# THE AGE OF BIG DATA

Asst. Prof. Natawut Nupairoj, Ph.D.

Mobile Application and System Services Research Group

Department of Computing Engineering

Chulalongkorn University

natawut.n@chula.ac.th

**Chula DataScience**

CHULA ƎNGINEERING
Foundation toward Innovation

**"Data is a new class of economic asset, like currency and gold"** - World Economic Forum

**"By 2020, total digital information will be 40ZB or 5.2TB for each human"** - IDC

**"Just 3% of all data is currently tagged and ready for manipulation, and only one sixth of this - 0.5% - is used for analysis"** - IDC

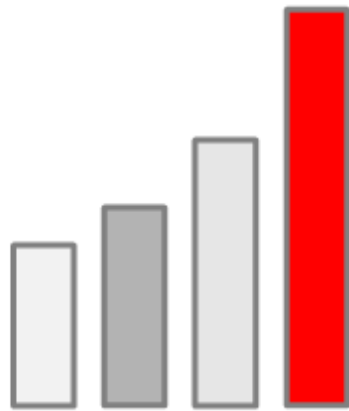B | KB | MB | GB | TB | PB | EB | ZB

# BIG DATA MOVEMENTS

Big Data Analytic worldwide market will reach $50 billion by 2018

Big data technology and services will grow worldwide at annual growth rate of 40% – about seven times that of the ICT market overall

There are at least 23 US University offering Master's Programs in Data Science

By 2018, there will be 140,000–190,000 data scientist job postings that go unfulfilled
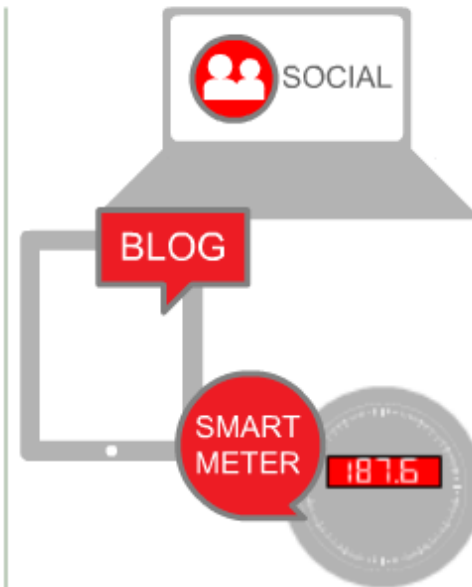
# "3VS" MODEL OF BIG DATA BY GARTNER (2012)



VOLUME | VELOCITY | VARIETY

# VOLUME

One-fifth of organizations store more than 1 petabyte of data

100TBs of data are uploaded to facebook each day

Walmart collects more than 2.5PBs from 1 million customers each hour

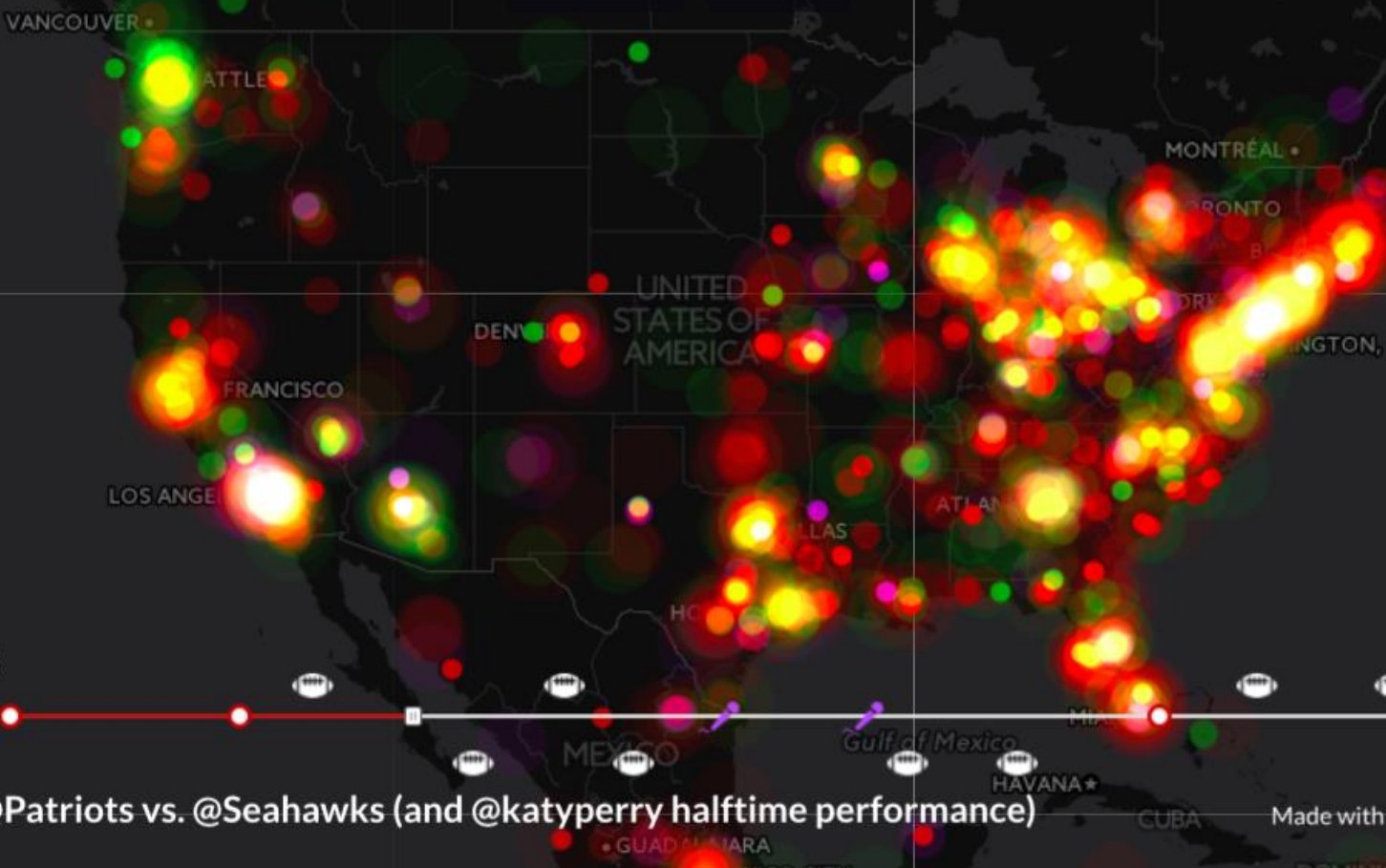Real-time data processing aka. Fast Data becomes common

