# Evolution to Open Source Data Management with Scale-out Storage & Processing

| Date | Paradigm | Processing Style/ Scale Out | Form Factor |
|------|----------|------------------------------|-------------|
| **90s** | • Reporting / Data Mining<br>• High Cost / Isolated use | • Batch– "sales reports"<br>• Sequential SQL queries<br><br>Scale → Multi-core | RDBMS |
| **2000s** | • Model-based discovery<br>• High Cost / Dept Use | • Batch-ie correlated buying pattern<br>• No SQL. parallel analysis<br>• Shared disk/memory<br><br>Scale → Node Node Node | No SQL RDBMS<br><br>Proprietary MPP/ DW Appliance |
| **Today** | • Unbounded Map Reduce Query<br>• Low Cost / Enterprise Use<br>• Arrival of vast amounts of unstructured data | • Real-time- ie recommend engine<br>• Process @ storage node<br>• Built-in data replication/reliability<br>• Shared nothing, in memory<br><br>Unlimited Linear Scale → Distributed node addition | Open Source SW coupled to commodity HW<br><br>Node Node Node<br><br>hadoop |

# Apache Hadoop Evolution

Source - Steven Nimmons
2/24-12

## 2006
- HDFS
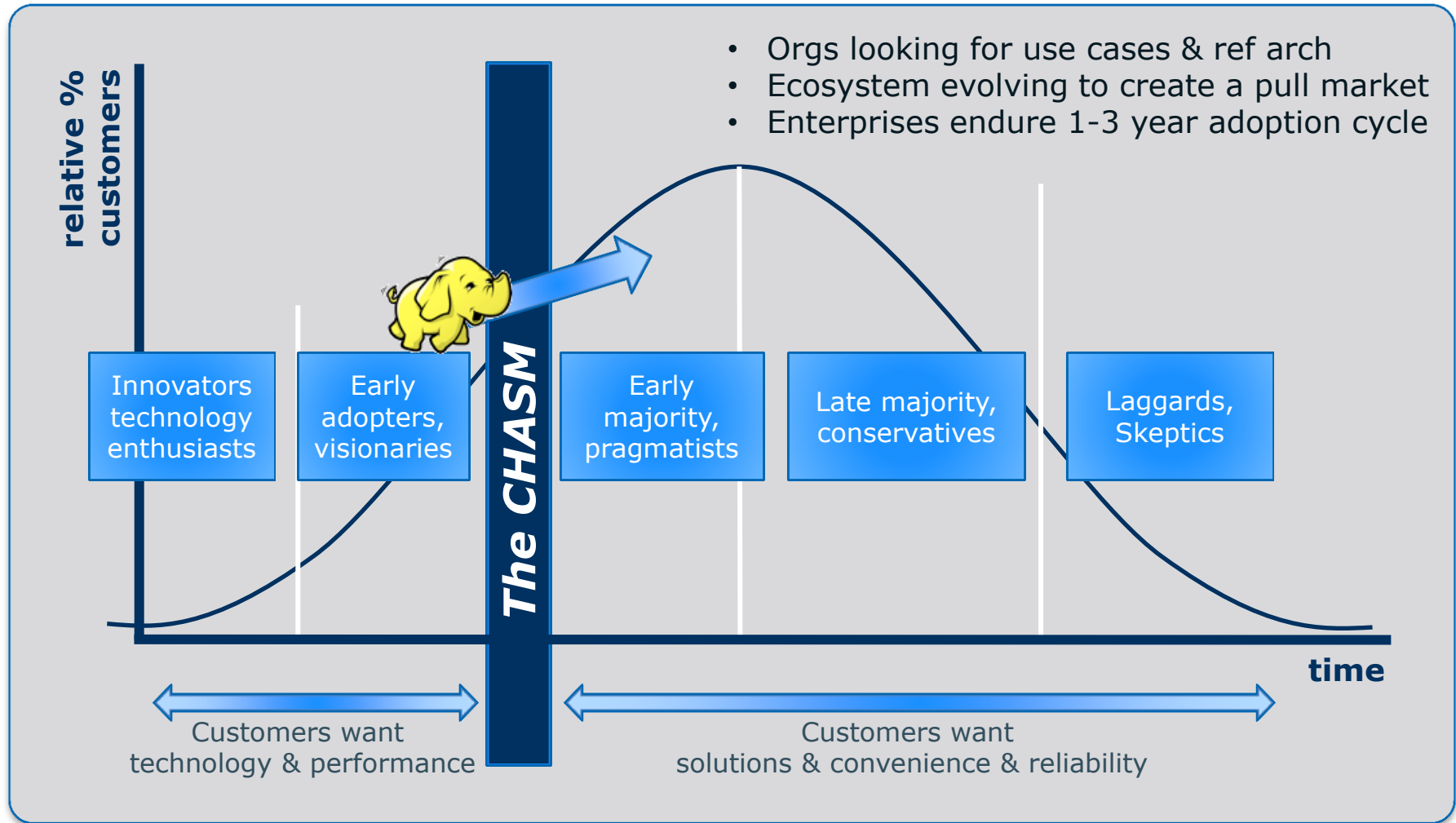- MapReduce

## 2008
- HBase
- ZooKeeper
- Pig
- Hive

## 2009-10
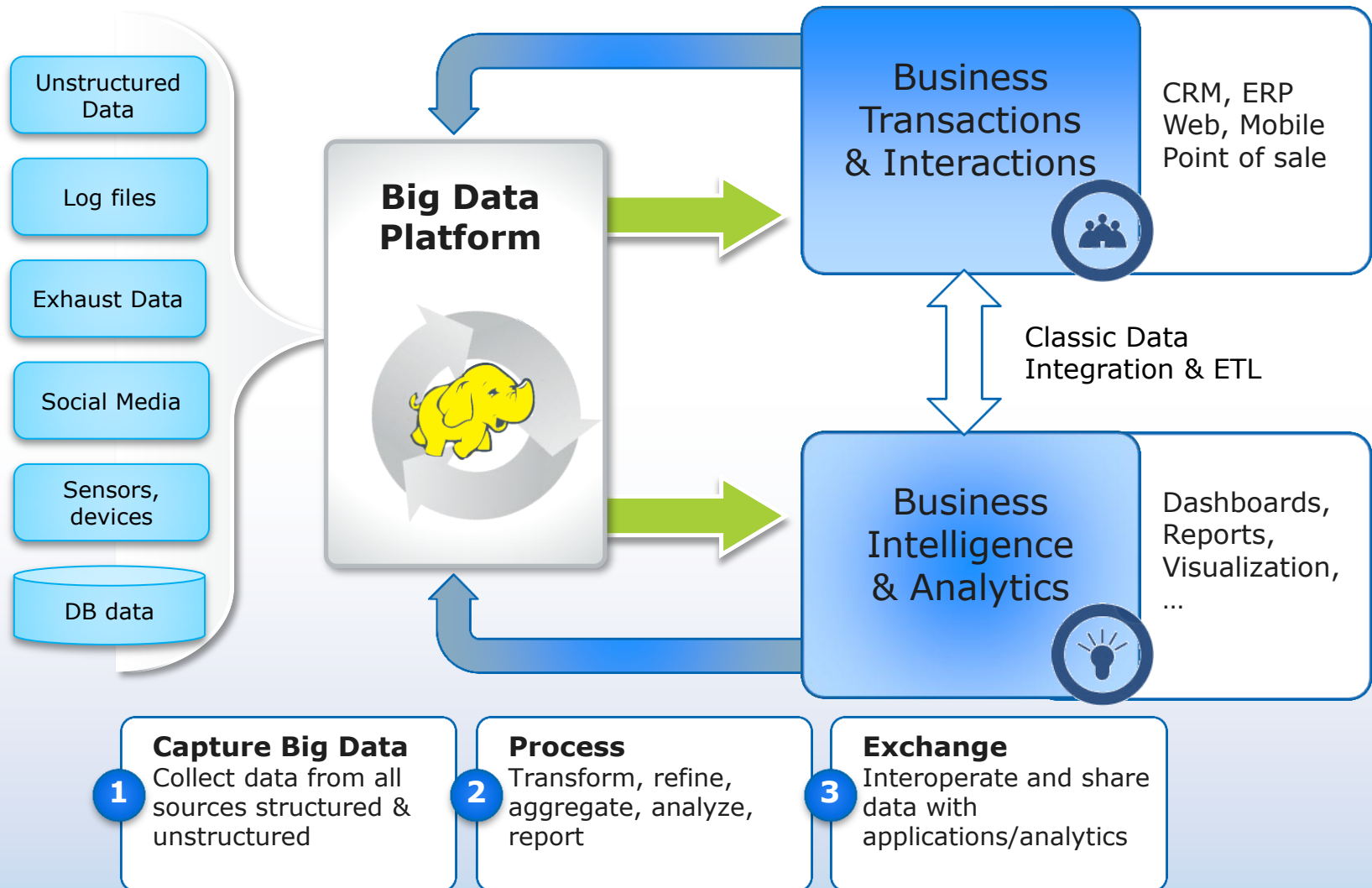- Flume
