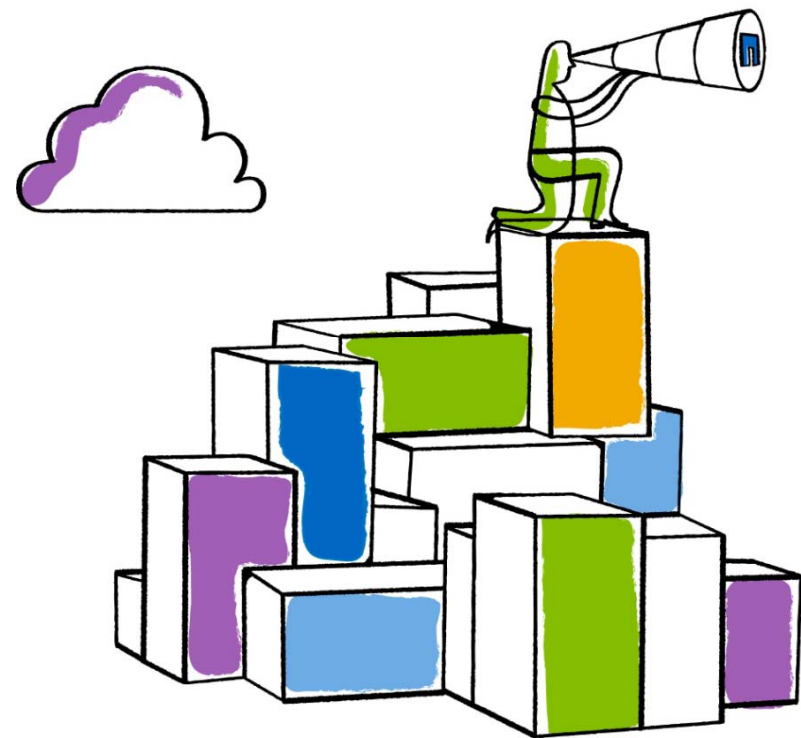# Big Data – Trends to Watch

Bill Peterson

NetApp

September, 2012

# HELLO
## My name is

**Bill Peterson**

@thebillp

# What I hope to accomplish today...

Think beyond structured and unstructured data

&

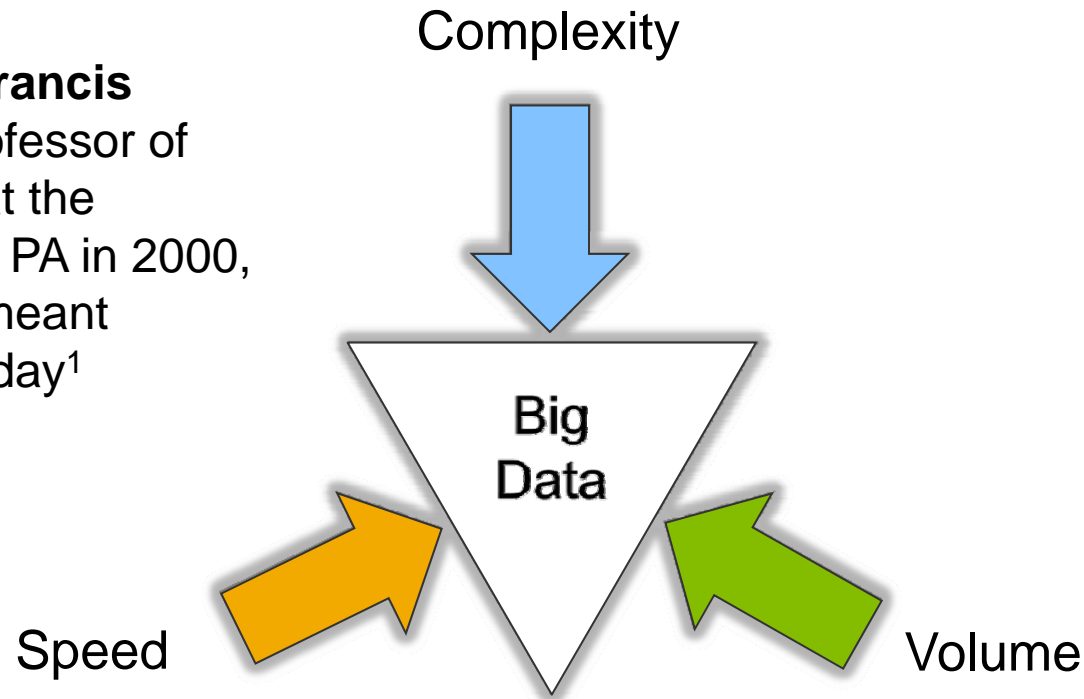Think beyond big data hype

...and avoid this.

