# Big Data – Data Science

**Srinivas Bhoosarapu**
**(TechnoEconomist, Former CISO and Chief FinTech & Innovation Officer)**

**$One Trillion Digital Economy**
**Viksit Bharat - 2047**

# Agenda

- **Big Data**

- **Data Science**

- **LLM**

- **Use cases and Case Studies**

- **Conclusion**

# Big Data – Data Science



*Large-Scale Data Management*

*Big Data Analytics*

*Data Science and Analytics*

- How to manage very large amounts of data and extract value and knowledge from them

# Introduction to Big Data

*What is Big Data?*

*What makes data, "Big" Data?*

# Big Data Definition

- No single standard definition…

  "***Big Data***" is data whose scale, diversity, and complexity require new architecture, techniques, algorithms, and analytics to manage it and extract value and hidden knowledge from it…

# What is Big Data?

- Massive sets of unstructured/semi-structured data from Web traffic, social media, sensors, etc
- Petabytes, exabytes of data
    - Volumes too great for typical DBMS
- Information from multiple internal and external sources:
    - Transactions
    - Social media
    - Enterprise content
    - Sensors
    - Mobile devices

- In the last minute there were …….
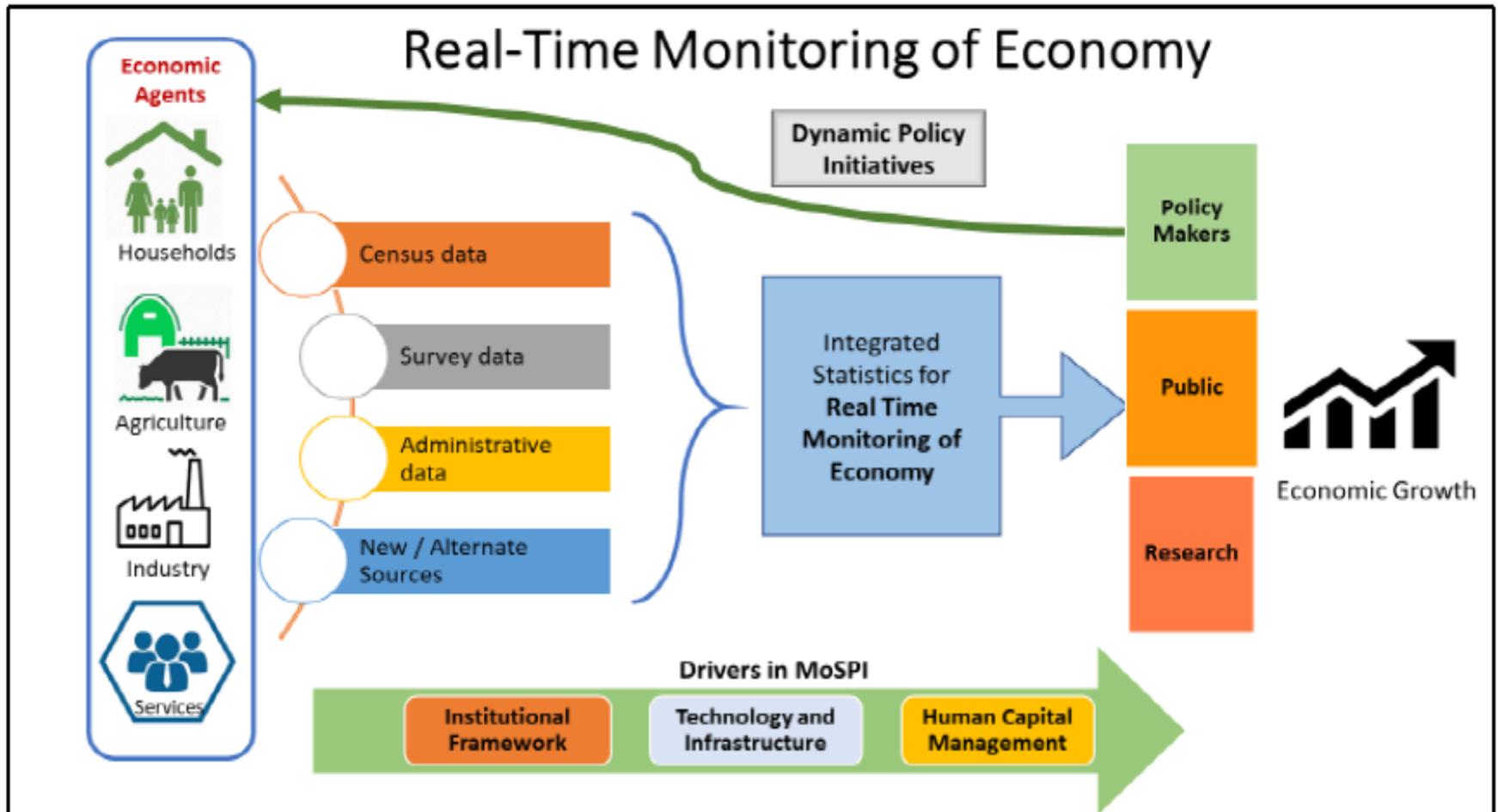
# What is Big Data?

- [What is Big Data](#)

- **204 million emails sent**
- **61,000 hours of music listened to on Pandora**
- **20 million photo views**

- **100,000 tweets**
- **6 million views and 277,000 Facebook Logins**
- **2+ million Google searches**
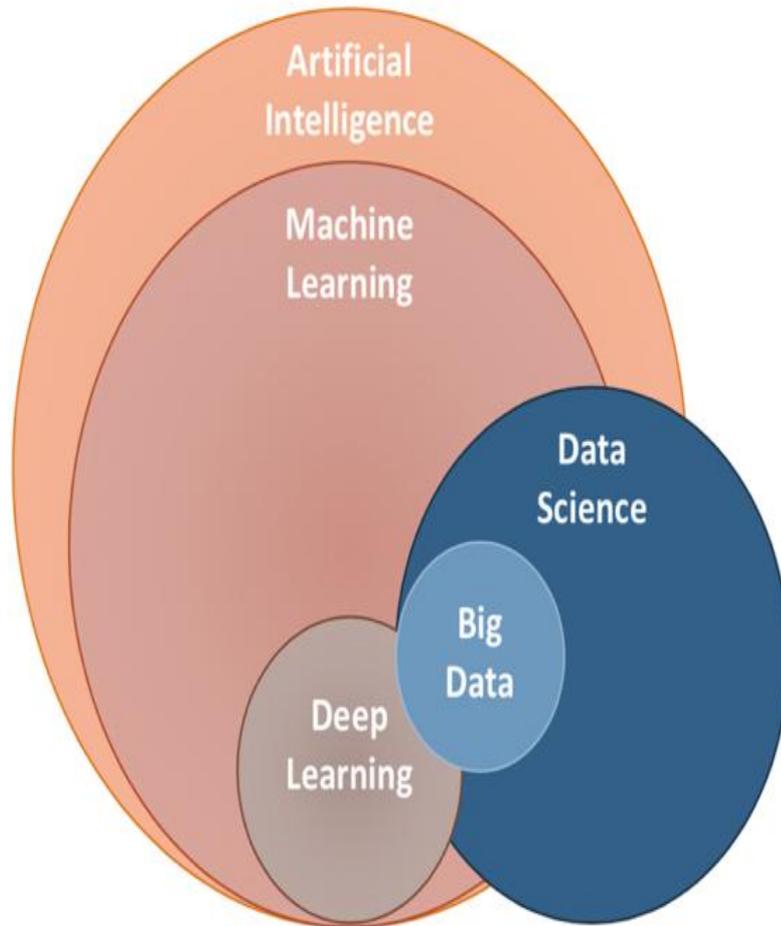- **3 million uploads on Flickr**

