# Apache Hadoop and the Big Data Opportunity in Banking

A webcast from Hortonworks & Tresata

# Presenters

- Arun C. Murthy
  - Co-founder of Hortonworks
  - Lead of NextGen MapReduce in Apache Hadoop
  - Long-time contributor and committer to Apache Hadoop

- Abhishek (Abhi) Mehta
  - Co-founder of Tresata
  - Creator of first Hadoop-powered big data & analytics platform for financial industry data

tresata

hortonworks

# Agenda

- What is Apache Hadoop?

- Creating Value from Big Data

- Tresata and Hadoop for Banking

- Future of Hadoop

# What is Apache Hadoop?

A set of **open source** projects owned by the Apache Foundation that transforms **commodity computers** and network into a **distributed service**
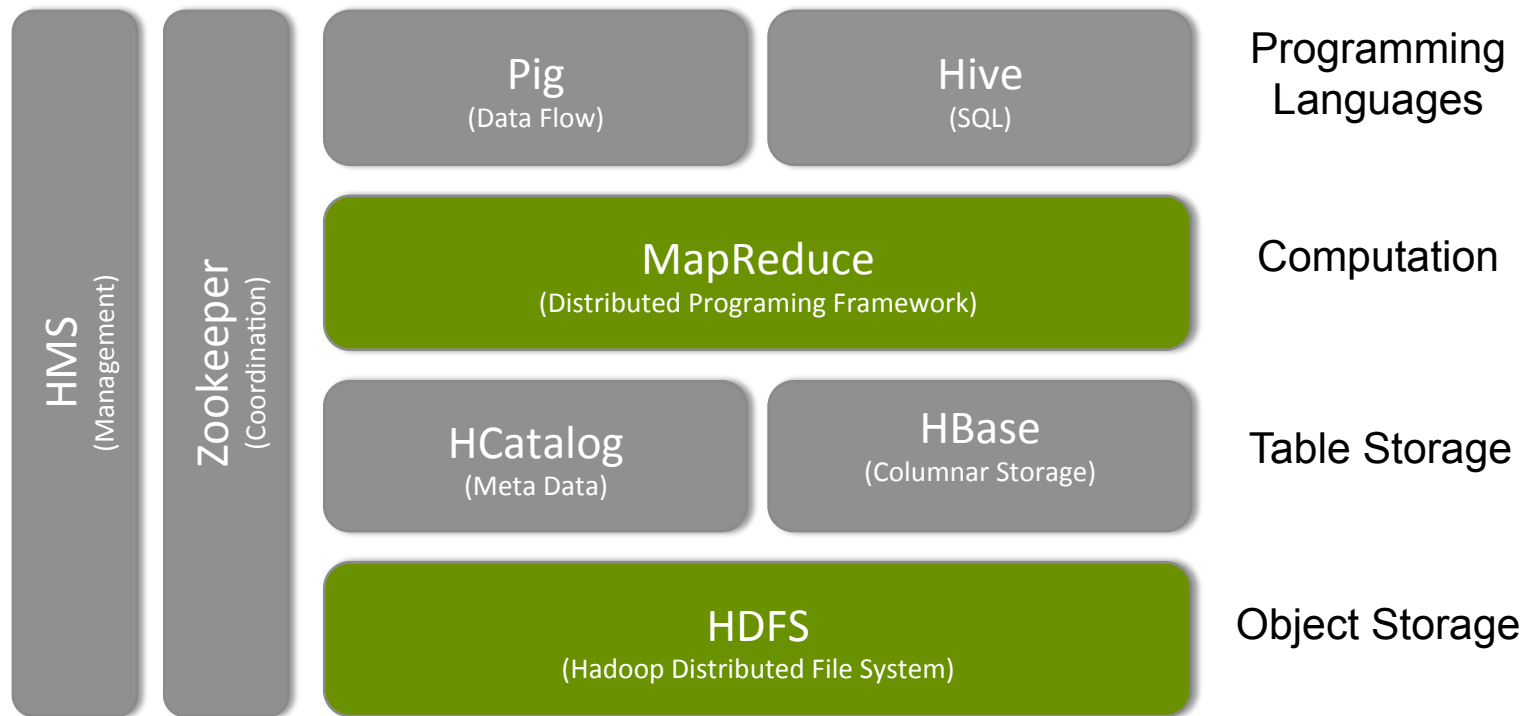


- HDFS – Stores petabytes of data reliably
- MapReduce – Allows huge distributed computations

