

Lecture 2: Introduction to the Big Data

Big Data System Design

Table of Contents

❖ Part 1

- Introduction to Big Data

❖ Part 2

- What is Big Data Analytics?

Part 1

INTRODUCTION TO BIG DATA

Intro to Big Data

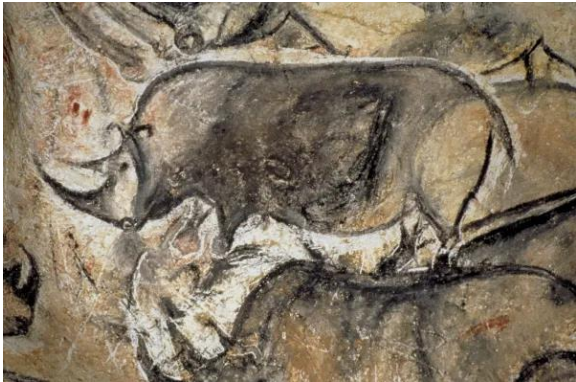
❖ What is data?

- Since birth, we are surrounded with data !



Intro to Big Data

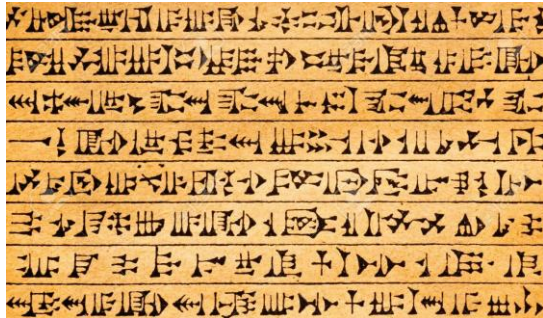
❖ What is data?



Intro to Big Data

❖ What is data?

- From the advent of written language, human observations have been recorded



Α α	Ι ι	Ρ ρ	Ω ω
Β β	Κ κ	Σ σ	Ψ ψ
Γ γ	Λ λ	Τ τ	ϛ ϛ
Δ δ	Μ μ	Υ υ	Ϛ Ϛ
Ε ε	Ν ν	Φ φ	Χ χ
Ζ ζ	Ξ ξ	Ψ ψ	Ϟ Ϟ
Η η	Ο ο	Ϝ Ϝ	ϟ ϟ
Θ θ	Π π	Ϟ Ϟ	Ϡ Ϡ



Intro to Big Data

❖ What is data?

- From the advent of written language, human observations have been recorded



Intro to Big Data

❖ What is data?

- The advent of computer technologies in 1950s, data most commonly refers to information that is transmitted or stored electronically
- The electronic sensors has additionally contributed to the volume and richness of recorded data



Intro to Big Data

❖ What is data?

Data signifies the documented result of
quantified or observed phenomena

Intro to Big Data

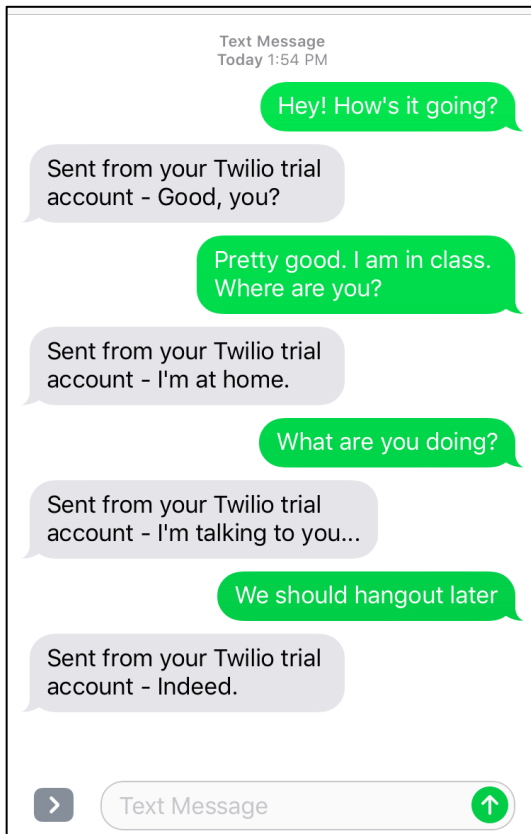
❖ Types of data

CUSTOMER

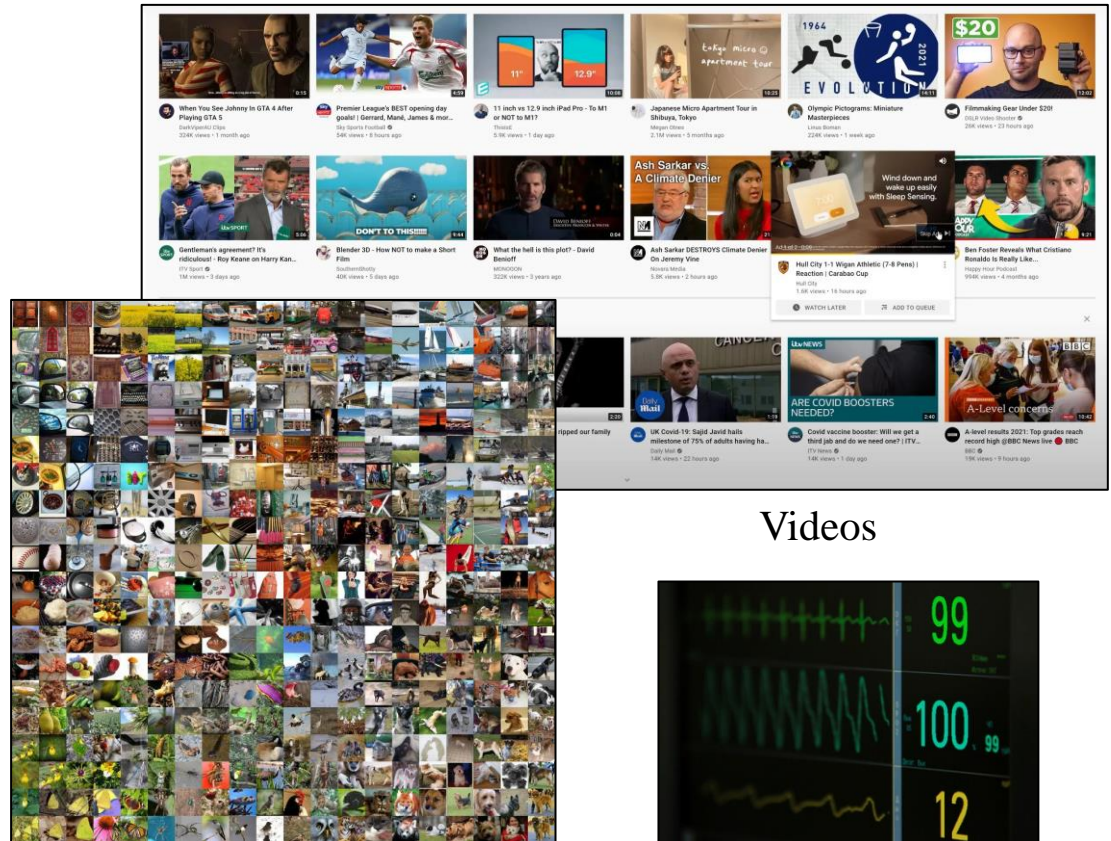
CUSTOMER_ID	LAST_NAME	FIRST_NAME	STREET	CITY	ZIP_CODE	COUNTRY
10302	Boucher	Leo	54, rue Royale	Nantes	44000	France
11244	Smith	Laurent	8489 Strong St	Las Vegas	83030	USA
11405	Han	James	636 St Kilda Road	Sydney	3004	Australia
11993	Mueller	Tomas	Berliner Weg 15	Tamm	71732	Germany
12111	Carter	Nataly	5 Tomahawk	Los Angeles	90006	USA
14121	Cortez	Nola	Av. Grande, 86	Madrid	28034	Spain
14400	Brown	Frank	165 S 7th St	Chester	33134	USA
14578	Wilson	Sarah	Seestreet #6101	Emory	1734	USA
14622	Jones	John	71 San Diego Ave	Arlington	69004	USA