- Avro
- Whirr
- Sqoop
- Mahout
- Oozie

## 2011-12
- HCatalog
- Bigtop
- Ambari
- Yarn

# Hadoop: What will it take to cross The Chasm?

relative % customers

- Orgs looking for use cases & ref arch
- Ecosystem evolving to create a pull market
- Enterprises endure 1-3 year adoption cycle

The CHASM

| Innovators technology enthusiasts | Early adopters, visionaries | Early majority, pragmatists | Late majority, conservatives | Laggards, Skeptics |

time

Customers want technology & performance

Customers want solutions & convenience & reliability

*Source: Geoffrey Moore - Crossing the Chasm*

# Enterprise Big Data Flows

**Intel IT Center**

**Big Data Platform**

Unstructured Data

Log files

Exhaust Data

Social Media

Sensors, devices

DB data

**Business Transactions & Interactions**

CRM, ERP
Web, Mobile
Point of sale

Classic Data Integration & ETL

**Business Intelligence & Analytics**

Dashboards, Reports, Visualization, …

**1  Capture Big Data**
Collect data from all sources structured & unstructured

**2  Process**
Transform, refine, aggregate, analyze, report

**3  Exchange**
Interoperate and share data with applications/analytics

**intel**

# What changes from POC to large clusters?

5-100 nodes
