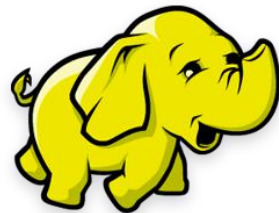# Big Data
*What it is and what is isn't...*

Vincent Staropoli

*CloudDevelop Conference*
*October, 23 2015 – Columbus, OH*

# About Me

# The Story of Hadoop
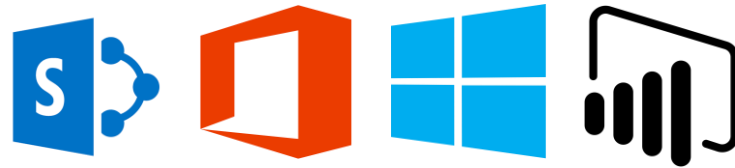
- **Google** and **Yahoo** got together, got drunk and made a baby elephant named **Hadoop**

- It was so big, they had to wrap a map around it's **flume** to reduce the **splunk** on **mahout**

- Personally, I'd rather learn **pig latin** from a guy named **ambari** than **sqoop** the **hive** of **hbase** that elephant leaves around

- But, when you're in a **storm** spinning a **yarn** like this one, there's only one place you wish you could go in a flash

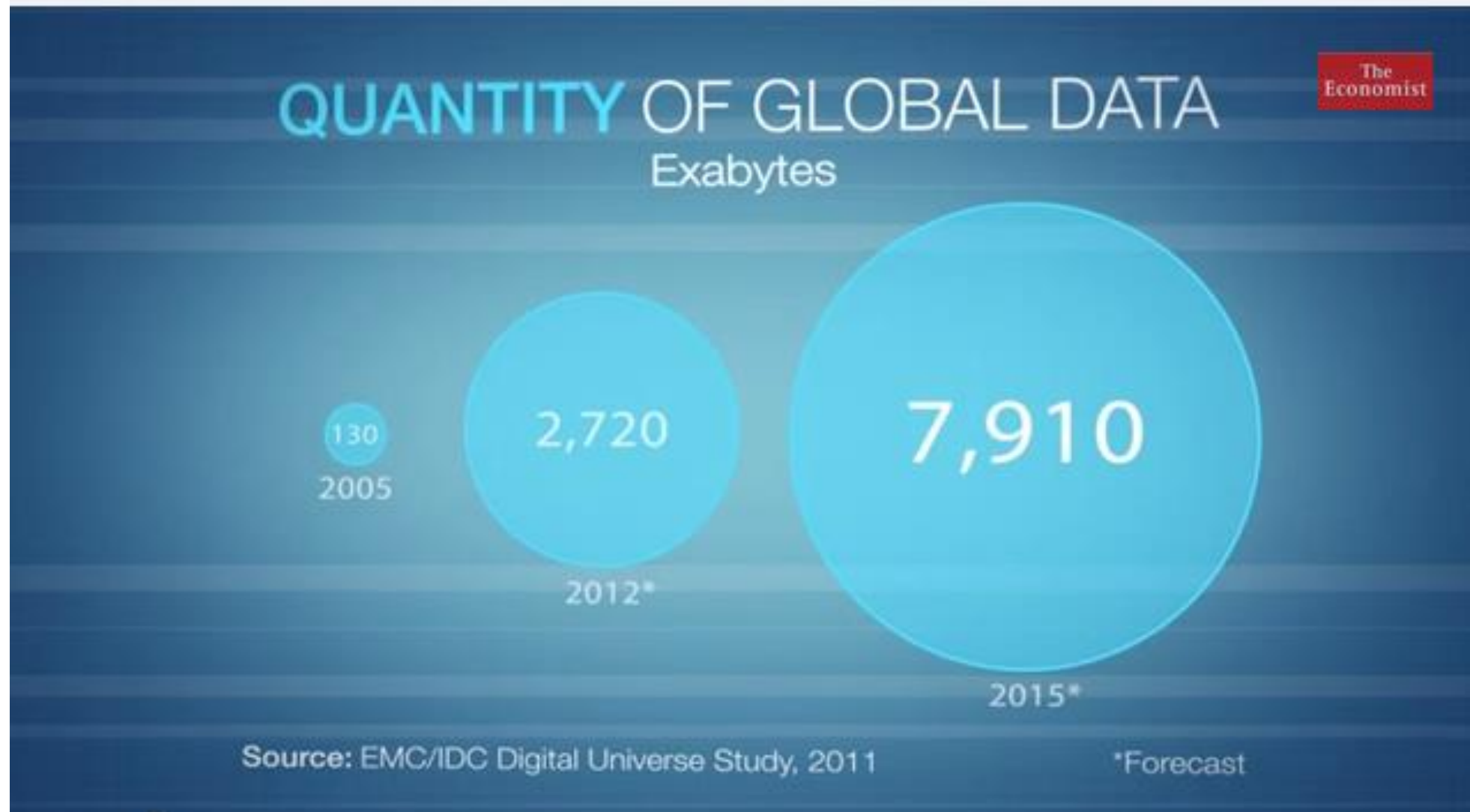- **The Cloud**

# Why Discuss Big Data?

- Not just a technology topic

- Adjective as well as a noun, occasionally a verb

- Tool in the toolkit for developers and businesses

- No magic, no silver bullet and won't change culture

- Changing *RAPIDLY*

# Our Journey...

- Data, data, data

- Systems & technology

- Data-driven

- "Welcome to real life"