Intro to Big Data

❖ Types of data



Texts



Images

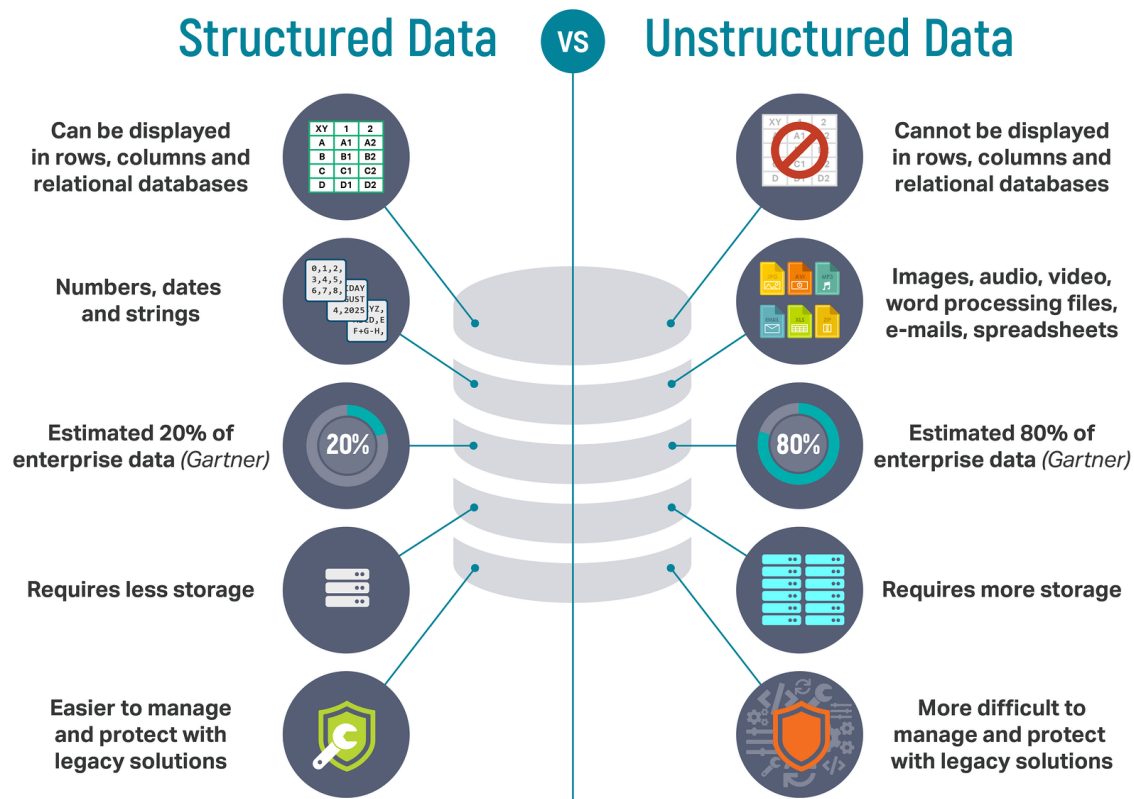
Videos

Sensors

Intro to Big Data

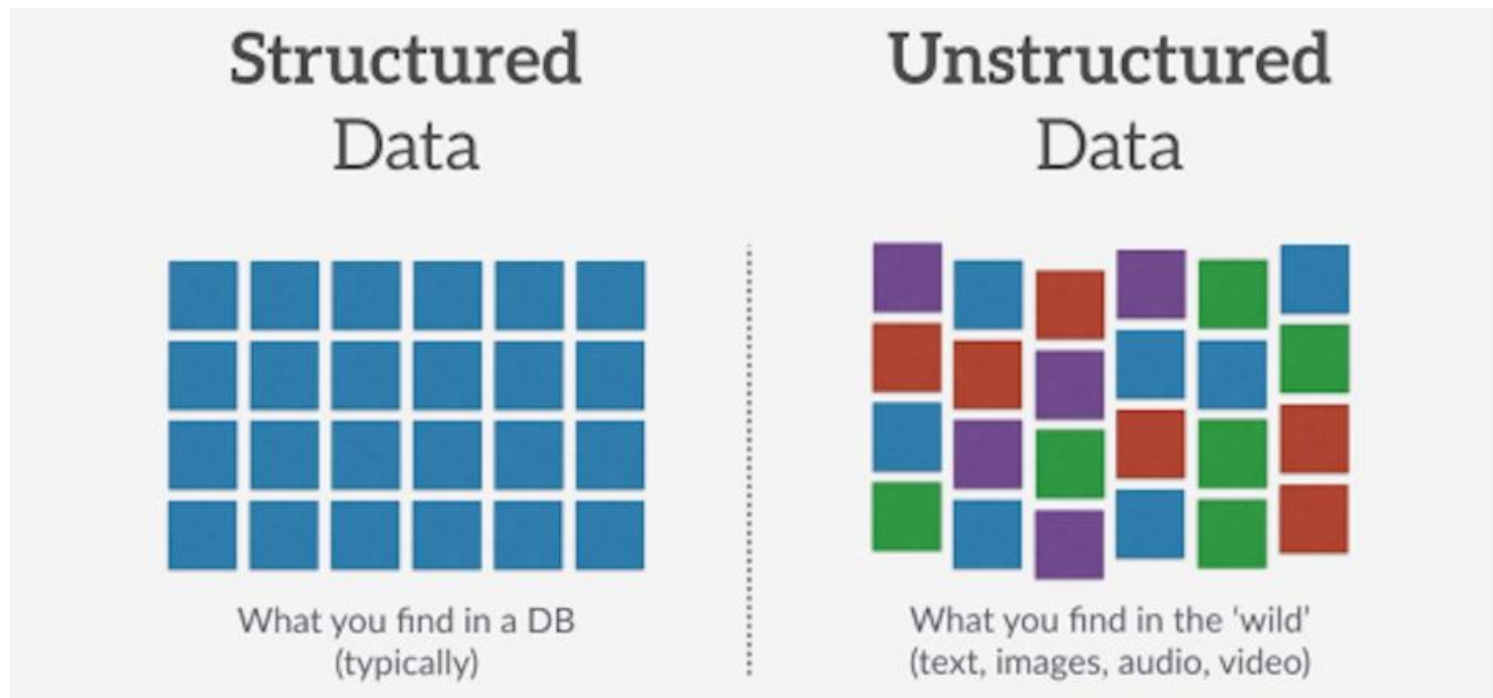
❖ Types of data

- Structured data vs. unstructured data



Intro to Big Data

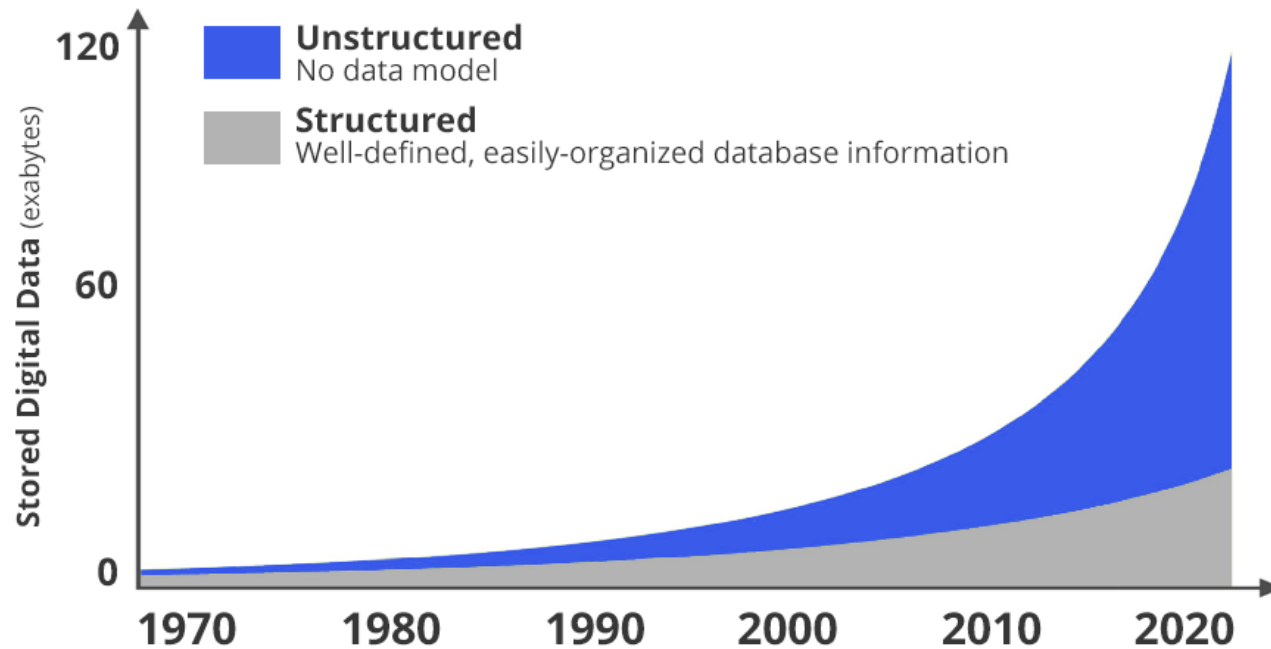
- ❖ Types of data
 - Structured data vs. unstructured data



