# Introduction to Big Data

Slides from Prof. Dr. Thoai Nam, Lecturer: Nguyen Quang Hung

High Performance Computing Lab (HPC Lab)

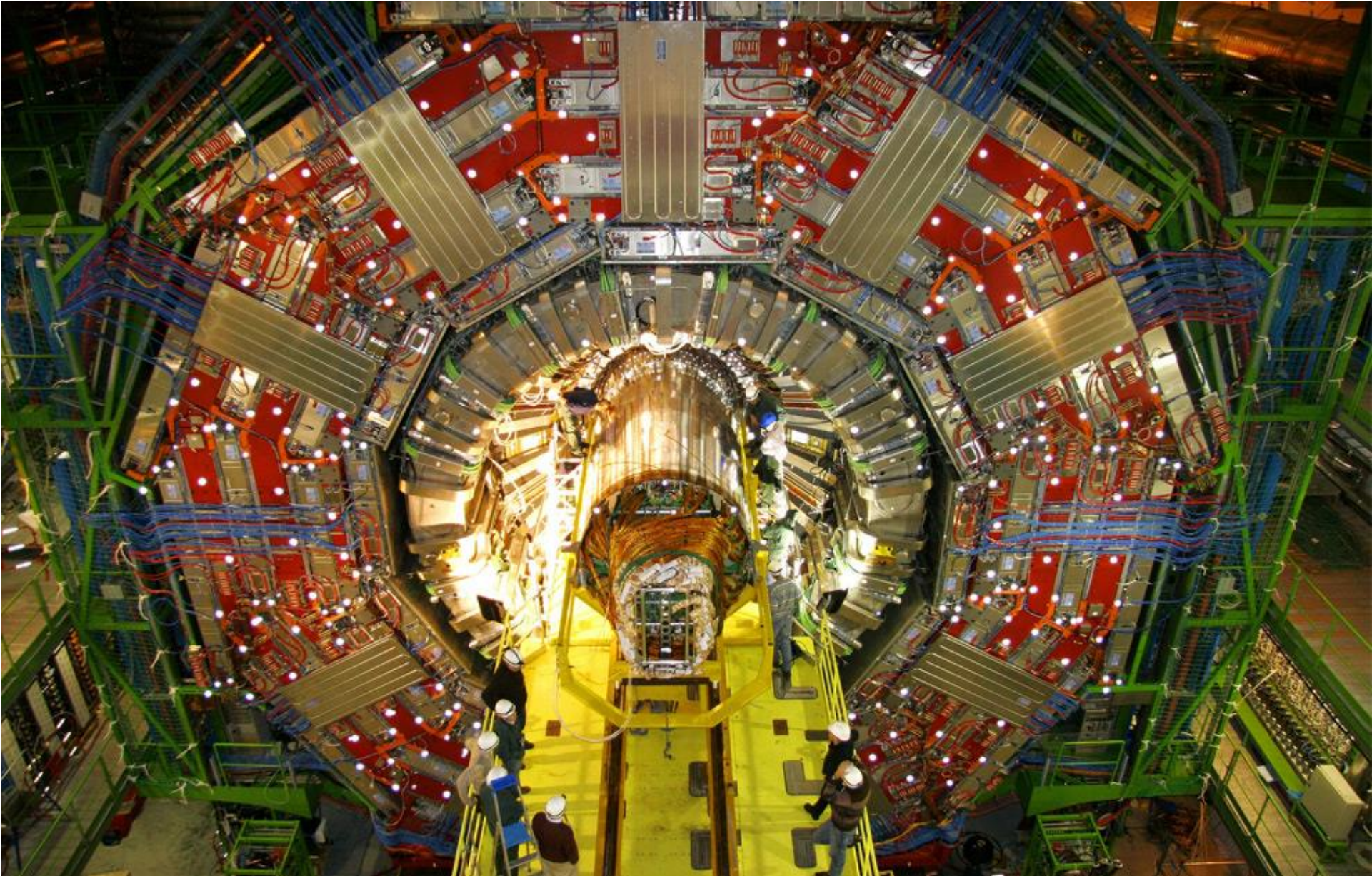Faculty of Computer Science and Technology