Large Hadron Collider has about 150 million sensors delivering data 40 million times per second

Google handles 2m search queries per minute

Superbowl 2015 final stats: 28.4M tweets with a peak of 395k tweets-per-minute

@PATRIOTS **14** **7** @SEAHAWKS

18:15    22:05

#SB49: @Patriots vs. @Seahawks (and @katyperry halftime performance)

Made with CartoDB

Chula DataScience

# VARIETY

The most important and most difficult to handle among 3Vs – unstructured data accounts more than 80% of corporate data

Unstructured data: customer feedback, social network comments, twitter messages, location information, email, photos, movies, documents, presentations, etc.

Patient information consists of medical records, medical images, clinical data, etc.
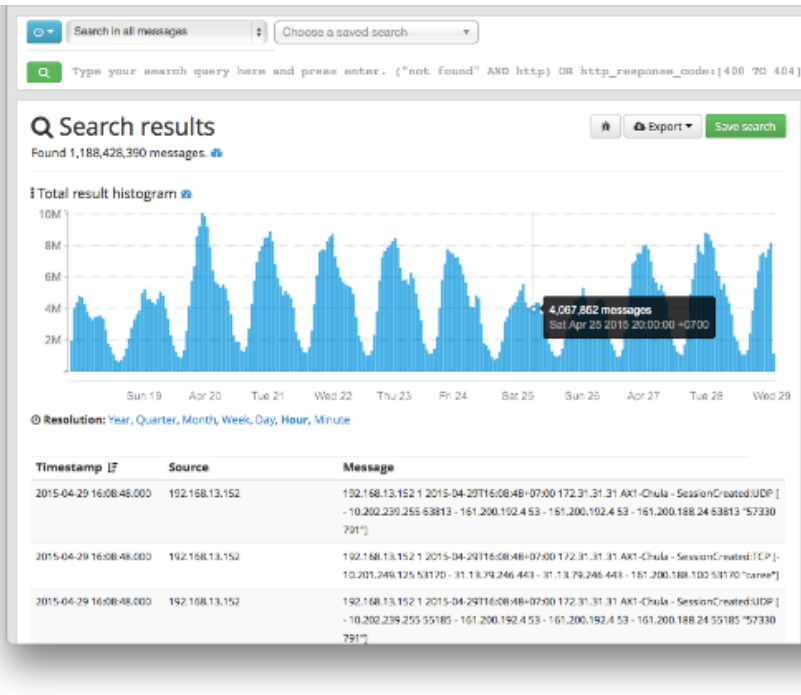
# SIMPLE EXAMPLE



**Chula IT collects various logs from 40 servers, specifically for computer crime act 2007**

Velocity = 2,000 events/sec

Volume = approx. 30+ GBs each day (or 3TBs for 3 months)

Variety = Different data sources, different data schema

# BIG DATA'S DRIVERS
# MOBILE PHONES AND PERSONAL DEVICES

World mobile penetration rate is more than 95%

Thailand's rate is 147% (smartphone = 49%)

Facebook has 798M daily active mobile users

Future: wearable computers

# BIG DATA'S DRIVERS
# INTERNET OF THINGS AND SENSOR NETWORKS



RFID, GPS, sensors, surveillance cameras, smart meters, appliances, medical devices, toys, etc.