Intro to Big Data

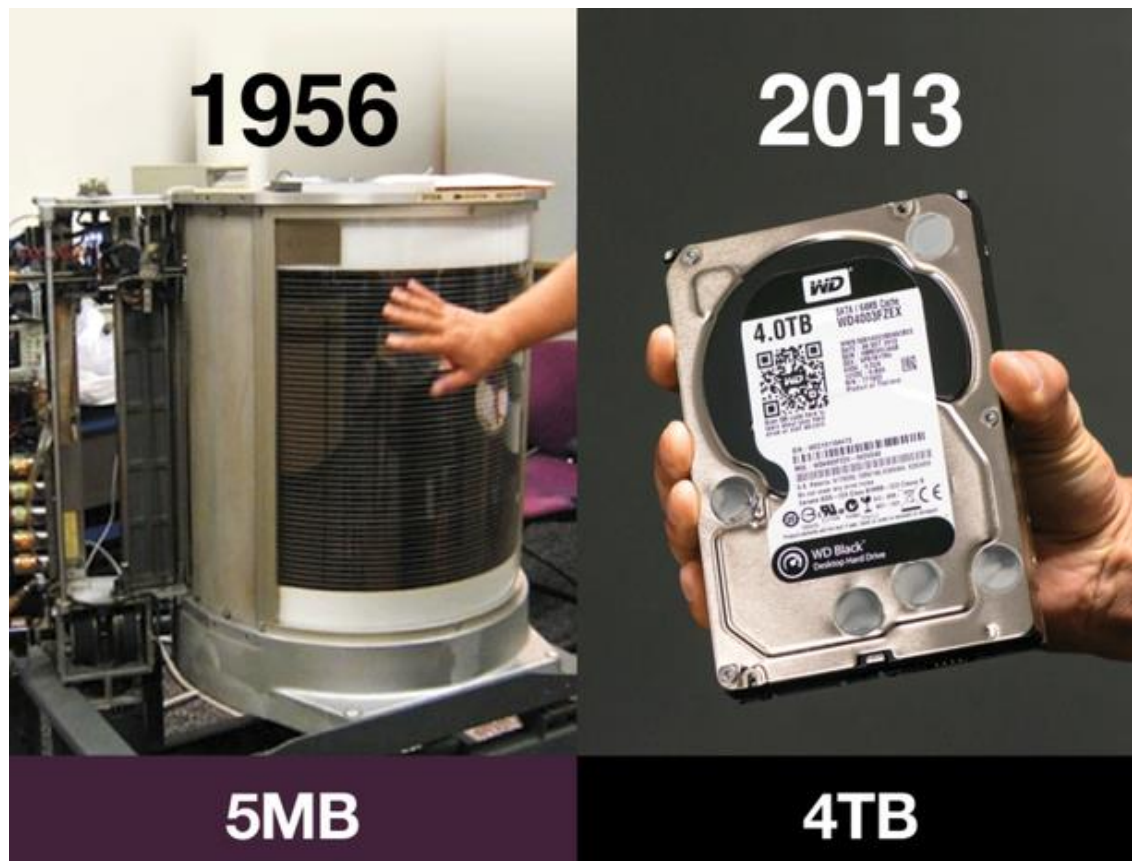
❖ Types of data

- Structured data vs. unstructured data



Intro to Big Data

- ❖ Reason for growth
 - Rapid advance in computer hardware



Intro to Big Data

- ❖ Reason for growth
 - Rapid advance in social networking



Intro to Big Data

❖ What is Big Data?

