



Education

# **NextGen Infrastructure for Big Data**

Anil Vasudeva, President & Chief Analyst, IMEX Research

Author: Anil Vasudeva, IMEX Research

- The material contained in this tutorial is copyrighted by the SNIA and author unless otherwise noted.
- Member companies and individual members may use this material in presentations and literature under the following conditions:
  - ◆ Any slide or slides used must be reproduced in their entirety without modification
  - ◆ The SNIA must be acknowledged as the source of any material used in the body of any document containing material from these presentations.
- This presentation is a project of the SNIA Education Committee.
- Neither the author nor the presenter is an attorney and nothing in this presentation is intended to be, or should be construed as legal advice or an opinion of counsel. If you need legal advice or a legal opinion please contact your attorney.
- The information presented herein represents the author's personal opinion and current understanding of the relevant issues involved. The author, the presenter, and the SNIA do not assume any responsibility or liability for damages arising out of any reliance on or use of this information.  
**NO WARRANTIES, EXPRESS OR IMPLIED. USE AT YOUR OWN RISK.**

## ➤ NextGen Infrastructure for Big Data

- ◆ This session will appeal to Business Planning, Marketing, Technology System Integrators and Data Center Managers seeking to understand the drivers behind the demand for and rise of Big Data.

### ◆ Abstract

- ◆ The internet has spawned an explosion in data growth in the form of data sets, called Big Data, that are so large they are difficult to store, manage and analyze using traditional RDBMS which are tuned for Online Transaction Processing (OLTP) only. Not only is this new data heavily unstructured, voluminous and streams rapidly and difficult to harness but even more importantly, the infrastructure cost of HW and SW required to crunch it using traditional RDBMS, to derive any analytics or business intelligence online (OLAP) from it, is prohibitive.
- ◆ To capitalize on the Big Data trend, a new breed of Big Data technologies (such as Hadoop and others) many companies have emerged which are leveraging new parallelized processing, commodity hardware, open source software and tools to capture and analyze these new data sets and provide a price/performance that is 10 times better than existing Database/Data Warehousing/Business Intelligence Systems.

### ◆ Learning Objectives

- ◆ The presentation will illustrate the existing operational challenges businesses face today using RDBMS systems despite using fast access in-memory and solid state storage technologies. It details how IT is harnessing the emergent Big Data to manage massive amounts of data and new techniques such as parallelization and virtualization to solve complex problems in order to empower businesses with knowledgeable decision-making.
- ◆ It lays out the rapidly evolving big data technology ecosystem - different big data technologies from Hadoop, Distributed File Systems, emerging NoSQL derivatives for implementation in private and hybrid cloud-based environments, Storage Infrastructure Requirements to Store, Access, Secure, Prepare for analytics and visualization of data while manipulating it rapidly to derive business intelligence online, to run businesses smartly.

# Big Data in IT Industry Roadmap

## IT Industry Roadmap

### Analytics – BI

#### Predictive Analytics - Unstructured Data

From Dashboards Visualization to Prediction Engines using Big Data.

### Cloudization

#### On-Premises > Private Clouds > Public Clouds

DC to Cloud-Aware Infrast. & Apps. Cascade migration to SPs/Public Clouds.

### Automation

#### Automatically Maintains Application SLAs

(Self-Configuration, Self-Healing<sup>©IMEX</sup>, Self-Acctg. Charges etc.)

### Virtualization

#### Pools Resources. Provisions, Optimizes, Monitors

Shuffles Resources to optimize Delivery of various Business Services

### Integration/Consolidation

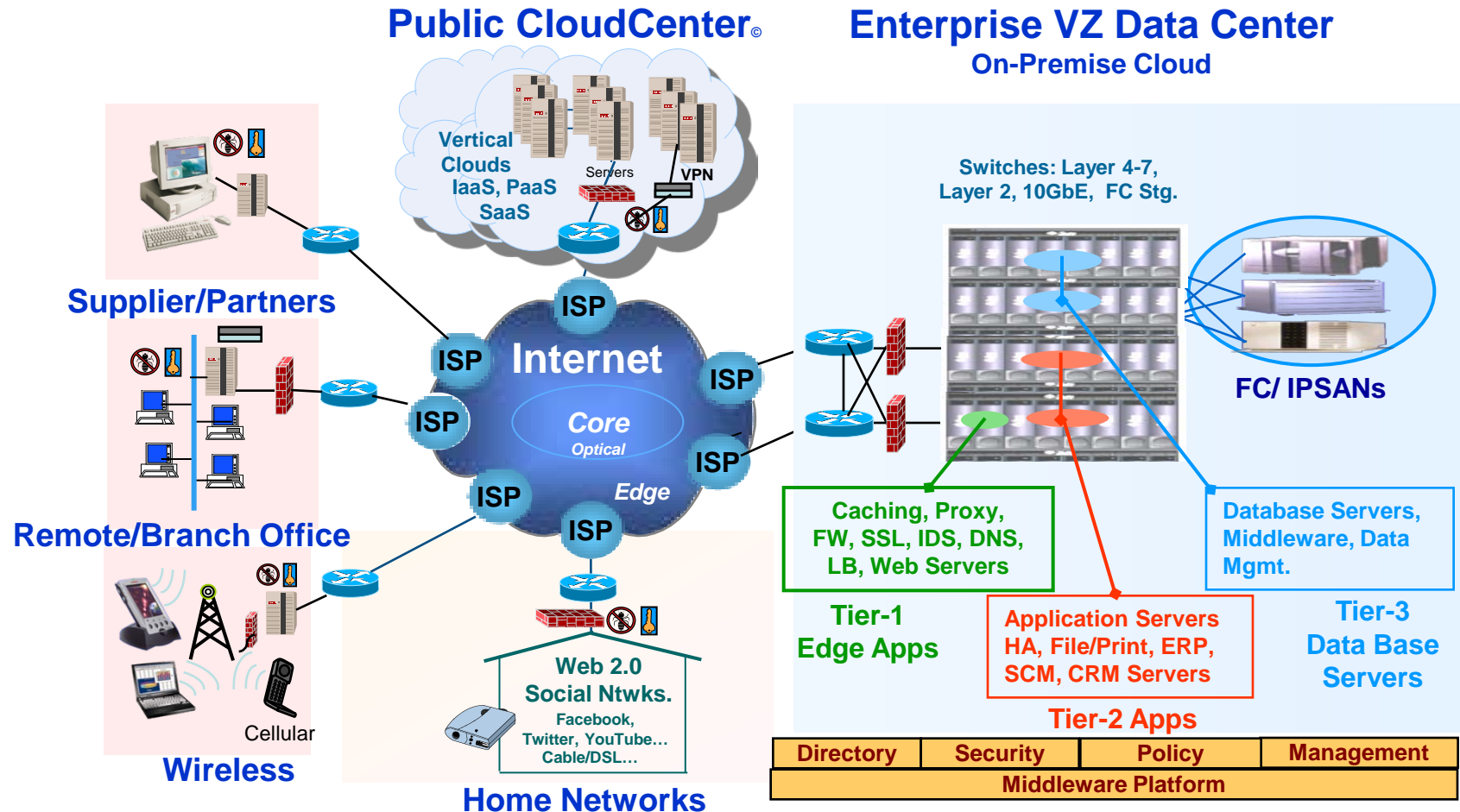
#### Integrate Physical Infrast./Blades to meet CAPSIMS<sup>©IMEX</sup>

Cost, Availability, Performance, Scalability, Inter-operability, Manageability & Security

### Standardization

#### Standard IT Infrastructure- Volume Economics HW/Syst SW

(Servers, Storage, Networking Devices, System Software (OS, MW & Data Mgmt. SW))



**Request for data from a remote client to a Data Center or Cloud crosses a myriad of systems and devices. Key is identifying bottlenecks & improving performance**

# Harnessing Big Data for Business Insights



Information is at the center of  
New Wave of opportunity

**44x**

as much Data and Content  
Over Coming Decade

2009  
800,000 petabytes

**Velocity**  
**Variety**  
**Volume**

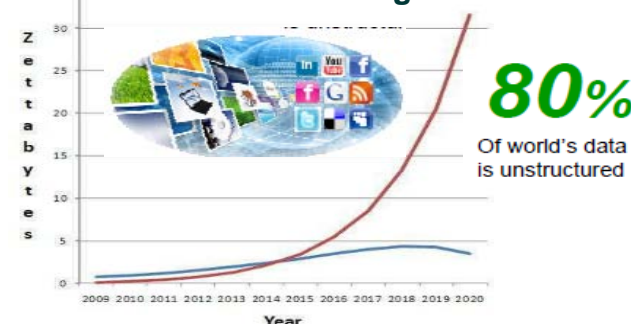


2020  
35 zettabytes

Majority of data growth is being  
driven by unstructured data  
and billions of large objects



80% of world's data is unstructured  
driven by rise in Mobility devices,  
collaboration machine generated data.



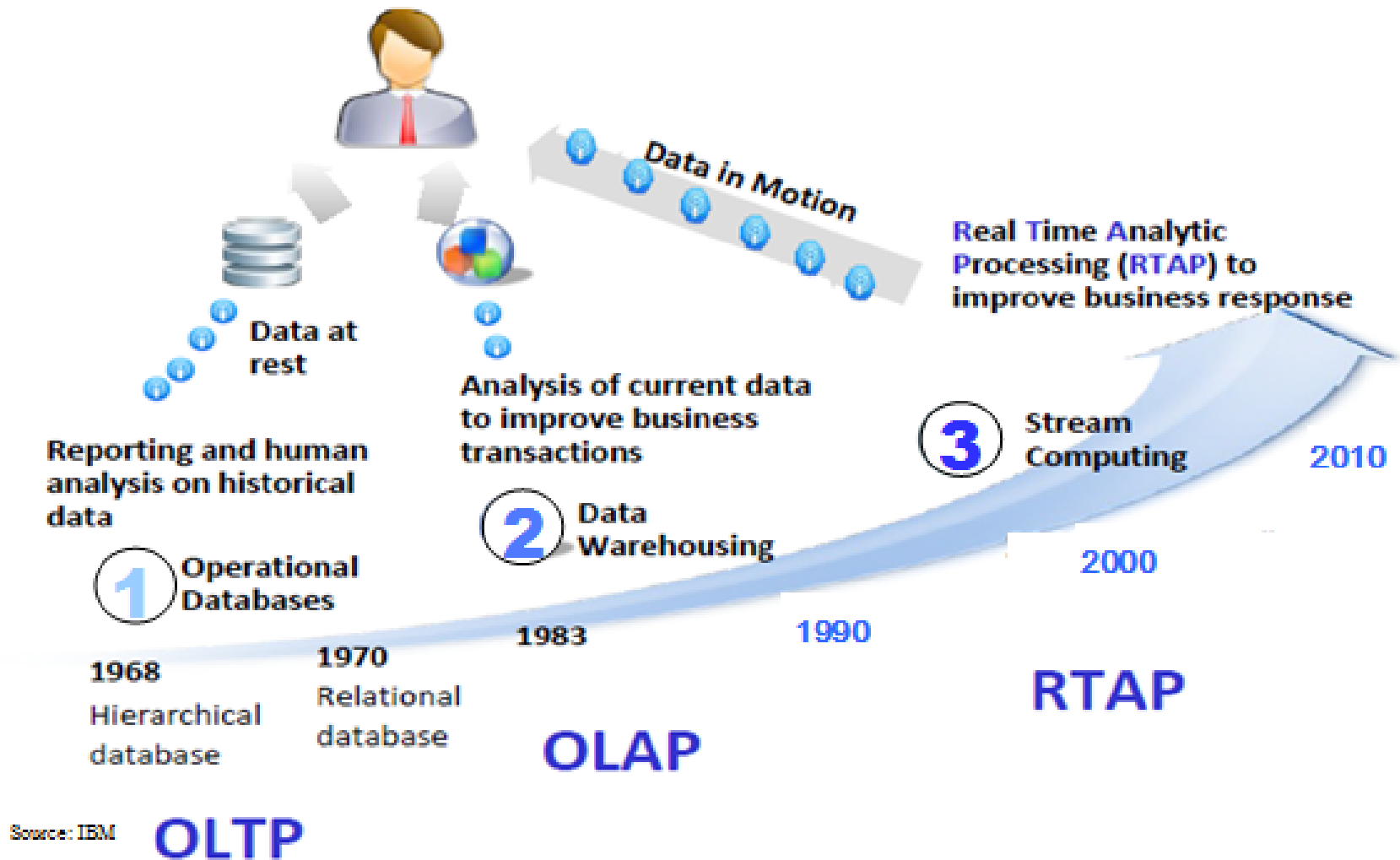


Unstructured Big Data can provide Next Gen Analytics to help businesses make informed, better decision in:

- Product Strategy
- Targeting Sales
- Just-In-Time Supply-Chain Economics
- Business Performance Optimization
- Predictive Analytics & Recommendations
- Country Resources Management



# Corporate Need: Real Time Analytics

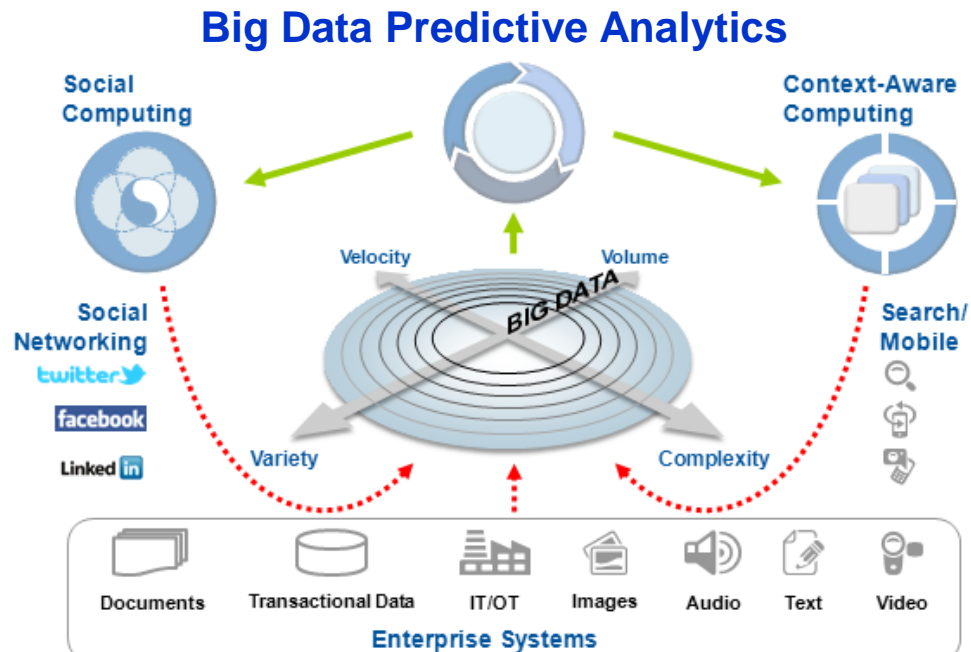




# Corporate Need: Business Insights

Item	Issue	Solution
<b>Store</b>	<b>Information Exploding</b> Volume: Digital Content doubling every 18 months. Velocity: >80% growth driven from unstructured data. Variety: sources of data changing	<b>A unified information/content storage methodology</b> that enables users to manage the volume, velocity and variety of information from multiple sources
<b>Manage</b>	<b>Complexity in "managing" information.</b> - Need to classify, synchronize, aggregate, integrate, share, transform, profile, move, cleanse, protect, retire	<b>A solution portfolio of tools and services</b> to manage all types of information in a hybrid storage environment
<b>Analyze</b>	<b>Current solutions limited to BI</b> tools focused on structured and lagging information	Build/buy packaged <b>Real-Time Predictive Analytical Solutions</b> for unstructured analytics tools
<b>Collaborate</b>	<b>Multiple access methods</b> needed to meet needs of a diverse audience.	<b>Centralized share, collaborate</b> and act on insights anytime, anywhere on any device.
<b>Model/ Adapt</b>	Ability to <b>understand how the information impacts the business.</b> How to transfer to action.	<b>Model Information on current operations w/potential strategy impact.</b> Leverage Tech. to adapt.

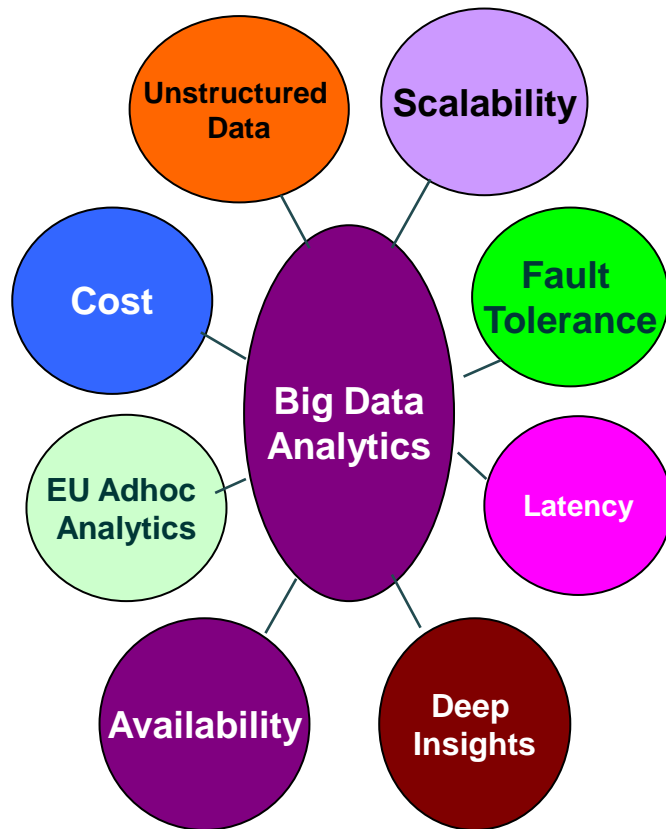
# Opportunity: Converting Big Data Deluge into Predictive Analytics & Insights



Personal Location Services Data Generated by	TB/Year
Navigation Devices	600
Navigation Apps on Phones	20
Smart Phone Opt-In Tracking	1000
Geo Targeted Ads	20
People Locator (Emergency Calls/Search..)	10
Location based Services (e.g.Games)	5
Other	45
<b>Total (Est.)</b>	<b>1700</b>

# Issues with Existing RDBMS

## Key Issues with RDBMS Technologies



### Handling Mixed Unstructured Data

- RDBMS don't handle non-tabular data  
(Notorious for doing a poor job on recursive data structure)

### Legacy Archaic Architecture

- RDBMS don't parallelize well to accommodate commodity HW clusters

### Speed

- Seek time of physical Storage has not kept pace with network speed improvements

### Scale

- Difficult to scale-out RDBMS efficiently – Clustering beyond few servers notoriously hard