212b sensors will be available with 30b IOT are expected to be connected to networks by 2020

# BIG DATA'S DRIVERS
# USER GENERATED CONTENTS AND CROWDSOURCING



Blogging, reviewing commenting, forum, digital video, podcasting, mobile phone photography, social networking, crowdsourcing, etc.

Highly influential to consumer behavior and also enable the study of consumer behavior

Generate lots of both structured and unstructured data
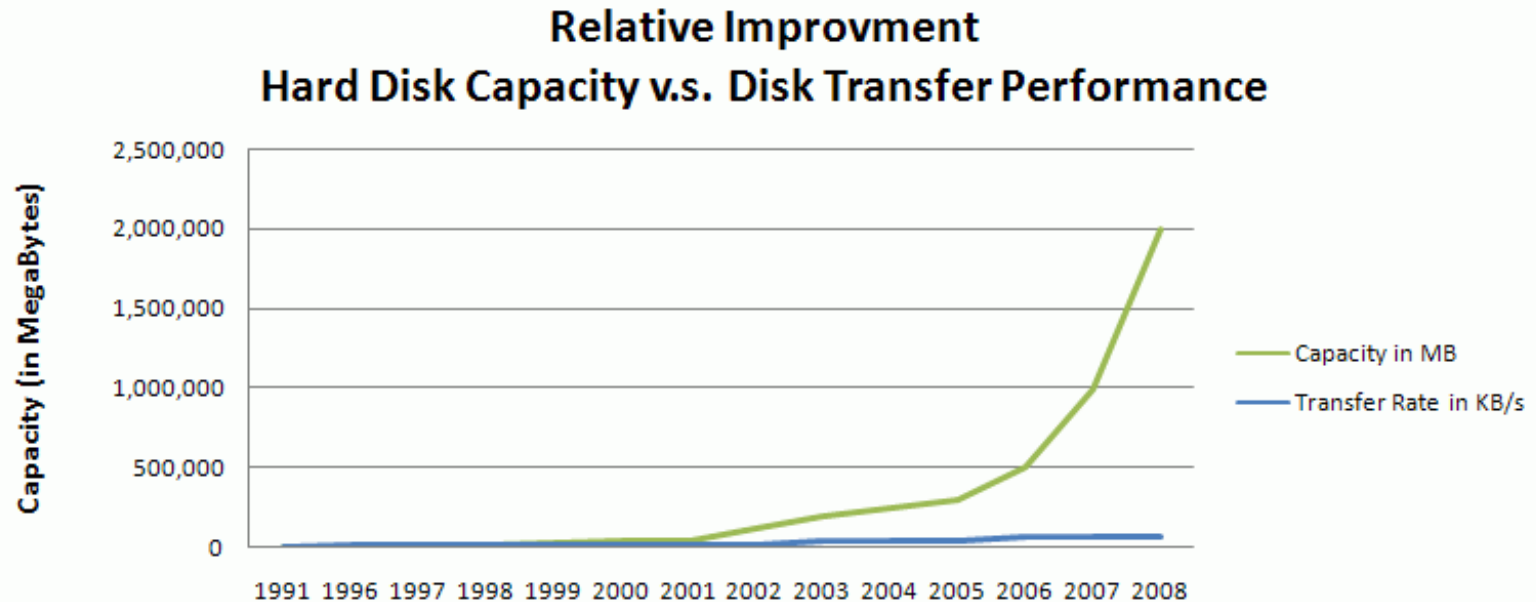
# BIG DATA'S DRIVERS
# CLOUD COMPUTING



Deliver computing services over a network

Evolution of technology, but revolution of economy

One of Big Data accelerators: significant big data sources and enabling platform for big data processing

# HOW CAN WE RESPOND TO THESE CHALLENGES?



## Relative Improvment
## Hard Disk Capacity v.s. Disk Transfer Performance

Source: http://wiki.r1soft.com/pages/viewpage.action?pageId=3016608

Disk data transfer improves very slowly

Increase throughput with **scale-out** – processing data on multiple disks in parallel

Opensource software framework inspired by Google Search Engine Architecture

Provide easy-to-program scale-out foundation for data-intensive applications on large clusters of commodity hardware

Hadoop File System (HDFS) has been widely used

Users: Yahoo!, Facebook, Amazon, eBay, American Airline, Apple, Google, HP, IBM, Microsoft, Netflix, New York Times, etc.