# What is Big Data?

- Companies leverage data to adapt products and services to:
  - Meet customer needs
  - Optimize operations
  - Optimize infrastructure
  - Find new sources of revenue
  - Can reveal more patterns and anomalies

- IBM estimates that by 2015 4.4 million jobs will be created globally to support big data
  - 1.9 million of these jobs will be in the United States
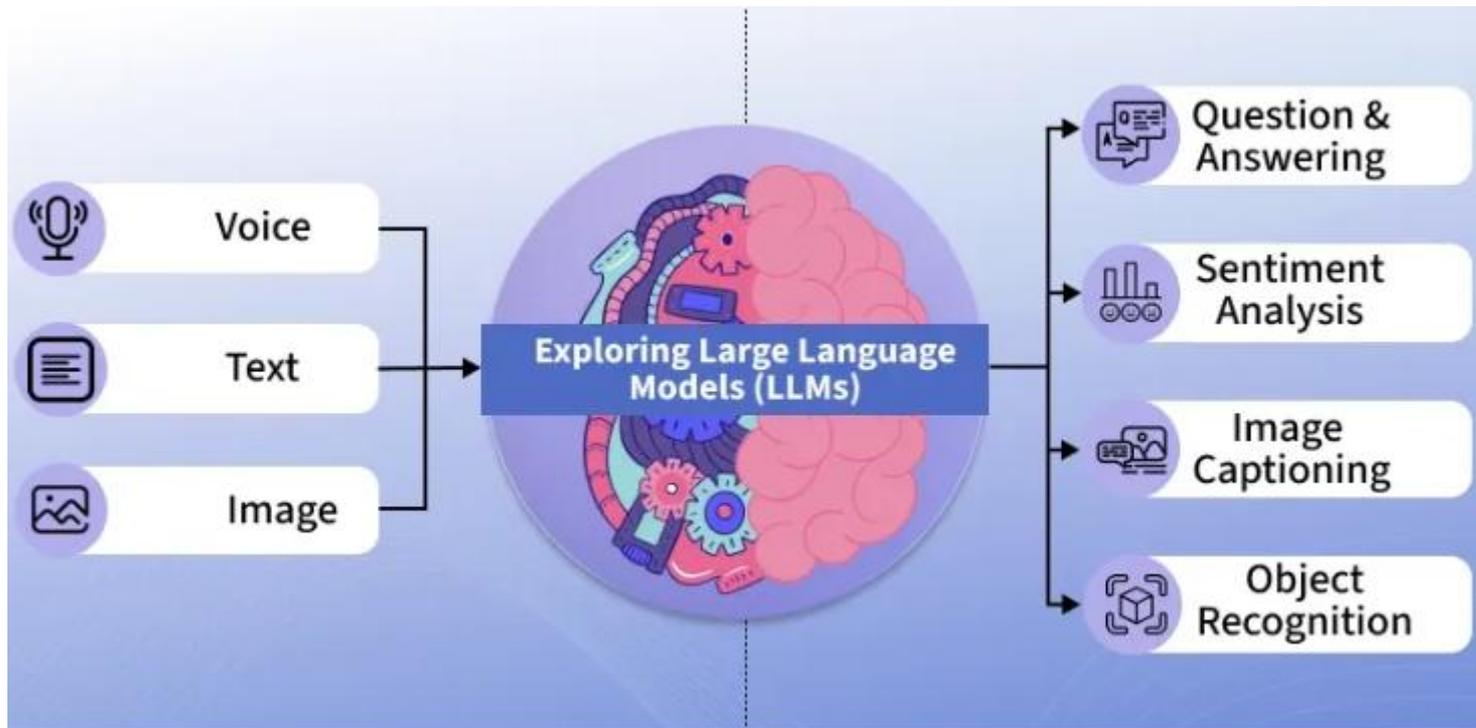
# Real Time Monitoring

# Data Science



Data science is an **interdisciplinary field** that combines scientific methods, algorithms, and systems to extract knowledge and actionable insights from both structured and unstructured data. It serves as the foundation for modern **Artificial Intelligence (AI)** and **machine learning**, allowing organizations to make informed, data-driven decisions rather than relying on intuition.

# LLM

Large Language Models (LLMs) are advanced AI systems built on deep neural networks designed to process, understand and generate human-like text
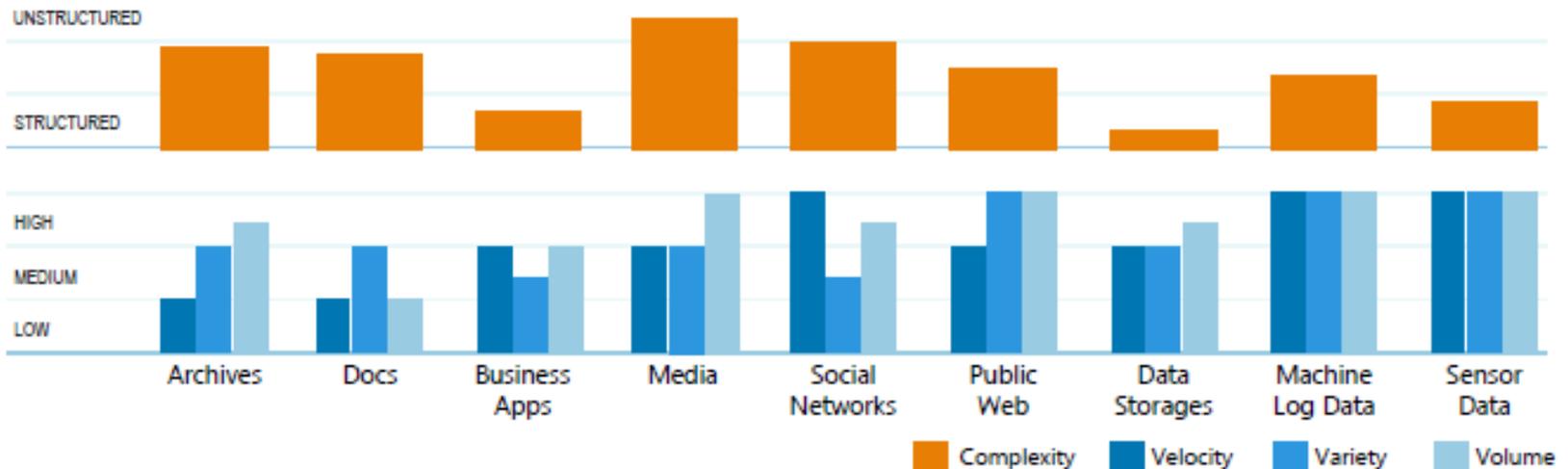
# LLM