# What is "Big Data"?

**"Big Data" refers to datasets whose volume, speed and complexity is beyond the ability of typical tools to capture, store, manage and analyze.**

Coined by **Francis Diebold**, professor of economics at the University of PA in 2000, when "Big" meant Gigabytes / day[1]

Complexity

Big Data

Speed

Volume

# Quantifying The Big Data Challenge

**NetApp**

## 60 Zettabytes

Estimated size of the digital universe in 2020

## 5 Billion
smart phones

## 30 Billion

pieces of new content to Facebook per month

Sensors
Video
Music
Location
Weblogs

**80%** of data is unstructured

## Growth Over the Next Decade:

Servers (Phys/VM):  10x
Data/Information:     50x
#Files:                      75x
IT Professionals:     <1.5x

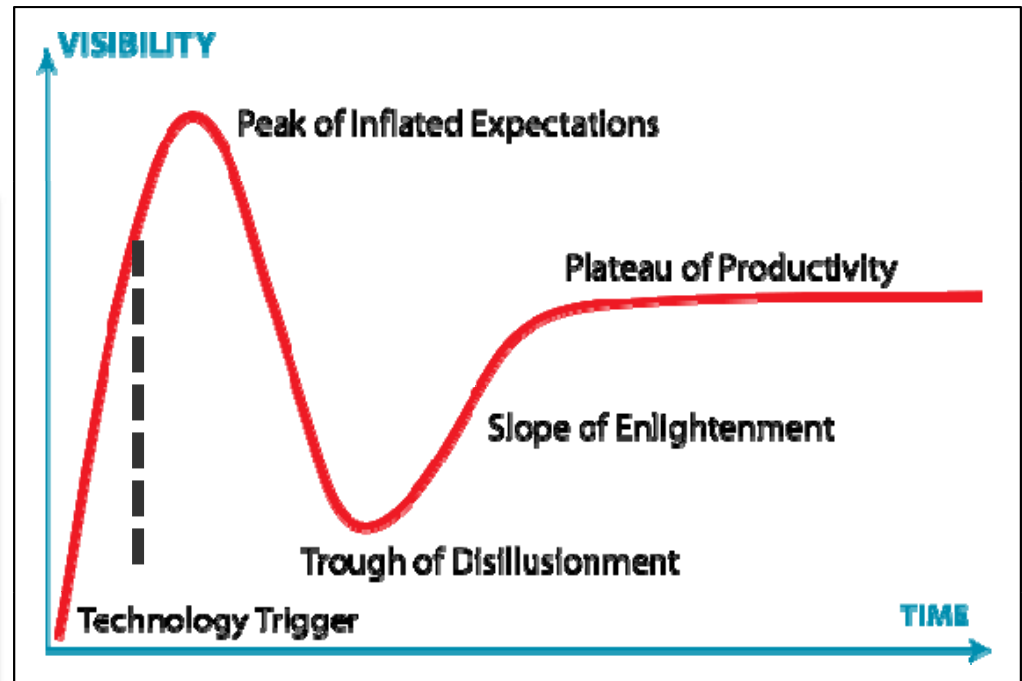**Source:** *Gantz, John and Reinsel, David, "Extracting Value from Chaos", IDC IVIEW, June 2011, page 4.*

**VISIBILITY**

Peak of Inflated Expectations

Plateau of Productivity

Slope of Enlightenment

Trough of Disillusionment

Technology Trigger

**TIME**

# The Big Data Push



High

Data Structure

Low

Performance

"Decision Support"
High Speed Analytics

NoSQL
Columnar
DBs

CS/DW

"Big Data"
High Bandwidth Throughput
Big Data Content

Sat Ground Stations

DVS

MV

Tech

HPC

Big Data Push

"Enterprise Applications"
Shared, Virtualized Infrastructure
Integrated Data Protection
Secure Multi-Tenancy