**Key Attributes**

- **Reliable and redundant** – Doesn't slow down or loose data even as hardware fails
- **Very powerful** – Harnesses huge clusters, supports best of breed analytics
- **Scalable** – scales linearly to handle "big data" volumes
- **Cost-effective** – runs on commodity machines & network
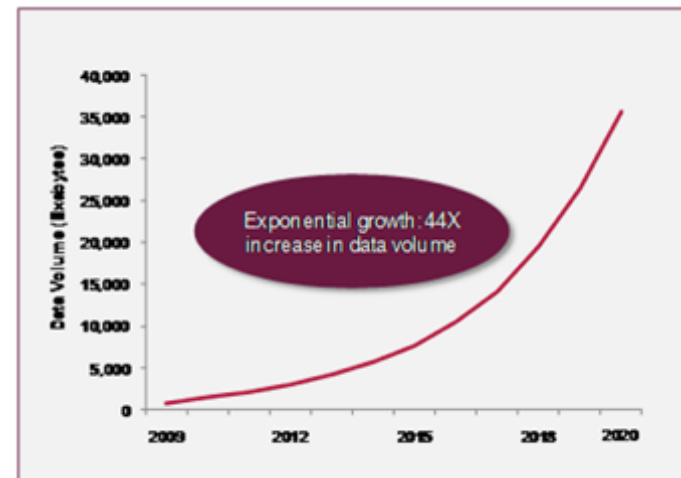- **Simple and flexible APIs** – enabling a large ecosystem of solution providers

tresata

hortonworks

# Core Apache Hadoop Projects

| HMS (Management) | Zookeeper (Coordination) | Pig (Data Flow) | Hive (SQL) | Programming Languages |
| | | MapReduce (Distributed Programing Framework) | | Computation |
| | | HCatalog (Meta Data) | HBase (Columnar Storage) | Table Storage |
| | | HDFS (Hadoop Distributed File System) | | Object Storage |

Core Apache Hadoop     Related Apache Projects

tresata

6

hortonworks

# Apache Hadoop: Why is it Transformational?

## Data Deluge (growth faster than Moore's law)

- Economist: Only 5% of generated data is structured

- Gartner: Data growth is the biggest data center hardware infrastructure challenge for large enterprises

- Forrester: Four Vs - Volume, Velocity, Variety, Variability

- Hundreds of **exabytes** of data per year!

Source: IDC Digital Universe Study, May 2011

# Apache Hadoop: Why is it Transformational?

New way of thinking about your data:

- Current

  - What data do I keep?
  - What reports do I run?
  - Sample – Store – Extrapolate → *What-If* scenarios → Variable *insights*

- New dawn

  - Store everything (viable and economical) – Process whenever!
  - Test every what-if situation, prove every hypothesis… do it in a timely manner!
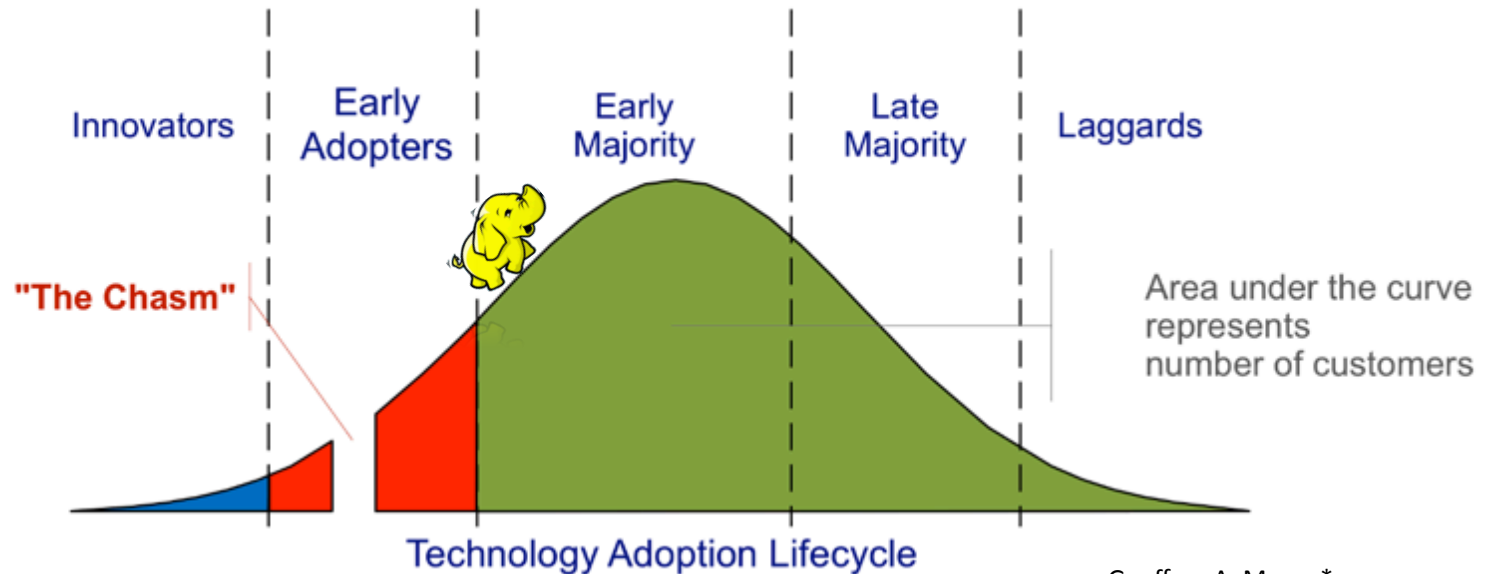  - No more down-sampled data