Products: IBM InfoSphere BigInsights, Google App Engine, Oracle Big Data Appliance, Microsoft HDInsight

In-Memory Data Processing from UC Berkeley

Extend MapReduce model to support batch executions, interactive queries, and stream processing

Support various languages (Java, Python, Scala, R) with built-in analytic libraries (machine learning, graph processing)

Strong and growing community

High performance, based on sorting benchmarks, Spark is 10x – 100x faster than Hadoop

# NOSQL – NOT ONLY SQL



Special DBMS for large data that does not require relational model e.g. unstructured data

Various types: Document Store, Graph, Key-Value store, etc.

Products: Parquet, Cassandra, HBASE, ElasticSearch, Accumulo, DynamoDB, Redis, Riak, CouchDB, MangoDB, Neo4j, etc.

CHULA ƎNGINEERING
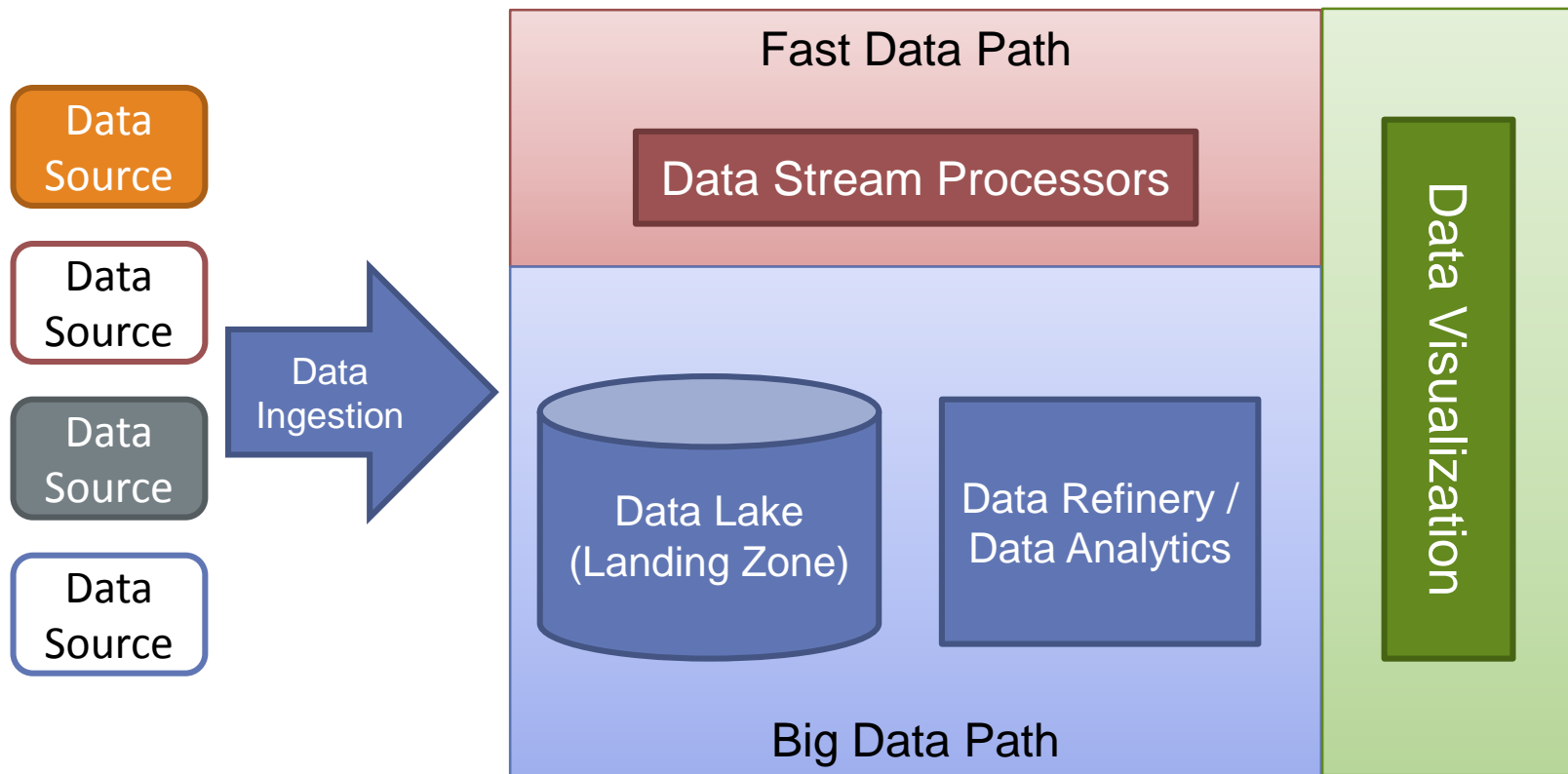Foundation toward Innovation

# PREDICTIVE ANALYTICS



Analyze current and historical data to automatically find patterns based on several techniques e.g. statistics, modeling, machine learning, data mining, time series analysis, deep learning, etc.