Tier 1 BP
OLTP

Tier 3
OLTP

App Dev

IT Infra

Web Infra

Content Repositories

Home Dirs

Large Block,
Sequential I/O
100s GB/sec

Small Block,
Random I/O
(100s KIOPS)

# What Does This Mean to You?



Information Becomes
a Propellant to the Organzation

Inflection
Point

Data Becomes a
Burden to IT Infrastructure

Business or Mission Velocity

2010

2020
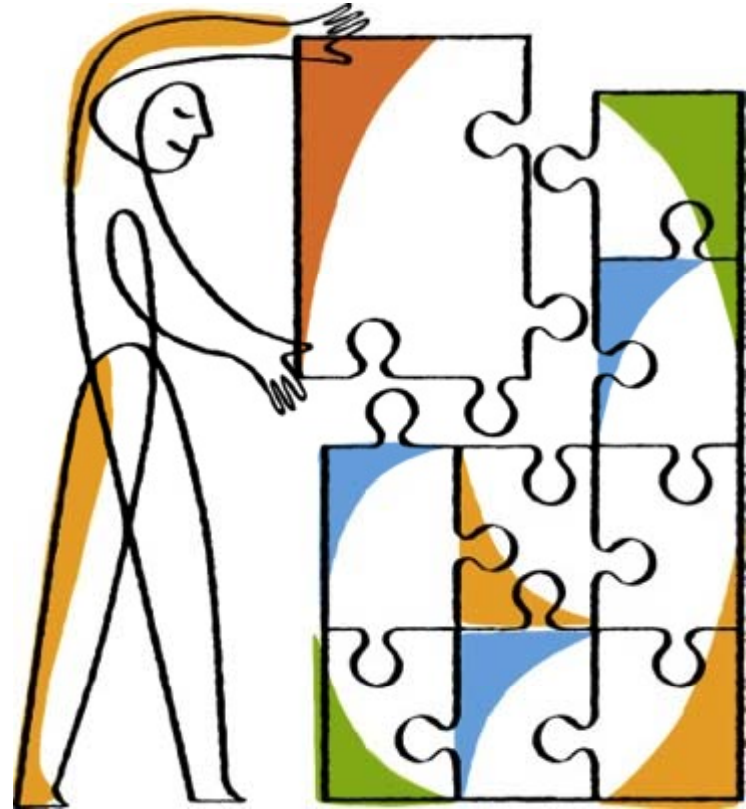
You are also at an Inflection Point:  You also have a
decision to make, as "business as usual" may not cut it!

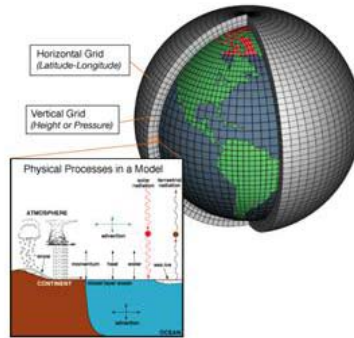# Dispelling the Misconceptions About Big Data

# Big Data Is NOT New



National Oceanographic and Atmospheric Administration

**3.5 billion** observations per day from NOAA sensors

"The Tank" 24-hour Ring Buffer

Data Assimilation System

6-Hour Update Cycle

Initial Conditions

Global Atmospheric Model (HIRAM)
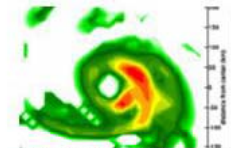
HPC Models

Other High Res Specialized Models:
- Hurricane (3 Km)
- Thunderstorms
- Tornados
- Fire Weather
- Ocean Models
- Volcanic Ash
- Etc.