### Integration

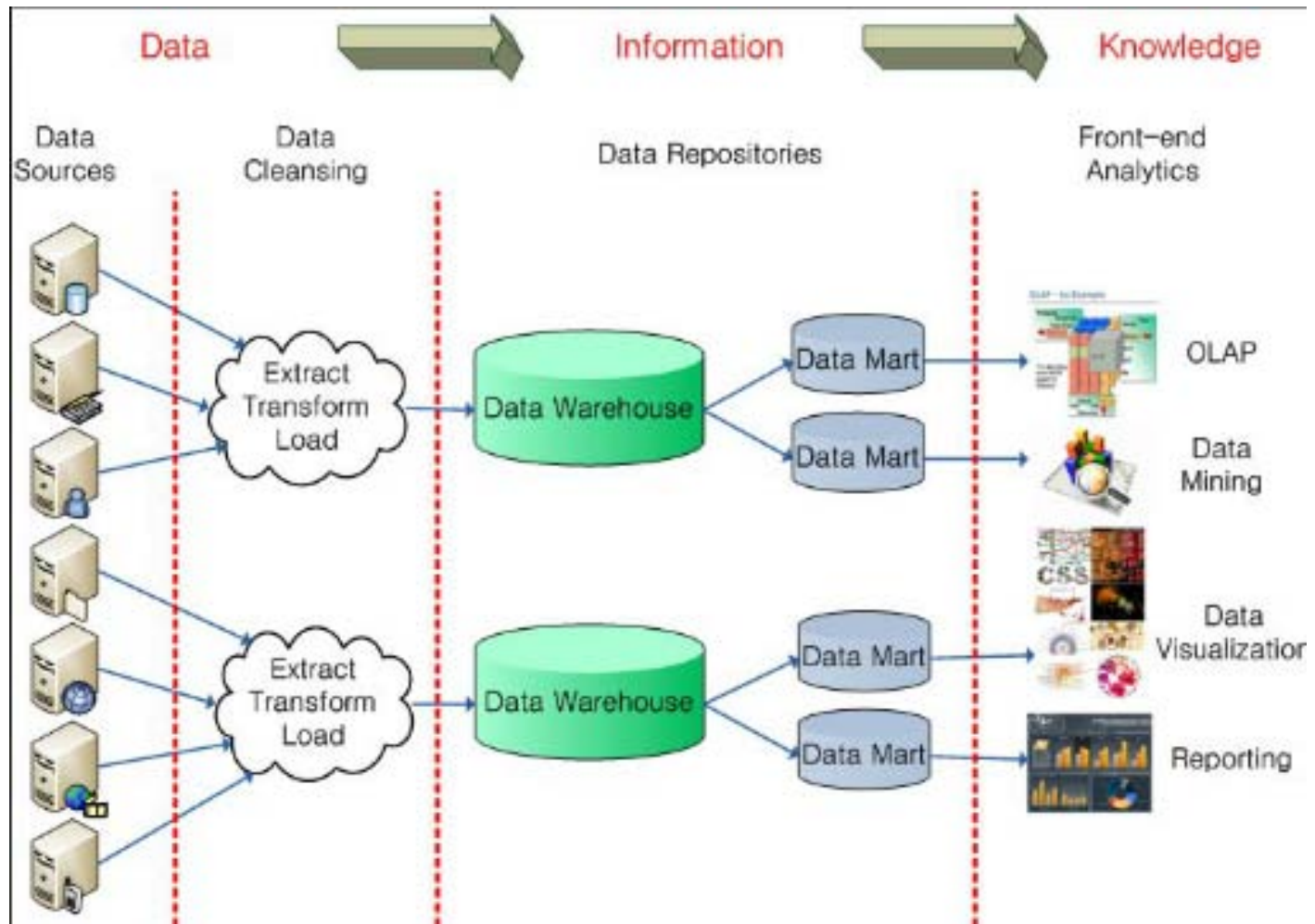
- Data processing tasks need to combine data from non-related sources, over a network

### Volume

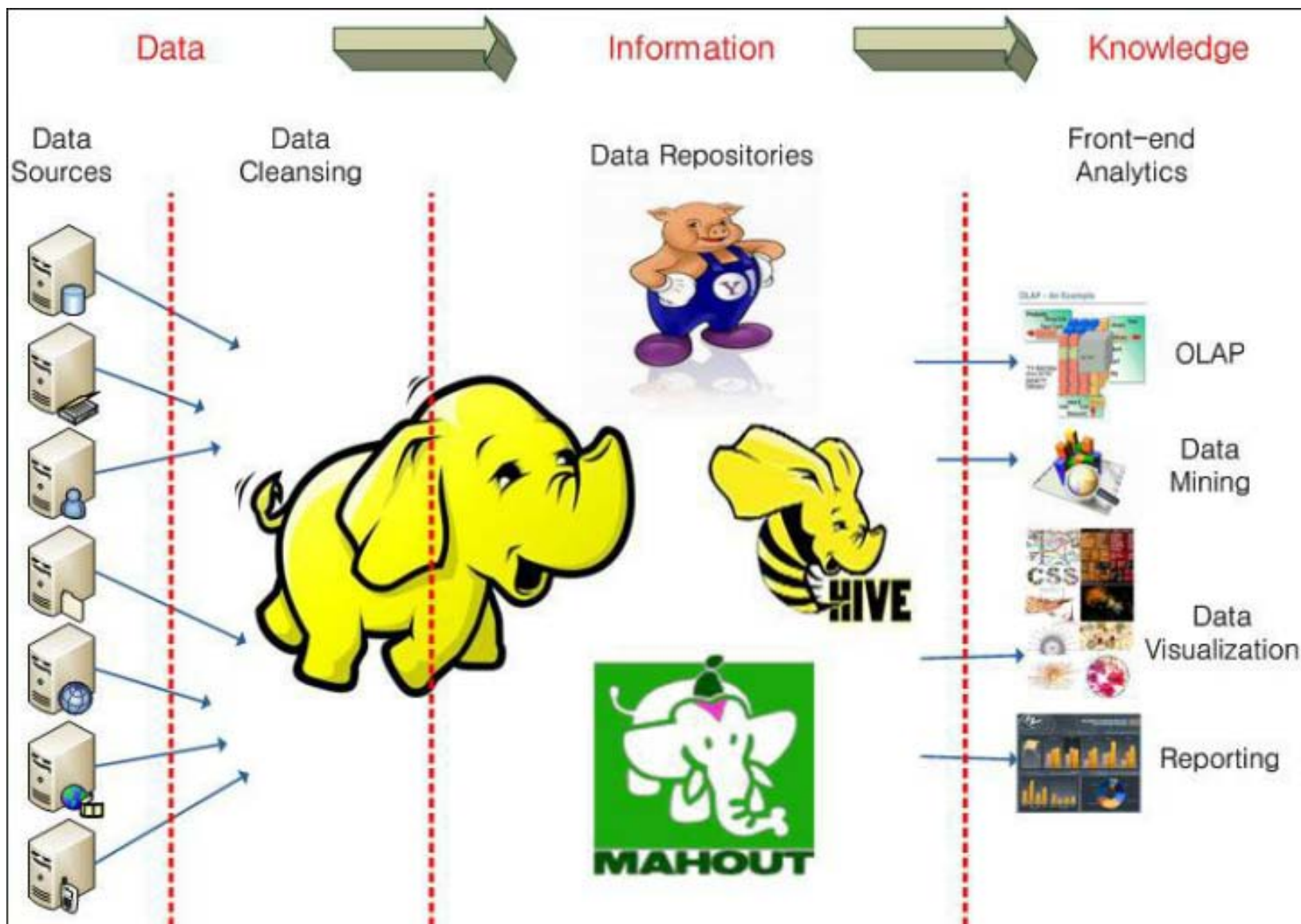
- Data volumes have grown from 10s GB >100s TB > PBs in recent years. Existing Tabular RDBMS can't handle such large DBs

# Issues with Existing RDBMS

## Present RDBMS struggling to Store & Analyze Big Data



# Big Data - Database Solutions



NextGen Infrastructure for Big Data

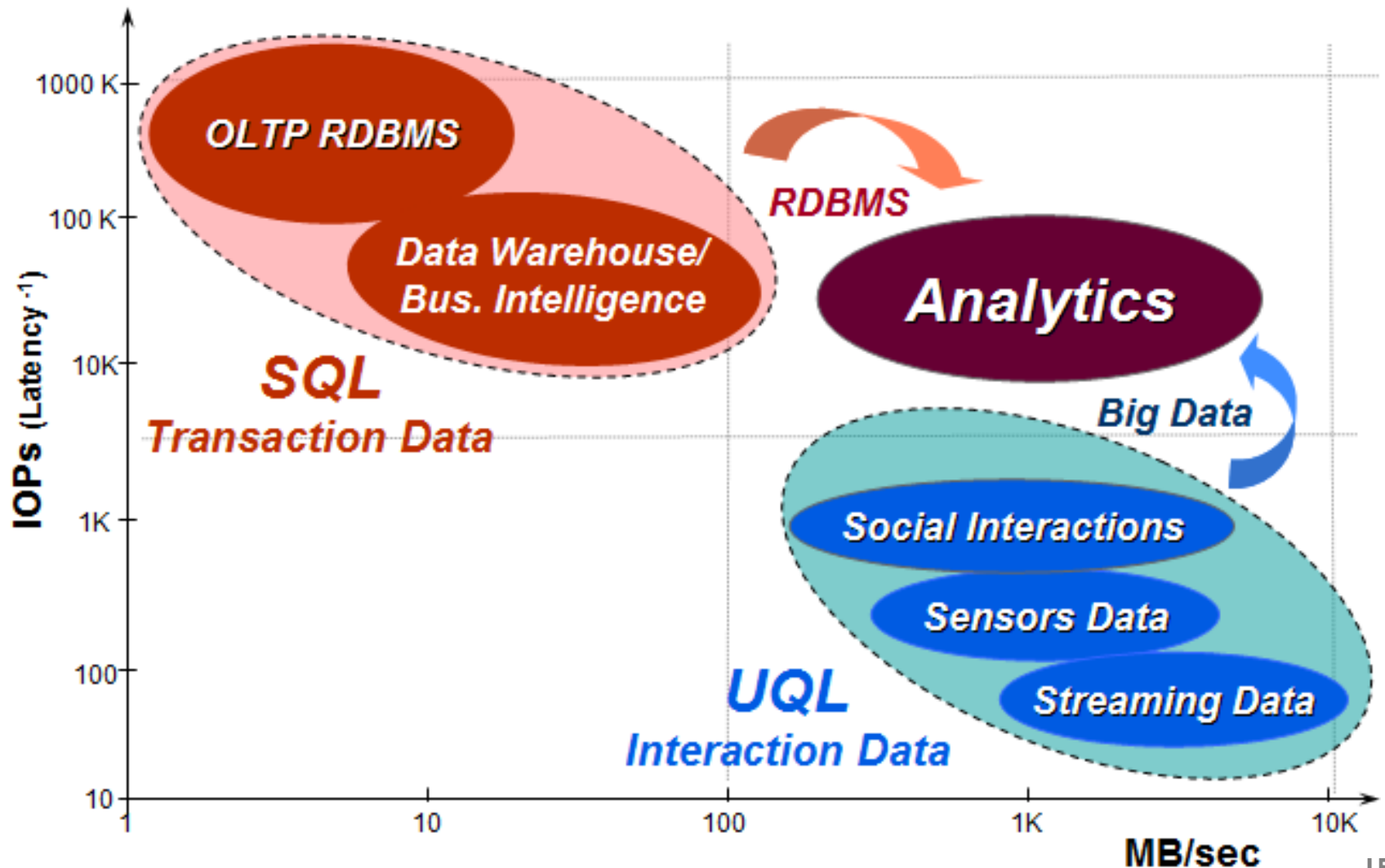
# Big Data – The New Face of DBs

## Big Data Paradigm - The New face of DB Systems

- Adopts Schema-Free Architecture
- Can do away with Legacy Relational DB Systems
  - Some data have sparse attributes, do not need relational property
- Key Oriented Queries
  - Some data stored/retrieved mainly by primary key, w/o complex joins
- Trade-off of Consistency, Availability & Partition Tolerance
- Scale Out, not up, - Online Load balancing cluster growth

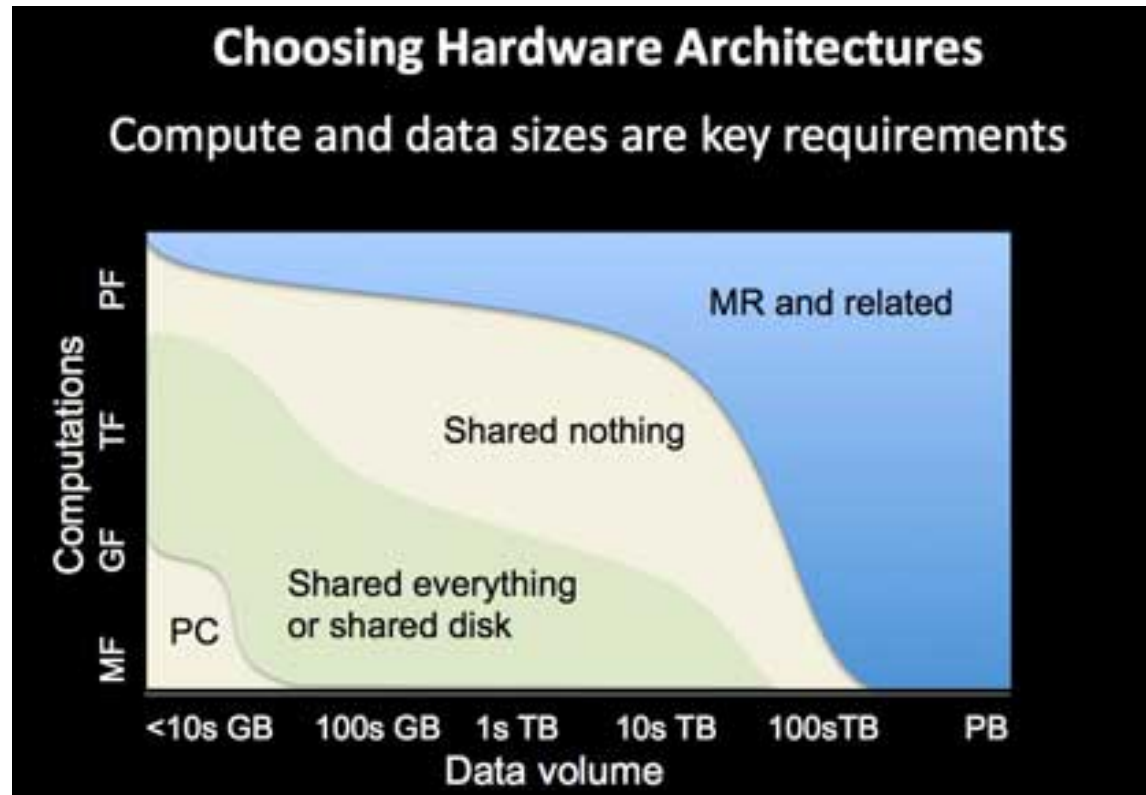


# Analytics – The Next Frontier in IT

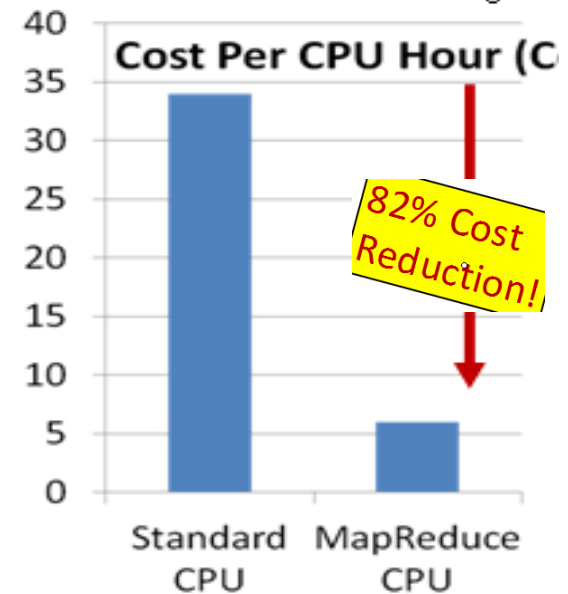
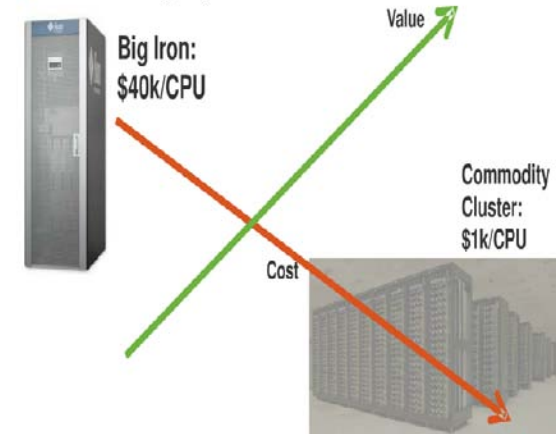




# Key Innovations: HW Technologies

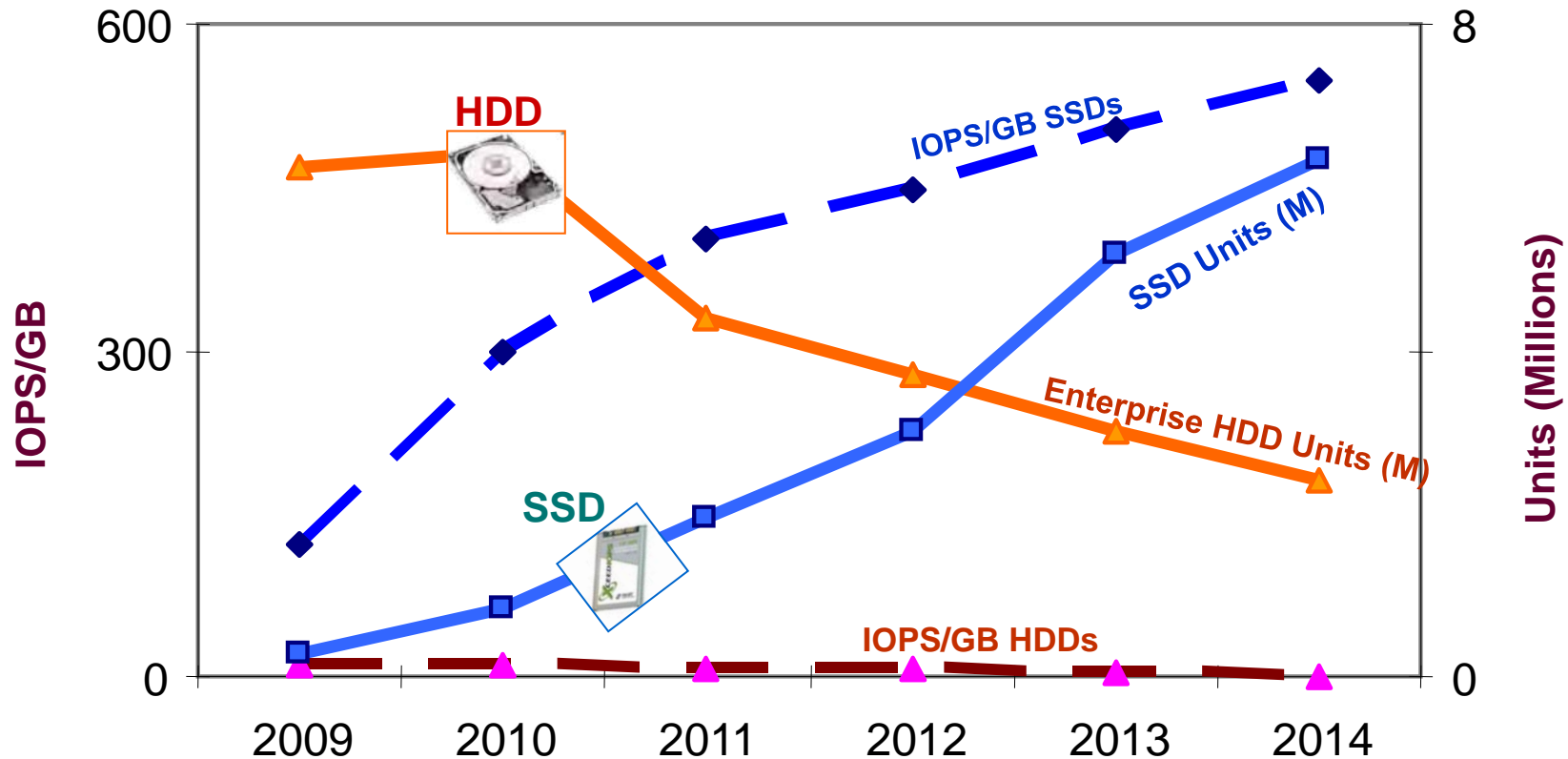


• Hardware cost halving every 18mo



# Key Innovations – Solid State Storage

## Storage - IOPS/GB & Price Erosion - HDD vs. SSDs



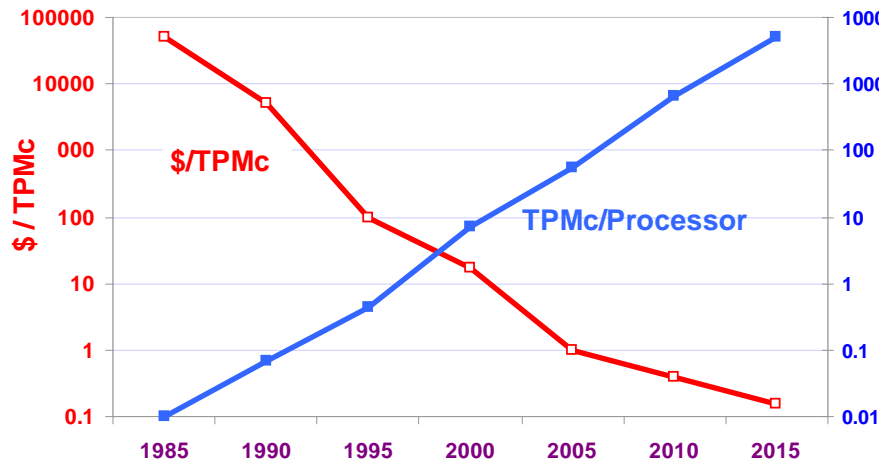
Note: 2U storage rack, • 2.5" HDD max cap = 400GB / 24 HDDs, de-stroked to 20%, • 2.5" SSD max cap = 800GB / 36 SSDs