# Data Getting Bigger

# Big Data Sources

"The biggest reason that investments in big data fail to pay off is that most companies don't do a good job with the information they already have."

-*"You May Not Need Big Data After All"*
https://hbr.org/2013/12/you-may-not-need-big-data-after-all

# Big Data Systems

- Framework of technologies to store and analyze data

- Exceeds the processing capacity of conventional database systems

- Characterized by **volume**, **velocity** and **variety**

- Significant effort to find the signal within the noise

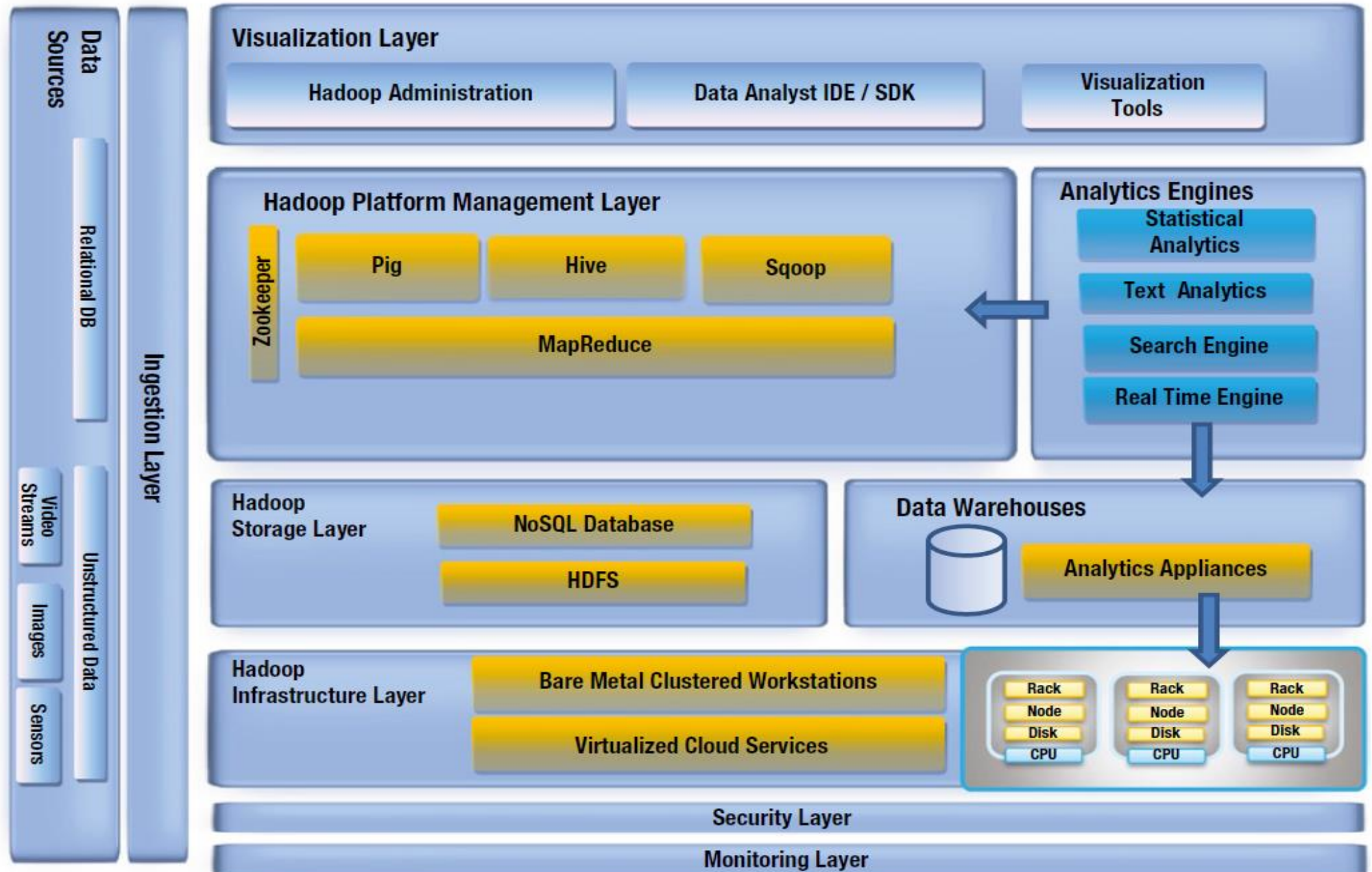# Big Data Systems

- Core Technologies
  - Hadoop
  - MapReduce
- New technologies to manage huge amounts of data
  - Store
  - Access
  - Analyze
- Monetize the benefits of owning huge amounts of data
  - Capability to process massive amounts of data
  - Efficient, cost-effective, and timely