**15 million** information products per day

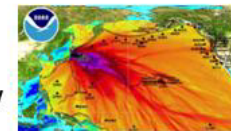Forecast & Warning Guidance to Public and Private Sector Forecasters

Horizontal Grid (Latitude-Longitude)

Vertical Grid (Height or Pressure)

Physical Processes in a Model

Thunderstorm Warnings

Hurricane Warnings

Flood Warnings

Fire Weather

Tsunami Energy Vectors)

National Weather Service Data Stream

30 PB of New Data Annually
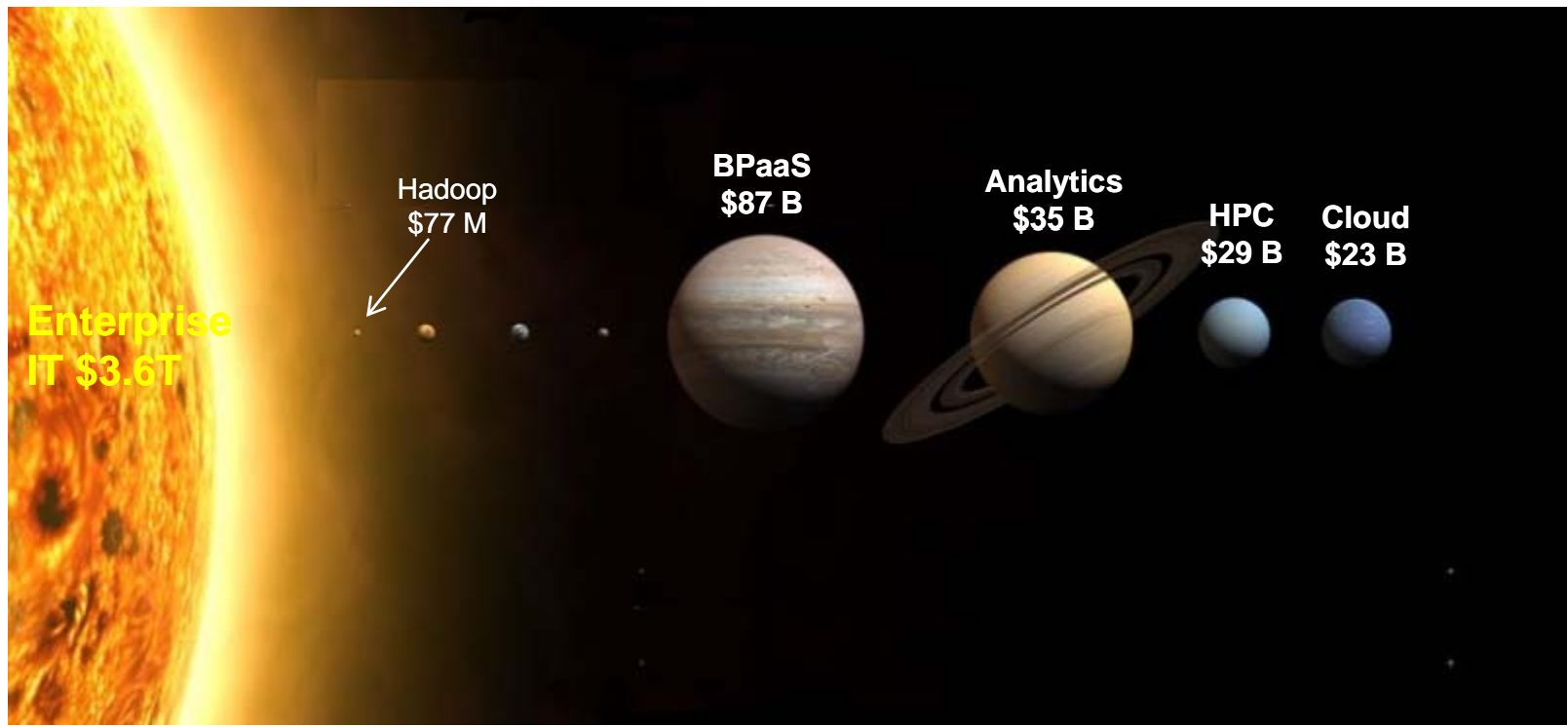
# Big Data = Big Analytics = Hadoop?

- That's What The Media Hype Implies, but it is NOT true!
- Traditional analytics (BI/DSS/DW) dominates the analytics market
- Like other technologies vying to gain broad adoption in Enterprise IT (e.g., Traditional Analytics, HPC & Cloud), it shows promise