**Key to Database performance are random IOPS. SSDs outshine HDD in IO price/performance – a major reason, besides better space and power, for their explosive growth.**

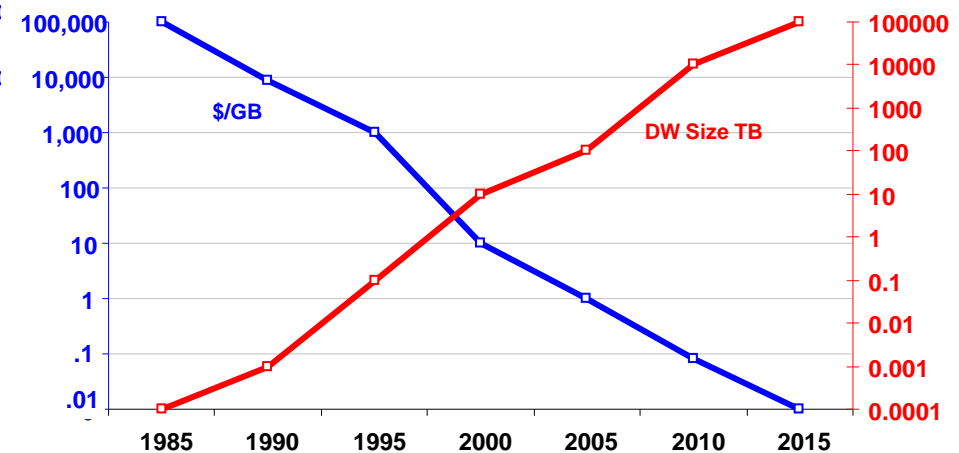
# Innovations – DB SW Technologies

Tech Innovation	1985	1990	1995	2000	2005	2010	2015
OLTP Transactions DB SW	Rows Locking	Optimizer	Parallel Query	Clustering	XML	Grid	Open Source / Hadoop
OLAP- Analytics DB SW	Indexing	Partitioning	Columnar	Materialized View	Bit Mapped Index	In-Memory	Query Binding
Hardware	32 bit	SMP	NUMA	64 bit	Multi-core/Blades	Flash	MPP
Big Data					Multi-core	Columnar In-Memory	MPP Visualization

## OLTP Database Innovation Progress

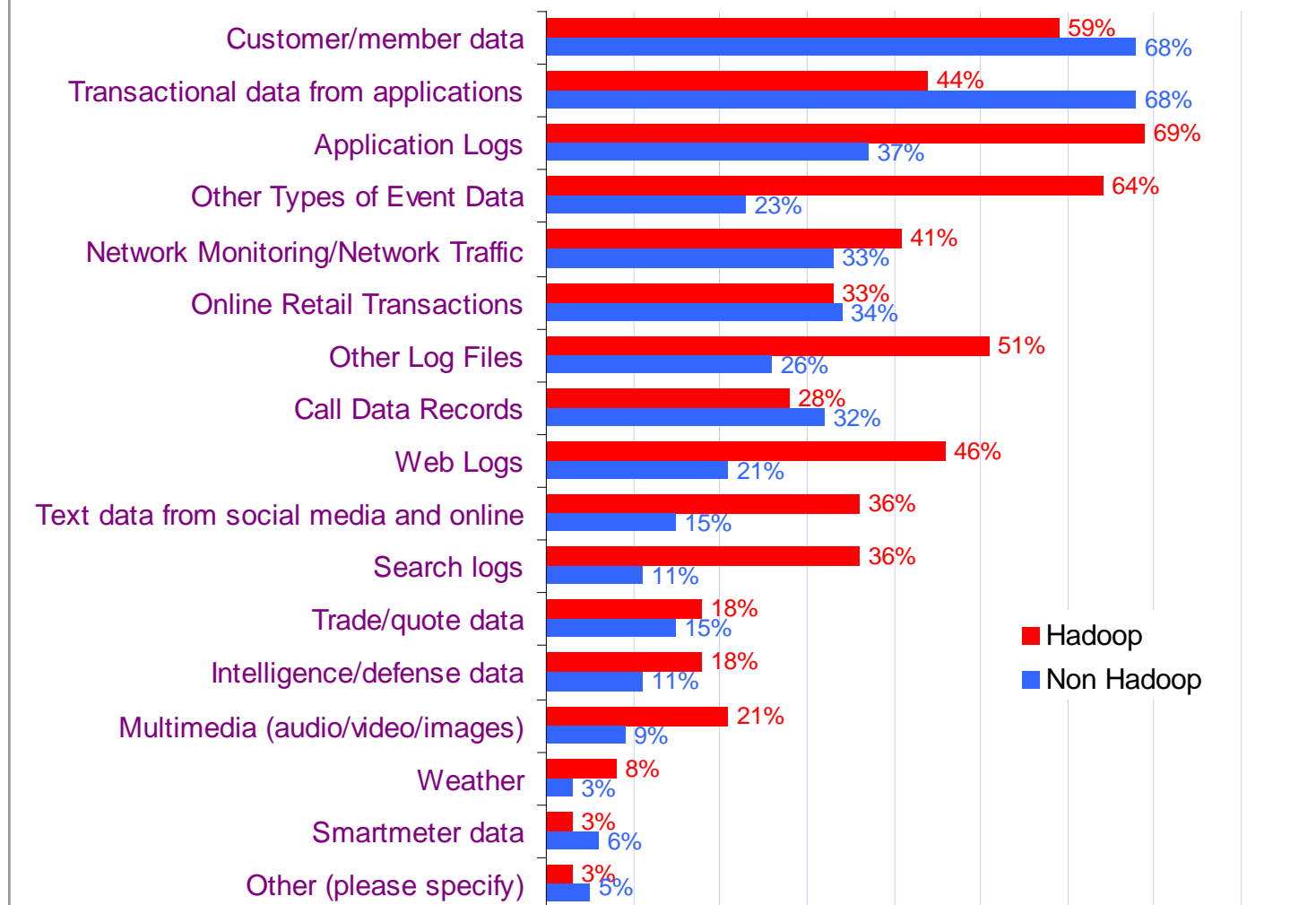


## Big Data: Analytics DB Technology Impact



# Big Data - Key Requirements

## Types of Data Organizations Analyze



# Big Data – Architectural Goals

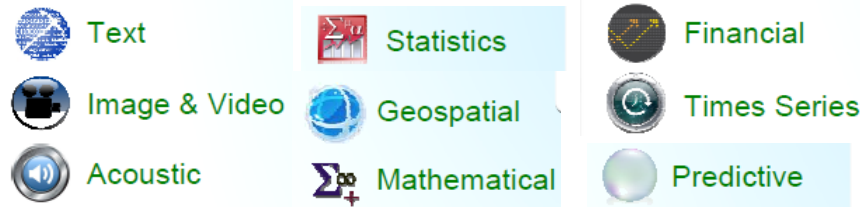
## Meet Requirements of V3



### Big Data Platform



## Meet Enterprise Criterion



## Analyze Data in Native Format

## **Unified system:** Pre-integrated for Ease of Installation and Management

- Platform – **Large Scale Indexing** Pre-integrated using Hadoop Foundation,
- Integrated Text Analytics - **Address Unstructured Data**
- **Usability** - User Friendly Admin Console including HDFS Explorer, Query Languages
- **Enterprise Class Features** – Provisioning, Storage, Scheduler, Advance Security
- Supports **search-centric, document-based XML data model**
  - store documents within a transactional repository.
- **Schema-Free:**
  - No advance knowledge of the document structure (its "schema") needed
  - Index words and values from each of the loaded documents together with its document structure.
- **Standard commodity hardware** leveraged

## Architectural

- **Shared-nothing clustered DB architecture**
  - programmable and extensible application servers.
- Support **massive scalability** to petabytes of source data
- Support **open-source XQuery- and XSLT-driven architecture**
- Simple to Deploy, Develop and Manage (**UI & Restful Interface**)
- Support **extreme mixed workloads** - a wide variety of data types including arbitrarily hierarchical data structures, images, waveforms, data logs etc.
- Support **thousands of geographically dispersed on--line users and programs** executing variety of requests from ad hoc queries to strategic analysis
- **Loading data before declaring or discovering its structure**
- **Load data in batch and streaming fashion**
- **Integrate data from multiple sources** during load process at very high rates
- **Spread I/O and data across instances**
- Provide consistent **performance with linear cost**
- Leverage **Open Source SW** Lo Costs, Multiple Sources, Hadoop Foundation Tools
- Connectivity with Oracle DB, Teradata Warehouse, JDBC Connectivity,



## Real Time Analytics Execution

- Execute “**streaming**” **analytic queries in real time** on incoming load data
- **Updating data in place** at full load speeds
- Scheduling and execution of **complex multi-hundred node workflows**
- Join a **billion row dimension table to a trillion row fact table** without pre-clustering the dimension table with the fact table

## Performance

- Analyze data, at very high rates >GB/sec
- **Predictable Sub-ms response time** for highly constrained standard SQL queries

## Availability

- Ability to **configure without any single point of failure**
- **Auto-Failover Extreme High Availability**
  - Automated failover and process continuation without operational interruption when processing nodes fail

# Big Data – Product Metrics Choices

<b>Big Data - Product Metrics</b>	<b>Data Set Size</b>	<b>PB</b>
		<b>TB</b>
		<b>GB</b>
	<b>Data Structure</b>	Transaction
		Machine
		Unstructured
		Other
	<b>Access/Use</b>	Transaction
		Search
		Analytics
	<b>Parallel Processing</b>	Appliance
		Cluster < 1K
		Cluster > 1K
	<b>Memory</b>	In-Memory
		Flash
	<b>DB Technique</b>	Columnar
		Zero Sharing
		No SQL
	<b>Data Cataloging SW</b>	Text
		Image
		Audio
		Video

# Advantage: Big Data Products

Characteristic	Legacy Paradigm	Big Data Paradigm
<b>Structure</b>	•Transactional/Corporate	•Unstructured/Derivative/Internet
<b>Mode</b>	•Data Collection	•Data Analysis
<b>Focus</b>	•Find Answers	•Find Questions
<b>Facility</b>	•Reportive / What Happened?	•Analytic / Why did it Happen? Predictive / What will Happen Next?
<b>Opportunity</b>	•Very Small Growth	•Massive Growth
<b>Players</b>	•Legacy Players	•Agile Start Ups, well funded
<b>Impact</b>	•Analyze Existing Businesses	•Create New Businesses

# Advantage: Big Data Products

Characteristic	Traditional RDBMS	Big Data/MapReduce
<b>Data Size</b>	•GB	•PB
<b>Access</b>	•Interactive	•Batch/Near Real-Time
<b>Latency</b>	•Low	•High
<b>Data Updates</b>	•Read & Write Many Times	•Write Once Read Many Times
<b>Schema/Structure</b>	•Static Schema	•Dynamic Schema
<b>Language</b>	•SQL	•UQL/Procedural (Java,C++..)
<b>Integrity</b>	•High	•Not 100%
<b>Works Well for</b>	•Process Intensive Jobs	•Data Intensive Jobs
<b>Works Well w Data Size</b>	•Gigabytes	•Petabytes
<b>Data/Processing Interactions</b>	•Low Latency/High BW – precursor to success. Ntwk. BW can be a bottleneck causing nodes to be idle	•Sends Code to Data, instead of Sending Data to other Nodes (Requiring Lower BW in Cluster)
<b>Fault Tolerance</b>	•Coordinating Processes with Node Failures – a challenge	•Fault Tolerant for HW/SW Failures
<b>Access</b>	•Interactive	•Batch/Near Real-time
<b>Scaling</b>	•Non-linear	•Linear
<b>Pgm-Distribution of Jobs</b>	•Difficult	•Simple & Effective

# Big Data Ecosystem

## Generation

### Data Class Types

#### Data Types

- Structured (Relational)
- Unstructured (Adhoc)

#### Data Class

- Human
- Machine

#### Data Velocity

- Batch
- Streaming

## Operational IT

### Store Access Prepare

#### Data Mgmt & Storage

- Store
- Secure
- Access
- Network

#### Engines

- Hadoop/MapReduce
- Apache Tools
- Cloudera/IBM/EMC ...
- Visualization ...

#### Prepare Data For

#### Analytics

- ETL / Data Integration
- Workflow Scheduler
- System Tools

## Analytics

### Analyze Visualize

#### Data Analytics

- Algorithmics
- Automation
- In Real Time

#### Business Analytics

- Visualization
- Interoperate with SQL- RDBMS
- BI/EDW

## Usage

### Analyze Business

#### Business Analysis

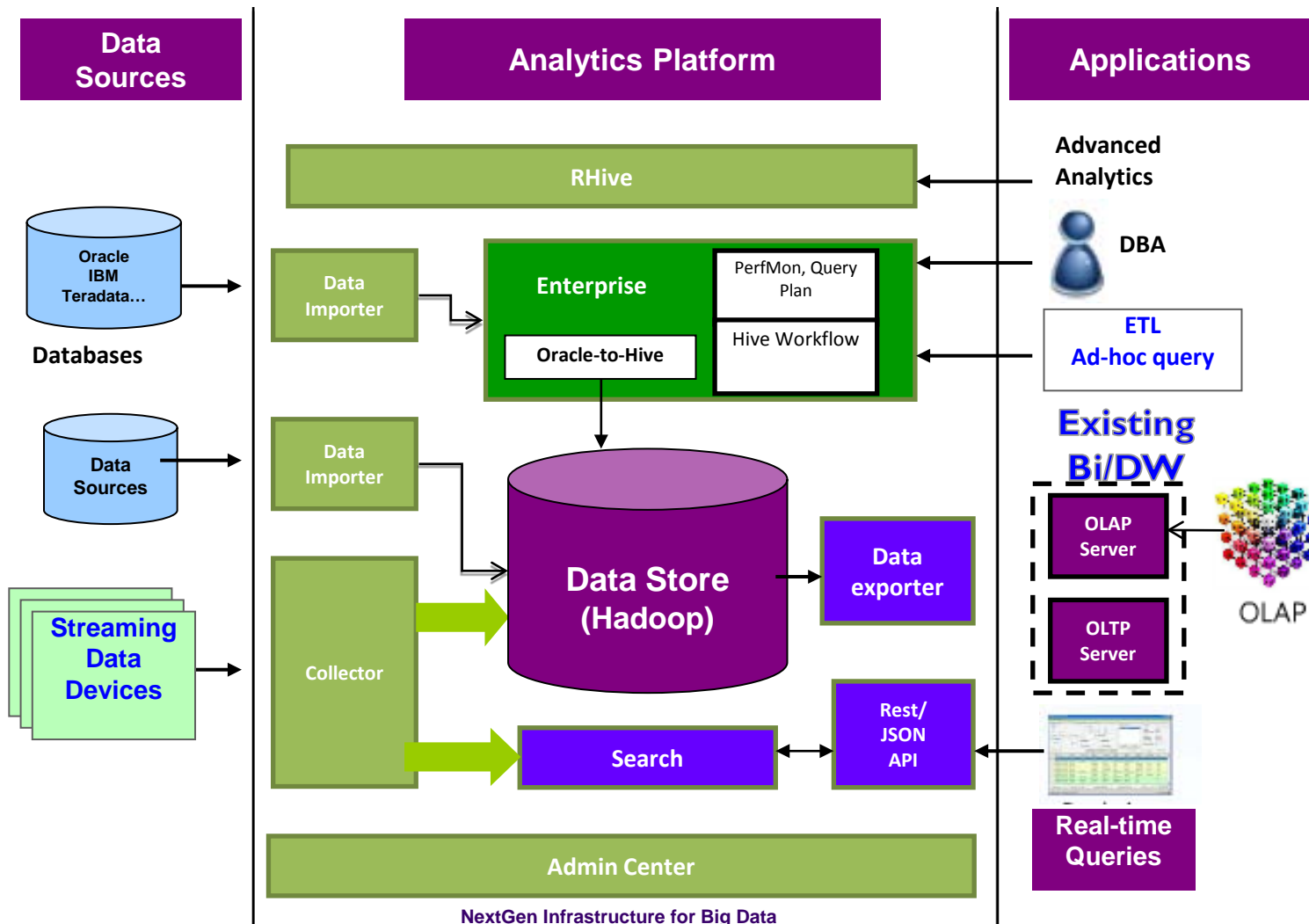
- Decision Support
- Just InTime
- Business Model

#### Business Use

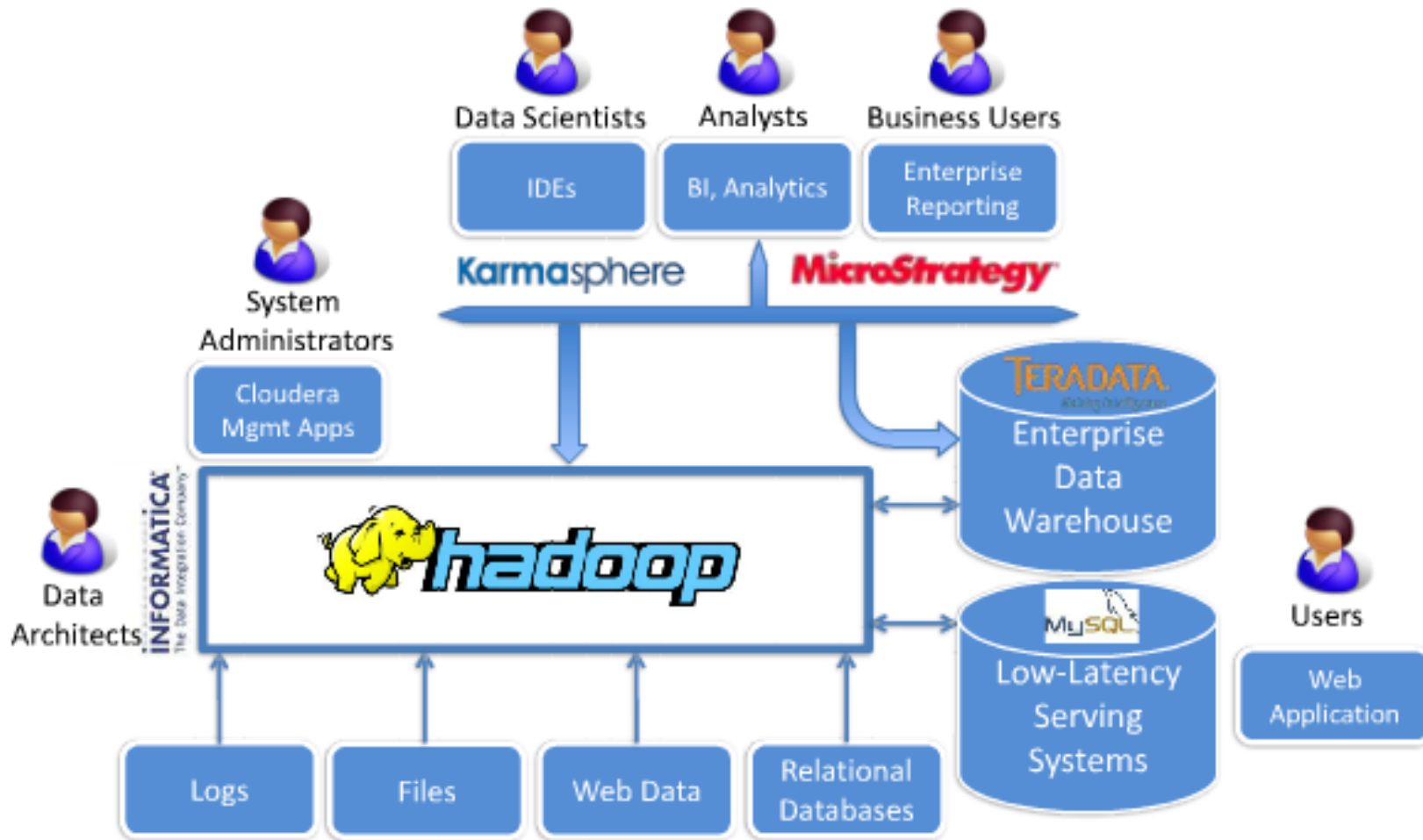
- Market Penetration
- Enhancements
- Cash Flow/ROI