# Big Data Systems

- Storage and infrastructure

- Platform management
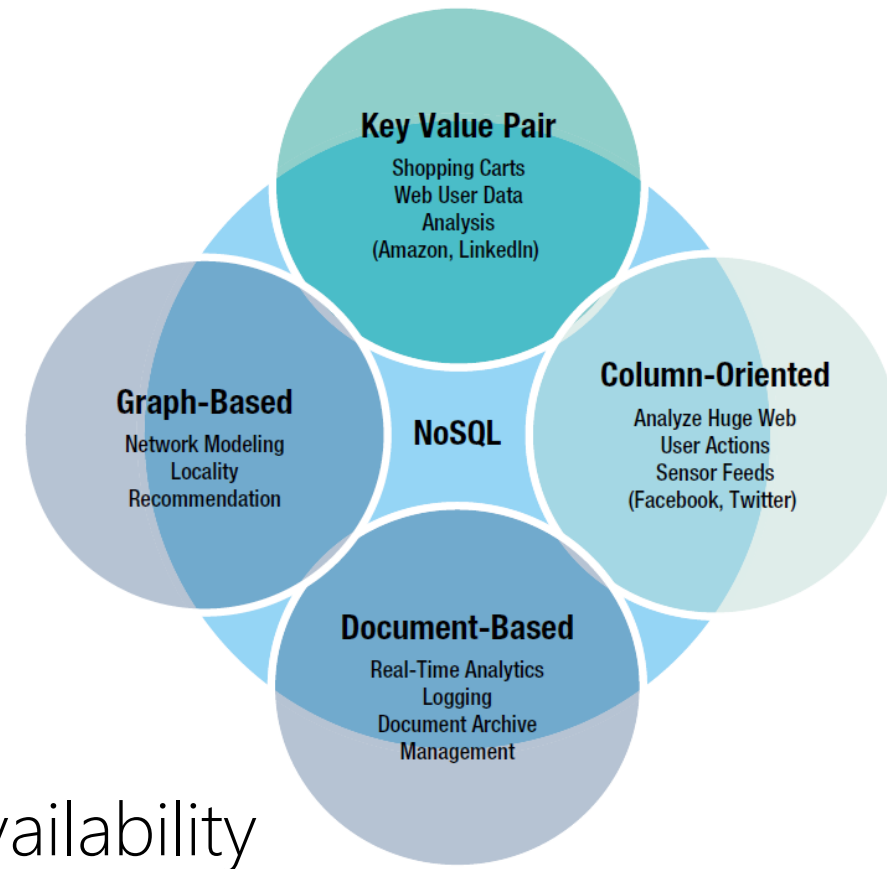
- Ingestion

- Visualization

# Big Data Systems
*Reference Architecture*

**Data Sources**

- Relational DB
- Video Streams
- Images
- Sensors

**Unstructured Data**

**Ingestion Layer**

## Visualization Layer

| Hadoop Administration | Data Analyst IDE / SDK | Visualization Tools |
|---|---|---|

## Hadoop Platform Management Layer

**Zookeeper**

| Pig | Hive | Sqoop |
|---|---|---|

**MapReduce**

## Analytics Engines

- Statistical Analytics
- Text Analytics
- Search Engine
- Real Time Engine

## Hadoop Storage Layer

**NoSQL Database**

**HDFS**

## Data Warehouses

**Analytics Appliances**

## Hadoop Infrastructure Layer

**Bare Metal Clustered Workstations**

**Virtualized Cloud Services**

| Rack | Rack | Rack |
|---|---|---|
| Node | Node | Node |
| Disk | Disk | Disk |
| CPU | CPU | CPU |

**Security Layer**

**Monitoring Layer**

# Big Data and NoSQL

- "Not Only SQL"

- Different solutions for different applications

- Must relax guarantees around consistency, availability and partition tolerance (the CAP Theorem)