- According to Seagate, the volume of data generated worldwide will increase from 33 in 2018 to about 175 zettabytes in 2025



Megabyte
1 million bytes
Capacity of a 3.5" floppy disk (remember those?).



Gigabyte
1000 megabytes = 1 billion bytes
Today, USB key drives typically hold single- or double-digit gigabytes.



Terabyte
1000 gigabytes = 1 trillion bytes
Today's large consumer hard drives hold single-digit terabytes.



Petabyte
1000 terabytes = 1 quadrillion (10^{15}) bytes
The information in every US academic research library represents about 2 petabytes of text.



Exabyte
1000 petabytes = 1 quintillion (10^{18}) bytes
Every word ever spoken by every person ever can be represented in about 5 exabytes of text.
In 2014, there was about 60 exabytes of global internet traffic each month.



Zettabyte
1000 exabytes = 1 sextillion (10^{21}) bytes
1.3 zettabytes will be transmitted over the internet in 2016.

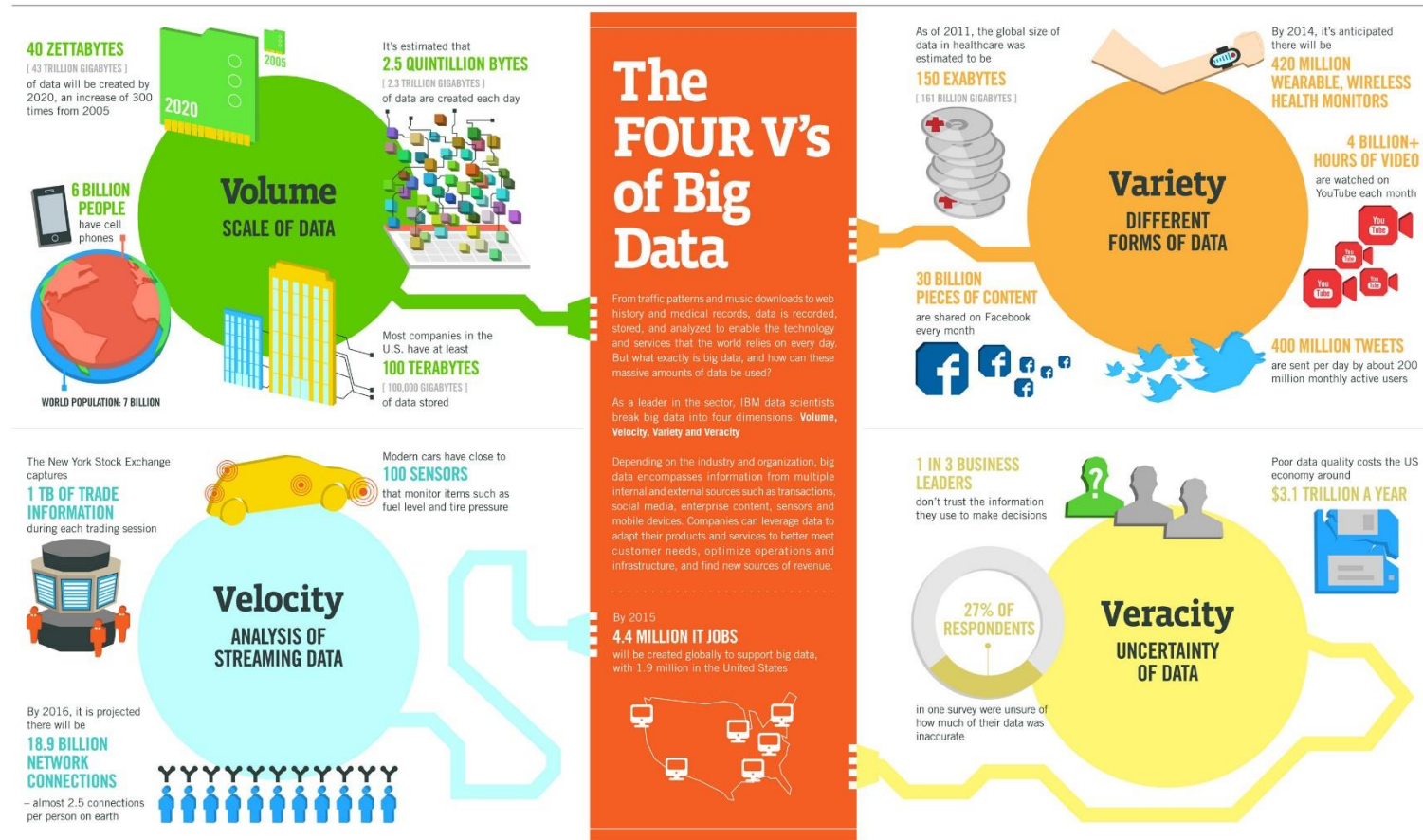
1 kilobyte	1,000,000,000,000,000,000
1 megabyte	1,000,000,000,000,000,000,000
1 gigabyte	1,000,000,000,000,000,000,000,000
1 terabyte	1,000,000,000,000,000,000,000,000,000
1 petabyte	1,000,000,000,000,000,000,000,000,000,000
1 exabyte	1,000,000,000,000,000,000,000,000,000,000,000
1 zettabyte	1,000,000,000,000,000,000,000,000,000,000,000,000



By 2019, global internet traffic will exceed 2 zettabytes per year.

Intro to Big Data

❖ What is Big Data?



IBM

Intro to Big Data

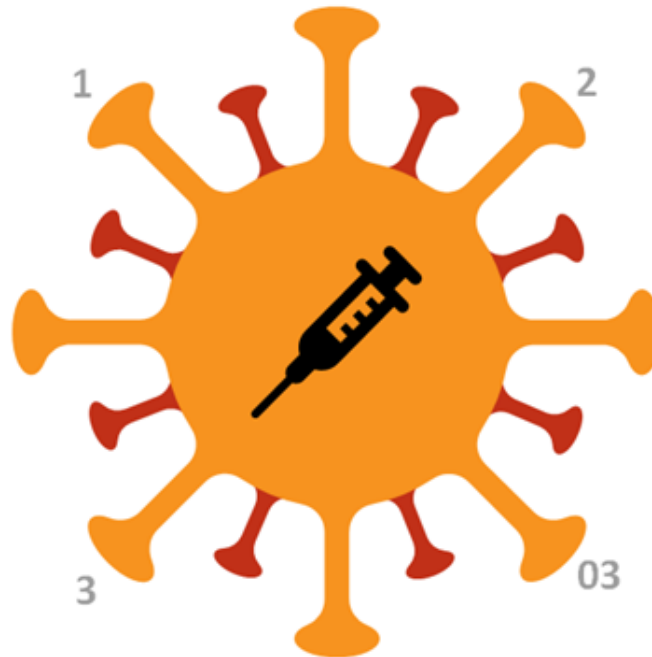
❖ Hospital Big Data Example

1. Volume

- Hospitals around the world generate a massive amount of data in the form of patient records and test results
- According to IBM, 2.314 Exabytes of medical data collected annually around the world