# Big Data Stack

## Merging Hadoop innovations into Nextgen DBMS

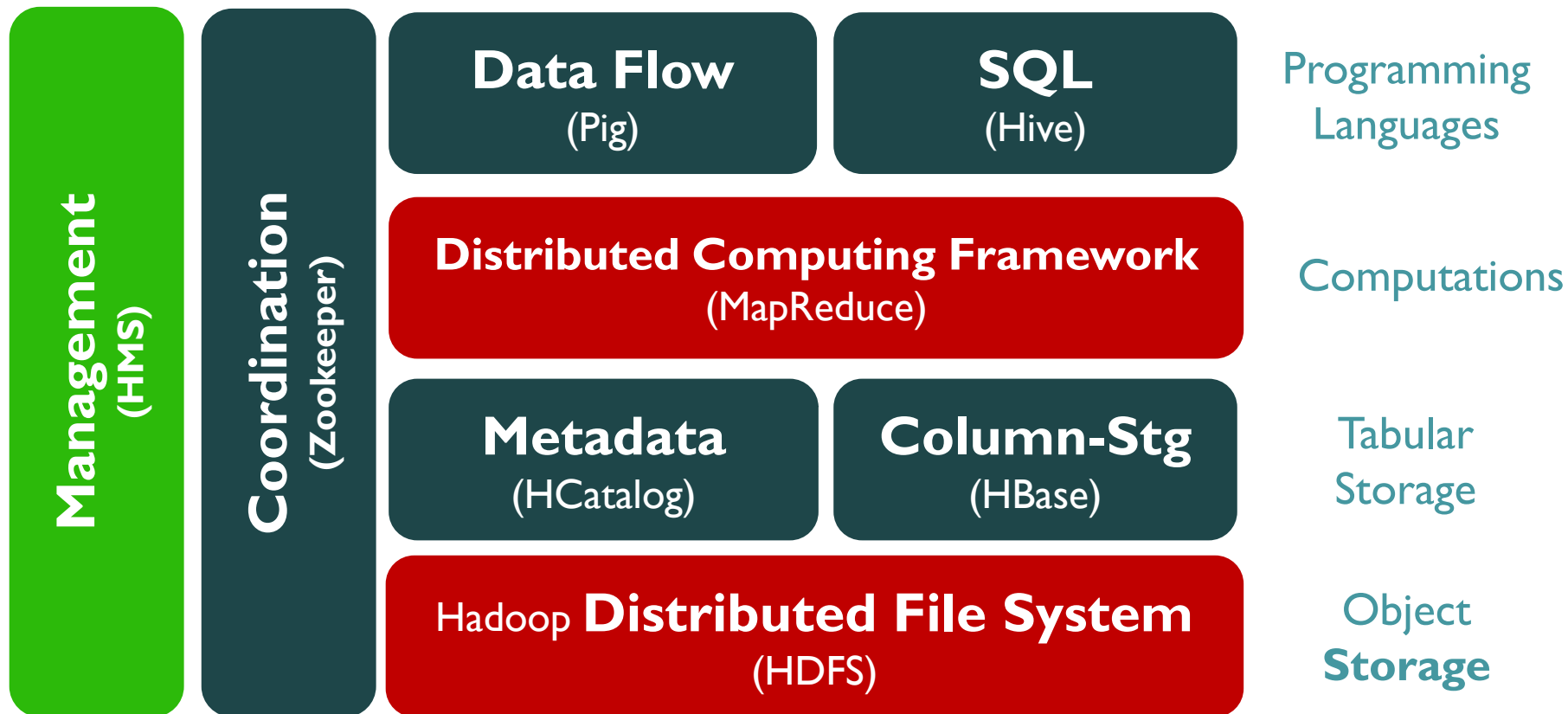


# Hadoop's Fit in Enterprise Stack



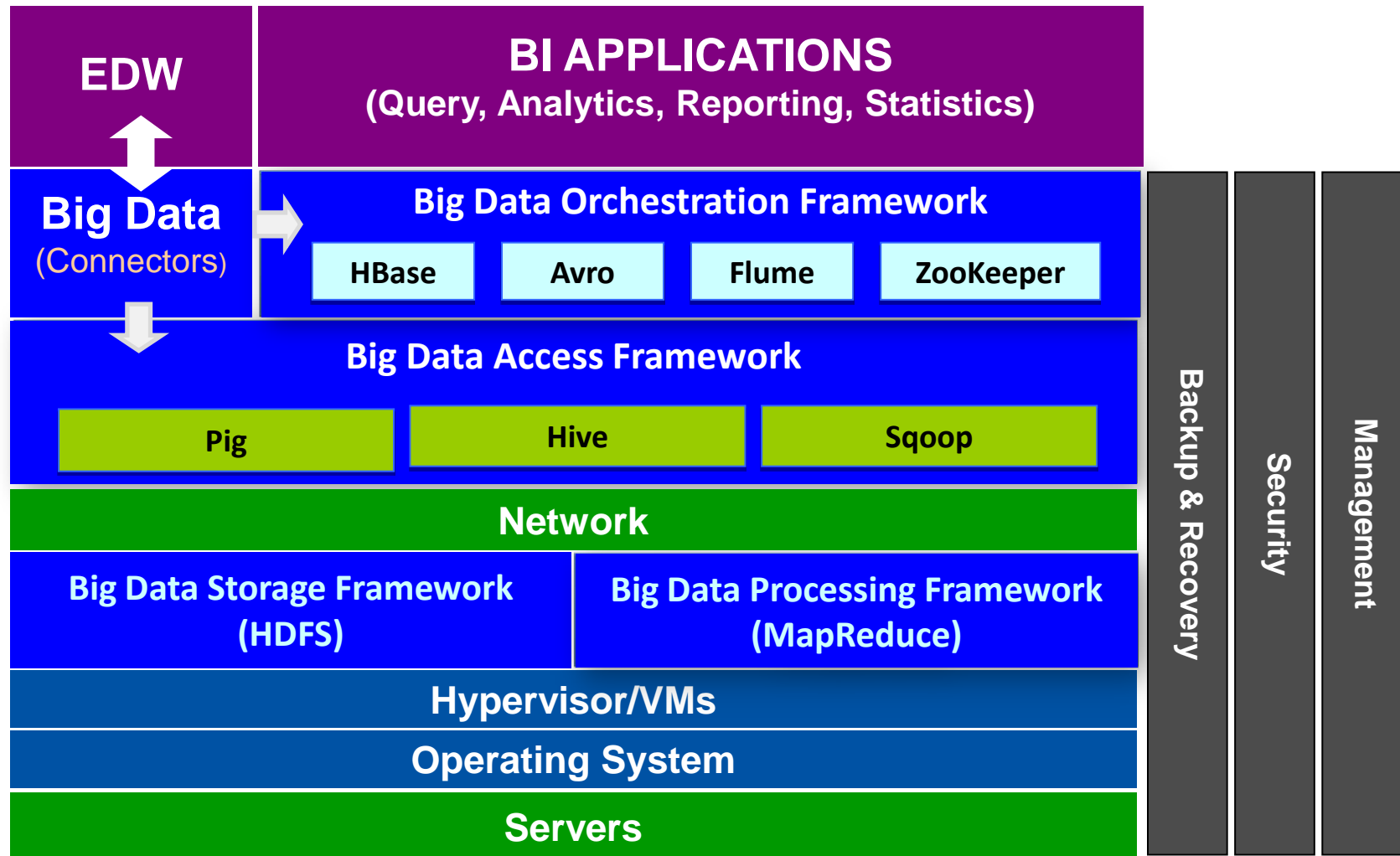


# Big Data - Hadoop Architecture

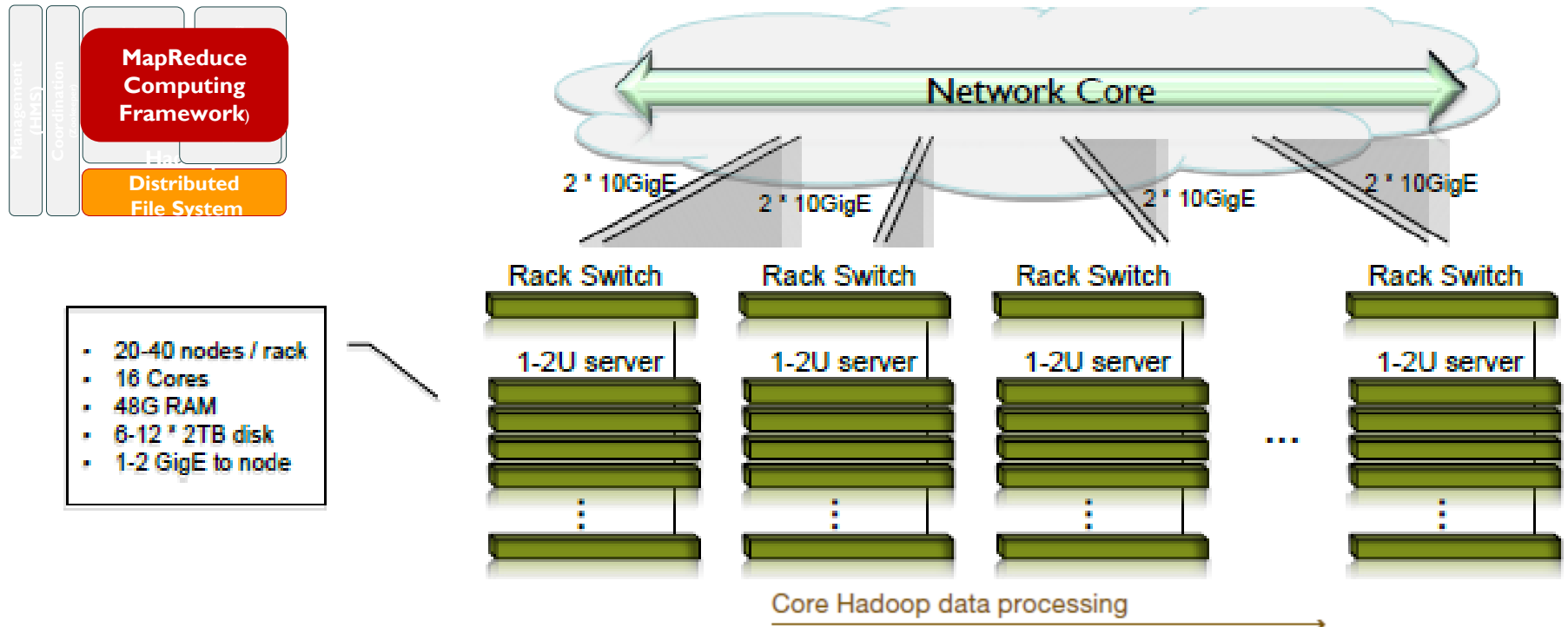


# Big Data Connectors to EDW/BI

## BI Framework - Interoperable with Enterprise Data Warehousing

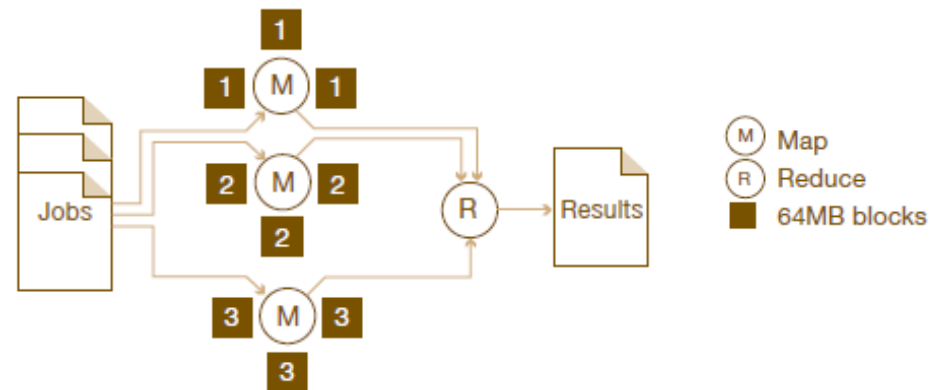


# Big Data Infrastructure – Map Reduce



## Map Reduce

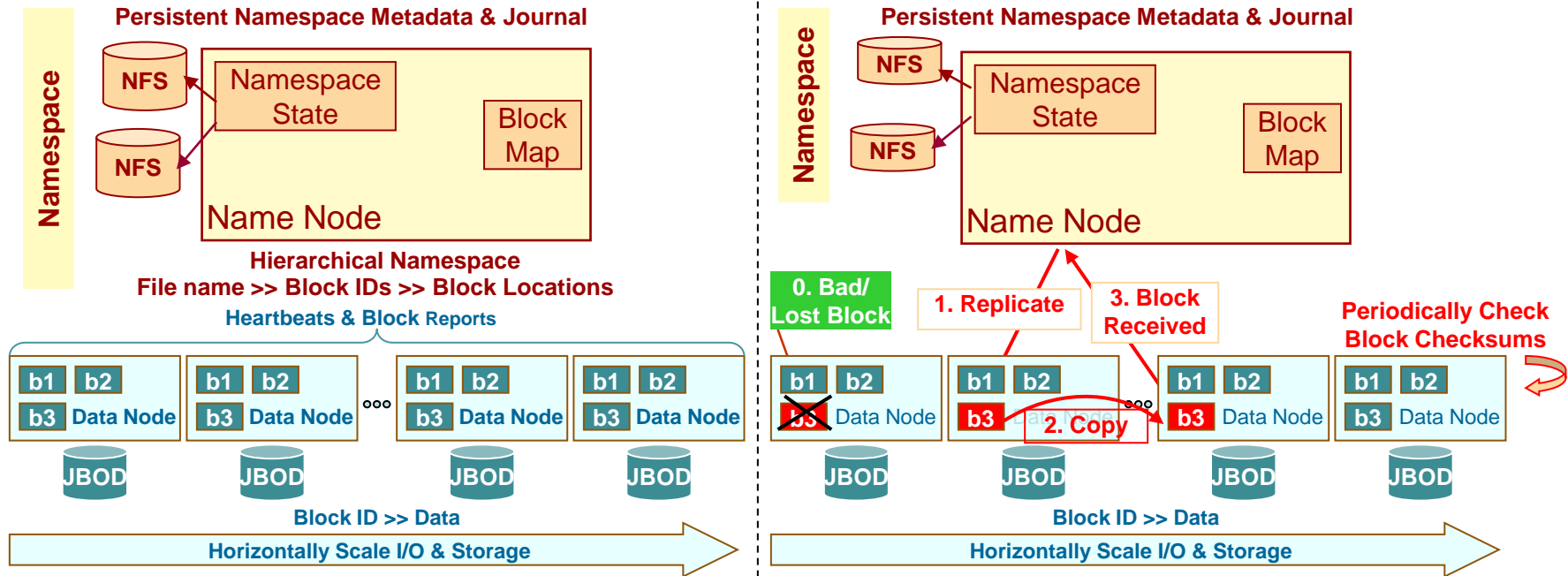
- A Distributed Computing Model
- Typical Pipeline:  
Input>**Map**>Shuffle/Sort>**Reduce**>Output
- Easy to Use , Developer writes few functions,  
Moves compute to Data
- Schedules work on HDFS node with data
- Scans through data, reducing seeks
- Automatic Reliability and re-execution on failure



# Big Data Infrastructure – HDFS

## HDFS Architecture

Actively Maintaining High Availability

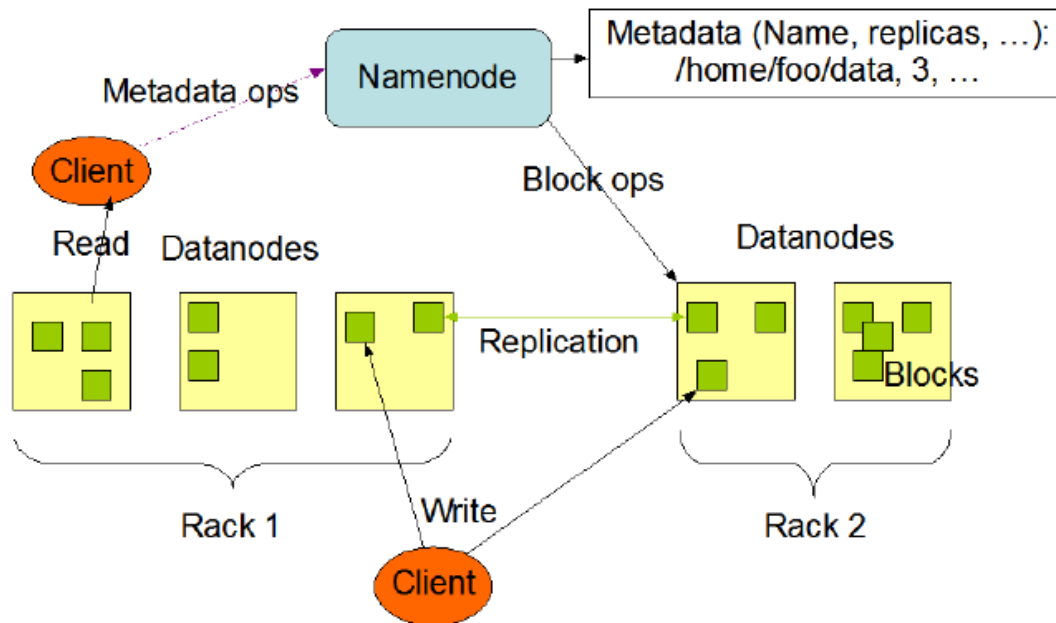


## HDFS

- Immutable File System – Read, Write, Sync/Flush – No random writes
- Storage Server used for Computation – Move Computation to Data
- Fault Tolerant & Easy Management – Built In Redundancy, Tolerates Disk & Node Failure, Auto-Managing addition/removal of nodes, One operator/8K nodes
- Not a SAN but high bandwidth network access to data via Ethernet
- Used typically to Solve problems not feasible with traditional systems: Large Storage Capacity >100PB raw, Large IO/computational BW >4K node/cluster, scale by adding commodity HW, Cost ~\$1.5/GB incl. MR cluster