- Likely have a combination of relational and NoSQL databases

**Key Value Pair**
Shopping Carts
Web User Data
Analysis
(Amazon, LinkedIn)

**Column-Oriented**
Analyze Huge Web
User Actions
Sensor Feeds
(Facebook, Twitter)

**Graph-Based**
Network Modeling
Locality
Recommendation

**NoSQL**

**Document-Based**
Real-Time Analytics
Logging
Document Archive
Management

# NoSQL Database Systems

| Key-Value Data Stores | Column-oriented Data Stores | Document Data Stores | Graph Data Stores |
|---|---|---|---|
| MEMCACHED, riak, Redis, V, membase | Amazon SimpleDB, Cassandra, H-BASE, Windows Azure, hadoop, HYPERTABLE | mongoDB, RAVEN DB, db, terrastore, Apache CouchDB relax | Neo4j the graph database, HyperGraphDB, AllegroGraph, InfiniteGraph The Distributed Graph Database, INFO GRID |

# MapReduce

- Map function
  - Applies a function on every key/value pair in the collection
  - Generates a new collection


- Reduce function
  - Works on the new generated collection
  - Applies an aggregate function to compute a final output

# MapReduce

Sample Data:
```
[{ "94303": "Tom"}, {"94303": "Jane"}, {"94301":
"Arun"}, {"94302": "Chen"}]
```