HCMC University of Technology

# What's Big Data?

"Data **too large & complex** to be effectively handled by standard database technologies currently founded in most organizations"

"Data whose **scale**, **diversity** and **complexity** require new **architectures**, **techniques**, **algorithms** and **analytics** to manage it and extract value and hidden knowledge from it"

# Data sources?

# CERN's Large Hydron Collider (LHC) generates 15 PB a year

# The Earthscope

- The Earthscope is the world's largest science project. Designed to track North America's geological evolution, this observatory records data over 3.8 million square miles, amassing 67 terabytes of data.

- It analyzes seismic slips in the San Andreas fault, sure, but also the plume of magma underneath Yellowstone and much, much more.

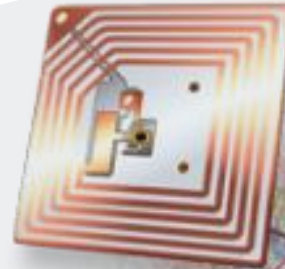- (http://www.msnbc.msn.com/id/44363598/ns /technology_and_science-future_of_technology/#.TmetOdQ--uI)



Annual budget: $25,000,000
Construction cost: $197,000,000
Staff: 110
Physical size: 3.8 million square miles
Scientific utility: 10
WIIFY: 10
Wow factor: 10

**12+ TBs**
of tweet data
every day

**? TBs** of
data every day

**25+ TBs** of
log data
every day

**30 billion** RFID
tags today
(1.3B in 2005)

**4.6
billion**
camera
phones
world wide

**100s of
millions
of GPS
enabled**
devices
sold
annually

**2+
billion**
people on
the Web
by end
2011

**76 million** smart
meters in 2009…
200M by 2014

**2015**

IoTs – Big Data - AI

100B sensors

**2005**

Mobile – Cloud

2.5B smart phones

**1980**

PC - Internet

1B users

# IoT and Services



Figure 4:
The Internet of Things and Services – Networking people, objects and systems

**Internet of People** $10^6$-$10^8$

Social Web

**CPS-platforms**

Smart Grid

Smart Factory

Smart Home

Smart Building

Business Web

**Internet of Things** $10^7$-$10^9$

**Internet of Services** $10^4$-$10^6$

Source: Bosch Software Innovations 2012

# Smart cities

# 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

- Data volume increases exponentially over time
- 33 ZB in 2018 to 175 ZB in 2025



**Big Data**

- 12. DomegemegrotteB $=10^{33}$ Byte
- 11. ShilentnoB $=10^{30}$ Byte
- 10. XenottaB $=10^{27}$ Byte
- 9. YottaB $=10^{24}$ Byte
- 8. ZettaB$=10^{21}$ Byte — 1 ZB =Global Internet Traffic 2016
- 7. ExaB$=10^{18}$ Byte — 5 EB =All of words ever spoken by one person
- 6. PetaB$=10^{15}$ Byte — 10 PB =All USA City Library
- 5. TeraB$=10^{12}$ Byte — 10 TB =One USA City Library
- 4. GB$=10^{9}$ Byte — 1 GB= One Library
- 3. MB$=10^{6}$ Byte — 1 MB= One Small Book
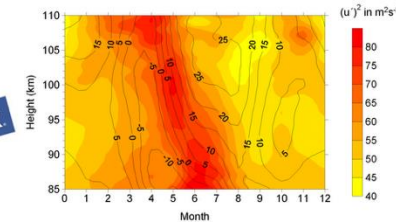- 2. KB$=10^{3}$ Byte — 1 KB= One Small Story
- 1. Bytes = 8 Bits — 1 Byte=One Character