# Hadoop Distributed File System

## HDFS Architecture

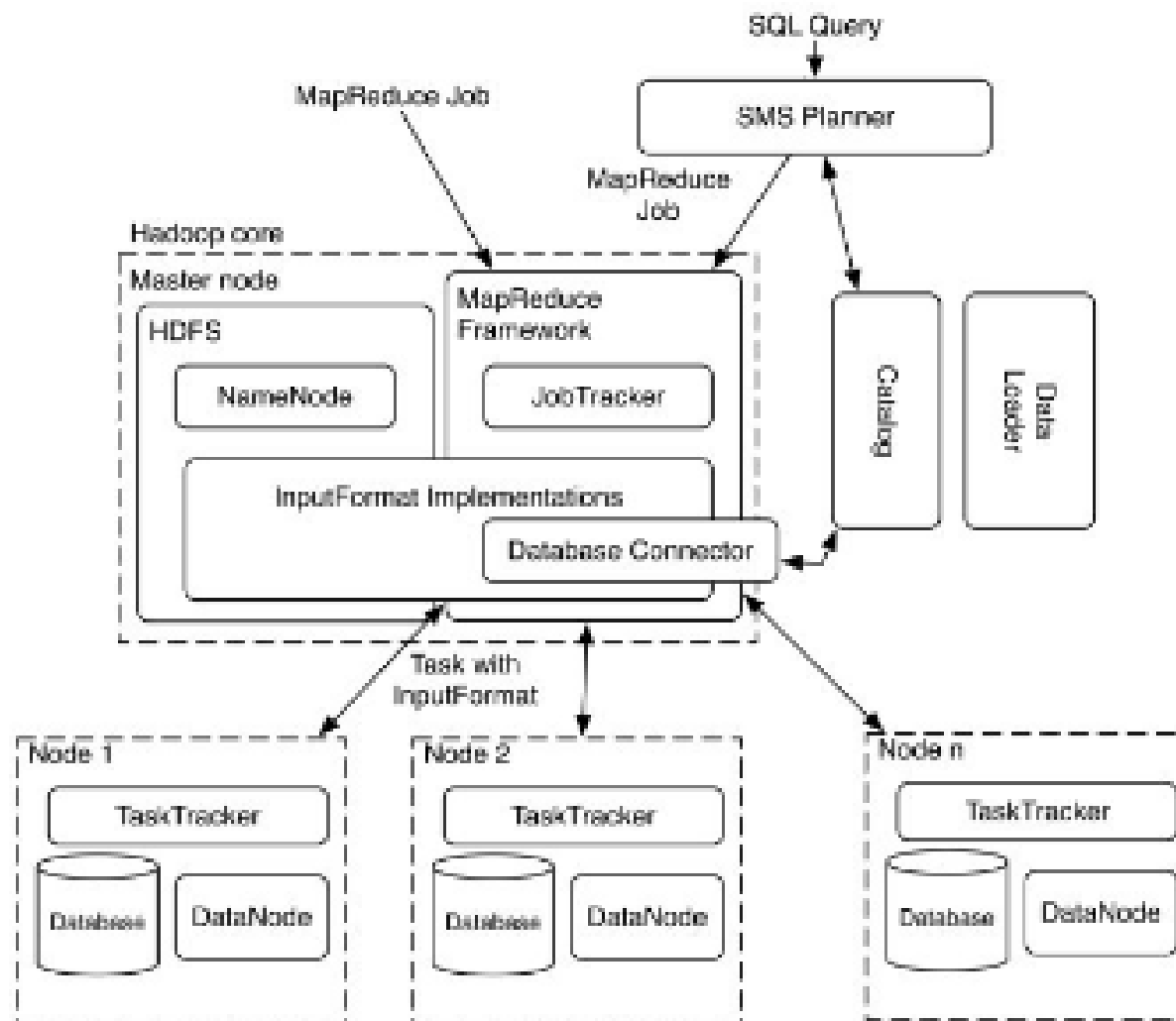


## HDFS Characteristics

- **Based on Google GFS** (Google File System)
- **Redundant Storage** for massive amounts of data
- **Data is distributed across all nodes** at load time – efficient MapReduce processing
- **Runs on commodity hardware** – assumes high failure rate for components
- **Works well with lots of large files**
- **Built around Write once – Read many times**
- **Large Streaming Reads** – Not random access
- **High Throughput** more important than low latency

# Hadoop Architecture - Overview

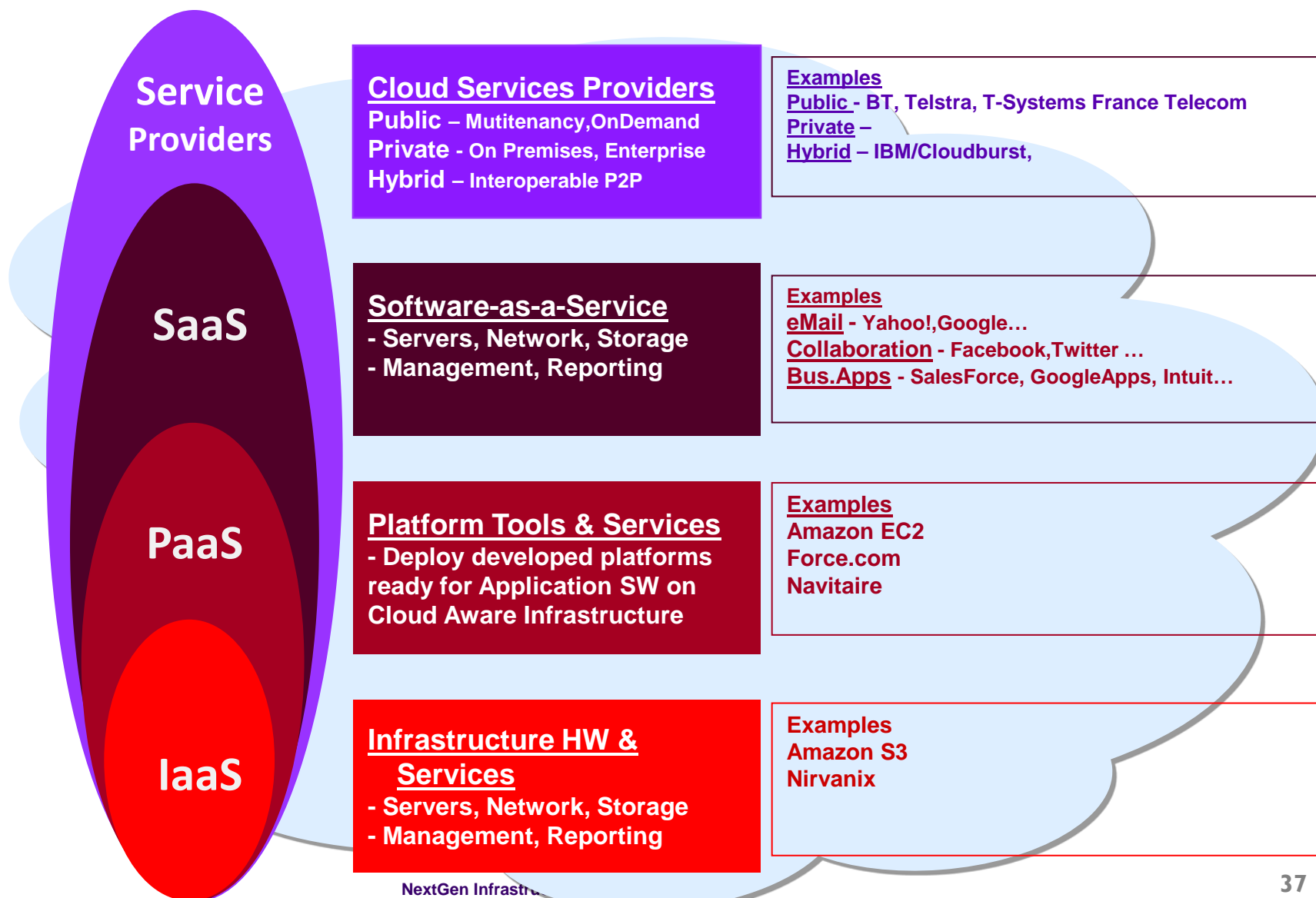
## Hadoop Data Processing Architecture



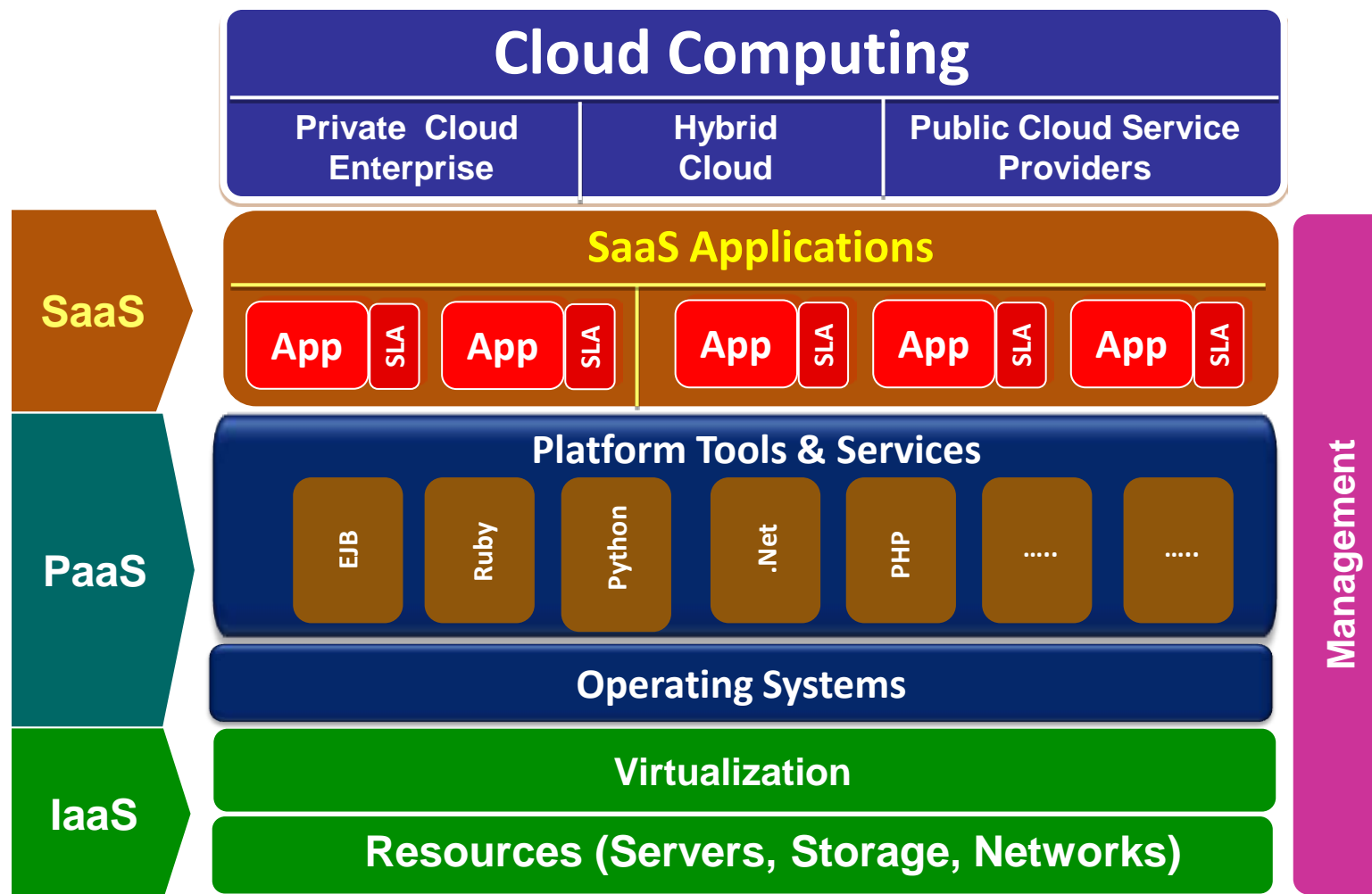
## Key Technologies Required for Big Data

- **Cloud Infrastructure**
- **Virtualization**
- **Networking**
- **Storage**
  - In-Memory Data Base (Solid State Memory)
  - Tiered Storage Software (Performance Enhancement)
  - Deduplication (Cost Reduction)
  - Data Protection (Back Up, Archive & Recovery)



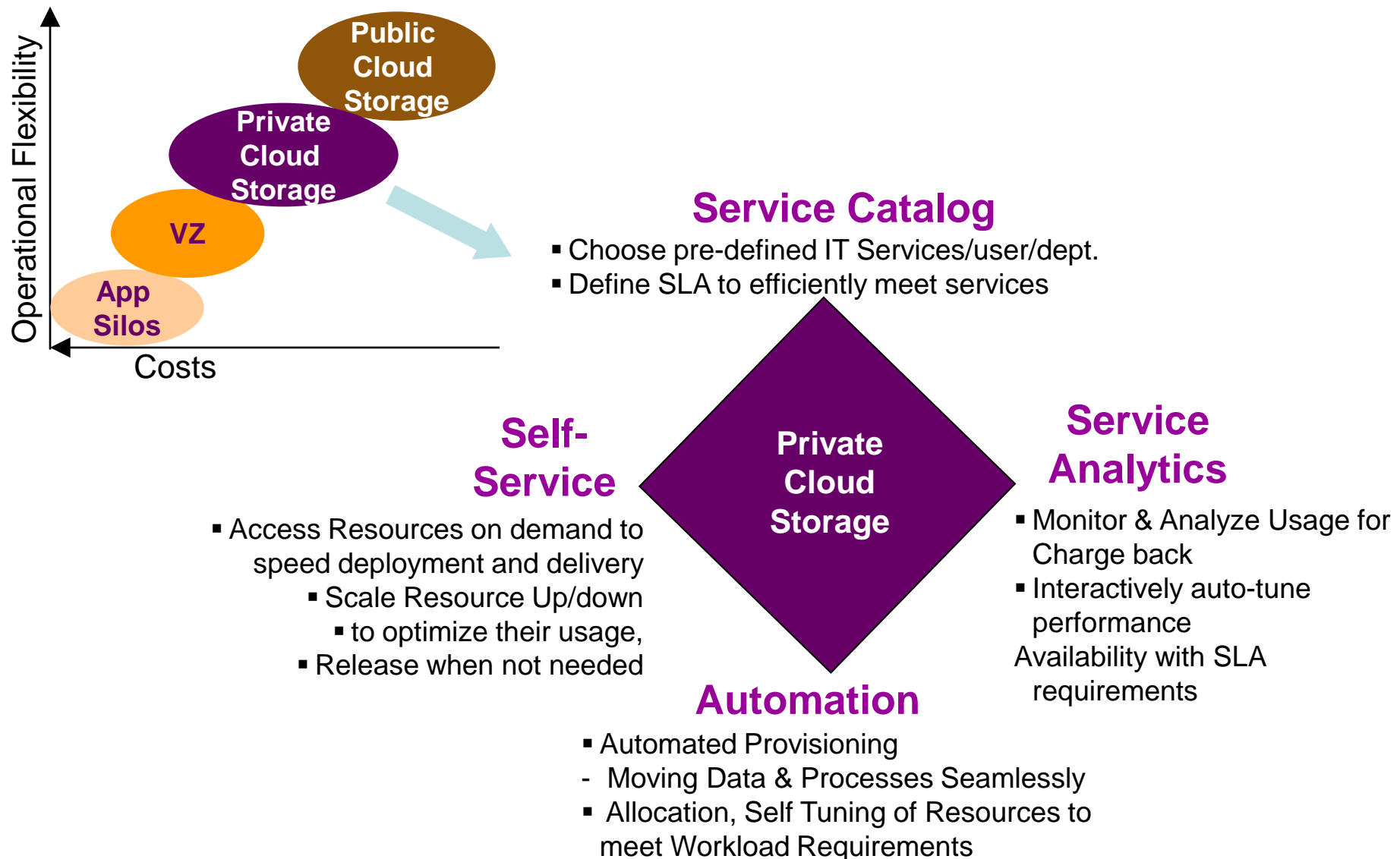


# Cloud Infrastructure for Big Data

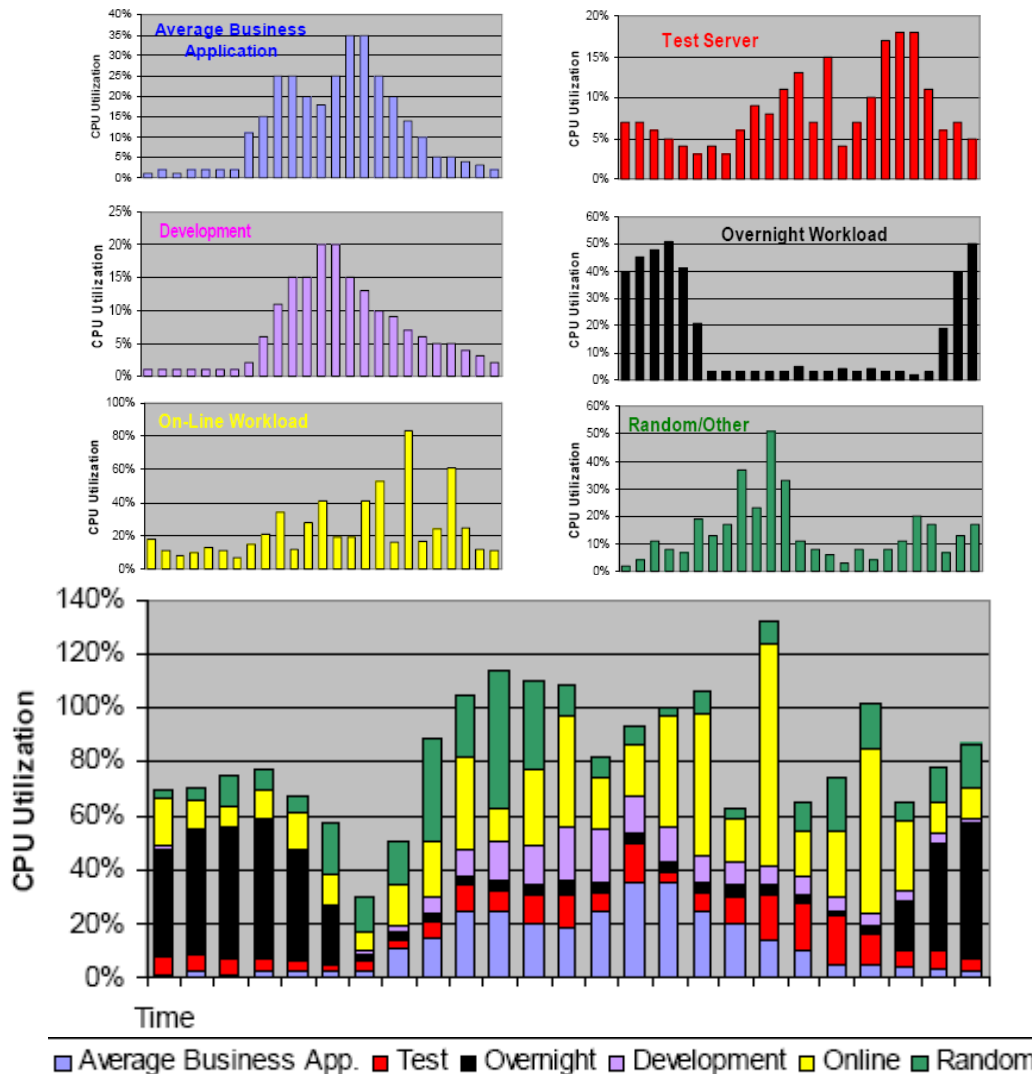


**Application's SLA dictates the Resources Required to meet specific requirements of Availability, Performance, Cost, Security, Manageability etc.**