Hadoop $77 M

BPaaS $87 B

Analytics $35 B

HPC $29 B
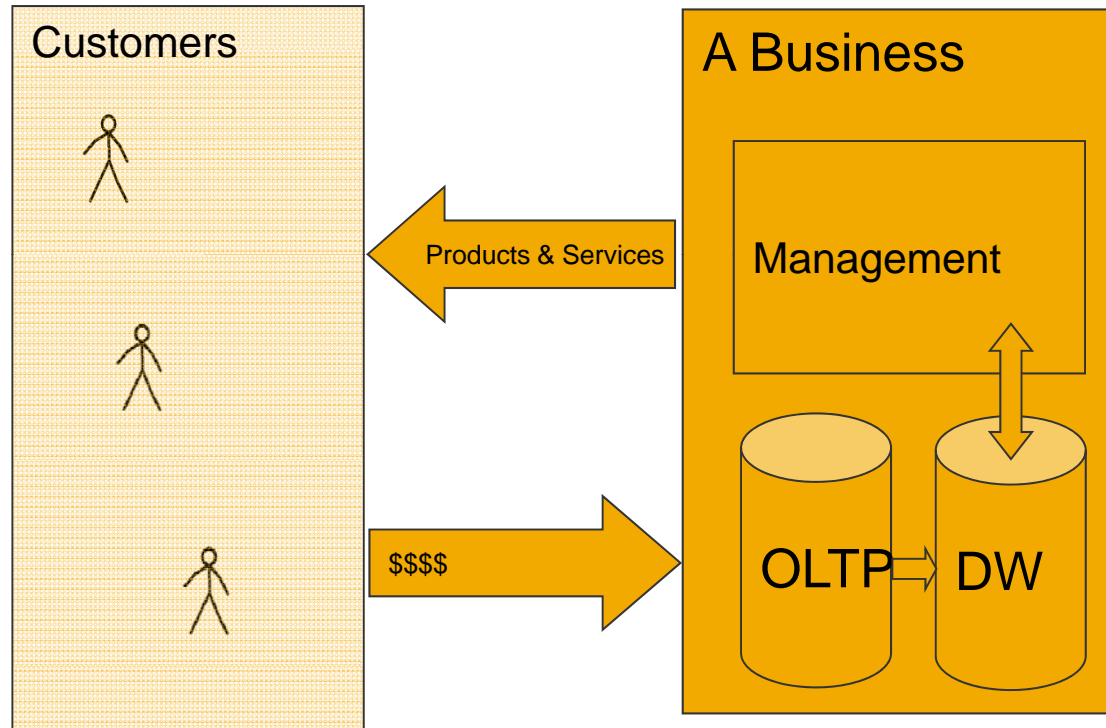
Cloud $23 B

Enterprise IT $3.6T

# Analytics

# Why Decision Support Systems are important?



DSS enables businesses to run "Closed Loop", ultimately improving their business through the use of feedback mechanisms.

# Big Analytics – An Emerging Market

**NoSQL / Column DBs**

aster data
*more data. big insights.*

Cassandra

COUCHBASE

HADAPT

hadoop

mongoDB

ParAccel

talend*
*open integration solutions*

VERTICA
An HP Company

VoltDB

**Legacy DBs**

TERADATA

ORACLE

NETEZZA
an IBM Company

GREENPLUM
A DIVISION OF EMC

**Middleware & Apps**

ATTIVIO
INSIGHT THAT MATTERS

DIGITAL REASONING

pentaho

tableau

INFORMATICA
The Data Integration Company

QUEST SOFTWARE

splunk>

**Open Source Distributors**

cloudera

Hortonworks

MAPR
TECHNOLOGIES
EASY. DEPENDABLE. FAST.

**Cloud & Cyber**

amazon
webservices

GOGRID

Google

KARMASPHERE

KEYW

tidemark

YAHOO!