Get the names of all those who reside in a particular zip code:
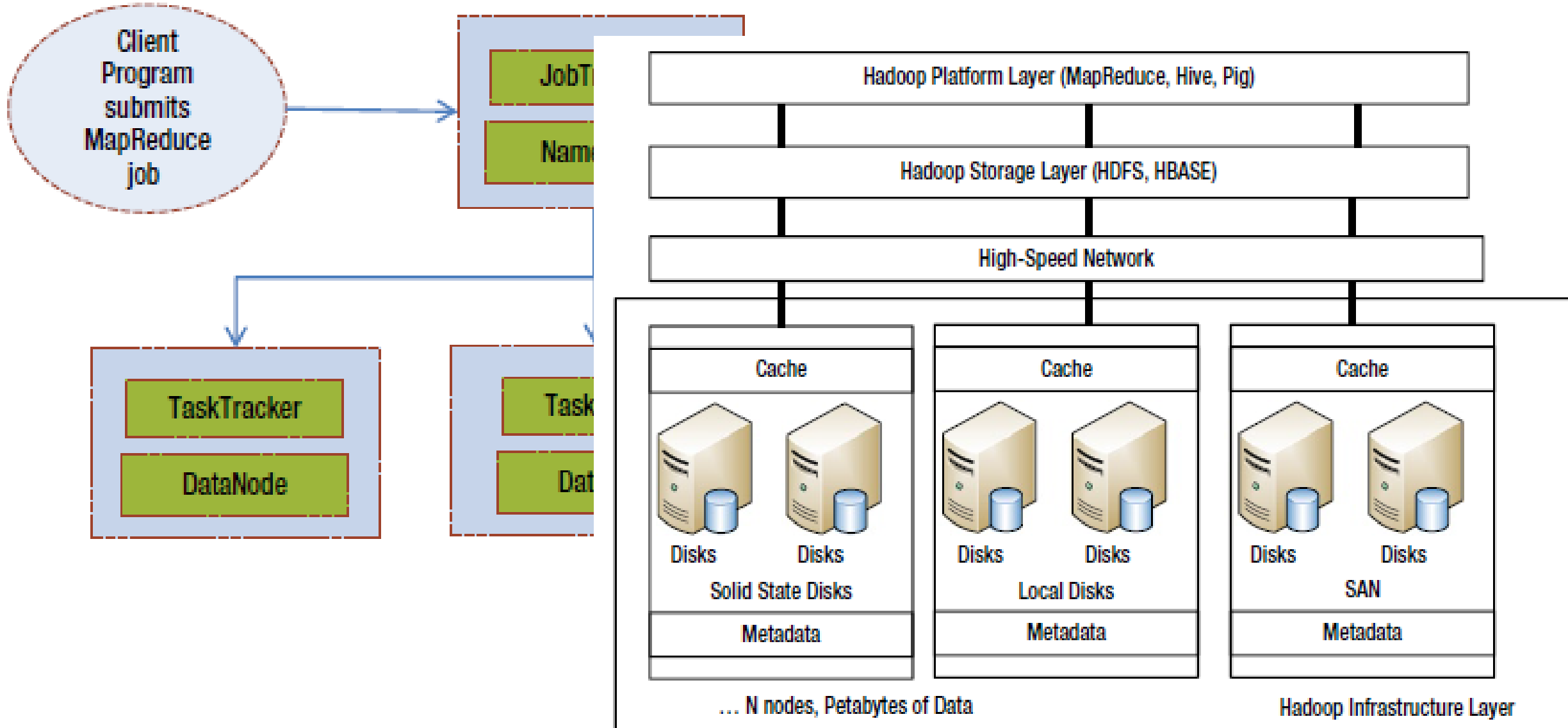```
[{"94303":["Tom", "Jane"]}, {"94301":["Arun"]},
{"94302":["Chen"]}]
```

Reduce function output to count the number of people by zip code:
```
[{"94303": 2}, {"94301": 1}, {"94302": 1}]
```

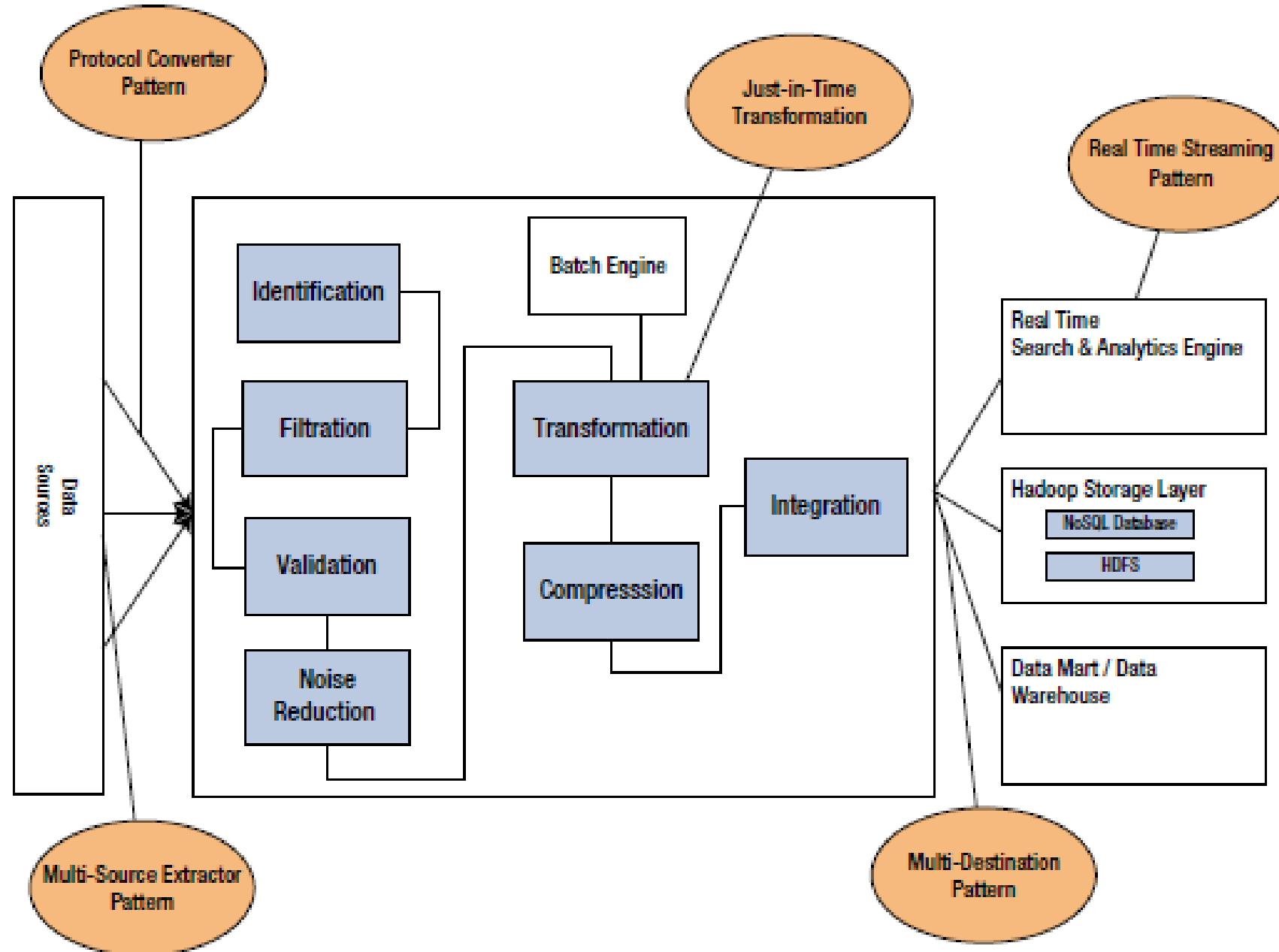# Hadoop Clustering
*Conceptual Architecture*

Client Program submits MapReduce job

JobT...

Nam...

TaskTracker

DataNode

Task...

Dat...

Hadoop Platform Layer (MapReduce, Hive, Pig)

Hadoop Storage Layer (HDFS, HBASE)

High-Speed Network

Cache

Disks    Disks
Solid State Disks
Metadata

Cache

Disks    Disks
Local Disks
Metadata

Cache

Disks    Disks
SAN
Metadata

... N nodes, Petabytes of Data

Hadoop Infrastructure Layer

# Hadoop Clustering

*Commodity Hardware*

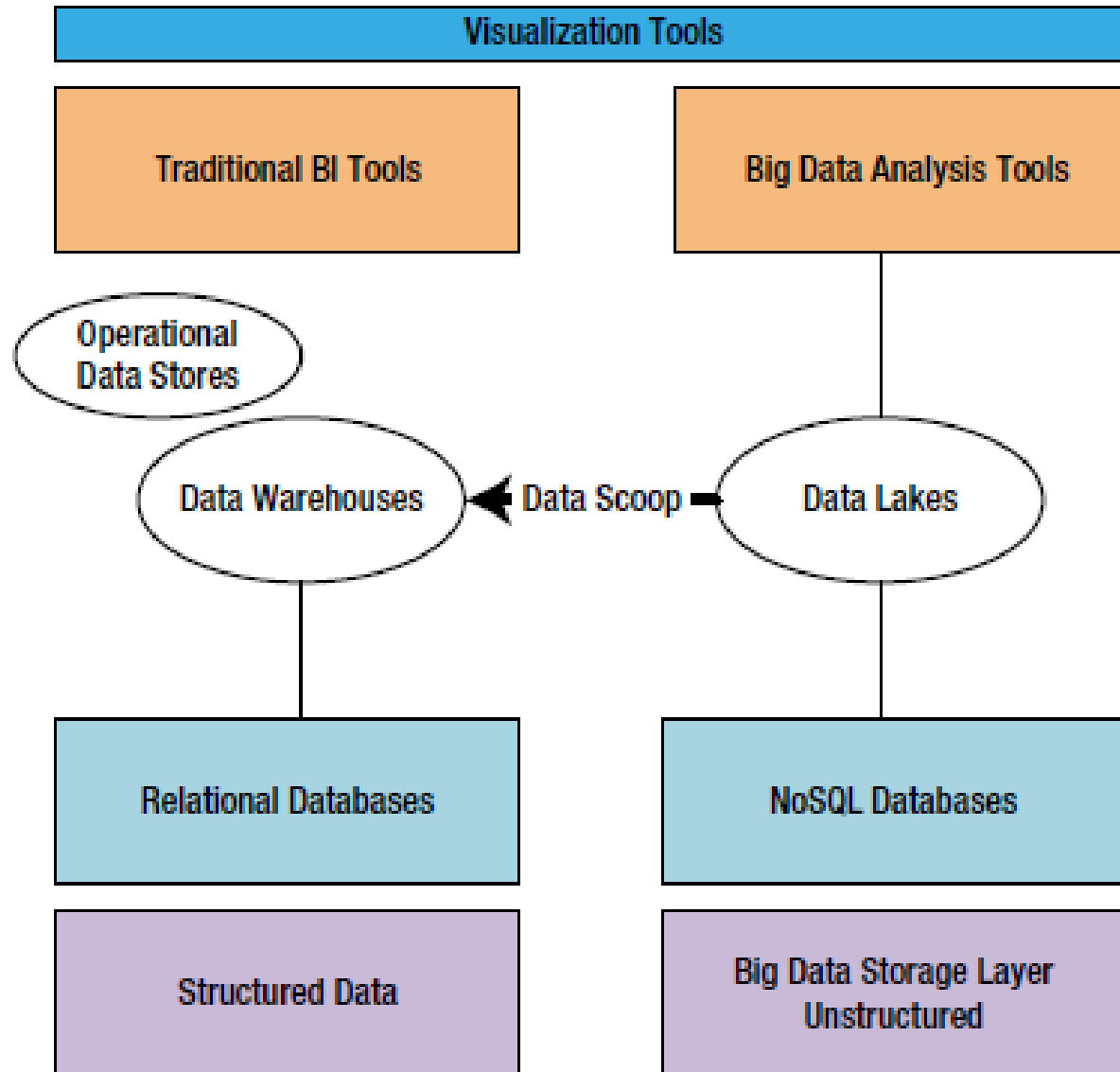| Entity | Configuration of Data Node |
|:---:|:---:|
| CPU | Two CPU sockets with six or eight cores, Intel Xeon processor E5-2600 series @ 2.9 GHz |