# Private Cloud Requirements for Big Data



# Virtualization: Workloads Consolidation

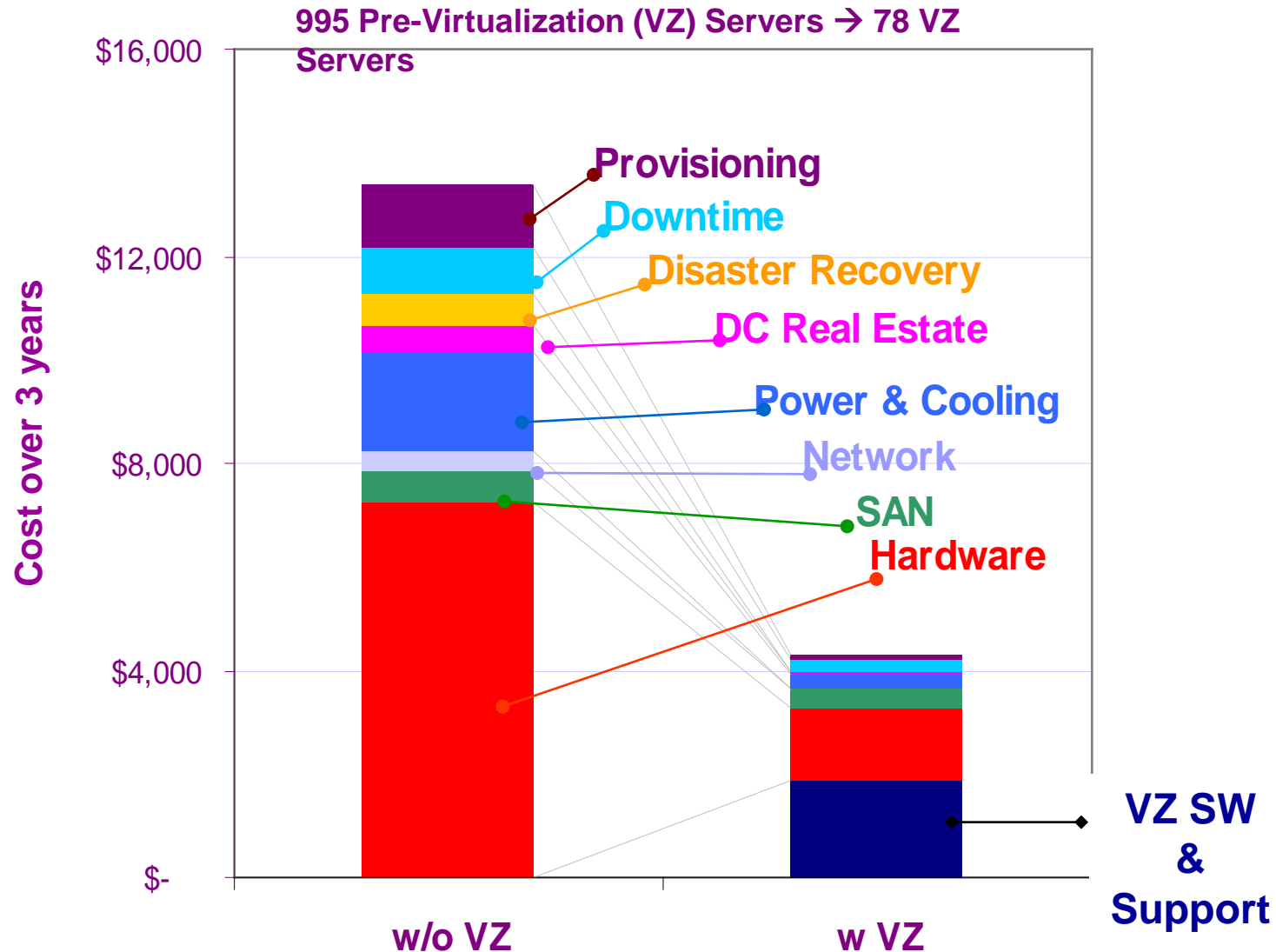


Source: Dan Olds & IMEX Research 2009

n Infrastructure for Big Data

- A single server 1.5x larger than standard 2-way server will handle consolidated load of 6 servers.
- VZ manages the workloads + important apps get the compute resources they need automatically w/o operator intervention.
- Physical consolidation of 15-20:1 is easily possible
- Reasonable goal for VZ x86 servers – 40-50% utilization on large systems (>4way), rising as dual/quad core processors becomes available
- Savings result in Real Estate, Power & Cooling, High Availability, Hardware, Management

# Virtualization: TCO Savings



<b>Storage Efficiency</b>	<b>Virtualization</b> <ul style="list-style-type: none"><li>▪ Mapping P &gt; V, VM Management</li></ul> <b>Performance</b> <ul style="list-style-type: none"><li>▪ In-Memory DB, Auto-Tiering-SSD/HDD</li></ul> <b>Costs Reduction</b> <ul style="list-style-type: none"><li>▪ Thin Provisioning</li><li>▪ Deduplication</li></ul> <b>Availability</b> <p>RAID/Auto recover HA, Snapshots, CDP, Cloning, DRS</p> <b>Security</b> <p>Encryption/DLP</p>
<b>Service Efficiency</b>	<b>Storage -as-a Service</b> <ul style="list-style-type: none"><li>▪ Service Catalogs by Workloads etc.</li><li>▪ Policy Infrastructure<ul style="list-style-type: none"><li>• Service Level Attributes</li><li>• Service Measurements</li></ul></li><li>▪ Performance Analytics<ul style="list-style-type: none"><li>• IOPS/Response Time, Bandwidth</li></ul></li><li>▪ Automation<ul style="list-style-type: none"><li>• Unified SAN/NAS Protocols</li><li>• Auto learning Workload Forensics</li><li>• Provisioning to Match Workloads</li><li>• Assured Auto recovery</li></ul></li></ul>

## Data Protection

Back Up/Archive/DR

RAID – 0,1,5,6,10

Virtual Tape

Replication

## Storage Efficiency

**Virtualization**

Thin Provisioning

Deduplication

Auto Tiering

MAID



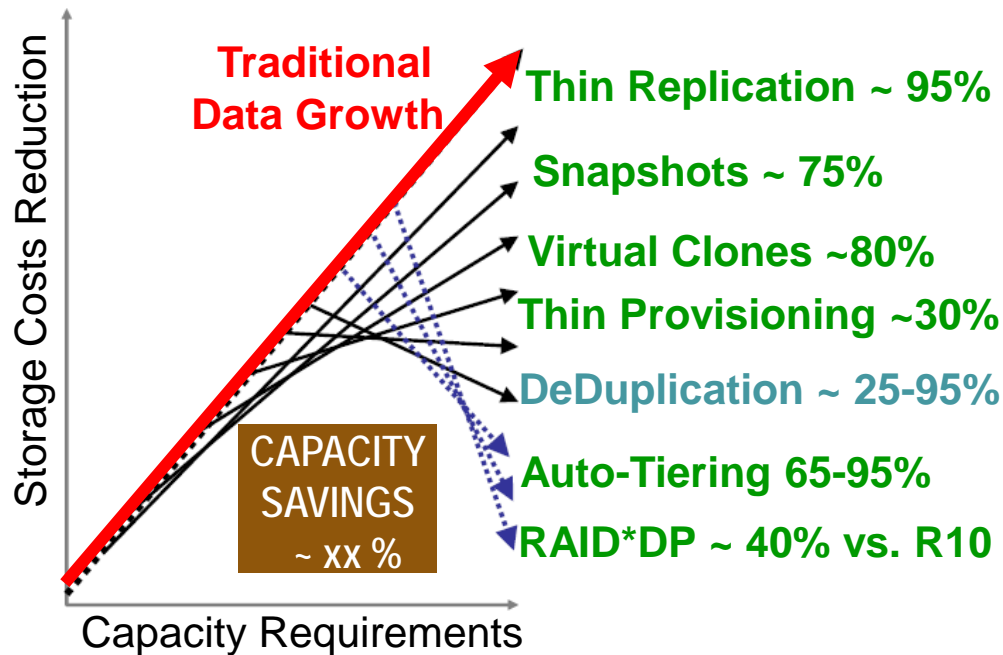
**Virtualization (VZ)  
requires Shared Storage for**

- VMotion
- Storage VMotion
- HA/DRS
- Fault Tolerance

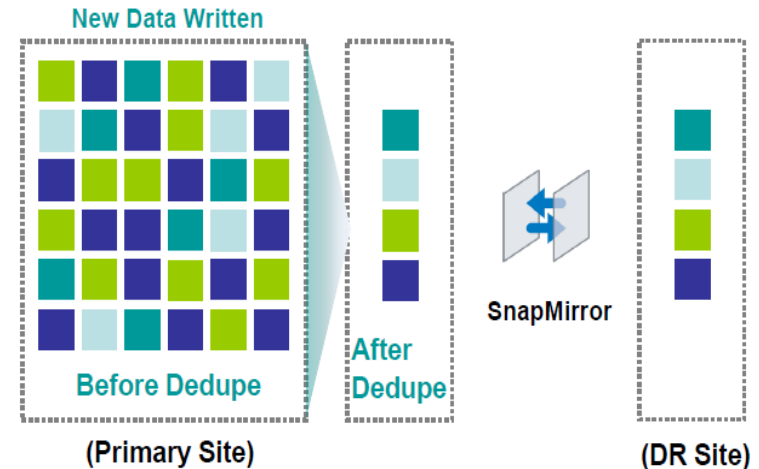
**Additional Capacity  
Consumed for**

- VZ snapshots,
- VM Kernel etc.

## Technologies Reducing Storage Costs

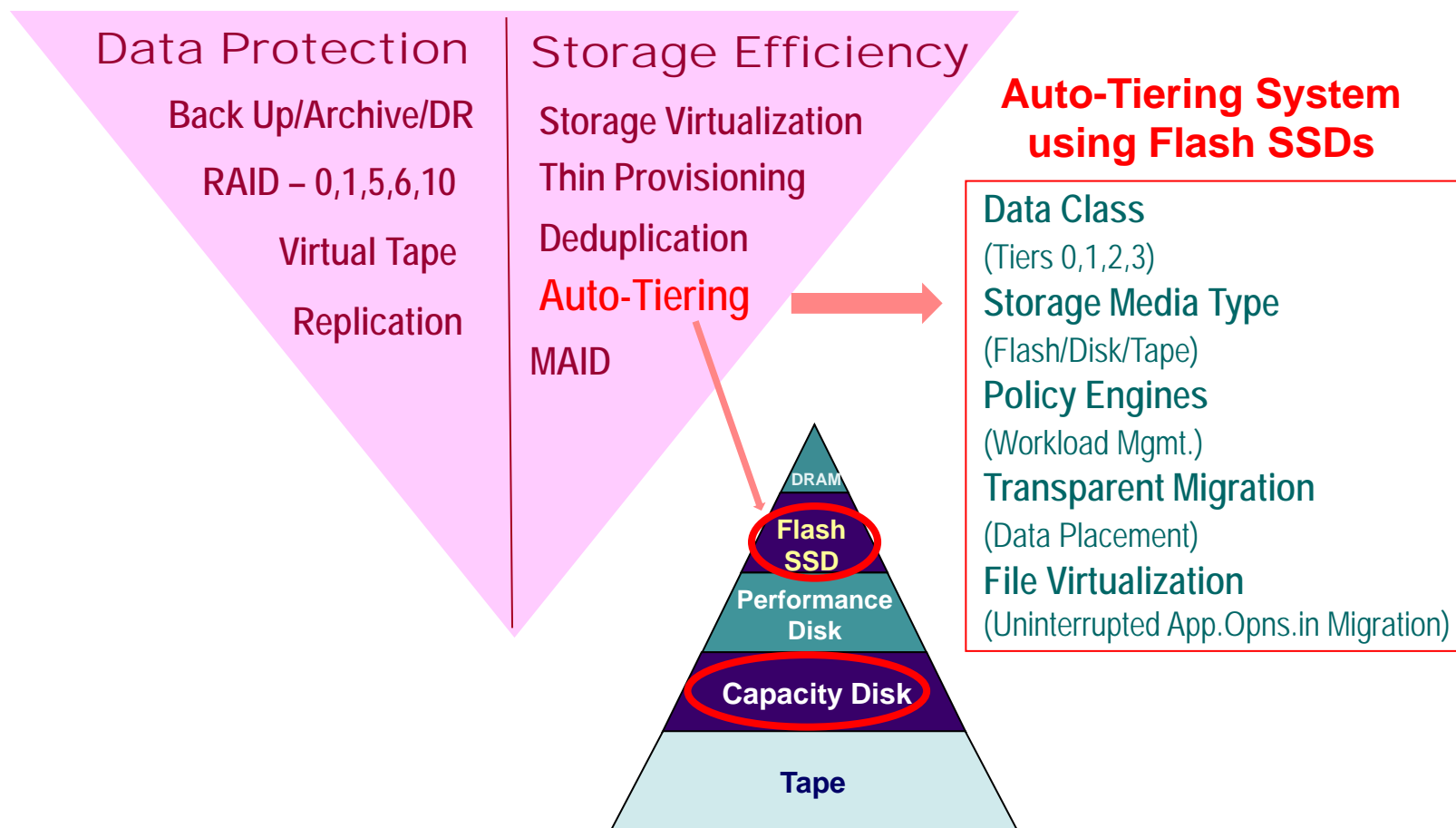


- Replicates only the deduplicated blocks
- Only unique data is replicated to the DR site



