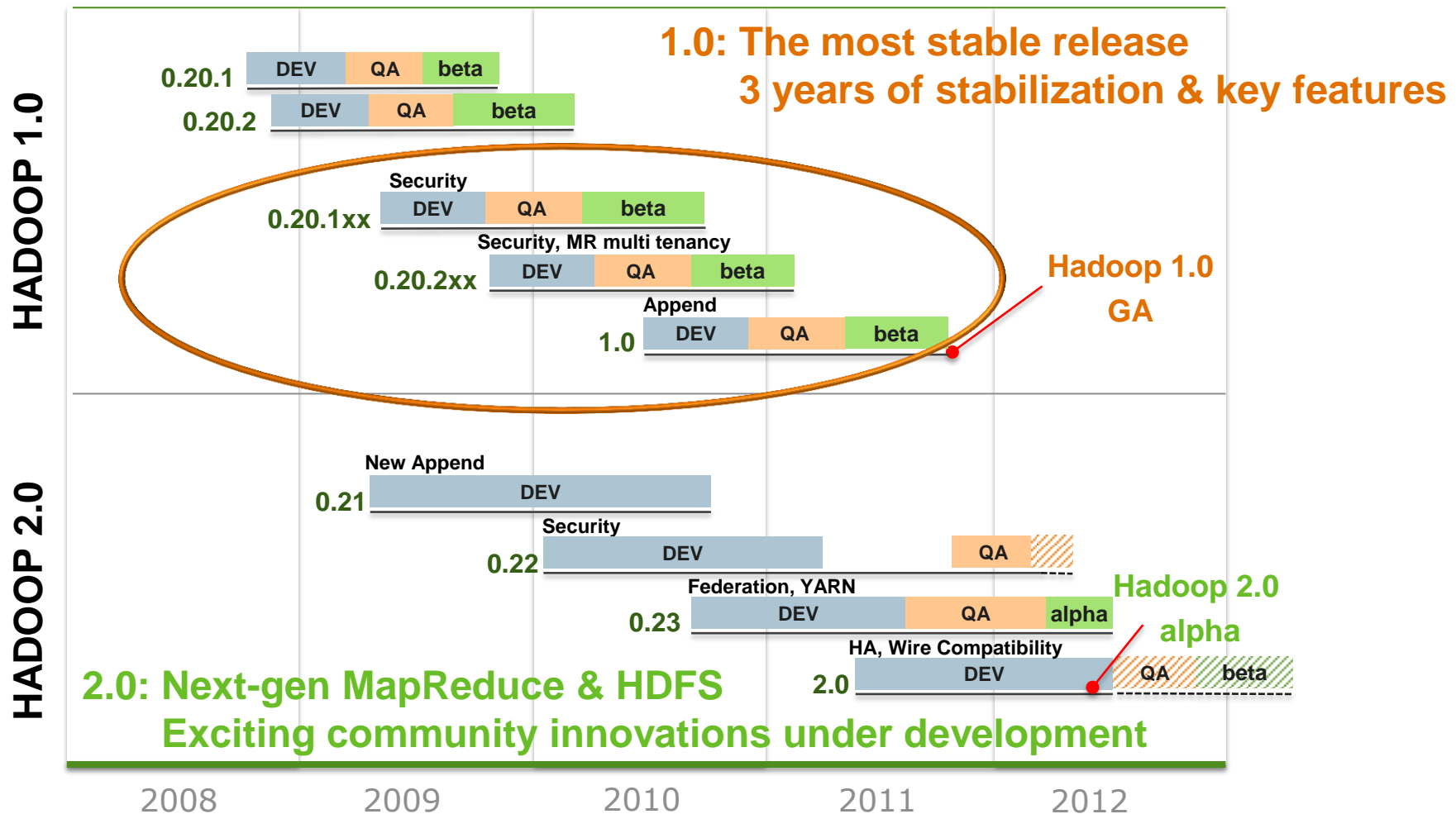
Open API allows use of Proven HA from multiple vendors

Minimized changes to clients and configuration

Auto-detects failures:

- Services, OS & Hardware

Timeline: Apache Hadoop 1.0 & 2.0



Hadoop 2.0 Innovations

- **Focus on Scale and Community Innovation**

- Designed to support 10,000+ computer clusters
- Extensible design to encourage innovation

- **YARN: Next Generation Execution**

- Improves MapReduce performance
- Supports **new frameworks** beyond MapReduce!
 - Low latency, Streaming, Services
 - Do more with a single Hadoop cluster

- **HDFS 2.0**

- Federation: Isolation & extensibility via multiple NameNodes
- Checkpoint support
- Performance improvement



Thank You!

Questions & Answers

Follow: @hortonworks