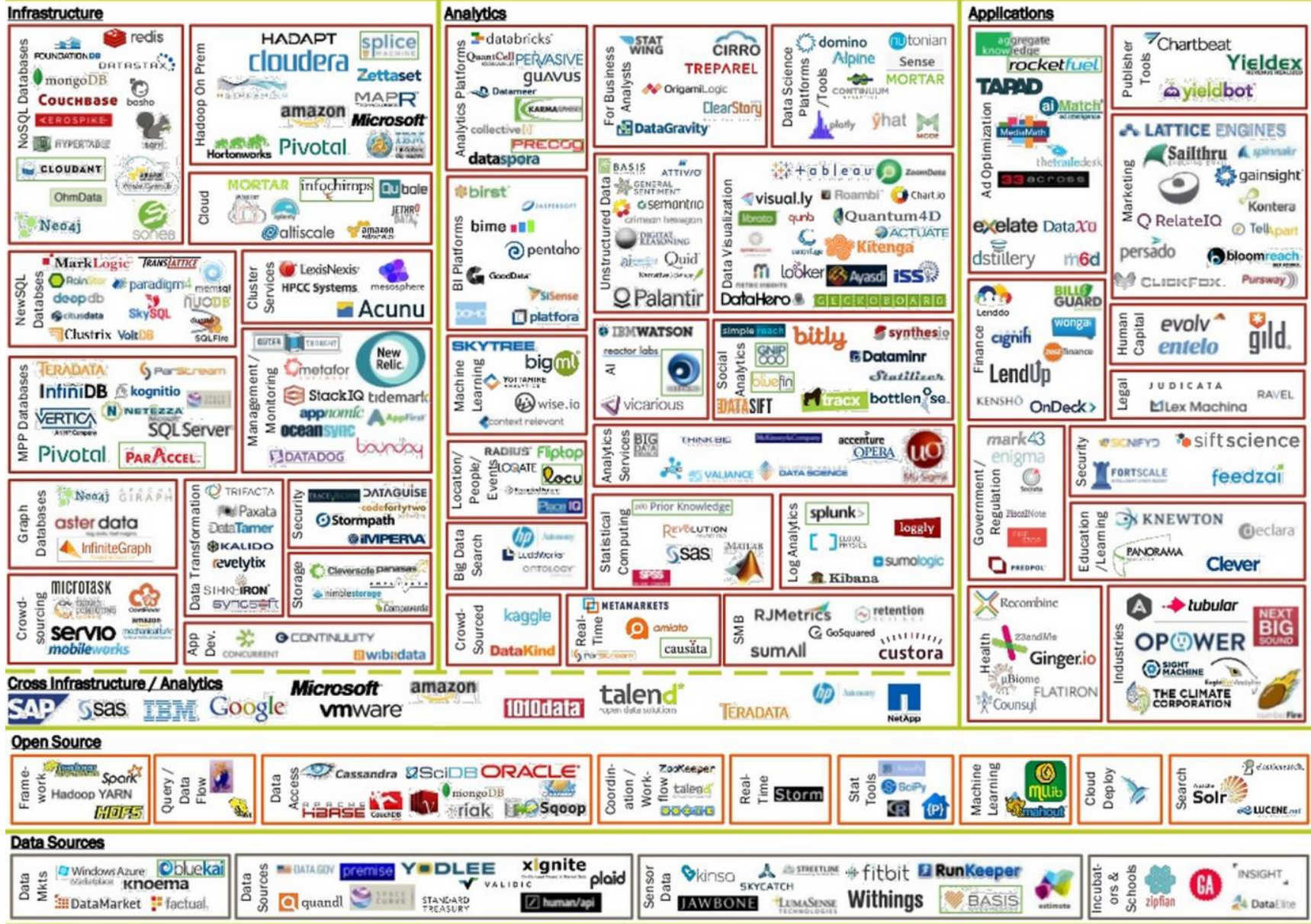Utilize other techniques e.g. text analytics, image processing, location analytics, etc.

Applications: Micro Customer Segmentation, Sentiment Analysis, Customer retention, Fraud detection, etc.