# Big Data Analytics

# Big Data Challenges



**Archives**
Scanned documents, statements, medical records, e-mails etc..

**Docs**
XLS, PDF, CSV, HTML, JSON etc.

**Business Apps**
CRM, ERP systems, HR, project management etc.

**Media**
Images, video, audio etc.

**Social Networks**
Twitter, Facebook, Google+, LinkedIn etc.

**Public Web**
Wikipedia, news, weather, public finance etc

**Data Storages**
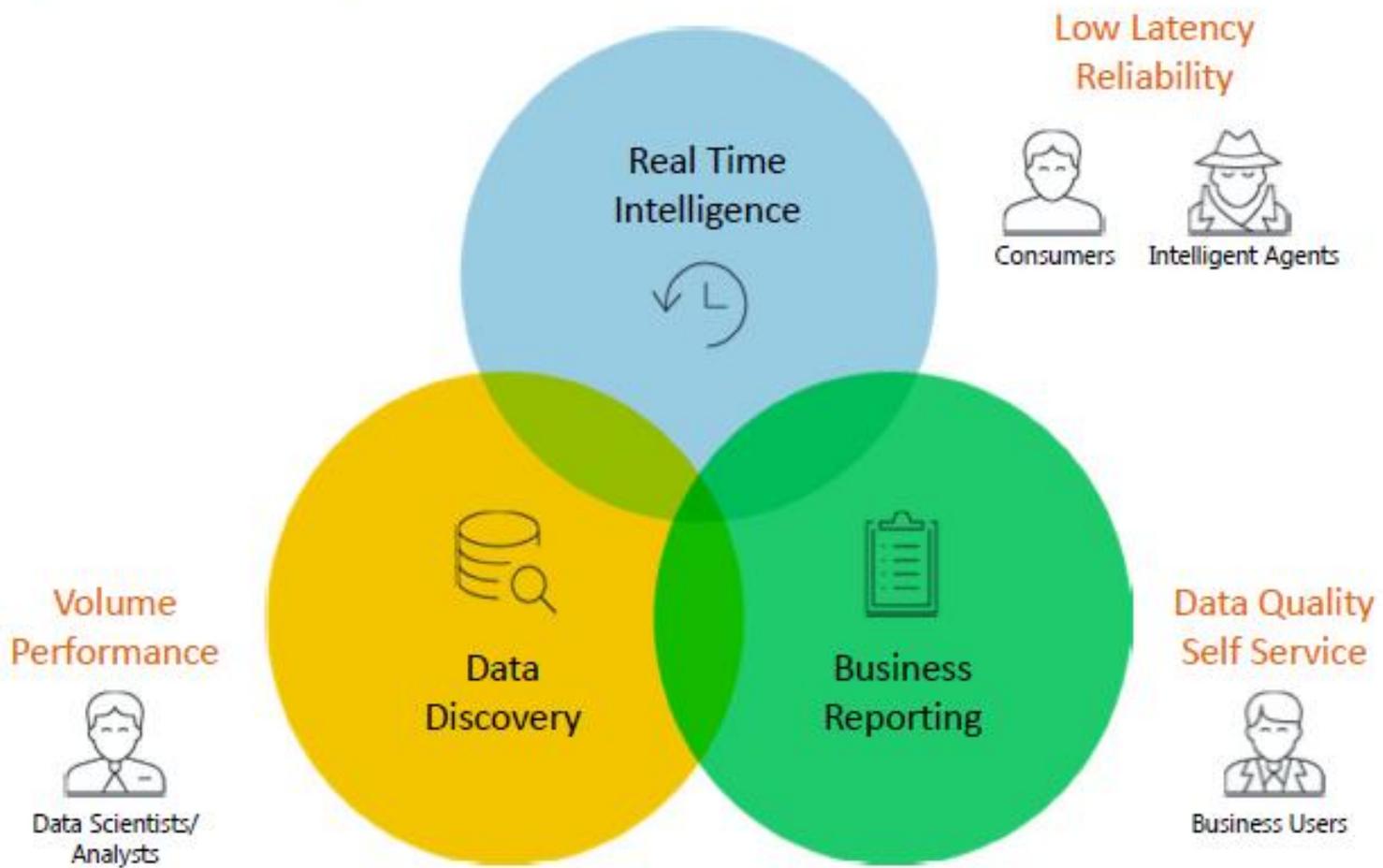RDBMS, NoSQL, Hadoop, file systems etc.

**Machine Log Data**
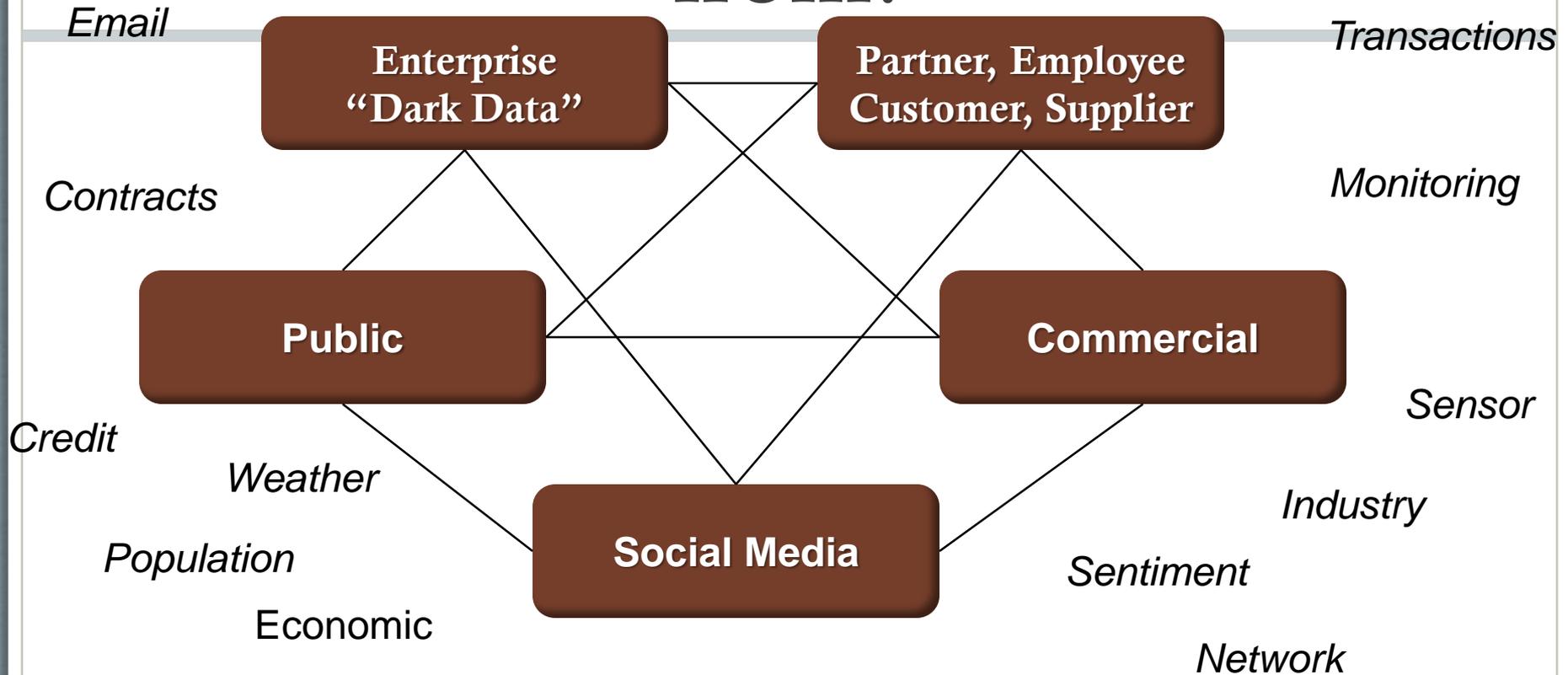Application logs, event logs, server data, CDRs, clickstream data etc.

**Sensor Data**
Smart electric meters, medical devices, car sensors, road cameras etc.

# Big Data Analytics Use Cases

# Where does Big Data come from?

Email

Transactions

**Enterprise "Dark Data"**

**Partner, Employee Customer, Supplier**

Contracts

Monitoring

**Public**

**Commercial**

Credit

Sensor

Weather

Industry

Population

**Social Media**

Sentiment

Economic

Network

# Characteristics of Big Data:
## 1-Scale (Volume)

- **Data Volume**
  - 44x increase from 2009 2020
  - From 0.8 zettabytes to 35zb
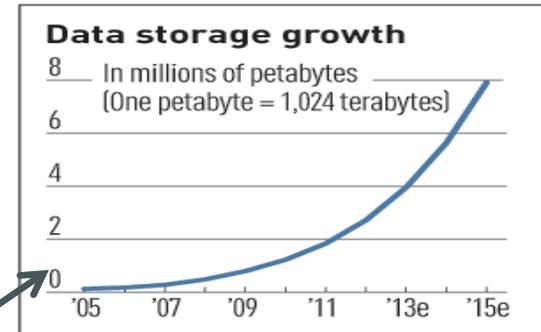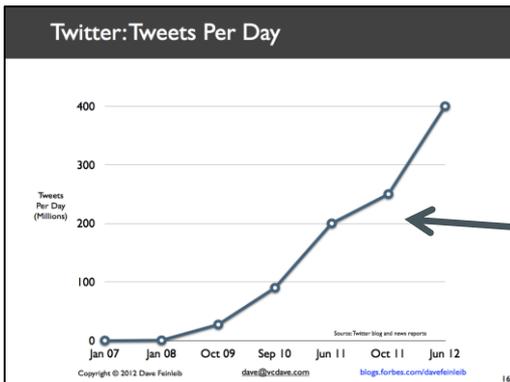
- Data volume is increasing exponentially

*Exponential increase in collected/generated data*

# Characteristics of Big Data:
## 2-Complexity (Variety)

- Various formats, types, and structures

- Text, numerical, images, audio, video, sequences, time series, social media data, multi-dim arrays, etc…

- Static data vs. streaming data

- A single application can be generating/collecting many types of data

To extract knowledge➔ all these types of data need to linked together

# Characteristics of Big Data:
## 3-Speed (Velocity)

- Data is begin generated fast and need to be processed fast

- Online Data Analytics

- Late decisions ➜ missing opportunities

- **Examples**
  - **E-Promotions:** Based on your current location, your purchase history, what you like ➜ send promotions right now for store next to you