# 5 Ways to Create Value from Big Data



**1** Create transparency

**2** Expose variability and enable experimentation

**3** Segment populations to customize actions

**4** Replace/support human decision-making with automated algorithms

**5** Innovate new business models, products, and services

Source: McKinsey & Company report. Big data: The next frontier for innovation, competition, and productivity. May 2011.

# Crossing the Chasm

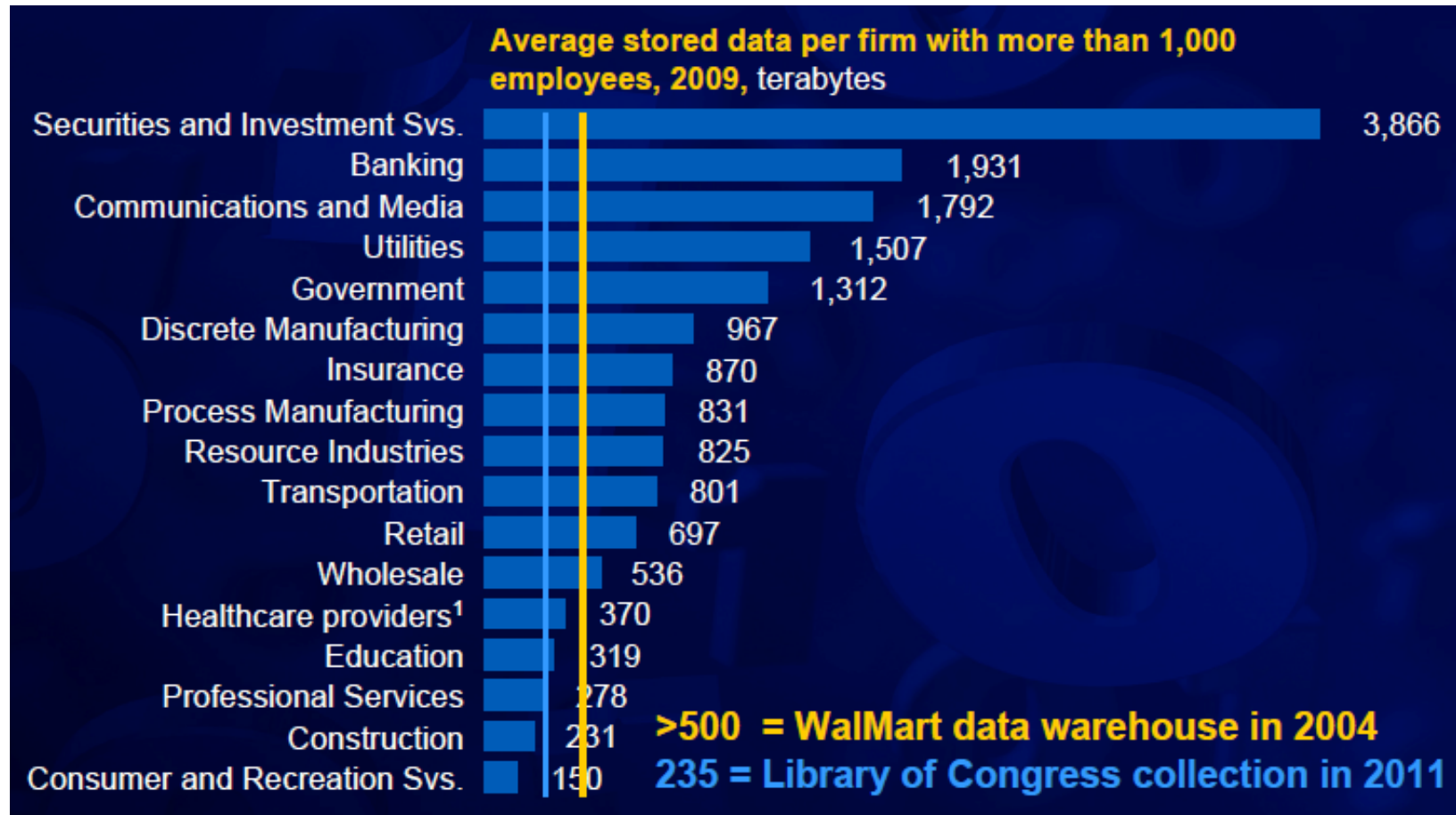Disruption: Data



**Innovators** | **Early Adopters** | **Early Majority** | **Late Majority** | **Laggards**

"The Chasm"

Area under the curve represents number of customers

Technology Adoption Lifecycle

Geoffrey A. Moore*

tresata

hortonworks

# Typical Applications & Early Adopters

analyzing web logs

data analytics

advertising optimization

**machine learning**

text mining

**web search**

mail anti-spam

content optimization

customer trend analysis

ad selection

**video & audio processing**

data mining

user interest prediction

social media

# Big Data Has Reached Every Market Sector

*Digital data is personal, everywhere, increasingly accessible, and will continue to grow exponentially*

**Average stored data per firm with more than 1,000 employees, 2009, terabytes**

| Sector | Terabytes |
|---|---|
| Securities and Investment Svs. | 3,866 |
| Banking | 1,931 |
| Communications and Media | 1,792 |
| Utilities | 1,507 |
| Government | 1,312 |