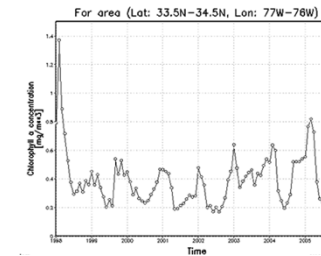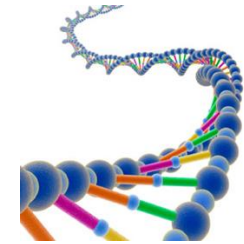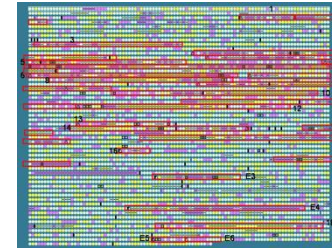Figure 1 – Annual Size of the Global Datasphere



**Annual Size of the Global Datasphere**

175 ZB

# Variety (Complexity)

- Various formats, types and structures
  - Numerical data, image data, audio, video, text, time series
  - Relational Data (Tables/Transaction/Legacy Data)
  - Text Data (Web)
  - Semi-structured Data (XML)
  - Graph Data
    - Social Network, Semantic Web (RDF), …

- A single application can be generating/collecting many types of data
  - Heterogeneous data
  - Complex data integration problem.

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

# A single view to the customer



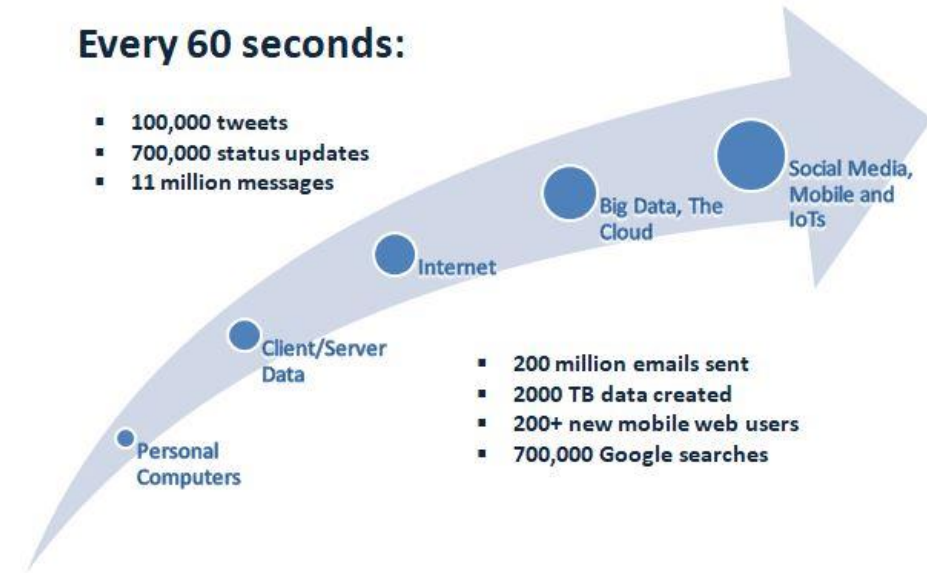a Single Customer View '*is an aggregated, consistent and holistic representation of the data known by an organization about its customers*'

# Velocity (Speed)

**Every 60 seconds:**

- 100,000 tweets
- 700,000 status updates
- 11 million messages

Social Media, Mobile and IoTs

Big Data, The Cloud

Internet

Client/Server Data

Personal Computers

- 200 million emails sent
- 2000 TB data created
- 200+ new mobile web users
- 700,000 Google searches

- 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

# Real-time/Fast 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.

# Veracity

Data quality

# Value



Example: US economy

Size of bubble indicates relative contribution to GDP

Big data: ease-of-capture index[1] (High / Low)

Big data: value potential index[1] (Low / High)