# PUTTING THEM ALL TOGETHER

BIG DATA LANDSCAPE, VERSION 3.0

© Matt Turck (@mattturck), Sutian Dong (@sutiandong) & FirstMark Capital (@firstmarkcap)

Chula DataScience

CHULA ENGINEERING
Foundation toward Innovation

# TYPICAL USE CASES

**Bigger / Faster / More Up-to-Date Data Warehouse**

**Product Recommendation**

**Social Listening**

**Fraud Detection and Risk Management**

**Micro Customer Segmentation**
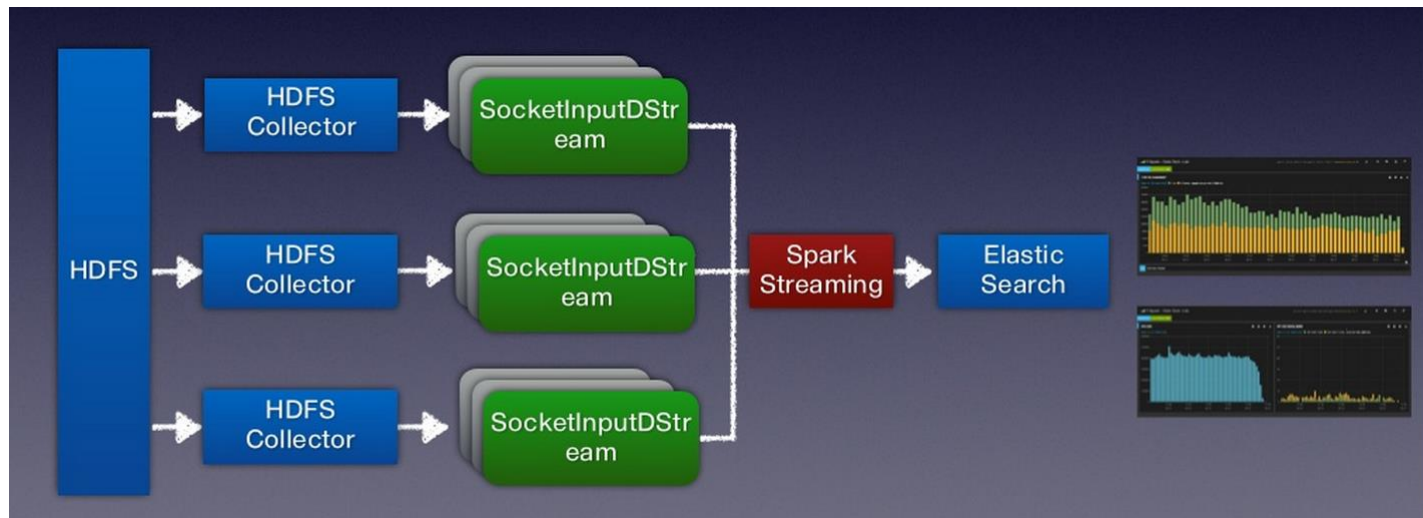
**Demand Sensing for Supply Chain**

**Bio Informatics**

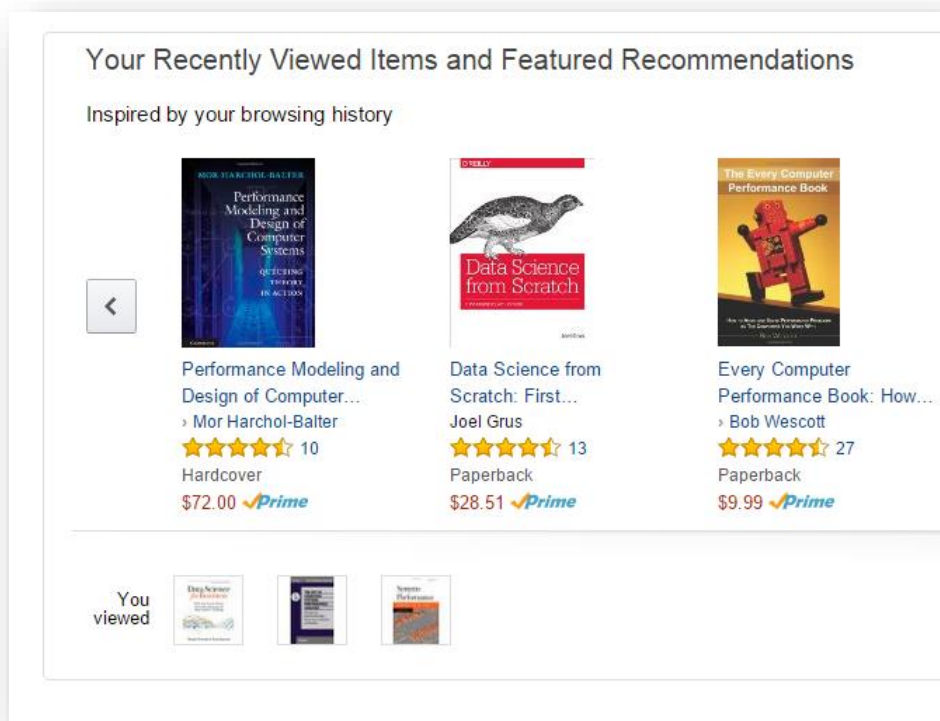# CASE STUDY:
# SK TELECOM'S USAGE PATTERN ANALYSIS