3. Velocity

- According to IBM, medical data is experiencing a 48 percent annual growth rate



2. Variety

- Hospital can collect medical records in variety form, such as structured and unstructured data
- It can be textual information, excel or images (e.g., X-Ray images)

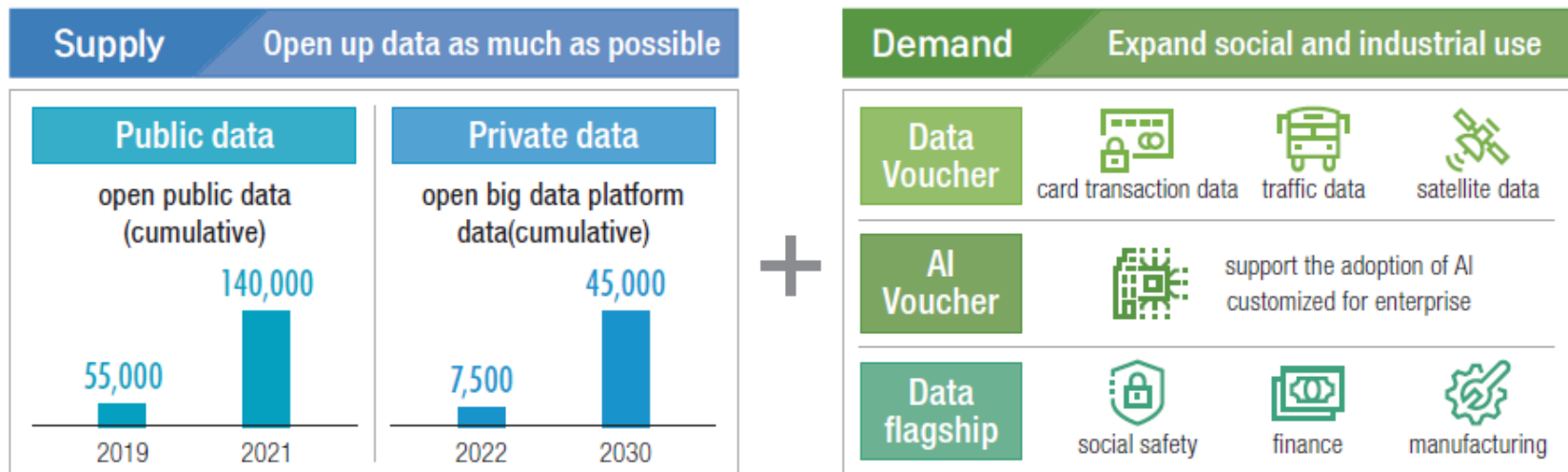
4. Veracity

- Since its healthcare field, the accuracy and trustworthiness of the data must be very high
- High accuracy in medical examination, prediction of disease

Intro to Big Data

❖ Big Data in South Korea

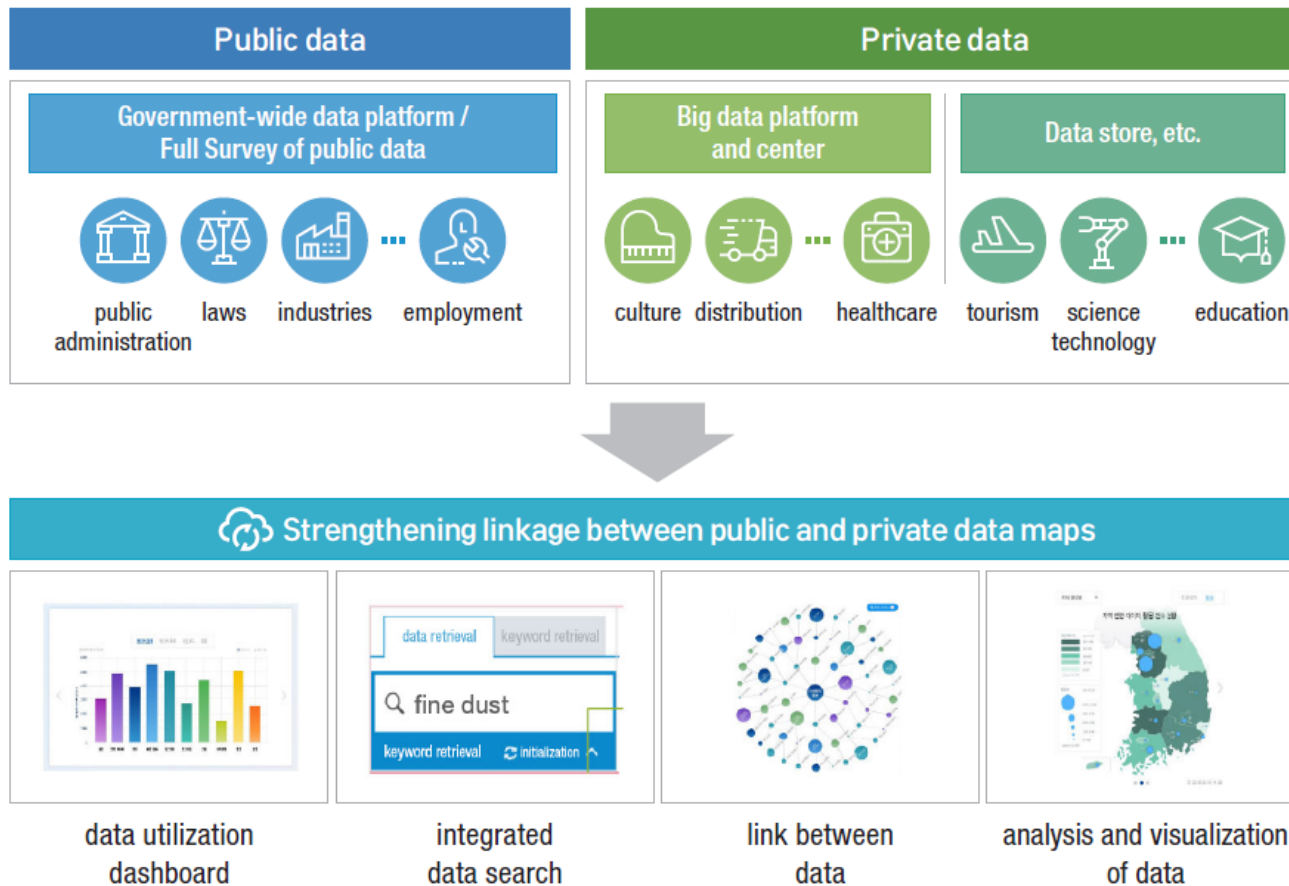
- Promotion of Opening Up Data and Reuse
 - Expanding the construction of AI learning data and securing of AI development infrastructure through the 'AI Hub' platform supply