| Discrete Manufacturing | 967 |
| Insurance | 870 |
| Process Manufacturing | 831 |
| Resource Industries | 825 |
| Transportation | 801 |
| Retail | 697 |
| Wholesale | 536 |
| Healthcare providers[1] | 370 |
| Education | 319 |
| Professional Services | 278 |
| Construction | 231 |
| Consumer and Recreation Svs. | 150 |

**>500 = WalMart data warehouse in 2004**
**235 = Library of Congress collection in 2011**

Source: McKinsey & Company report. Big data: The next frontier for innovation, competition, and productivity. May 2011.

tresata

hortonworks

# Big Data Value Creation Opportunities

| **Financial Services** | **Healthcare** |
|---|---|
| • Detect fraud<br>• Model and manage risk<br>• Improve debt recovery rates<br>• Personalize banking/insurance products | • Optimal treatment pathways<br>• Remote patient monitoring<br>• Predictive modeling for new drugs<br>• Personalized medicine |
| **Retail** | **Web / Social / Mobile** |
| • In-store behavior analysis<br>• Cross selling<br>• Optimize pricing, placement, design<br>• Optimize inventory and distribution | • Location-based marketing<br>• Social segmentation<br>• Sentiment analysis<br>• Price comparison services |
| **Manufacturing** | **Government** |
| • Design to value<br>• Crowd-sourcing<br>• "Digital factory" for lean manufacturing<br>• Improve service via product sensor data | • Reduce fraud<br>• Segment populations, customize action<br>• Support open data initiatives<br>• Automate decision making |