  - **Healthcare monitoring:** sensors monitoring your activities and body ➜ any abnormal measurements require immediate reaction

# Big Data: 3V's



BIG DATA?

VOLUME
Large amounts of data.

VELOCITY
Needs to be analyzed quickly.

VARIETY
Different types of structured and unstructured data.



Complexity

Big Data

Speed          Volume

## Big Data = Transactions + Interactions + Observations



**BIG DATA**

| Petabytes | Sensors / RFID / Devices | | User Generated Content |
| | Mobile Web | Sentiment | Social Interactions & Feeds |
| | User Click Stream | | Spatial & GPS Coordinates |

**WEB**

| Terabytes | Web logs | A/B testing | External Demographics |
| | Offer history | Dynamic Pricing | Business Data Feeds |
| | | Affiliate Networks | HD Video, Audio, Images |

**CRM**

| Gigabytes | | Segmentation | Search Marketing | |
| | | Offer details | Behavioral Targeting | Speech to Text |
| | ERP Purchase detail | Customer Touches | | Product/Service Logs |
| Megabytes | Purchase record | Support Contacts | Dynamic Funnels | SMS/MMS |
| | Payment record | | | |

Increasing Data Variety and Complexity

**Source**: Contents of above graphic created in partnership with Teradata, Inc.

# Some Make it 4V's

| Volume | Velocity | Variety | Veracity* |
|---|---|---|---|
| **Data at Rest** | **Data in Motion** | **Data in Many Forms** | **Data in Doubt** |
| Terabytes to exabytes of existing data to process | Streaming data, milliseconds to seconds to respond | Structured, unstructured, text, multimedia | Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations |

# Volume



**40 ZETTABYTES**
[ 43 TRILLION GIGABYTES ]
of data will be created by 2020, an increase of 300 times from 2005

2005
2020

**6 BILLION PEOPLE**
have cell phones

**WORLD POPULATION: 7 BILLION**

**Volume**
**SCALE OF DATA**

It's estimated that
**2.5 QUINTILLION BYTES**
[ 2.3 TRILLION GIGABYTES ]
of data are created each day

Most companies in the U.S. have at least
**100 TERABYTES**
[ 100,000 GIGABYTES ]
of data stored