**Integration Services**

BAYSWATER
BUSINESS SOLUTIONS THAT EMPOWER

Datameer

lunexa
illuminate your data

platfora

ZALONI

**Compute**

hp

DELL

sgi

StackIQ

SUPERMICRO

**Network**

CISCO

ARISTA

Mellanox
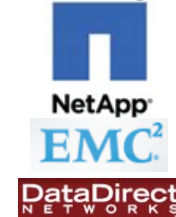TECHNOLOGIES
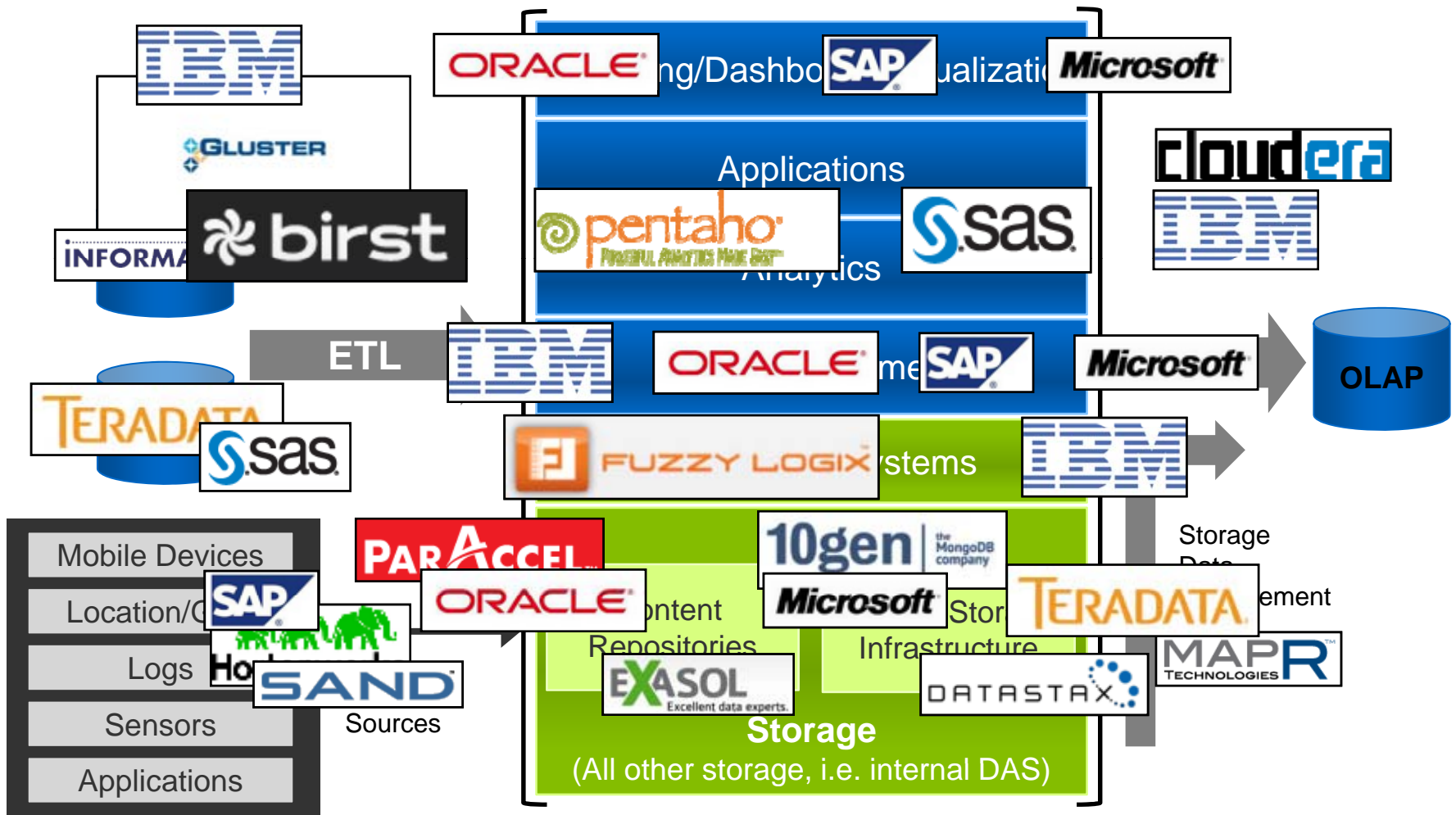
SOLARFLARE
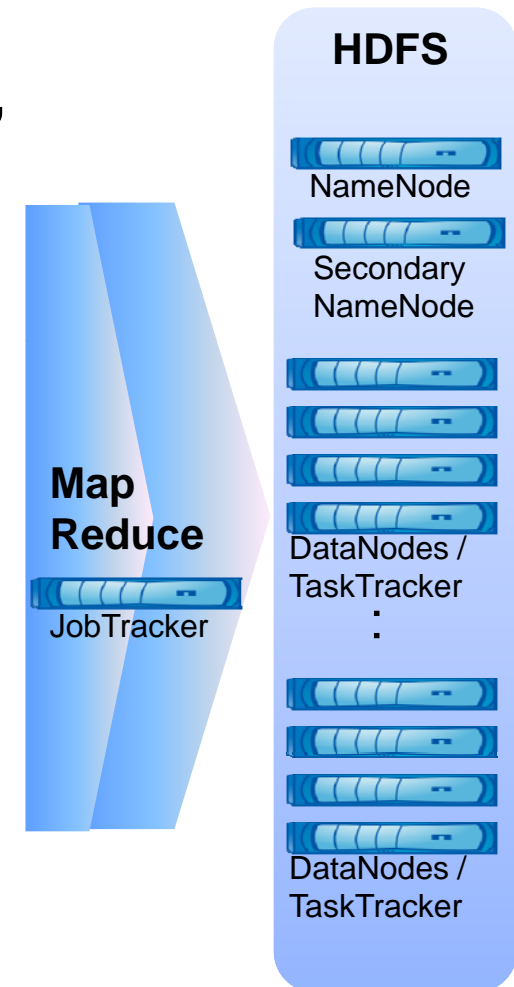
**Storage**