Intro to Big Data

❖ Big Data in South Korea

- Strengthening Linkage between Public/Private Data Map



Part 2

WHAT IS BIG DATA ANALYTICS

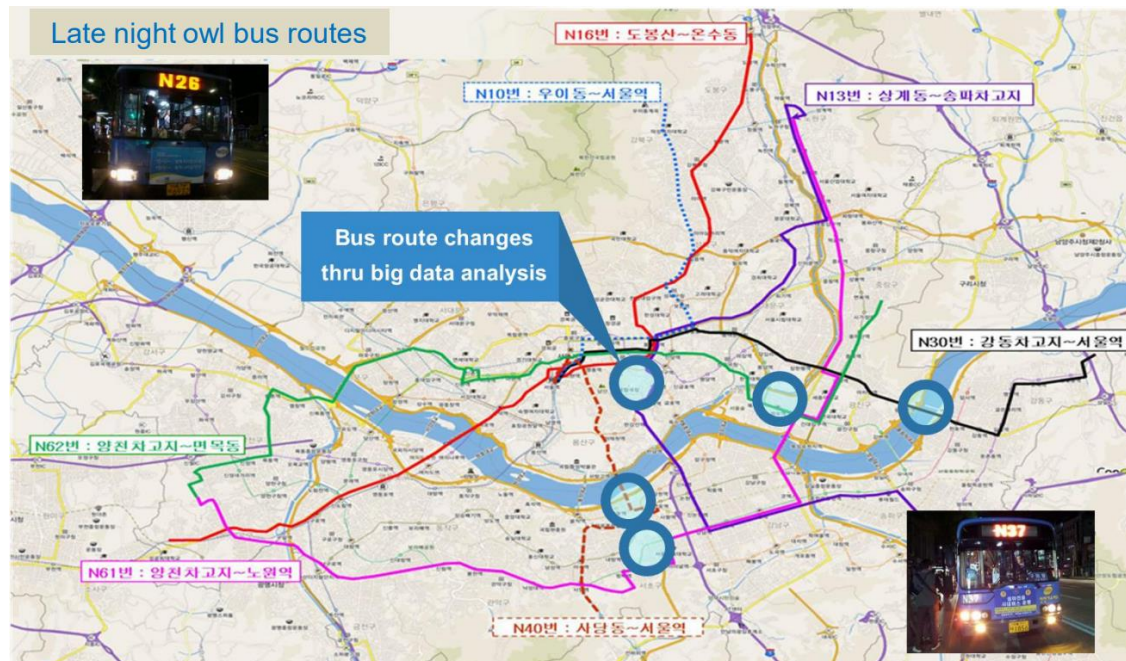
What is Big Data Analytics?

- ❖ Big Data analytics is a process used to extract meaningful insights
 - hidden patterns
 - unknown correlations
 - market trends
 - customer preferences
- ❖ Big Data analytics provides various advantages
 - It can be used for better decision making, preventing fraudulent activities, reduce cost among other things.

What is Big Data Analytics?

❖ Example of Big Data Analytics (The OWL Service)

- Taxi at night is expensive and difficult to catch
- Through a partnership with Korea Telecom, Seoul Government gained access to anonymized mobile communication data
 - 3 billion mobile call logs, 5 million taxi ride data



What is Big Data Analytics?

❖ Example of Big Data Analytics (The OWL Service)

Big Data Analytics for Bus Route Optimization



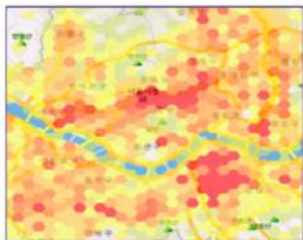
1. Data collection and analysis



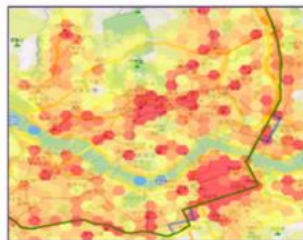
2. Layer modelling based on ride locations



3. Hexagon mapping



4. Floating population density analysis



5. Bus route optimization with floating population



6. Dispatch timetable adjusted accordingly

Impact

After three months of operating two routes

- Covers 42% of Seoul residents
- 7,900 passengers per day
- 2.3 million less car trips per year
- \$13 million fare savings
- 500 metric tons reduction in greenhouse gas emission per year
- A service satisfaction score of 82 points (74.3 points for standard buses)

What is Big Data Analytics?

❖ Example of Big Data Analytics (POSCO)

- POSCO is one of the largest hot rolling plant in the world
- POSCO reduced energy input by 2% and save 1 billion won annually
 - Collecting and analyzing manufacturing environment data through sensors in factory
 - Maintaining the optimal working conditions through AI



What is Big Data Analytics?

- ❖ Example of Big Data Analytics (Siemens Amberg Factory)
 - The most successful case of smart factory using big data
 - Produce prototypes in virtual environment
 - The reasons of defects: materials, man(worker), machine, method
 - Automation 75%, Defects: 0.0009%, Energy: 30%, Time: 50%



What is Big Data Analytics?

❖ Example of Big Data Analytics (Netflix)

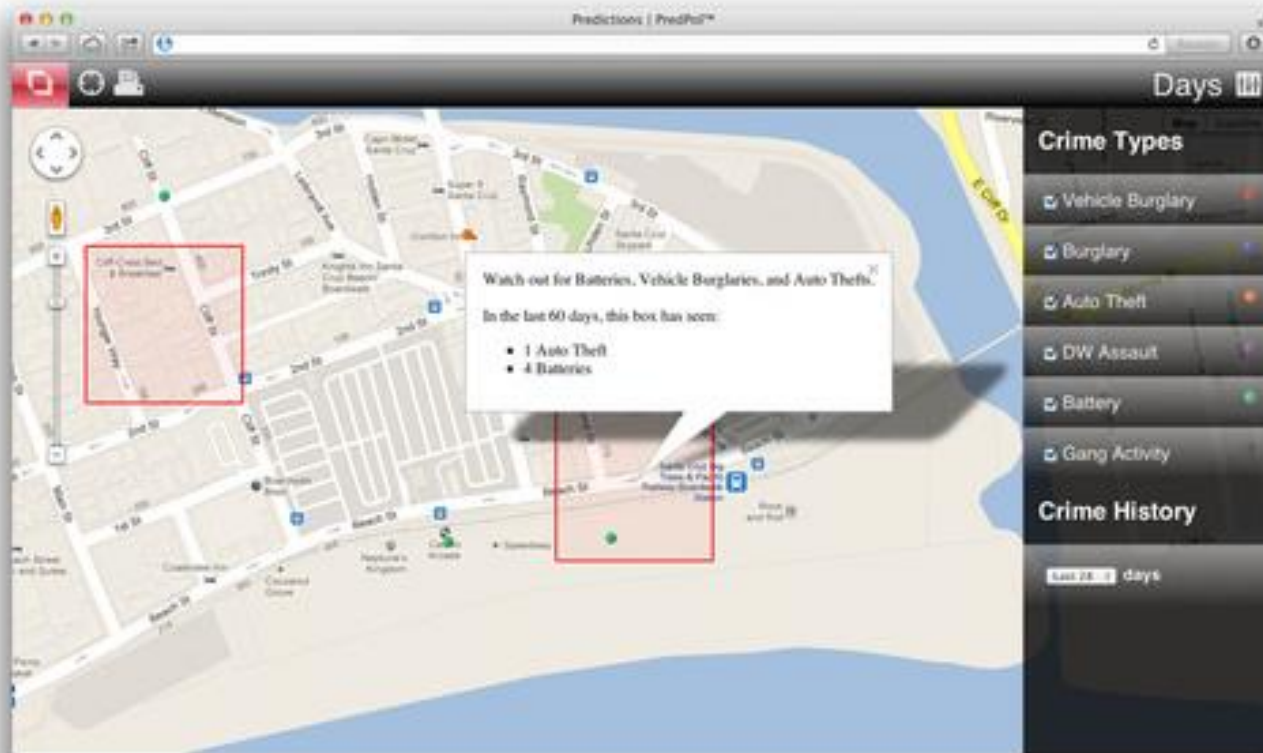
- With 115 million subscribers, Netflix collect a huge amount of data
 - Ratings, watch history, searchers and others
- Recommend the next movie you should watch or smart advertising