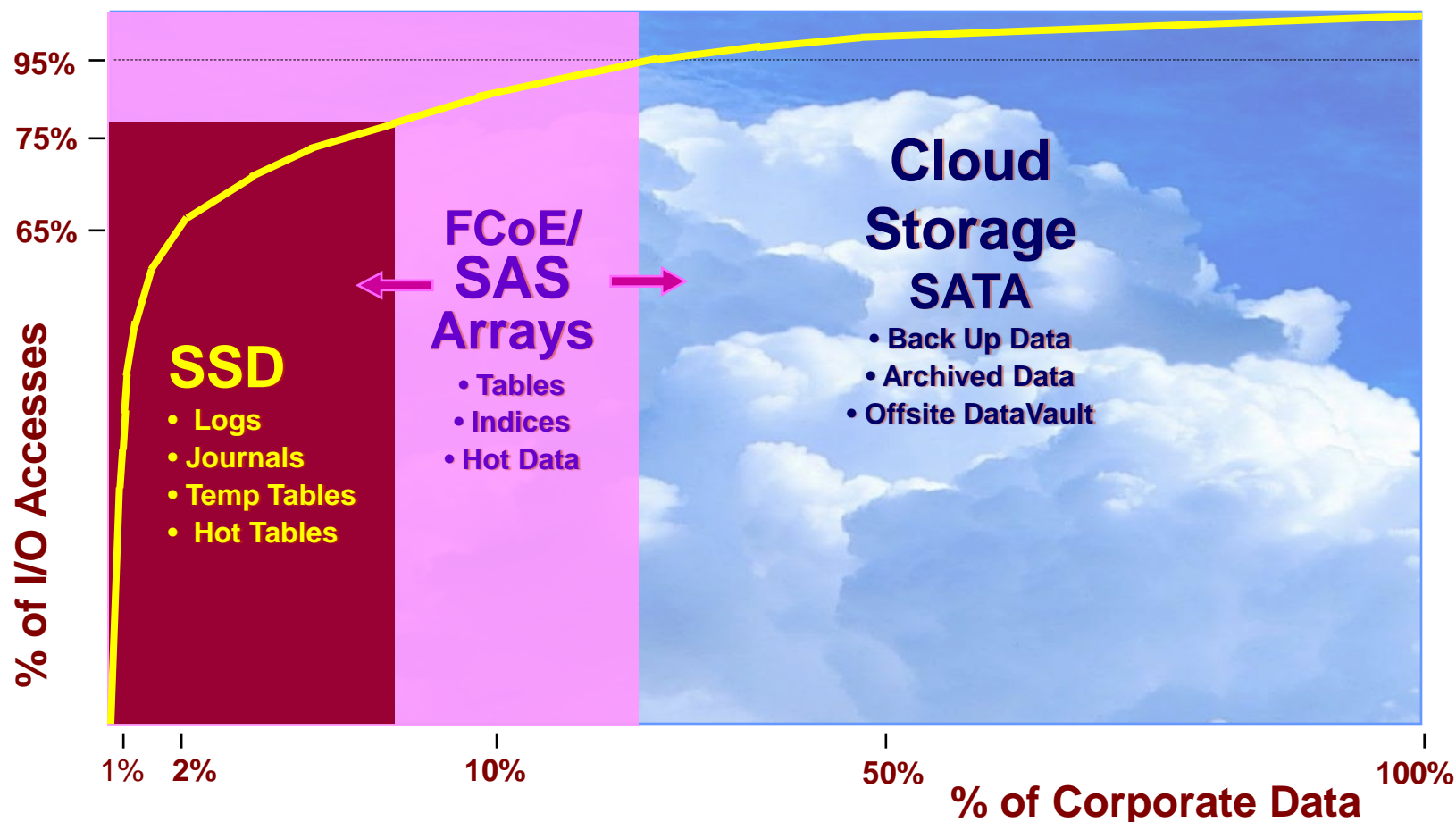


# Storage Architecture Impacting Big Data

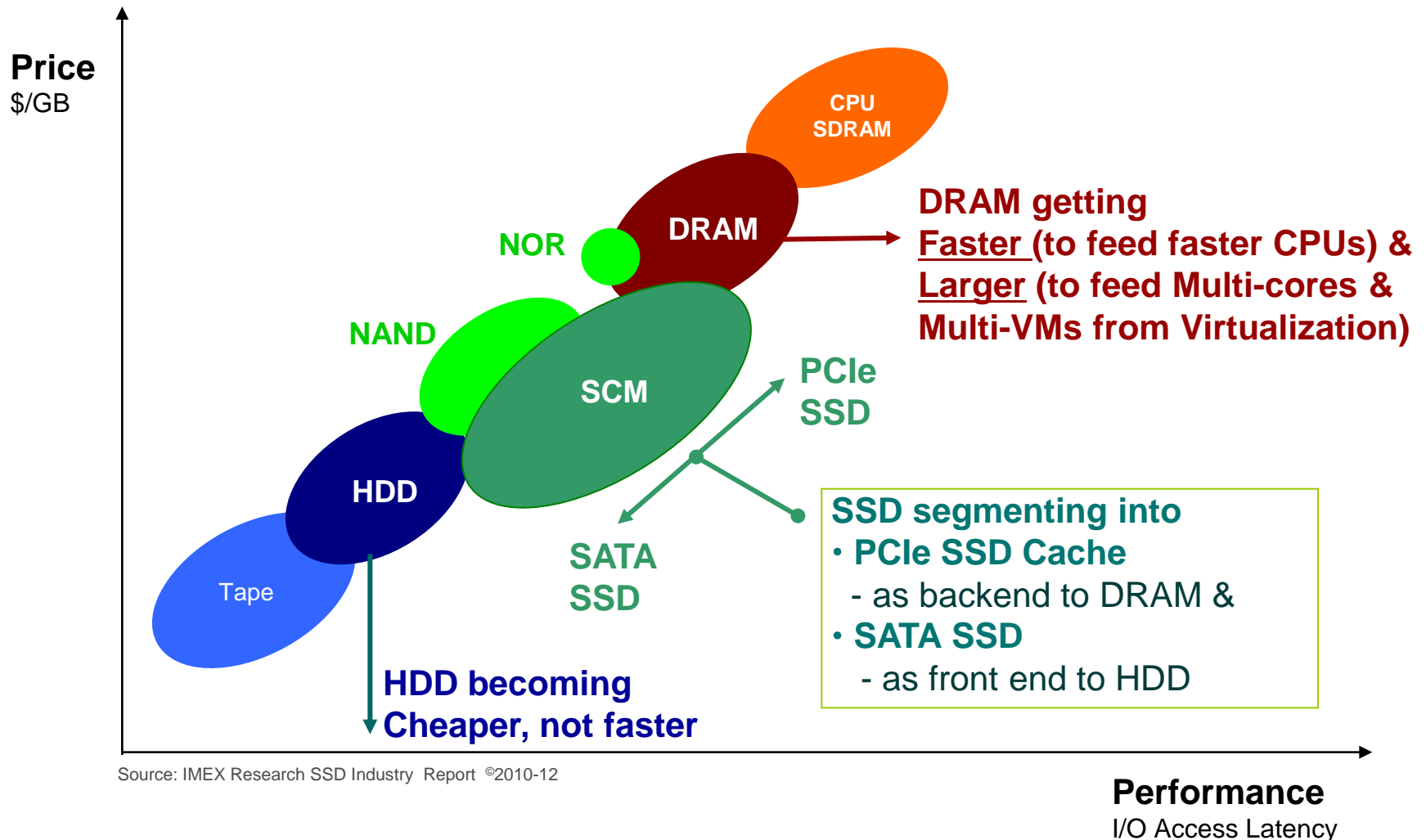


Source: IMEX Research SSD Industry Report ©2011

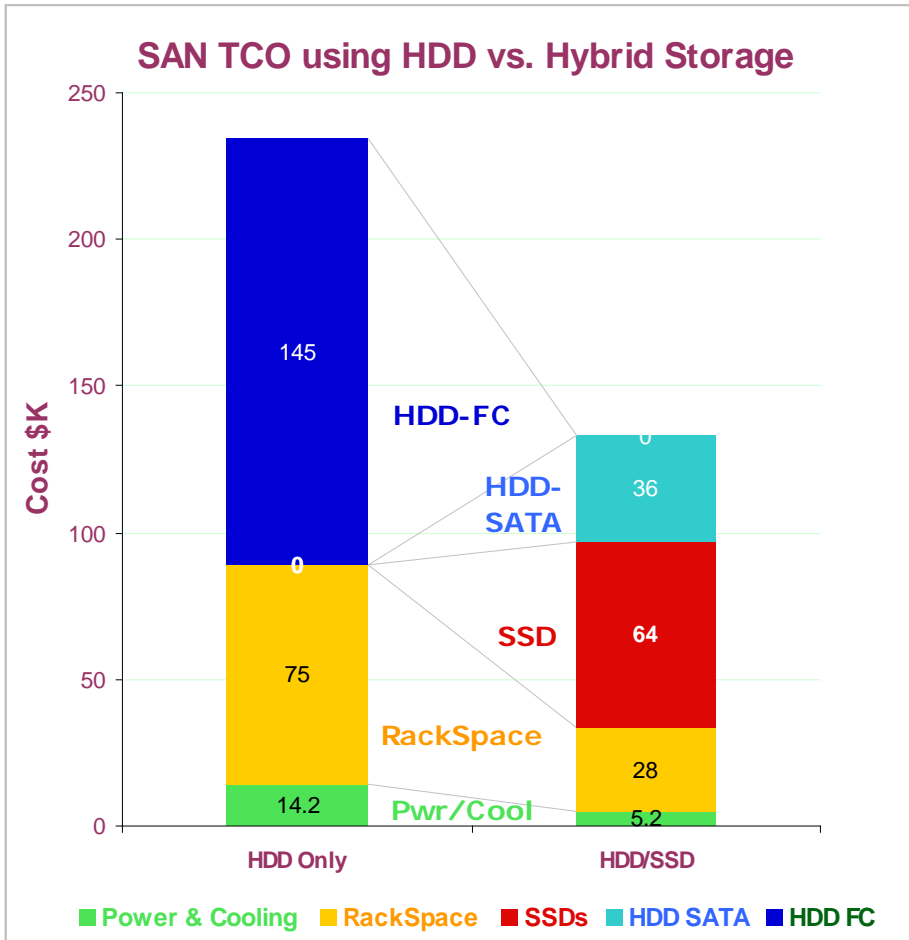
## I/O Access Frequency vs. Percent of Corporate Data



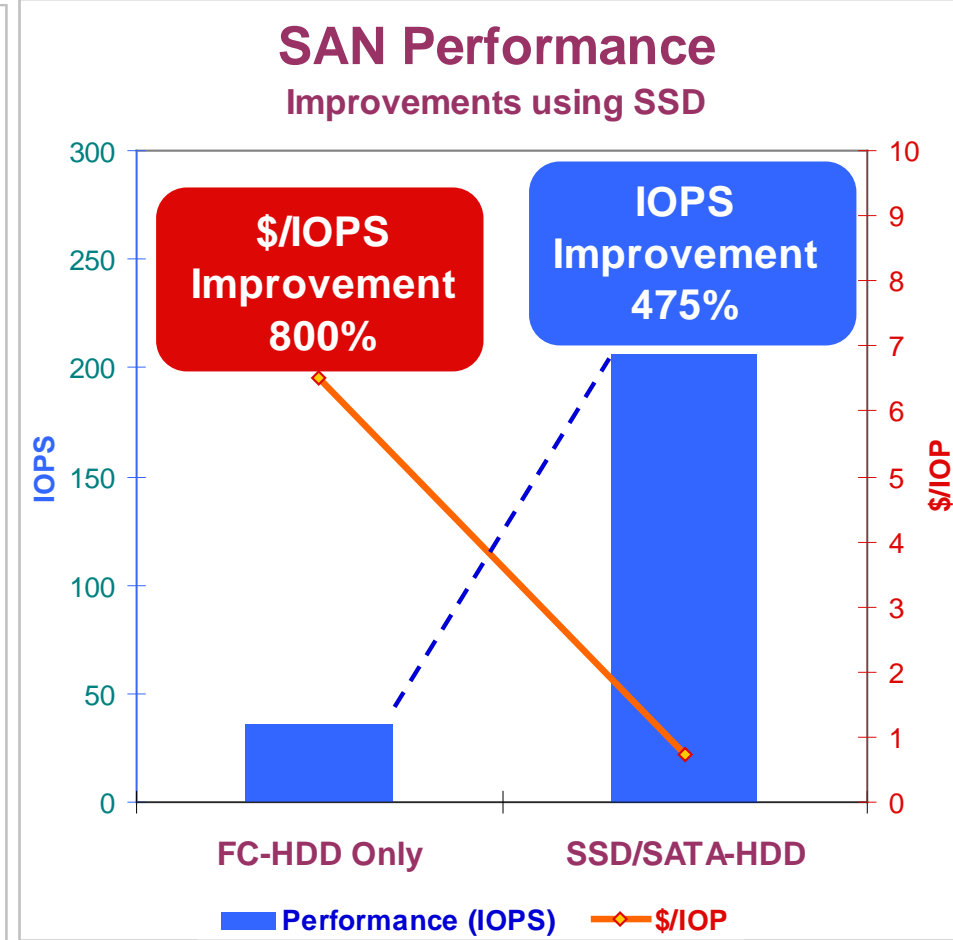
# SSD Storage: Filling Price/Perf.Gaps



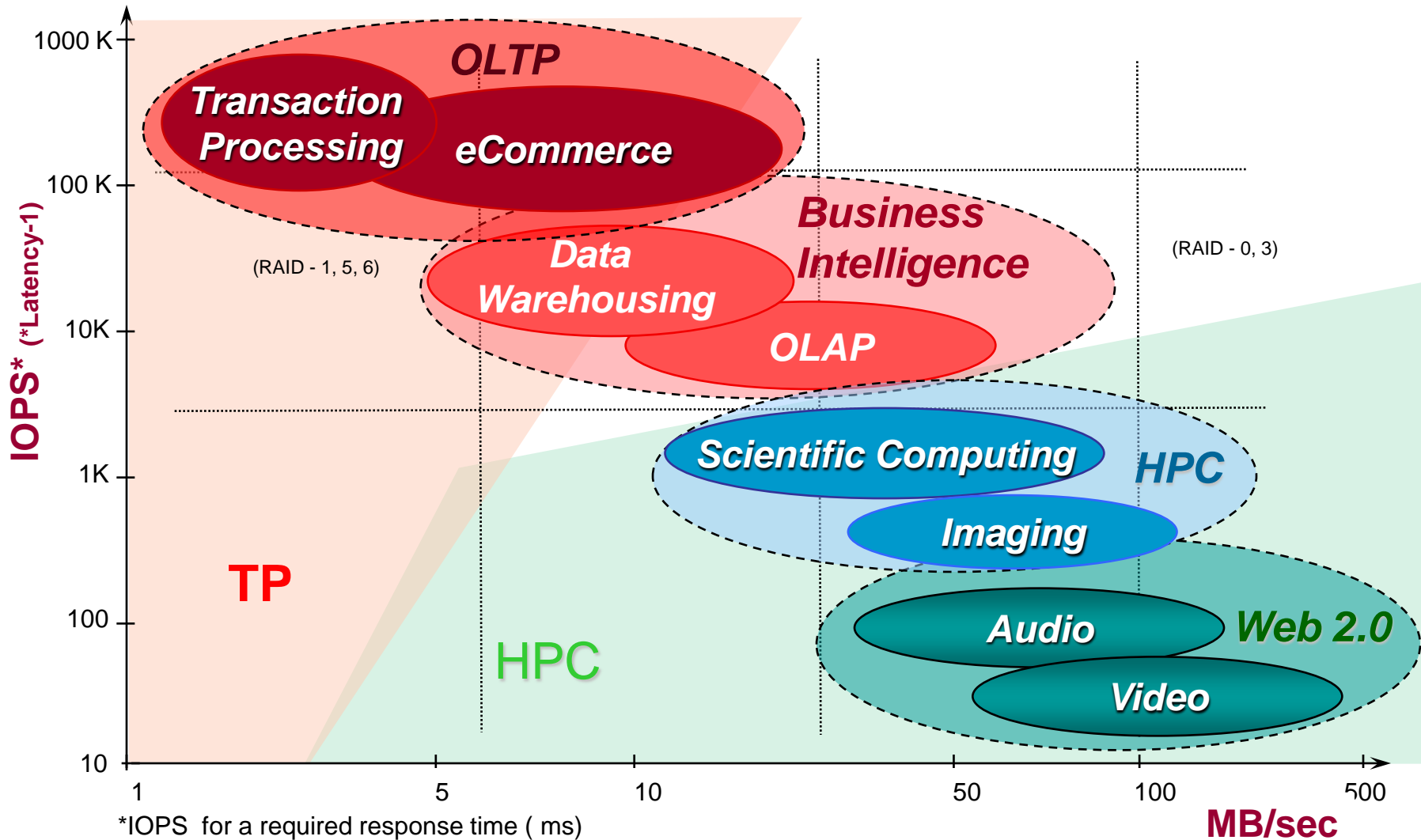
# SSD Storage - Performance & TCO



Source: IMEX Research SSD Industry Report ©2011



# Workloads Characterization



## Storage performance, management and costs are big issues in running Databases

---

- **Data Warehousing Workloads are I/O intensive**
  - Predominantly read based with low hit ratios on buffer pools
  - High concurrent sequential and random read levels
    - ✓ Sequential Reads requires high level of I/O Bandwidth (MB/sec)
    - ✓ Random Reads require high IOPS)
  - Write rates driven by life cycle management and sort operations
- **OLTP Workloads are strongly random I/O intensive**
  - Random I/O is more dominant
    - ✓ Read/write ratios of 80/20 are most common but can be 50/50
    - ✓ Can be difficult to build out test systems with sufficient I/O characteristics
- **Batch Workloads are more write intensive**
  - Sequential Writes requires high level of I/O Bandwidth (MB/sec)
- **Backup & Recovery times are critical for these workloads**
  - Backup operations drive high level of sequential IO
  - Recovery operation drives high levels of random I/O

## Goals & Implementation

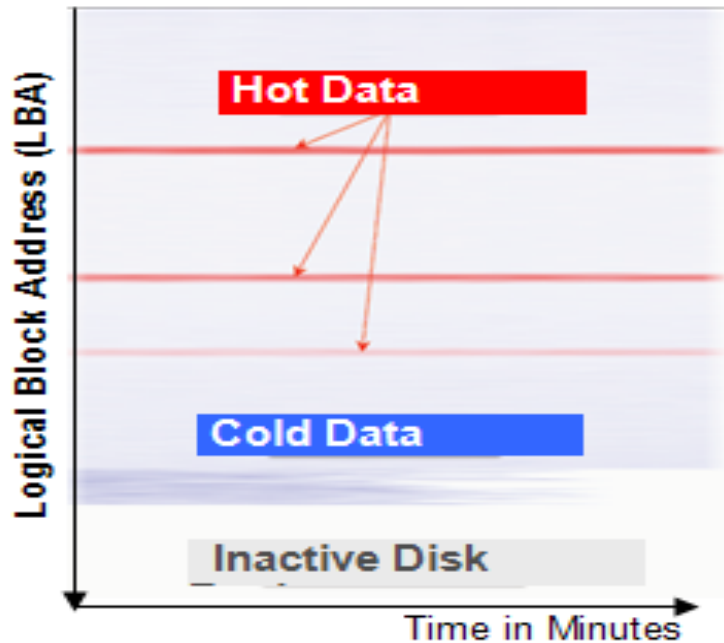
- ◆ Establish **Goals for SLAs** (Performance/Cost/Availability), BC/DR (RPO/RTO) & Compliance
- ◆ **Increase Performance for DB, OLTP and OLAP Apps:**
  - › Random I/O > 20x , Sequential I/O Bandwidth > 5x
  - › Remove Stale data from Production Resources to improve performance
- ◆ Use **Partitioning Software to Classify Data**
  - › By Frequency of Access (Recent Usage) and
  - › Capacity (by percent of total Data) using general guidelines as:
  - › Hyperactive (1%), Active (5%), Less Active (20%), Historical (74%)

## Implementation

- ◆ **Optimize Tiering** by Classifying Hot & Cold Data
  - › Improve Query Performance by reducing number of I/Os
  - › Reduce number of Disks Needed by 25-50% using advance compression software achieving 2-4x compression
- ◆ **Match Data Classification vs. Tiered Devices** accordingly
  - › Flash, High Perf Disk, Low Cost Capacity Disk, Online Lowest Cost Archival Disk/Tape
- ◆ **Balance Cost vs. Performance** of Flash
  - › More Data in Flash > Higher Cache Hit Ratio > Improved Data Performance
- ◆ **Create and Auto-Manage Tiering** (Monitoring, Migrations, Placements) without manual intervention



# Best Practices: I/O Forensics in Storage-Tiering

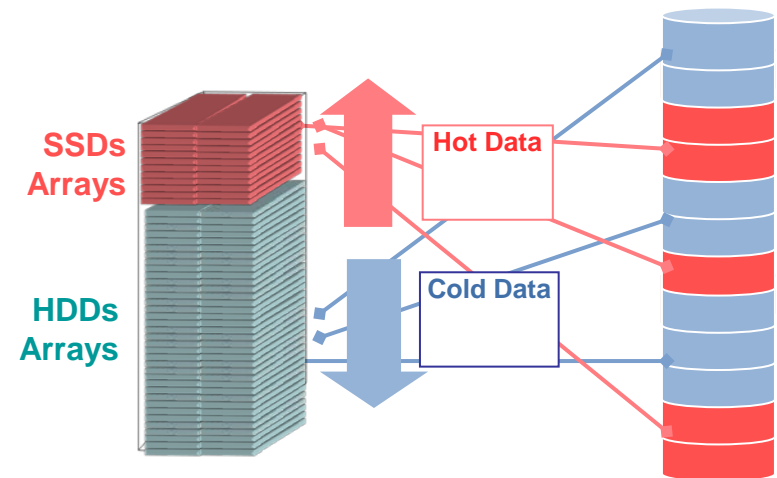


## Storage-Tiered Virtualization

Storage-Tiering at LBA/Sub-LUN Level

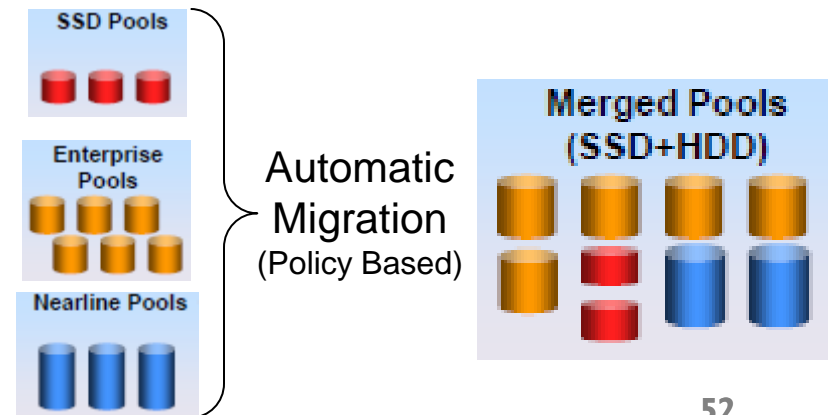
Physical Storage

Logical Volume

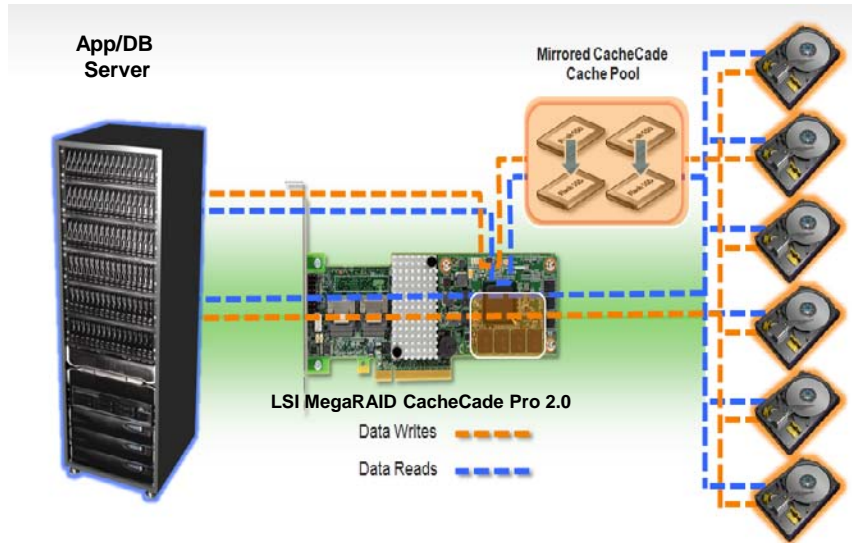


## LBA Monitoring and Tiered Placement

- Every workload has unique I/O access signature
- Historical performance data for a LUN can identify performance skews & hot data regions by LBAs

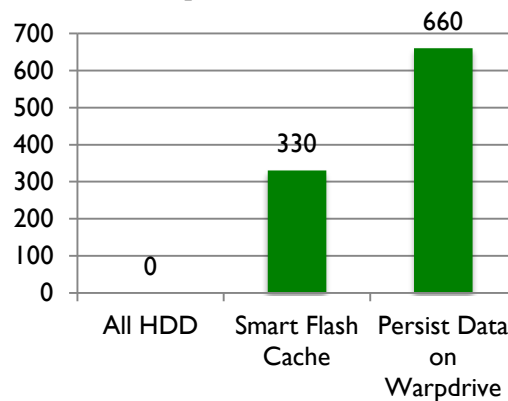


# Best Practices: Cached Storage

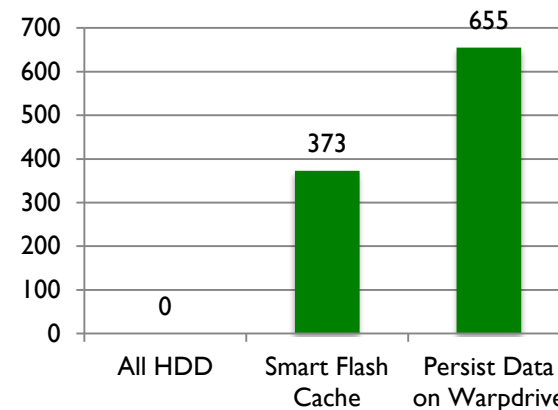


Application	Improvement over Cached vs.HDD only
<b>Oracle OLTP Benchmarks</b>	681%
<b>SQL Server OLTP Benchmark</b>	1251%
<b>Neoload (Web Server Simulation)</b>	533%
<b>SysBench (MySQL OLTP Server)</b>	150%