- Utilities
- Natural resources
- Manufacturing
- Health care providers
- Computers and other electronic products
- Information
- Finance and insurance
- Professional services
- Transportation and warehousing
- Accommodation and food
- Real estate
- Management of companies
- Construction
- Wholesale trade
- Administrative services
- Retail trade
- Other services
- Educational services
- Arts and entertainment
- Government
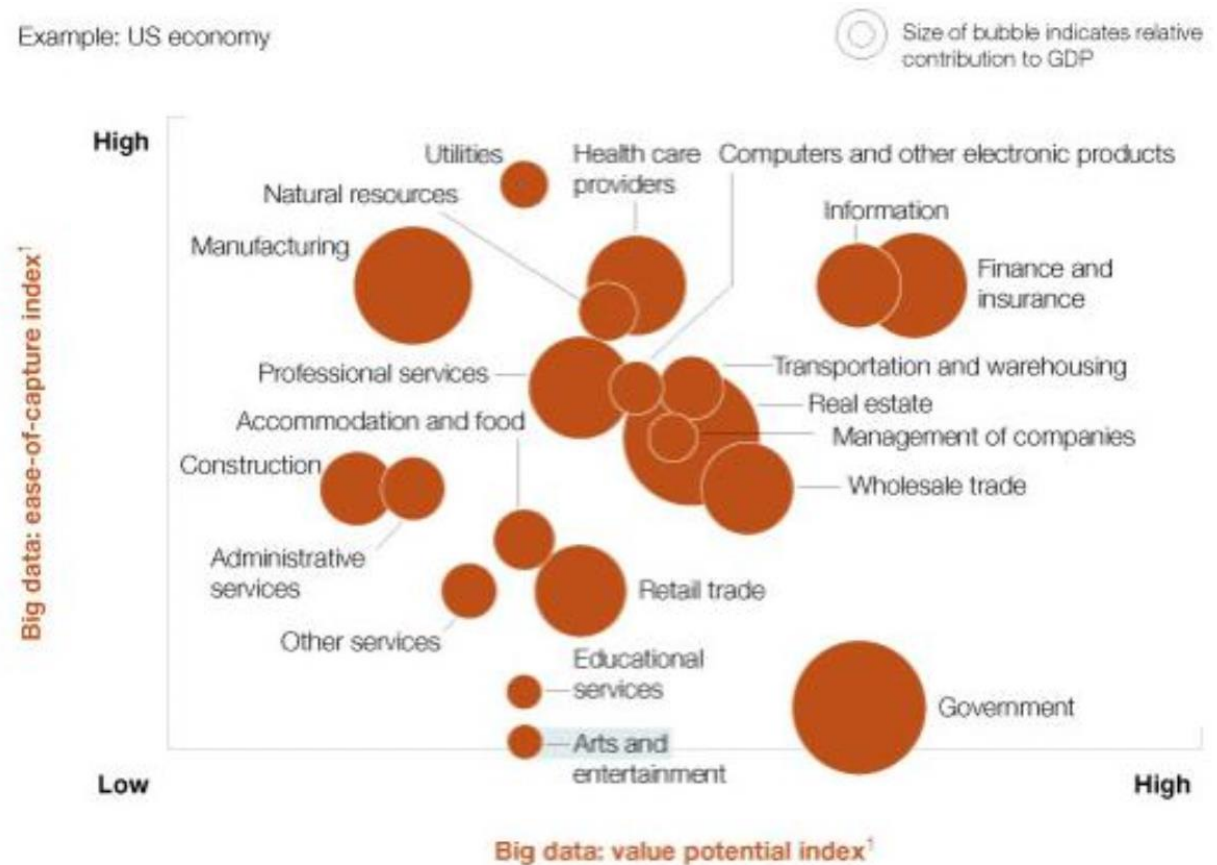
➢ Translate data into business advantage.

# Answering tough questions

- **Problem**
  - sales for lollipops are going down
- **Data**
  - all sales data by customer, region, time, ...
- **Information**
  - lollipops bought by people older than 25
  - (but eaten by people younger than 10)
- **Knowledge**
  - moms believe: lollipops = bad teeth
- **Value**
  - dentists advertise your lollipops

# Why is this difficult?

- You need more data than your data warehouse
  - you need more data that you have
  - logs, Twitter feeds, blogs, customer surveys, ...
- You need to ask the right questions
  - data alone is silent
- You need technology and organization that help you concentrate on asking the right questions.

# What is Big Data?

- Three alternative perspectives
  - Philosophical
  - Business
  - Technical
- (Ultimately, it is a buzz word for everybody.)