What is Big Data Analytics?

❖ Example of Big Data Analytics (Predpol)

- Projection of areas where criminal activity is most likely
- Reduced crime rates in Los Angeles, US



What is Big Data Analytics?

❖ Example of Big Data Analytics

- Drug data reveal sneaky side effects

The screenshot shows the top of the Nature website with a dark red header. The 'nature' logo is on the left, and a search bar is on the right. Below the header is a navigation bar with links like 'Home', 'News & Comment', 'Research', etc. The main content area has a breadcrumb trail: 'News & Comment > News > 2019 > May > Article'. The article title 'Drug data reveal sneaky side effects' is prominently displayed, followed by a subtitle 'Mining of surveillance data highlights thousands of previously unknown consequences when drugs are taken together.' and the author 'Heidi Ledford'. A date '14 March 2012' is shown. There is a 'Rights & Permissions' button. The article text begins with 'An algorithm designed by US scientists to trawl through a plethora of drug interactions has yielded thousands of previously unknown side effects caused by taking drugs in combination.' A quote from the study's lead author, Russ Altman, is also present. On the right side, there is a 'nature briefing' section with a smartphone image and a 'Sign up' button, and a 'Listen' section with a large red 'n' and headphones icon.

nature International weekly journal of science

Search [Advanced search](#)

[Home](#) | [News & Comment](#) | [Research](#) | [Careers & Jobs](#) | [Current Issue](#) | [Archive](#) | [Audio & Video](#) | [For Authors](#)

[News & Comment](#) > [News](#) > [2019](#) > [May](#) > [Article](#)

NATURE | NEWS

Drug data reveal sneaky side effects

Mining of surveillance data highlights thousands of previously unknown consequences when drugs are taken together.

Heidi Ledford

14 March 2012

[Rights & Permissions](#)

An algorithm designed by US scientists to trawl through a plethora of drug interactions has yielded thousands of previously unknown side effects caused by taking drugs in combination.

The work, published today in *Science Translational Medicine*¹, provides a way to sort through the hundreds of thousands of 'adverse events' reported to the US Food and Drug Administration (FDA) each year. "It's a step in the direction of a complete catalogue of drug-drug interactions," says the study's lead author, Russ Altman, a bioengineer at Stanford University in California.

nature briefing

What matters in science — and why — free in your inbox every weekday.

[Sign up](#)

Listen

What is Big Data Analytics?

❖ Example of Big Data Analytics

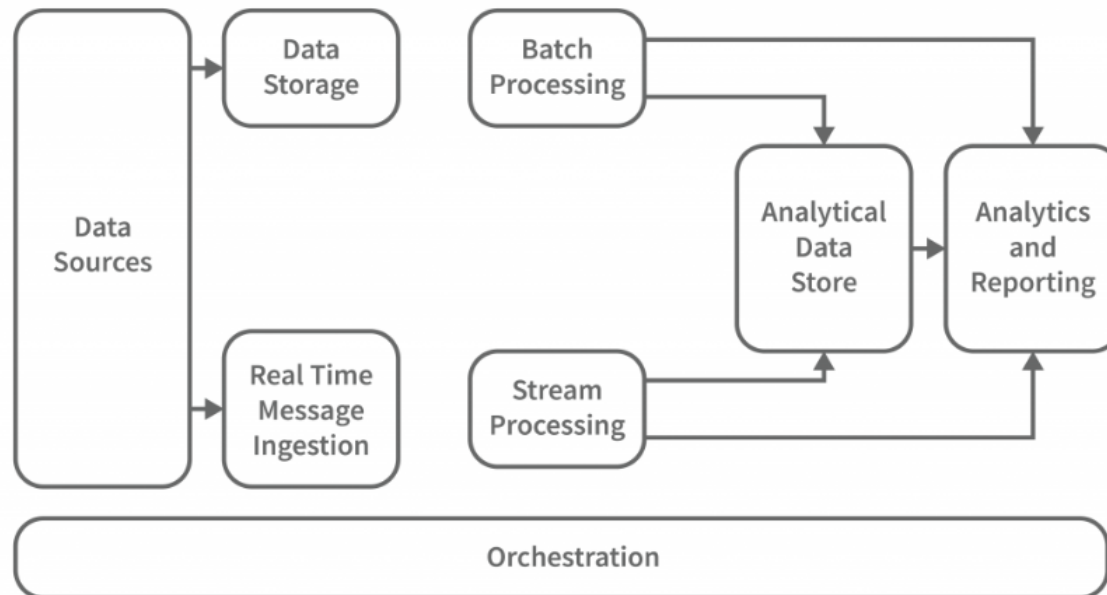
- Shoplifting detection using artificial intelligence (AI)



How to design Big Data System?

❖ Big Data Frameworks (required)

- Processing tremendous data
- Real-time data processing
- Low cost
- Fault Tolerance



How to design Big Data System?

❖ Framework

- Software environment
- Reuse of the designs and implementations of software functions
- Support the development of new applications or solutions



How to design Big Data System?