## Response Time

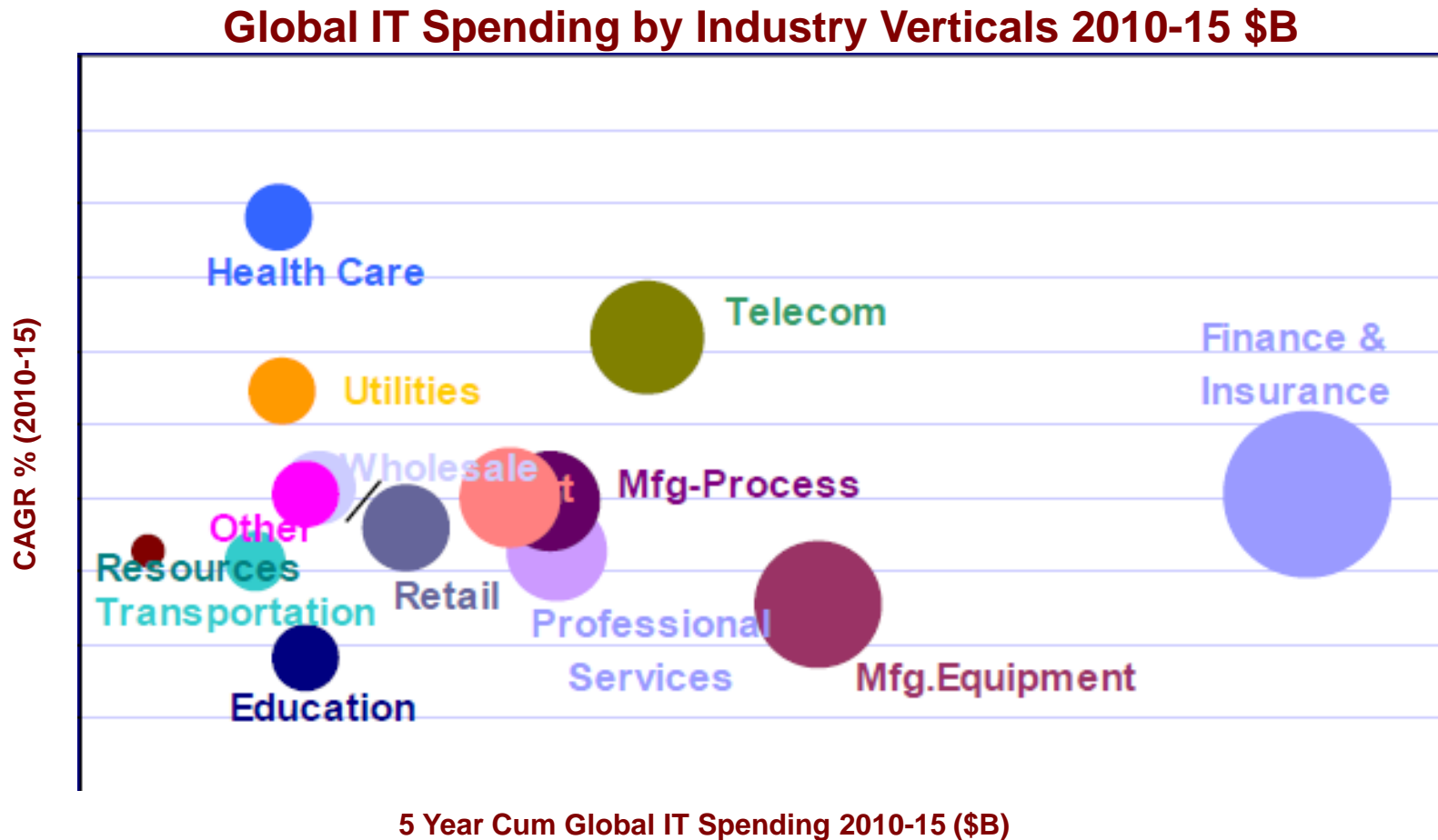


## TPS



# Big Data Targets: Analytics

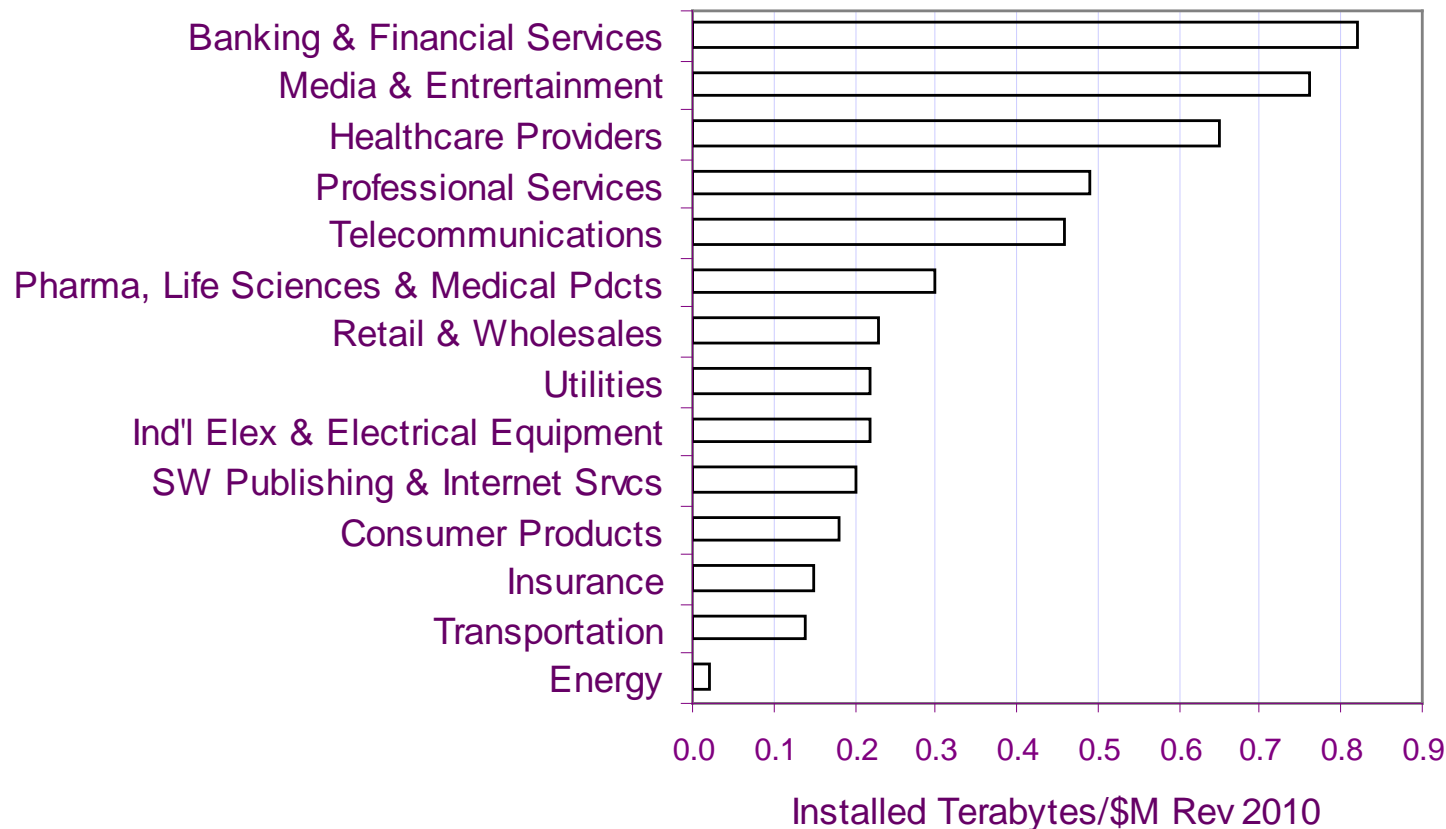
## Key Industries Benefitting from Big Data Analytics



# Big Data Targets – Storage Infrastructure

## Value Potential of Using Big Data by Data Intensive Verticals

### Data Intensity by Industry Vertical

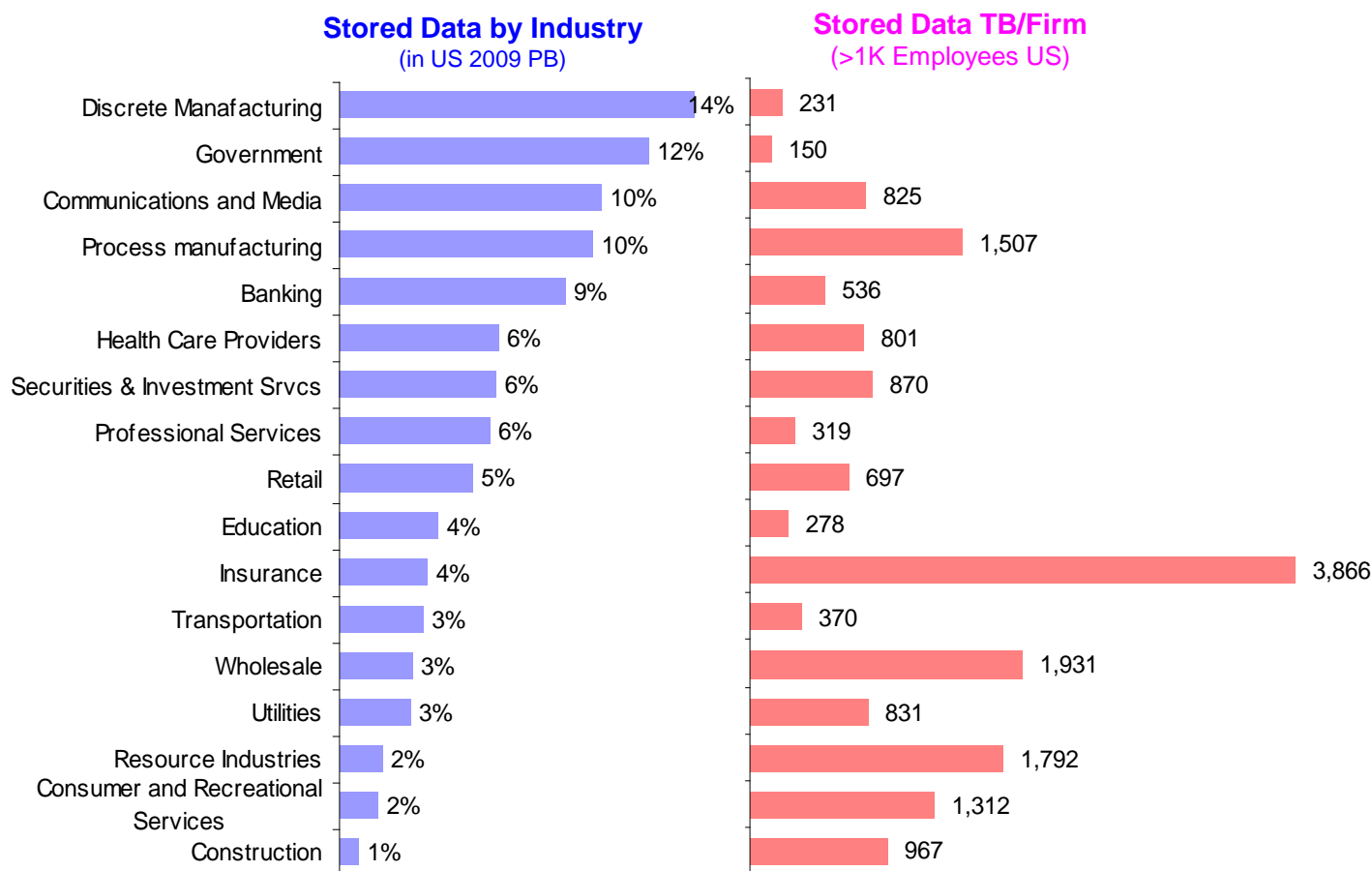


# Big Data Targets: Storage Infrastructure

## Data Stored by Large US Enterprises

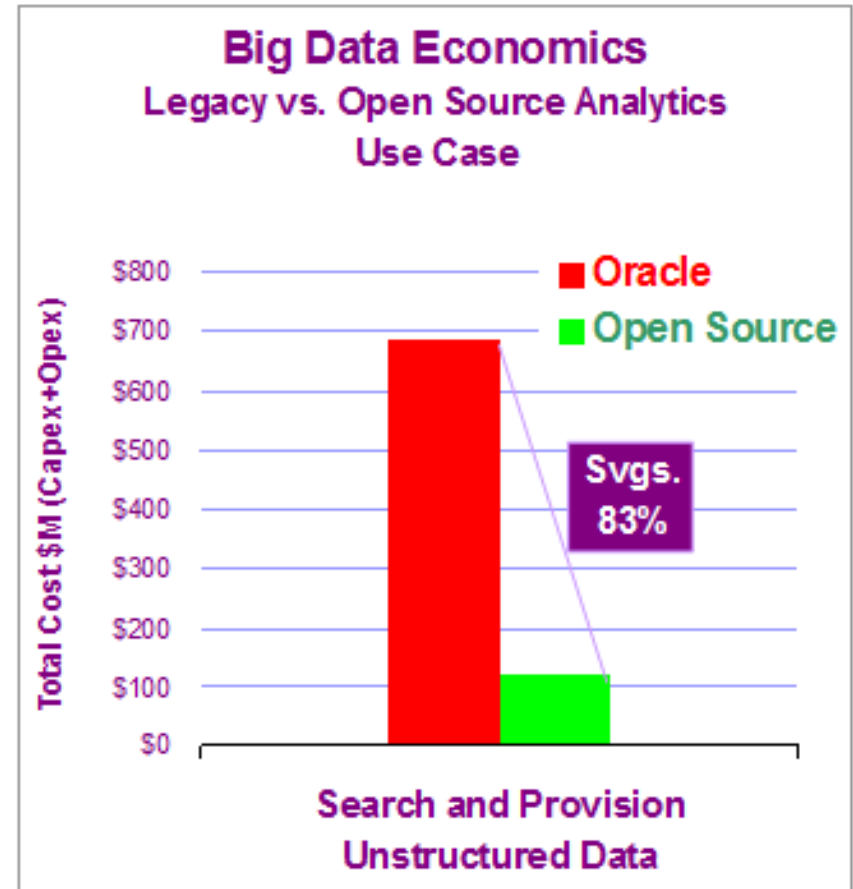
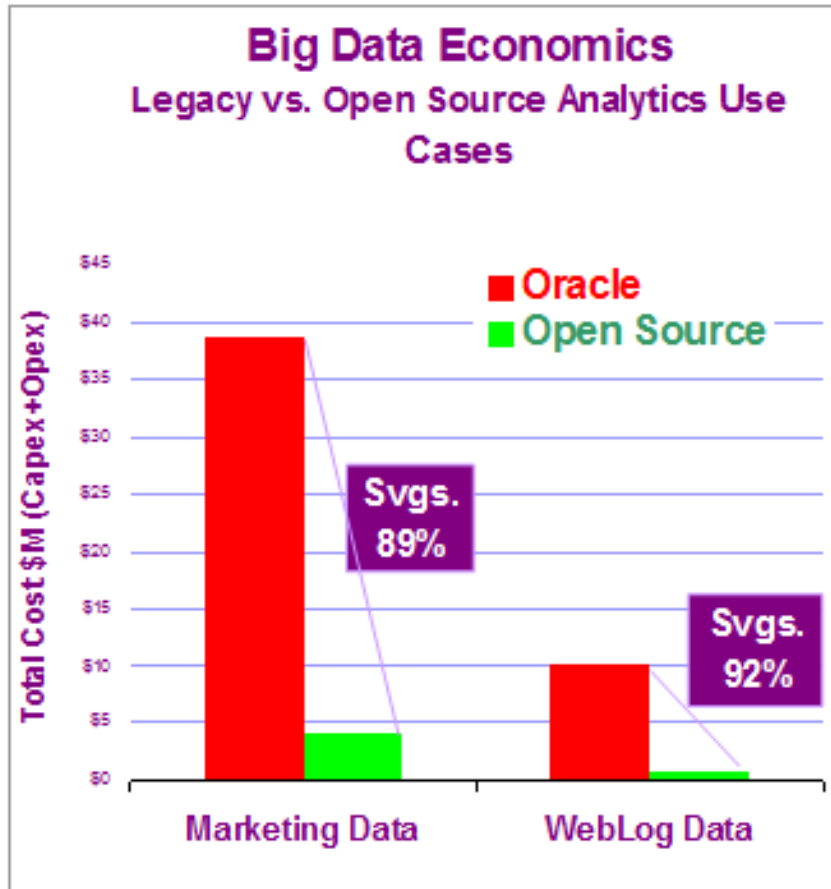
### Big Data Storage Potential

Data Stored by Large US Enterprises



# Big Data Targets: Savings w Open Source

## Legacy BI vs. Open Source Big Data Analytics



# Rise of Big Data Adoption

2006

YAHOO!



2007



last.fm

2008

Google able grape

ImageShack Cascading

TBM facebook

ENORMO Every property has secrets AQ

krugle rockspace

Lookery Control freaks welcome!

The New York Times Jodel

Zvents FORMATION SCIENCES

News Corporation

Cornell University Computing and Information Science Visible MEASURES

LOTAME NetSeer

parc veeva

2009

AOL

amazon web services

deepdyve

cooliris

eyealike

TEXTMAP

PRG College of Technology

iterend

tailsweep

hulu

RapLeaf

USCMS

Ning quxntcast

cloudera pressflip

detikSearch

WorldLingo

Systems@ETH Zürich

VK SOLUTIONS

TERRACON

HOSTING HABITAT

HOLA

Terrier

adknowledge

stampede

2010

SAMSUNG rubicon

BERKELEY LAB VISIBLE TECHNOLOGIES

APOLLO GROUP ADSDAQ

rackspace HOSTING RapLeaf

wordnik MODIGEN

COMSCORE trulia

Accela Forward3D

LinkedIn Microsoft

Infochimps Pham 2Pham

ADMELD gumgum

Pronux The Datagraph Blog

NETFLIX mobileanalytics.tv

markt24.de twitter

media6degrees BEEBLER

SLC Security eBay



## ➤ **Big Data creating paradigm shift in IT Industry**

- ◆ Leverage the opportunity to optimize your computing infrastructure with Big Data Infrastructure after making a due diligence in selection of vendors/products, industry testing and interoperability.
- ◆ Apply best storage technologies listed in this presentation and elsewhere

## ➤ **Optimize Big Data Analytics for Query Response Time vs. # of Users**

- ◆ Improving Query Response time for a given number of users (IOPs) or Serving more users (IOPS) for a given query response time

## ➤ **Select Automated Storage Management Software**

### ◆ **Data Forensics and Tiered Placement**

- Every workload has unique I/O access signature
- Historical performance data for a LUN can identify performance skews & hot data regions by LBAs. Non-disruptively migrate hot data using auto-tiering Software

## ➤ **Optimize Infrastructure to meet needs of Applications/SLA**

### ◆ **Performance Economics/Benefits**

- Typically 4-8% of data becomes a candidate and when migrated for higher performance tiering can provide response time reduction of ~65% at peak loads. Many industry Verticals and Applications will benefit using Big Data



**Many thanks to the following individuals  
for their contributions to this tutorial.**

**Source: IMEX Research**

**Joseph White  
Anil Vasudeva**

Send any questions or comments on this presentation to SNIA: [tracktutorials@snia.org](mailto:tracktutorials@snia.org)