| Memory | 48 GBs ( 6X8 GBs 1.35v 1333 MHz DIMMs) or 96 GBs (6x16 GBs 1.35v 1333 MHz DIMMs) |
| Disk | 10-12, 1-3 TB SATA drives |
| Network | 1x dual port 10 GbE NIC, or 1x quad port 1 GbE NIC |

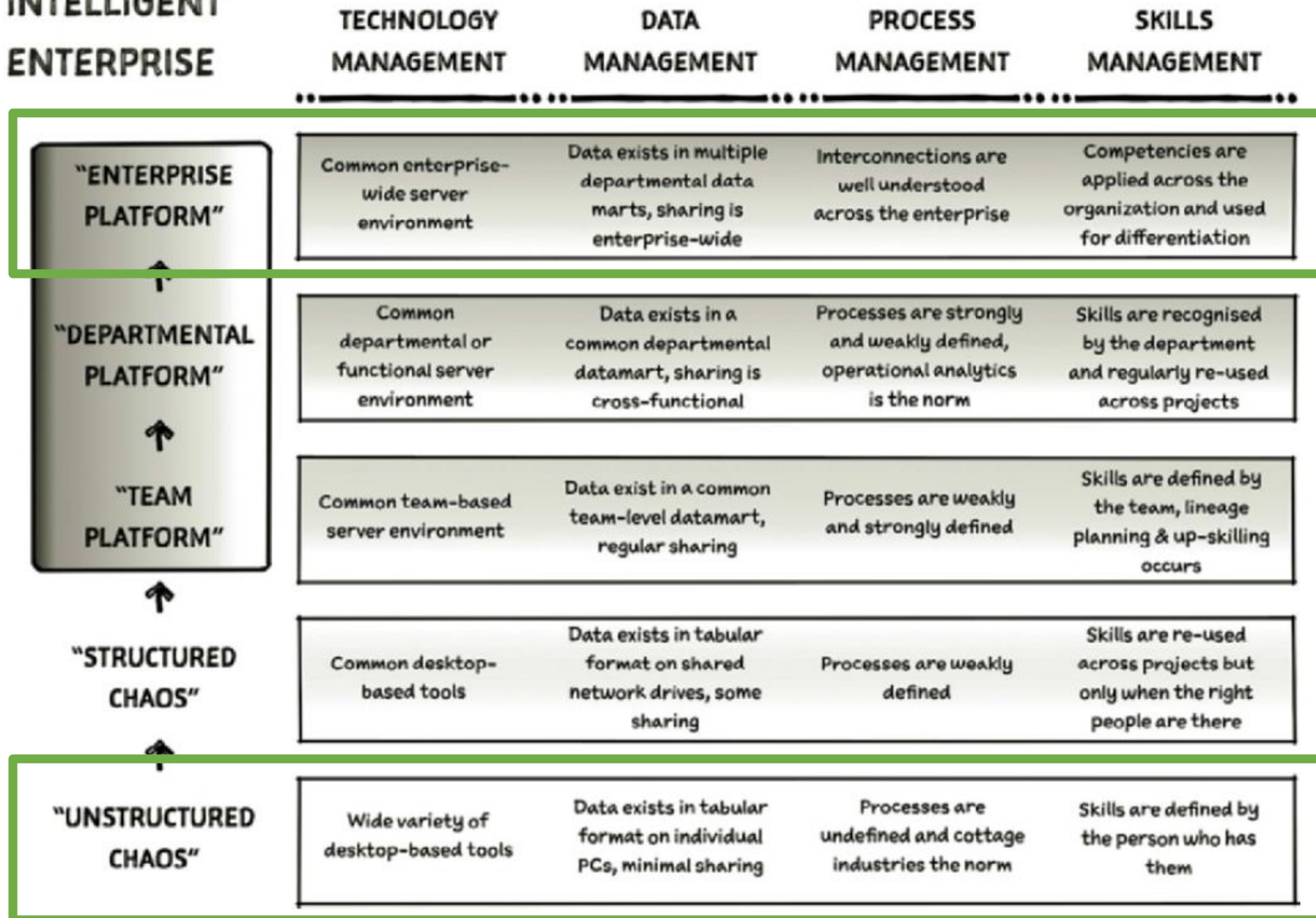# Hadoop Clustering

*Ingestion Engine*
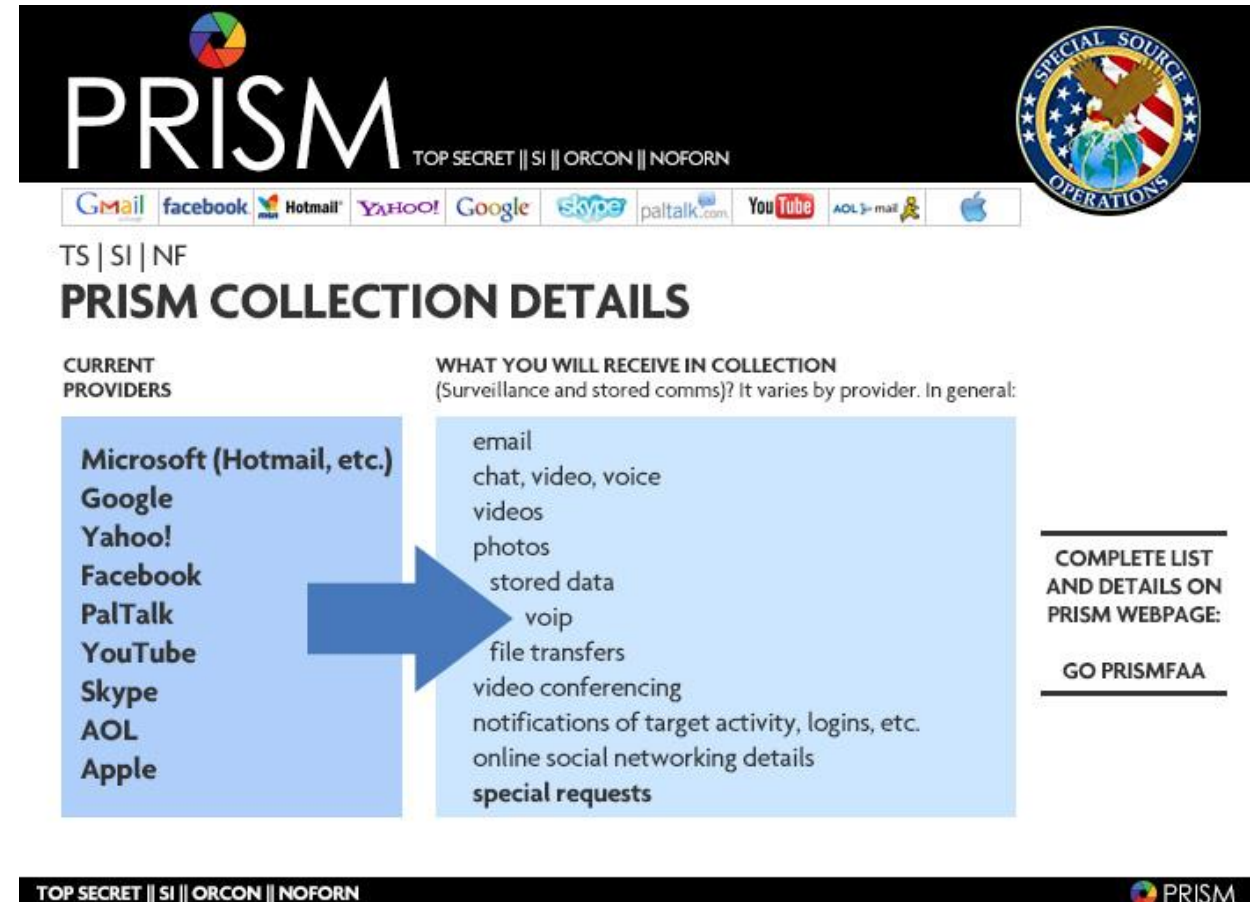
# Hadoop Clustering

*Scalability*

# Data Driven Organizations

- Establish one undisputed source

- Near-real-time feedback

- Articulate business rules (and regularly update them)

- High-quality coaching to employees

# the INTELLIGENT ENTERPRISE

| | TECHNOLOGY MANAGEMENT | DATA MANAGEMENT | PROCESS MANAGEMENT | SKILLS MANAGEMENT |
|---|---|---|---|---|
| "ENTERPRISE PLATFORM" | Common enterprise-wide server environment | Data exists in multiple departmental data marts, sharing is enterprise-wide | Interconnections are well understood across the enterprise | Competencies are applied across the organization and used for differentiation |