**"Small cluster"**

4000 node
**"Hadoop at Scale"**
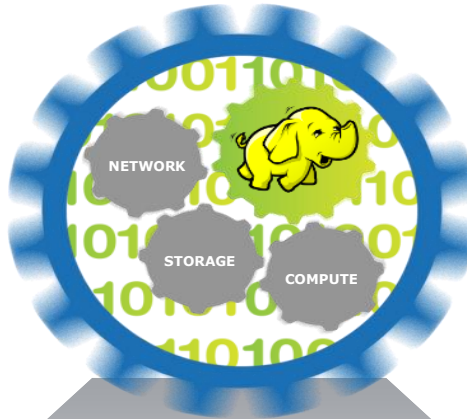
**Cluster Size**

Node    Node    Node

hadoop

- Staff & consultants are dominant costs
- Redundant networks, hardware reliability features save human capital & support
- Need to focus on simplicity

- Hardware + Power + Hosting are dominant costs
- Hardware Optimization
- Failures are inevitable, Hadoop software handles this
- Hadoop operations expertise

# Optimizing Hadoop Deployments
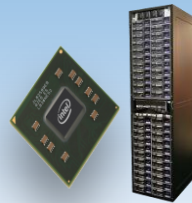
Address Potential Deployment Bottlenecks

NETWORK

STORAGE

COMPUTE

| Benchmark Tuning | Security & APIs | Compute | Disk Write/memory | Fast Fabric |
|---|---|---|---|---|
| Hi-tune Hi-Bench | Encryption | Instruction Sets | SSDs Non-volatile memory | 10GbE |

(intel)

# Talk to an Expert: Question & Answer

**Today's Experts:**

• Eric Baldeschwieler, CTO, Hortonworks  - @JERIC14

• Avik Dey, Director, Hadoop Services, Intel - @AvikonHadoop

**Submit your questions:**

• Ask questions at anytime by pressing the Question tab at the top of the player.

**Download today's content:**

• Located under the attachment tab at the top of the player

**More information:**

• www.intel.com/bigdata