tresata

hortonworks
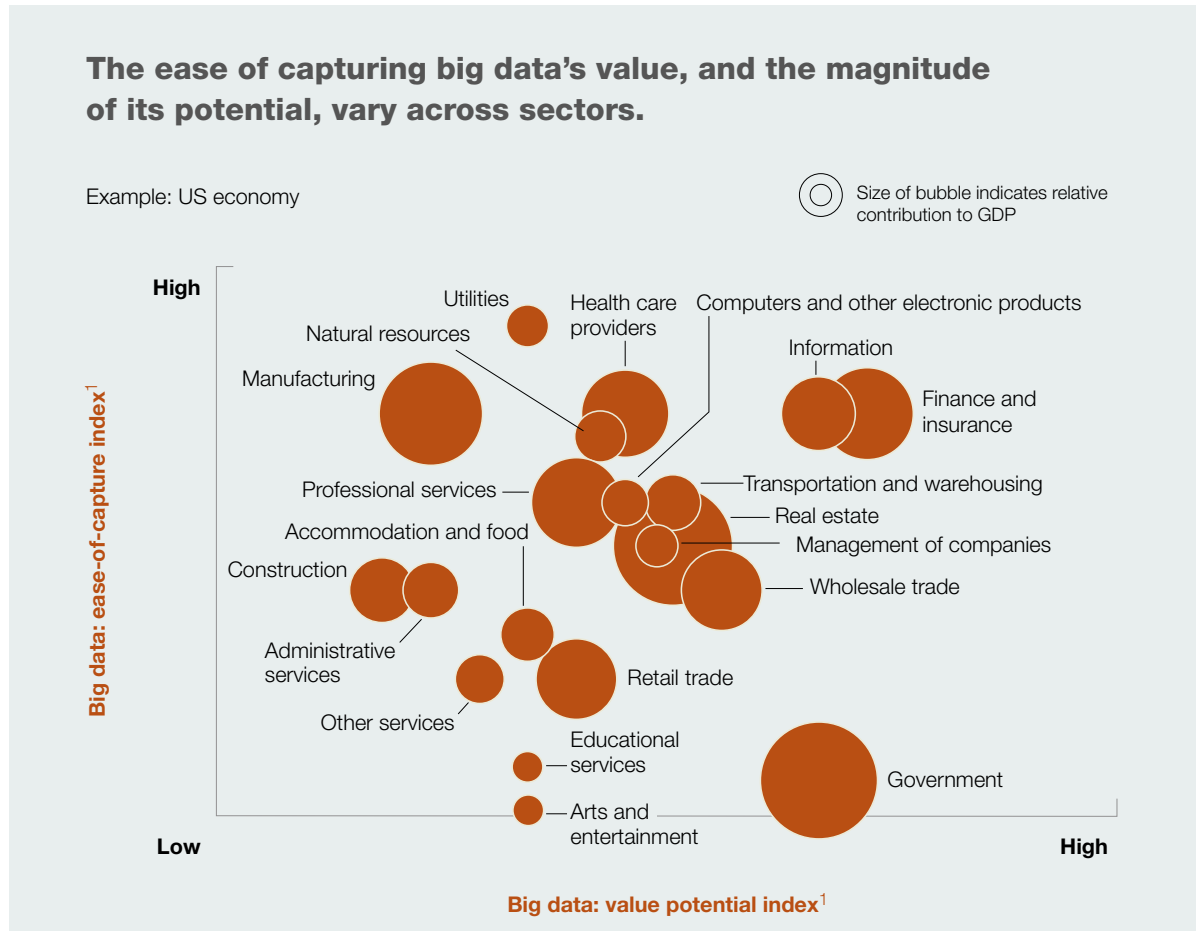
# Why can you bank on Hortonworks?

- Architects of Hadoop since big bang, circa 2006

- Real world experience supporting the largest Hadoop install in the world:
  - 50,000 node footprint
  - Over 200PB of data
  - 24x7 service
  - Billions of ad dollars

- ***We've taken the 3am calls to fix stuff when it breaks!***

**tresata**

**hortonworks**

# Tresata = Hadoop for Banking

# What We Believe

1.  **Data will reboot** the financial service industry

2.  Data growth is **viral**…existing tools can't keep up

3.  The **economics** to store, process, analyze and visualize **all** of your data makes it a **'no-brainer'** to do so

4.  **Big Data capabilities** are needed to address **business problems**

# and What McKinsey Said



The ease of capturing big data's value, and the magnitude of its potential, vary across sectors.

Example: US economy

Size of bubble indicates relative contribution to GDP

Big data: ease-of-capture index[1]

High / Low

Big data: value potential index[1]

Low / High

Labels: Utilities, Health care providers, Computers and other electronic products, Natural resources, Information, Manufacturing, Finance and insurance, Professional services, Transportation and warehousing, Accommodation and food, Real estate, Construction, Management of companies, Wholesale trade, Administrative services, Retail trade, Other services, Educational services, Government, Arts and entertainment