NetApp

EMC²

DataDirect
NETWORKS

# Analytics & Enterprise Apps Environment
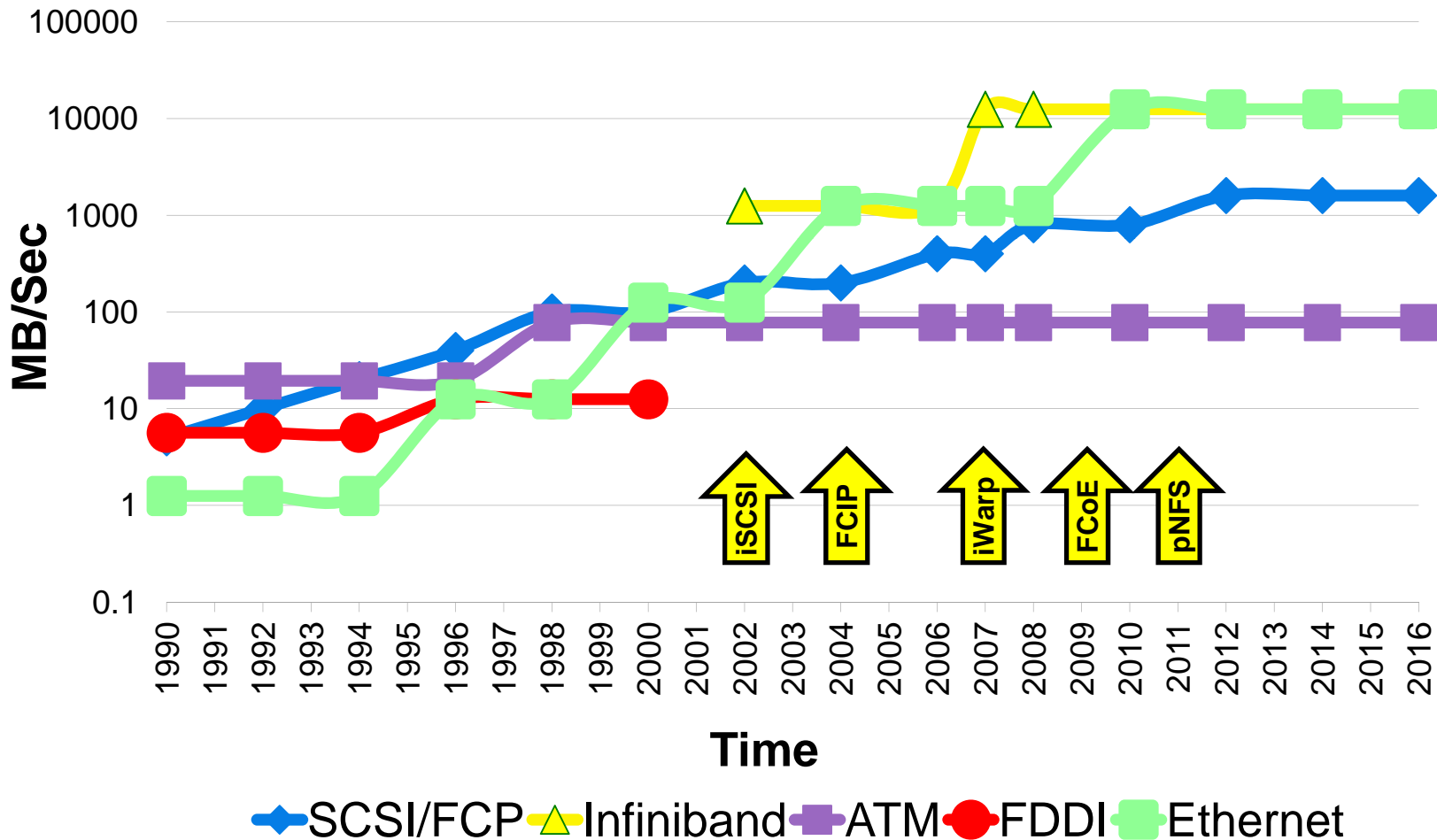
# What Does Hadoop Look Like Today?

- Runs on a collection of cheap, commodity servers, in a distributed, shared nothing architecture

- Two key components

  - HDFS

    - Hadoop Distributed File System

  - MapReduce

    - Programming model for processing and generating large datasets

**HDFS**

NameNode

Secondary NameNode

DataNodes / TaskTracker

:

DataNodes / TaskTracker

**Map Reduce**

JobTracker

# Ethernet's Relentless March

Data will be growing by 50x, but bandwidth only by 10x!

# Why Should You Care?

It's the Value of your data

**THOMSON REUTERS**

5 Billion Records
Anywhere, Anytime
Faster time to market
50% Increase in Revenue

**Sprint**

Over 1PB of data
Growth of 175% YOY
90 days of data within
24 hours of a failure

- **Top line revenue**
  - Leverage their data assets into business advantage

- **Bottom Line savings**
  - Lower the cost of compliance
  - Manage ever growing data efficiently

# AutoSupport: Hadoop Use Case at NetApp

- "Call-home" service for all NetApp® systems
- Foundation of NetApp proactive support strategies
- Machine-generated data doubles every 16 months

| CHALLENGE | NETAPP SOLUTION | BENEFITS |
|-----------|-----------------|----------|
| **4 weeks to run a query on 24 billion unstructured records** | **10-node Hadoop Cluster w/ shared Storage** | **Time reduced from 4 weeks to 10.5 hours** |