Process usage data from 28 millions subscribers: 40TB/day – 15PB total

Must process data with 530MB/sec or 1 million records/sec

Use Hadoop, Spark, and ElasticSearch to provide mobile usage pattern analytics with low latency ad-hoc query (< 2 secs)

# CASE STUDY:
# AMAZON'S RECOMMENDATION ENGINE



Amazon mines data from 152 million customers to suggest products to customers

Use Hadoop + DynamoDB to perform collaborative filtering, click-stream analysis, historical purchase data analytics

Other similar offerings: Facebook, LinkedIn, Netflix

# USHAHIDI

2007
Kenya

2010
Haiti
Chile
Washington DC
Russia

2011
Christchurch
Middle East
India
Japan
Australia
US
Macedonia

2012
Balkans

2014
Kenya

**Chula DataScience**

CHULA ƎNGINEERING
Foundation toward Innovation

# CASE STUDY:
# JP MORGAN
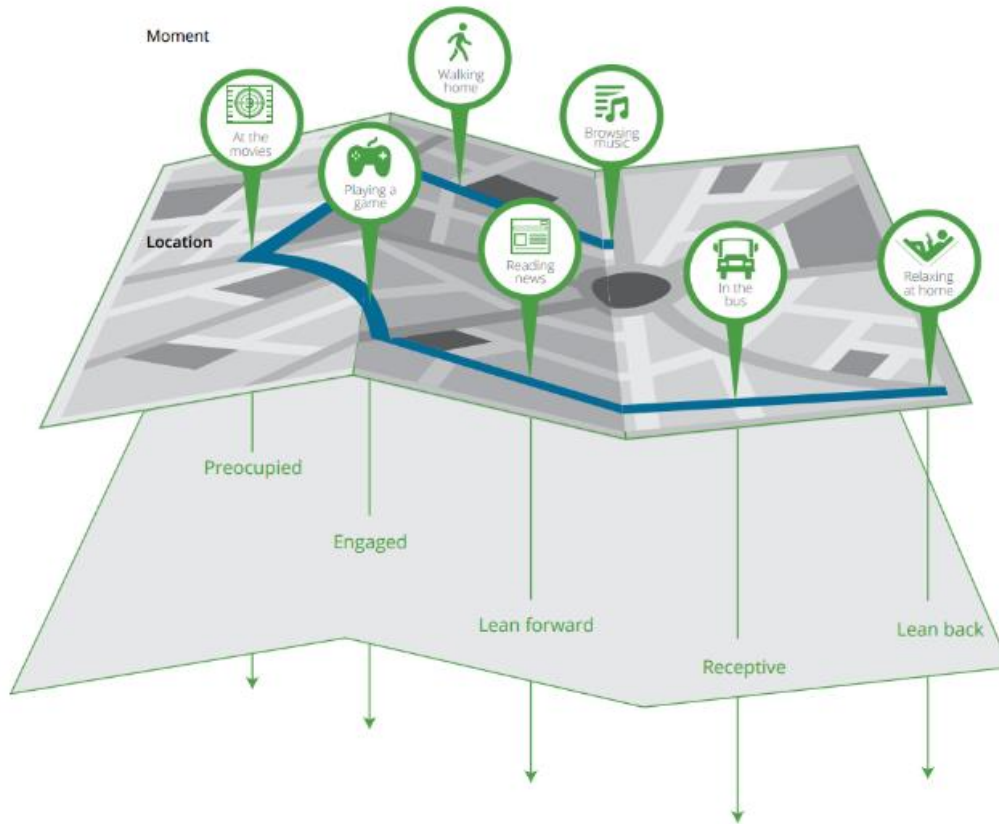


In 2013, JP Morgan Chase & Co use Big Data to aggregate all available information about a single customer

Data included monthly balances, credit card transactions, credit bureau data, demographic data

This allowed bank to offer lower interest rates by reducing credit card fraud

# CASE STUDY:
# INMOBI'S TARGETED MARKETING



User behaviour changes dramatically across work, home, commute, and other location contexts

Geo context targeting: create customer micro segmentation from customer's location activities, time of day, and app being used

Analyze 250b monthly activities of 750m customers for 65m POI in 50+ countries - - 500TB/day