Source: mcksinsey global institute october 2011

17

# The Opportunity

1.  **1-5% of data** in a financial institution **is analyzed**

2.  **Not all data is stored**

3.  **Top-down macros** cannot be implemented or acted on

4.  **New approaches to data & analytics** are essential

Source:  tresata research

**tresata**

**hortonworks**

# The Business Problems

1. **Storage & retrieval –** archive data on disk

2. **Consumer behavior analysis –** as it happens

3. **Modeling & Analytics** – model entire data sets…

4. **Single View Of Customer** – structured & unstructured data

tresata

hortonworks

# Hadoop in Banking Today

1. **Early adopter**, like any other technology trend

   a. Interest is **global**, not just limited to the US

   b. **Approved for use** by most technology teams & architects

   c. **Proof of concepts** ongoing at most financial services institutions

2. **Broad agreement** that:

   a. **Hadoop** will be the **Big Data operating system**

   b. A hadoop powered **Data processing & analytics platform** will spur rapid adoption

   c. Need to apply to **business problems** with revenue/profit impact

# Elephants in the Room

1. **Resources**

   a.  **Talent** – how many MapReduce developers can we get

   b.  **Applications** – data, analytics and business problems are the same, why should each institution build their own hadoop applications to run the same processes

   c.  **Essentials** – how to manage security, provisioning, & performance

2. **Implementation**

   a.  **Integration** – fit with existing technology infrastructures & business processes

   b.  **Support** – experience with managing thousands of nodes/ servers

   c.  **Open Source** – 'out of the box' applicability

tresata

hortonworks

# How Tresata Changes the Game…

1. **Our Application**

   a. **FS for Hadoop** – fully built on hadoop.  Store, process, analyze and visualize  leveraging the full power of hadoop

   b. **Processing Pipeline** – automated ingestion, cleaning, de-duping, matching engine built for financial data

   c. **Analytics** – massively parallel scoring and algorithm containers codified by business problem


2. **Tame the elephants** (making it work at scale)

# Tresata Case Study

## A. Client Business Problem

 i. Problem – Process data and score for >**30 MM client** applications

 ii. Data Sources – **23** separate data sources, **multiple time series**

 iii. Raw Variables – **100 variables** per client per data source

 iv. Current state - **expensive** legacy platform, algorithms developed on **sub-samples**, **unable to scale** algorithms to full data set

## B. Tresata Solution

 i. Data Engine – **automated** data import, cleaning, matching, scoring

 ii. Compute Engine – algorithms process & score >30MM in **minutes**

 iii. Integration – work with existing tools and processes

 iv. Scalable deployment – **Big Data as a Service** delivery

tresata

hortonworks

# Tresata Big Data Application

# Tresata Analytics



Sacramento Metro Area Tresata Equity Indicator - Individual Property Level

# Tresata + Hortonworks

1. **Commitment to train**

   a. Series of webinars

   b. Hadoop training tailored for financial services

   c. Dedicated training programs (tailored for client needs)

2. **Commitment to support**

   a. Production scale support model (Tresata certified for Financial Svcs)

   b. Meet and/or exceed industry standards on distro

**tresata**

**hortonworks**

# Where Hadoop is Going

# The Future

- Hadoop is going mainstream

  – Amazon, Microsoft, Oracle, EMC, NetApp, etc.

- Apache Hadoop is covering new ground (Hadoop .Next)

  – Much of the development led by Hortonworks

  – More scale, more performance

  – Other paradigms than MapReduce for data processing

  – Enhanced operability and management (Ambari)

  – Metadata management (Hcatalog)

tresata

hortonworks

# Technology Roadmap

| | |
|---|---|
| **Hadoop.Now – Making Apache Hadoop Accessible**<br>• Release the most widely deployed version of Hadoop ever (0.20.205)<br>• Release directly usable code via Apache (RPMs, .debs…)<br>• Frequent sustaining releases off of the stable branches | **Q4 2011** |
| **Hadoop.Next – Next Generation Apache Hadoop**<br>• Address key product gaps (HBase support, HA, Management…)<br>• Enable community & partner innovation via modular architecture & open APIs<br>• Work with community to define integrated stack | **2012**<br>(Alphas starting late 2011) |

tresata

hortonworks

# Next Steps

- Engage with Hortonworks & Tresata
  - www.hortonworks.com
  - www.tresata.com

- Additional Webcast Series
  - Reference Architecture for Hadoop in Financial Services - Nov 15
  - Hadoop in Financial Services DEEP DIVE - Dec 6
  - www.hortonworks.com/webcasts

- Hortonworks party @ Hadoop World
  - Tuesday Nov 8 at 7pm
  - Inc Lounge at the Time Hotel
  - RSVP: hortonworks.eventbrite.com

tresata

hortonworks

# Questions?