| **Impossible to run a query: 240 billion unstructured records** | | **Previously impossible, now achievable in just 18 hours** |

**"NetApp ASUP is a mission-critical application"**

# Analytics of Tomorrow

- Traditional & Big Analytics side-by-side for years to come

- Hadoop moves to shared, virtualized infrastructure, for better efficiency and ease of management:
  - Hadoop remains logically distributed, shared nothing, but runs on a virtualized shared everything architecture (e.g., FlexPod for Vmware + eSeries)
  - Same as above, except Hadoop becomes logically shared everything, as HDFS is replaced by a parallel file system (e.g., Lustre Cluster, StorNext or GPFS)

- Enterprise class resiliency (no SPoF) and reliability with HPC-like performance (no need for triplicas)

- Use of a single copy of data for the map phase (higher storage utilization)

- Natural intersection with Cloud (Analytics as a Service)

# Summary

- Despite the hype, **Big Data is not new and is more than just analytics!** (Many agencies and private companies have struggled with Big Data for decades)

- **Analytics:** Traditional BI/DSS analytics still dominate. Importance of newer NoSQL & Columnar DB applications, enabled by MapReduce will grow with the growth of multi-structured data

- Big Data applications, such as Hadoop, will need to adopt shared, virtualized infrastructure (and its management benefits) if they are to be widely adopted by Enterprise IT

questions

**YOU'VE GOT**

?

**I'VE GOT**
rambling responses that sound like

answ

@thebillp or william.peterson@netapp.com

ers