- Volume
  - Petabytes, exabytes of data
  - Volumes too great for typical DBMS

# Volume - Bytes Defined

| | Managerial Definition | Exact Amount | To Put It in Perspective |
|---|---|---|---|
| **1 Terabyte (TB)** | One trillion bytes | $2^{40}$ bytes | Printed collection of the Library of Congress = 20 TB |
| **1 Petabyte (PB)** | One quadrillion bytes | $2^{50}$ bytes | eBay data warehouse (2010) = 10 PBC. Monash, "eBay Followup—Greenplum Out, Teradata > 10 Petabytes, Hadoop Has Some Value, and More," October 6, 2010. Note eBay plans to increase this value 2.5 times by the end of 2011. |
| **1 Exabyte (EB)** | One quintillion bytes | $2^{60}$ bytes | |
| **1 Zettabyte (ZB)** | One sextillion bytes | $2^{70}$ bytes | Amount of data consumed by U.S. households in 2008 = 3.6 ZB |

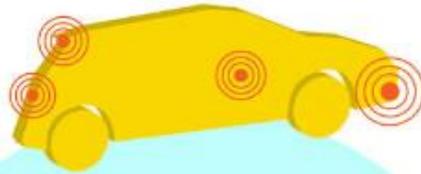Megabyte:  $2^{20}$ bytes or, loosely, one million bytes

Gigabyte:   $2^{30}$ bytes or, loosely one billion bytes

# Velocity

The New York Stock Exchange captures **1 TB OF TRADE INFORMATION** during each trading session

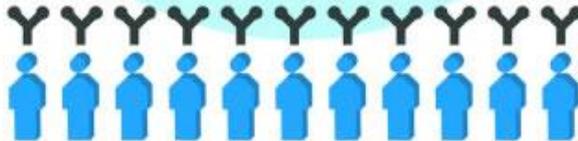Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure

**Velocity**

**ANALYSIS OF STREAMING DATA**

By 2016, it is projected there will be **18.9 BILLION NETWORK CONNECTIONS** – almost 2.5 connections per person on earth

- Velocity
  - Massive amount of streaming data

# Variety



As of 2011, the global size of data in healthcare was estimated to be

**150 EXABYTES**

[ 161 BILLION GIGABYTES ]

By 2014, it's anticipated there will be

**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO**

are watched on YouTube each month

**Variety**

**DIFFERENT FORMS OF DATA**

**30 BILLION PIECES OF CONTENT**

are shared on Facebook every month

**400 MILLION TWEETS**

are sent per day by about 200 million monthly active users

- Variety
  - Massive sets of unstructured/semi-structured data from Web traffic, social media, sensors, and so on

# Which source of data represents the most immediate opportunity?



| Existing underutilized "dark data" | More detail from customers, suppliers, etc. | Social media content | Commercially available data | Publicly available data |
|---|---|---|---|---|
| 38% | 38% | 16% | 4% | 4% |

Source: *Getting Value from Big Data*, Gartner Webinar, May 2012

# Big Data Opportunities

**Making better informed decisions**
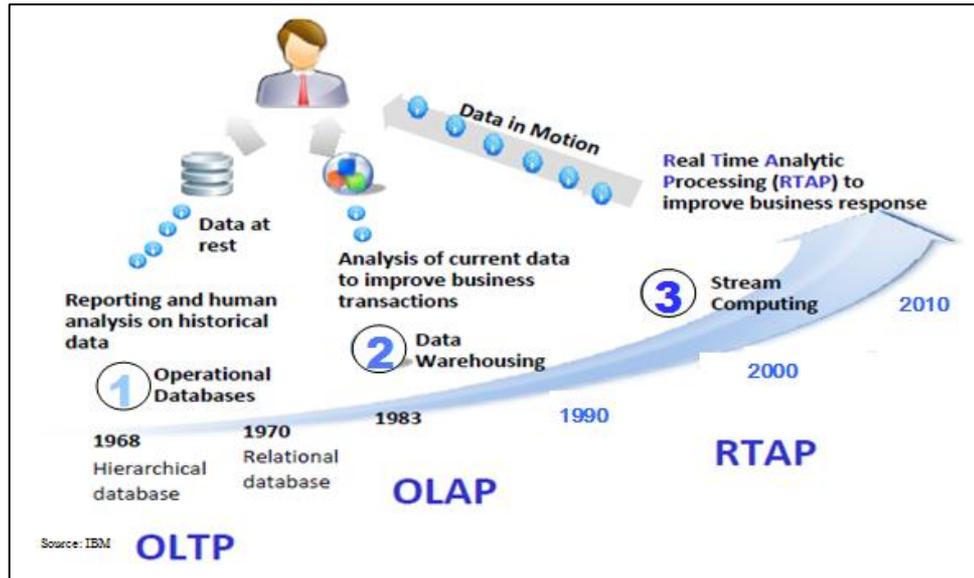e.g. strategies, recommendations

**Discovering hidden insights**
e.g. anomalies forensics, patterns, trends

**Automating business processes**
e.g. complex events, translation

# Harnessing Big Data



- **OLTP:** Online Transaction Processing   (DBMSs)

- **OLAP:** Online Analytical Processing   (Data Warehousing)

- **RTAP:** Real-Time Analytics Processing  (Big Data Architecture & technology)

# Who's Generating Big Data



**Social media and networks**
(all of us are generating data)

**Scientific instruments**
(collecting all sorts of data)

**Mobile devices**
(tracking all objects all the time)

**Sensor technology and networks**
(measuring all kinds of data)

- The progress and innovation is no longer hindered by the ability to collect data

- But, by the ability to manage, analyze, summarize, visualize, and discover knowledge from the collected data in a timely manner and in a scalable fashion

# The Model Has Changed…

- **The Model of Generating/Consuming Data has Changed**

**Old Model:** Few companies are generating data, all others are consuming data



**New Model:** all of us are generating data, and all of us are consuming data

# What's driving Big Data



COMPLEXITY

HIGH

Predictive Analytics
and Data Mining

Business
Intelligence

LOW                    BUSINESS VALUE                    HIGH

- Optimizations and predictive analytics
- Complex statistical analysis
- All types of data, and many sources
- Very large datasets
- More of a real-time

- Ad-hoc querying and reporting
- Data mining techniques
- Structured data, typical sources
- Small to mid-size datasets