| "DEPARTMENTAL PLATFORM" | Common departmental or functional server environment | Data exists in a common departmental datamart, sharing is cross-functional | Processes are strongly and weakly defined, operational analytics is the norm | Skills are recognised by the department and regularly re-used across projects |
| "TEAM PLATFORM" | Common team-based server environment | Data exist in a common team-level datamart, regular sharing | Processes are weakly and strongly defined | Skills are defined by the team, lineage planning & up-skilling occurs |
| "STRUCTURED CHAOS" | Common desktop-based tools | Data exists in tabular format on shared network drives, some sharing | Processes are weakly defined | Skills are re-used across projects but only when the right people are there |
| "UNSTRUCTURED CHAOS" | Wide variety of desktop-based tools | Data exists in tabular format on individual PCs, minimal sharing | Processes are undefined and cottage industries the norm | Skills are defined by the person who has them |

# Big Data Systems

- Why is this so powerful?
  - A more complete picture
  - Machine learning, data mining
  - Patterns, patterns, patterns
  - Tradeoffs?

# How the Cubs Emerged From the Stone Age

Chicago has worked hard at modernizing every facet of its operation—from the scoreboard to its data gathering—and the club is now benefitting.



"These were not revolutionary advances within the industry. But for the Cubs, **it felt like the space age**. The team's previous information hub was a lone secretary who kept player contracts in file cabinets."

# Making the Case for the 'Long-Tail' of Big Data



"Instead, the **future of data management lies in "data curation,"** which he describes as being "aimed directly at the 'long tail'."

"The long tail refers to the hundreds or thousands of data silos not captured within the traditional data warehouse, and **which can only be captured and integrated at scale by applying automation and machine-learning** based on statistical patterns."

Do you consider analytics in your applications?

Are you helping create a data driven culture?