Chula **DataScience**

# CASE STUDY:
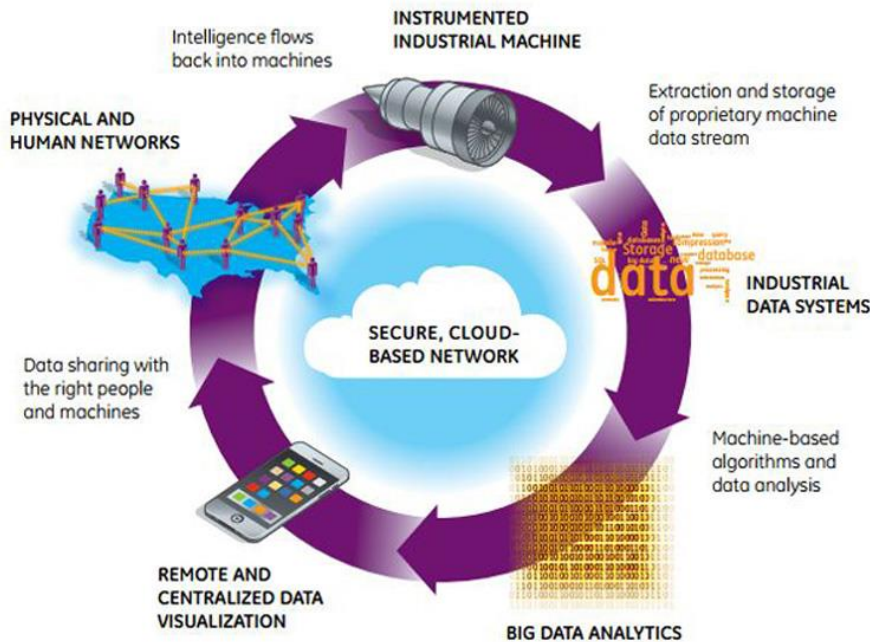# UBER'S DYNAMIC PRICING FARES



Uber's entire business model is based on the very Big Data principle of crowd sourcing

"dynamic pricing" fares are calculated automatically, using GPS, street data, demand forecast, and predictive algorithms

Due to traffic conditions in New York on New Year's Eve 2011, the fare of journey of one mile rose from $27 to $135

CHULA ƩNGINEERING
Foundation toward Innovation

# CASE STUDY:
# GE'S SMART MACHINES



GE has launched Industrial Internet initiative

Jet engine has 20 sensors generating 5,000 data samples per second

Data can be used for fuel efficiency and service improvements

"In the future it's going to be digital. By the time the plane lands, we'll know exactly what the plane needs."

# CASE STUDY:
# SOCIAL NETWORK-BASED THAI MONITOR CORPUS



Cooperation with Department of Linguistics

Construct Thai Monitor Corpus by extracting data from social networks e.g. Twitter, forums, blogs, websites, etc.

Provide corpus repository with temporal analytic capabilities that can keep expanding

Target: 500M words

# OTHER USE CASES

Automatic – Data Processing for IoT

Baidu – Business insights data analytics

Cisco – Entity resolution for Whois clean-up

eBay – Log Transaction Analytics

ESRI – Geospatial and temporal analysis

Goldman Sachs – Data curation processing

NBC Universal – Multimedia data distribution prediction

Novartis – Genomic data analytic

Salesforce – Recommendation engine

Shopify – E-commerce transaction data analytics

Tresata – Anti-Money Laundering real-time graph analytic

Yahoo – Real-time personalization engine

CHULA ƐNGINEERING
Foundation toward Innovation

**"Big data is about having the technology and people with the appropriate analysis skills to allow firms to make sense of huge volumes of data in an affordable manner."**

*Source: Forrester Research, 2012*

# DATA SCIENTIST IS A TEAM SPORT



Data Scientist

Data Processing

Statistical Research

Domain Knowledge

Computer Science

Math & Statistics

Machine Learning

CHULA ƐNGINEERING
Foundation toward Innovation

# OUR CAPABILITIES

## Big Data Architecture (6+ faculties)

- Map-Reduce and Hadoop
- In-Memory Processing
- NoSQL
- Data Ingestion
- Social Network Information Retrieval
- Internet of Things
- Embedded System
- Wireless Sensor Network
- Cloud Computing
- Mobile Computing

# OUR CAPABILITIES

## Big Data Analytics (7+ faculties)

- Machine Learning
- Clustering / K-Mean
- Data Mining
- Neuron Networks
- Social Network Analytics
- Text Analytics
- Image Processing
- Time-Series Analytics
- Location Analytics
- Data Visualization

CHULA ƎNGINEERING
Foundation toward Innovation

# COME ALONG DATA-DRIVEN ECONOMY

In July 2014, the European Commission outlined a new strategy on Big Data, supporting and accelerating the transition towards a data-driven economy in Europe

In Feb 2015, The White House appointed the first US chief data scientist

As of today, US Government's open data publishes more than 150,000 datasets to the public

# FINAL THOUGHTS

**We are definitely in the age of Big Data**
**(for at least 2 years)**

**Big Data is more than just "Big" data**

**Next Battlefield … Data-Driven Economy**

**Are you ready?**

CHULA ƎNGINEERING
Foundation toward Innovation