# Value of Big Data Analytics

- Big data is more real-time in nature than traditional DW applications

- Traditional DW architectures (e.g. Exadata, Teradata) are not well-suited for big data apps
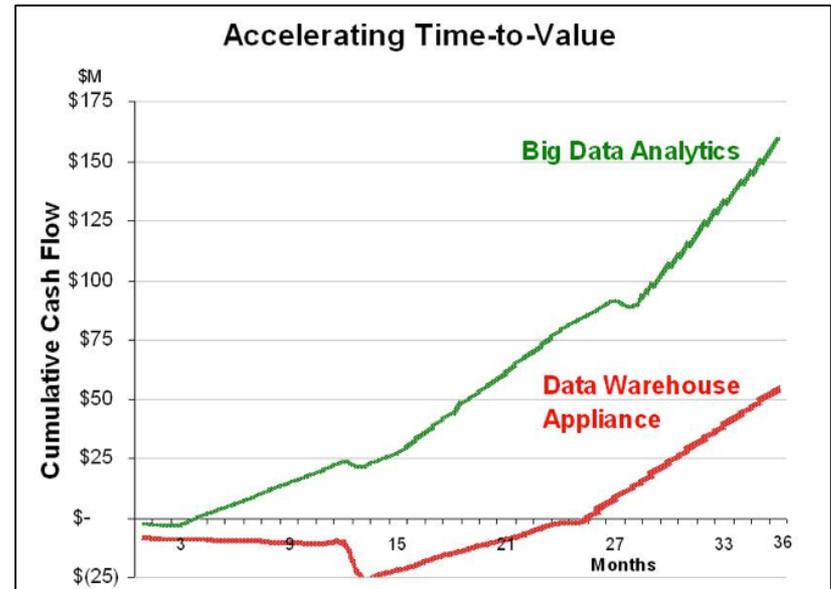
- Shared nothing, massively parallel processing, scale out architectures are well-suited for big data apps



Accelerating Time-to-Value

Cumulative Cash Flow ($M) vs Months. Big Data Analytics (green) rises faster than Data Warehouse Appliance (red).

# Data Discovery: Non-Relational Architecture

# Business Reporting: Hybrid Architecture

# Lambda Architecture



*Source:*

# Architectural Decisions

Architecture Drivers:

- **Volume (45 TB)**
- Sources (Semi-structured - JSON)
- **Throughput (> 20K/sec)**
- Latency (1 hour)
- **Extensibility (Custom tags)**
- Data Quality (Not critical)

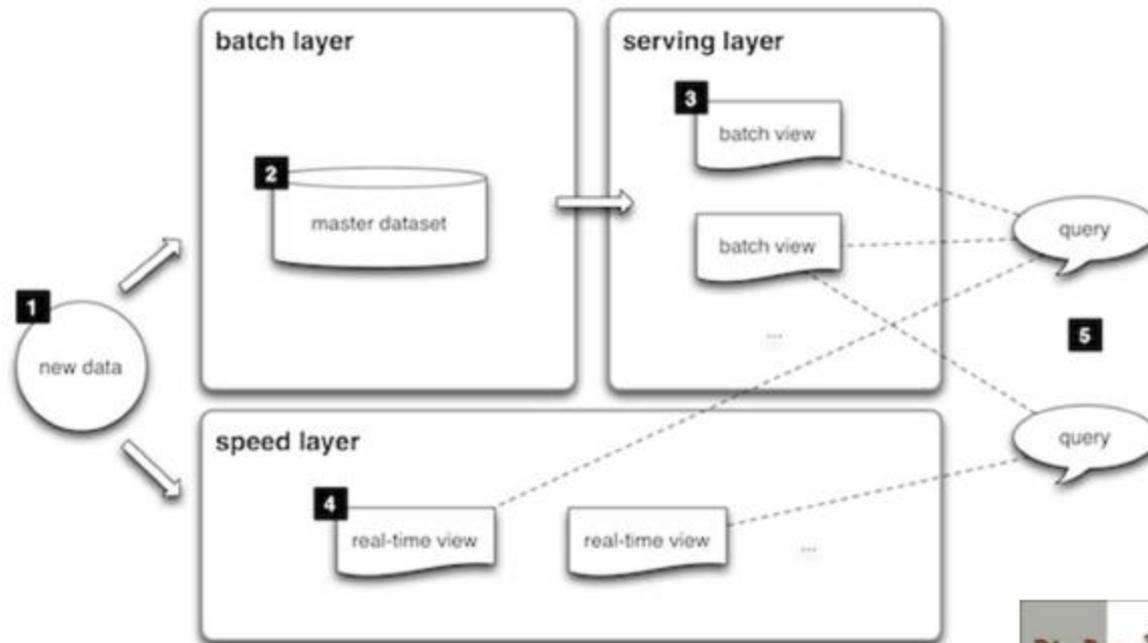- Reliability (24/7)
- Security (Multitenancy)
- **Self-Service (Canned reports, Data science)**
- Cost (The less the better ☺)
- Constraints (Public Cloud)

Trade-off:

| | Extended Relational | Non-Relational |
|---|---|---|
| Volume/Scalability | +/- | + |
| Throughput | + | + |
| Self-Service | + | +/- |
| Extensibility | - | + |

✓ **Non-Relational Architecture**
✓ Reporting via **Materialized View** pattern

# Solution Architecture

**Technologies:**
- Amazon S3
- Flume
- Hadoop/HDFS, MapReduce
- HBase
- Oozie
- Hive

# Tips for Designing Big Data Solutions

- ❑ Understand data users and sources
- ❑ Discover architecture drivers
- ❑ Select proper reference architecture
- ❑ Do trade-off analysis, address cons
- ❑ Map reference architecture to technology stack
- ❑ Prototype, re-evaluate architecture
- ❑ Estimate implementation efforts
- ❑ Set up devops practices from the very beginning
- ❑ Advance in solution development through "small wins"
- ❑ Be ready for changes, big data technologies are evolving rapidly

# India Stack

| | | |
|---|---|---|
| **CONSENT LAYER** | Provides a modern privacy data sharing framework | Open Personal Data Store |
| **CASHLESS LAYER** | Game changing electronic payment systems and transition to cashless economy | IMPS, AEPS, APB, and UPI |
| **PAPERLESS LAYER** | Rapidly growing base of paperless systems with billions of artifacts | e-KYC, E-sign, Digital Locker |
| **PRESENCE-LESS LAYER** | Unique digital biometric identity with open access of nearly a Billion users | Authentication |

# India Stack

# Identifying Insurance Fraud

- ## Opportunity
  - Save and make money by reducing fraudulent auto insurance claims

- ## Data & Analytics
  - Predictive analytics against years of historical claims and coverage data
  - Text mining adjuster reports for hidden clues, e.g. missing facts, inconsistencies, changed stories

- ## Results
  - Improved success rate in pursuing fraudulent claims from 50% to 88%; reduced fraudulent claim investigation time by 95%
  - Marketing to individuals with low propensity for fraud

# Identifying Insurance Fraud

**INFINITY.**

*What  **"dark data" is just laying around that can transform business processes?*

**Operational **data** that is not being used. Consulting and market research company Gartner Inc. describes **dark data** as "information assets that organizations collect, process and store in the course of their regular business activity, but generally fail to use for other purposes."