# Philosophical

- What is more valuable, if you had to pick one?
  - experience or intelligence?

- Traditional (computer) science: **logic!** [intelligence]
  - understand the problem, build model / algorithm
  - answer question from implementation of model

- New science: **statistics!** [experience]
  - collect data
  - answer question from data (what did others do?)

# Data Science, 4<sup>th</sup> Paradigm

- New approach to do science
  - Step 1: Collect data
  - Step 2: Generate Hypotheses
  - Step 3: Validate Hypotheses
  - Step4: (Goto Step 1 or 2)

- Why is this a good approach?
  - it can be automated: no thinking, less error
- Why is this a bad approach?
  - how do you debug without a ground truth?

# Is bigger = smarter?

- **Yes!**
  - tolerate errors
  - discover the long tail and corner cases
  - machine learning works much better
- **But!**
  - more data, more error (e.g., semantic heterogeneity)
  - with enough data you can prove anything
  - still need humans to ask right questions

# What is Big Data?

- Business Perspective
  - it is a new business model
- People pay with data
  - e.g. Facebook, Google, Twitter:
    - use service, give data
    - Google sells your data to advertisers • (you pay advertisers indirectly)
  - e.g, Amazone
    - pay service + give data
    - sells data and uses data to improve service

# Business Perspective

- Bank
  - keeps your money securely (kind of...)
  - puts your money at work (lends it to others), interest
  - you keep ownership of money and take it when needed
- Databank
  - keeps your data securely (kind of...)
  - puts your data at work: interest or better service
  - (you keep ownership of data: hopefully to come)

# Technical Perspective (?)

- You collect all data
  - the more the better -> statistical relevance, long tail
  - keeping all is cheaper than deciding what to keep
- You decide independently what to do with data
  - run experiments on data when question arises
- Huge difference to traditional information systems
  - design upfront what data to keep and why!!!
  - (e.g., waterfall model of software engineering!)

# Big data value chain (1)

| Generation | Acquisition | Storage | Analysis |

- ## Generation
  - ○ Passive recording
    - ▪ Typical structured data
    - ▪ Bank trading transactions, shopping records, government sector archives
  - ○ Active generation
    - ▪ Semi-structured or unstructured data
    - ▪ User-generated content, e.g., social networks
  - ○ Automatic production
    - ▪ Location-aware, context-dependent, highly mobile data
    - ▪ Sensor-based Internet-enabled devices.

# Big data value chain (2)

Generation › Acquisition › Storage › Analysis

- Acquisition
  - Collection
    - Pull-based, e.g., web crawler
    - Push-based, e.g., video surveillance, click stream
  - Transmission
    - Transfer to data center over high capacity links
  - Preprocessing
    - Integration, cleaning, redundancy elimination.

# Big data value chain (3)

| Generation | Acquisition | Storage | Analysis |
|---|---|---|---|

- ## Storage
  - o Storage infrastructure
    - Storage technology, e.g., HDD, SSD
    - Networking architecture, e.g., DAS, NAS, SAN
  - o Data management
    - File systems (HDFS), key-value stores (Memcached), column-oriented databases (Cassandra), document databases (MongoDB)
  - o Programming models
    - Map-Reduce, stream processing, graph processing.

# Big data value chain (4)

| Generation | Acquisition | Storage | Analysis |
|---|---|---|---|

- Analysis
  - Objectives
    - Descriptive analytics, predictive analytics, prescriptive analytics
  - Methods
    - Statistical analysis, data mining, text mining, network and graph data mining
    - Clustering, classification and regression, association analysis
  - Diverse domains call for customized techniques.

# Big data challenges

- Technology and infrastructure
  - New architectures, programming paradigms and techniques are needed
- Data management and analysis
  - New emphasis on "data"
  - Data science.

# The bottleneck

- Processors process data

- Hard drives store data

- We need to transfer data from the disk to the processor.

# The solution

- Transfer the processing power to the data

- Multiple distributed disks
  - Each one holding a portion of a large dataset

- Process in parallel different file portions from different disks.