# Quality Improvement

- **Opportunity**
  - Move from manual to automated inspection of burger buns production to ensure and improve quality

- **Data & Analytics**
  - Photo-analyze over 1000 buns-per-minute for color, shape and seed distribution
  - Continually adjust ovens and process automatically

# Quality Improvement

- Result
  - Eliminate 1000s of pounds of wasted product per year; speed production; save energy; Reduce manual labor costs

*Is the company using all of its "senses" to observe, measure and optimize business processes?*

# Improving Corporate Image

- Opportunity
  - Improve reputation, brand and buzz by tapping social media

- Data & Analytics
  - Continually scanning twitterverse for mentions of their business
  - Integrating tweeters with their robust customer management system
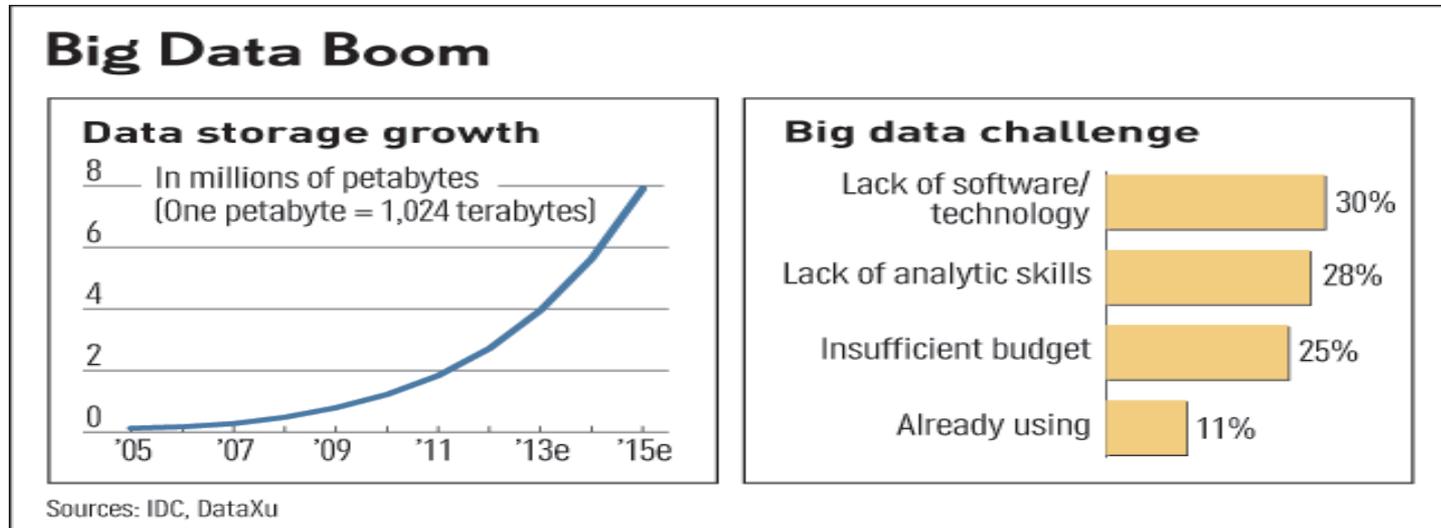
# Improving Corporate Image

- Results
  - Saw tweet from a top customer lamenting late flight—no time to dine at Morton's
  - Tuxedo-clad waiter waiting for him when he landed with a bag containing his favorite steak, prepared the way he normally likes it with all the fixin's

*How can the company listen, analyze and respond in real-time?*

# Challenges in Handling Big Data



**Big Data Boom**

Data storage growth
8 — In millions of petabytes
(One petabyte = 1,024 terabytes)
6
4
2
0
'05 '07 '09 '11 '13e '15e

Big data challenge
Lack of software/technology — 30%
Lack of analytic skills — 28%
Insufficient budget — 25%
Already using — 11%

Sources: IDC, DataXu

- **The Bottleneck is in technology**
  - New architecture, algorithms, techniques are needed

- **Also in technical skills**
  - Experts in using the new technology and dealing with big data

# What Technology Do We Have

# For Big Data ??

# Data Mining Techniques...

- Classification [Predictive]

- Clustering [Descriptive]

- Association Rule Discovery [Descriptive]

- Sequential Pattern Discovery [Descriptive]

- Regression [Predictive]

- Deviation Detection [Predictive]

# Big Data Landscape

## Vertical Apps
PREDICTIVE POLICING
bloomreach. GET FOUND.
MYRRIX

## Log Data Apps
splunk> loggly sumologic

## Ad/Media Apps
rocketfuel
bluefin
Media Science
TURN
collective [i]
Recorded Future
LuckySort
DataXu
Data. Insight. Action.

## Data As A Service
factual.
GNIP DATASIFT Windows Azure Marketplace
INRIX LexisNexis SPACE CURVE
kaggle
knoema beta
LOQATE Everything Location

## Business Intelligence
ORACLE | Hyperion
SAP Business Objects RJMetrics
Microsoft | Business Intelligence
IBM COGNOS birst
Autonomy MicroStrategy
QlikView bime DOMO
Chart.io GoodData

## Analytics and Visualization
tableau Palantir
OPERA metaLayer
METAMARKETS dataspora centrifuge
TERADATA ASTER
SAS TIBCO KARMASPHERE
panopticon Real-Time Visual Data Analysis
Datameer pentaho
platfora ClearStory CIRRO
alteryx visual.ly AYATA

## Analytics Infrastructure
Hortonworks VERTICA An HP Company MAPR TECHNOLOGIES
cloudera INFOBRIGHT
ParAccel
EMC² GREENPLUM
NETEZZA kognitio
DATASTAX EXASOL calpont

## Operational Infrastructure
COUCHBASE 10gen the MongoDB company
TERADATA HADAPT
TERRACOTTA VoltDB
MarkLogic INFORMATICA

## Infrastructure As A Service
amazon web services
Windows Azure
infochimps
Google BigQuery

## Structured Databases
ORACLE MySQL
Microsoft SQL Server PostgreSQL
IBM DB2 SYBASE
memsql

## Technologies
hadoop
hadoop Map Reduce
mahout
APACHE HBASE
Cassandra

50

dave@vcdave.com
blogs.forbes.com/davefeinleib

# Big Data Technology



**Big Data:** The Moving Parts

Increasing Age & Maturity

**Fast Data:** Hadoop, Vertica, MapReduce, Esper, kdb, Greenplum, ETL, ECL, Netezza, Teradata

**Big Analytics:** Hive, SciPy, Mahout, MATLAB, Revolution R, SPSS, AMPL, SAS

**Deep Insight:** unsupervised learning, social media analytics, sentiment analysis, predictive modeling, BPO, BI, network analysis, visualization, simulation

**Business Objectives:** mass customization of services, quicker response to market trends, identifying real-time cost optimizations, faster, more accurate decision making, better and more holistic R&D, autonomic supply chain management

From http://blogs.zdnet.com/Hinchcliffe

the growth of data will be exponential for the foreseeable future

terabytes | petabytes | exabytes | zettabytes

the amount of data stored by the average company today

51

# Uncertainty of Data

**1 IN 3 BUSINESS LEADERS**

don't trust the information they use to make decisions

**27% OF RESPONDENTS**

in one survey were unsure of how much of their data was inaccurate

**Veracity**

**UNCERTAINTY OF DATA**

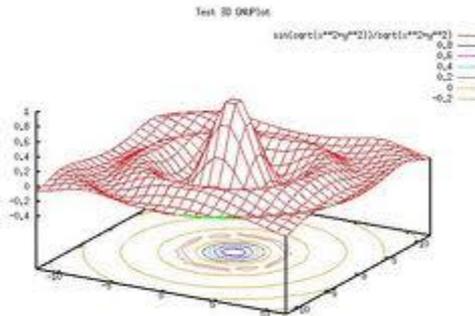Poor data quality costs the US economy around

**$3.1 TRILLION A YEAR**

# Analytics Models

# Descriptive Analytics

- Descriptive analytics, such as reporting/OLAP, dashboards, and data visualization, have been widely used for some time.
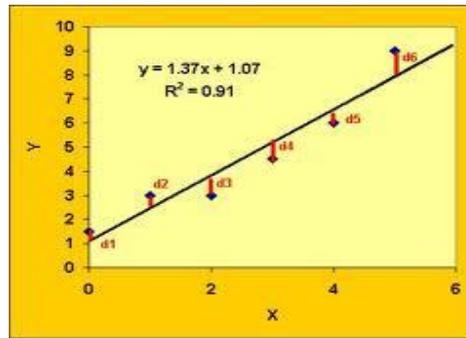- They are the core of traditional BI.
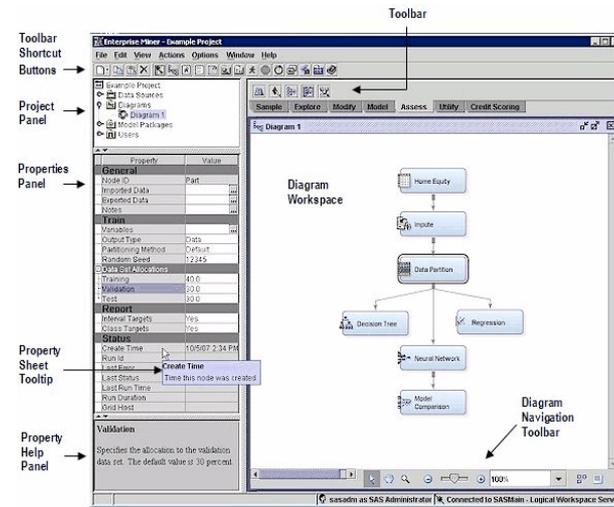


## What has occurred?

Descriptive analytics, such as data visualization, is important in helping users interpret the output from predictive and predictive analytics.

# Predictive Analytics

- Algorithms for predictive analytics, such as regression analysis, machine learning, and neural networks, have also been around for some time.
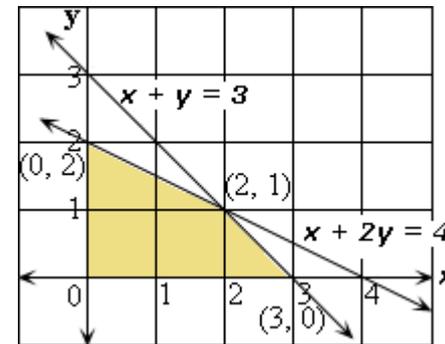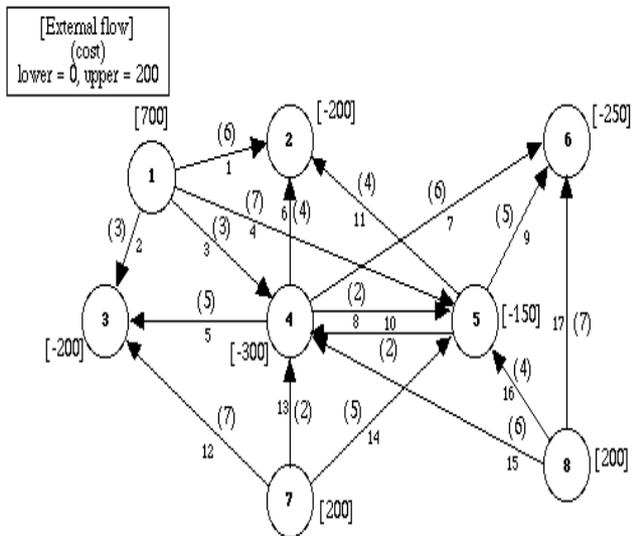


*What will occur?*

- Marketing is the target for many predictive analytics applications.
- Descriptive analytics, such as data visualization, is important in helping users interpret the output from predictive and prescriptive analytics.

# Prescriptive Analytics

- Prescriptive analytics are often referred to as advanced analytics.
- Often for the allocation of scarce resources
- Optimization

*What should occur?*

Prescriptive analytics can benefit healthcare strategic planning by using analytics to leverage operational and usage data combined with data of external factors such as economic data, population demographic trends and population health trends, to more accurately plan for future capital investments such as new facilities and equipment utilization as well as understand the trade-offs between adding additional beds and expanding an existing facility versus building a new one.

# What You Will Learn…

- We focus on *Hadoop/MapReduce technology*

- **Learn the platform (how it is designed and works)**
  - How big data are managed in a scalable, efficient way

- **Learn writing Hadoop jobs in different languages**
  - Programming Languages: Java, C, Python
  - High-Level Languages: Apache Pig, Hive

- **Learn advanced analytics tools on top of Hadoop**
  - RHadoop: Statistical tools for managing big data
  - Mahout: Data mining and machine learning tools over big data

- **Learn state-of-art technology from recent research papers**
  - Optimizations, indexing techniques, and other extensions to Hadoop

# Course Output: What You Will Learn…

- We focus on *Hadoop/MapReduce technology*

- **Learn the platform (how it is designed and works)**
  - How big data are managed in a scalable, efficient way

- **Learn writing Hadoop jobs in different languages**
  - Programming Languages: Java, C, Python
  - High-Level Languages: Apache Pig, Hive

- **Learn advanced analytics tools on top of Hadoop**
  - RHadoop: Statistical tools for managing big data
  - Mahout: Analytics and data mining tools over big data

- **Learn state-of-art technology from recent research papers**
  - Optimizations, indexing techniques, and other extensions to Hadoop