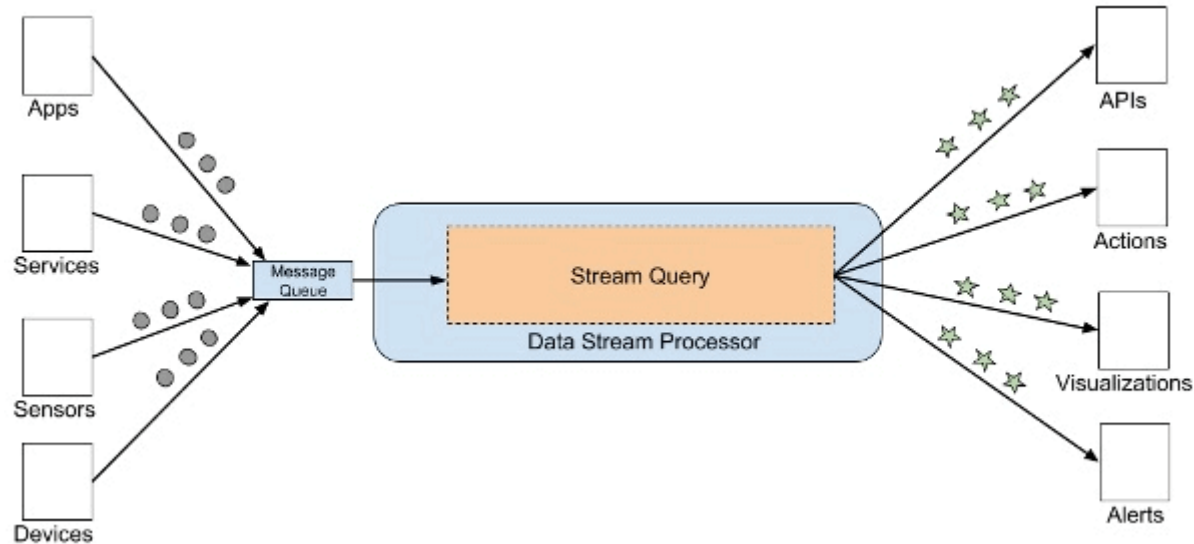
❖ Real-time data processing (Why?)

- Numerous data sources: SNS, IoT, Smartphone, Sensors, RFID, Log, etc
- Real-time analysis → Pattern identification → Decision or New BM
- Data integration and analysis techniques (required)
 - Data filtering and cleaning
 - Data regularization and normalization
 - Outlier detection
 - Data interpolation
 - Data integration
 - Data analysis
 - Data visualization, etc
- Data security & privacy (required)

How to design Big Data System?

❖ Streaming data

- Data generated continuously from data sources
- Typically, data records (in kilobytes) are transferred simultaneously
- E.g., Customer logs, transactions, SNS, Stock market, Video, GPS, etc

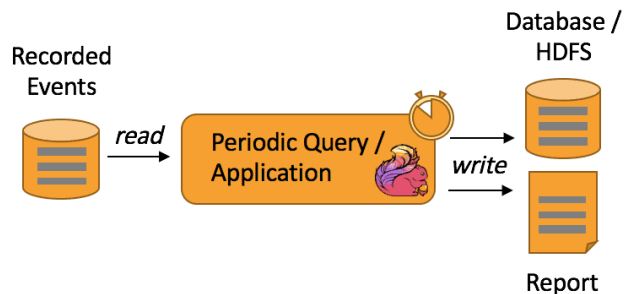


How to design Big Data System?

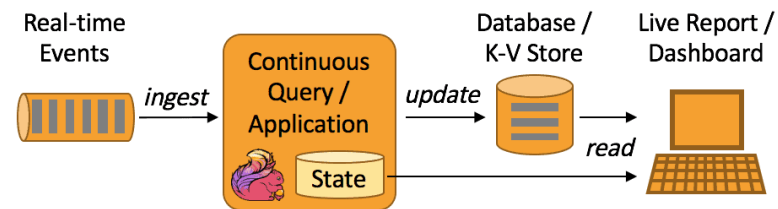
❖ Big Data System

- Distributed and Parallelized System for managing large amounts of data
- Real-time data processing and batch data processing
- Functions: Data collection, management, transmission and analysis

Batch processing



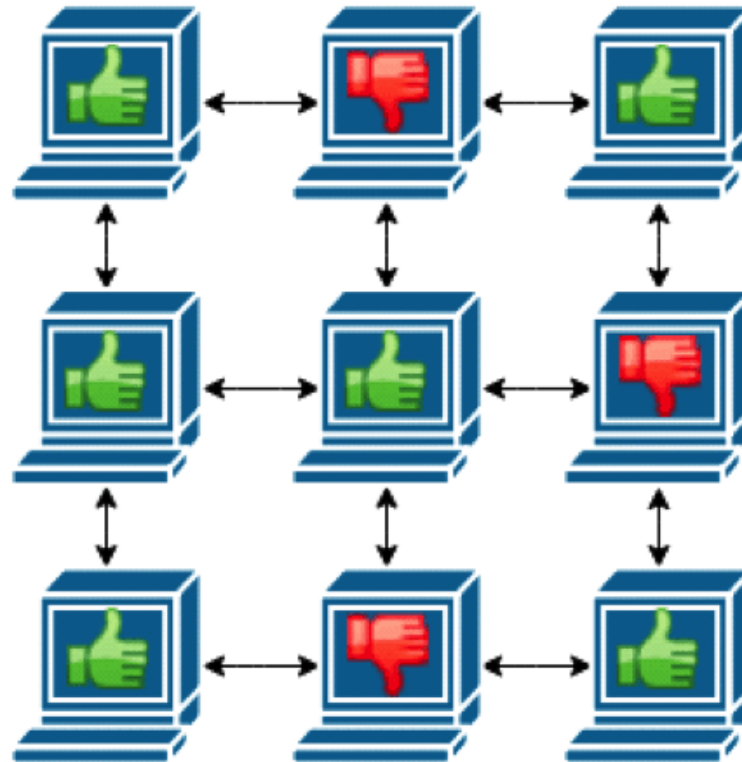
Real-time processing



How to design Big Data System?

❖ Fault tolerance system

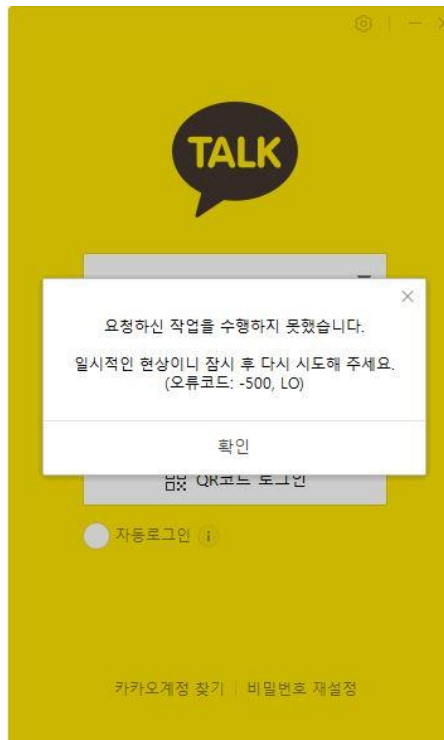
- System that performs tasks even when faults occur



How to design Big Data System?

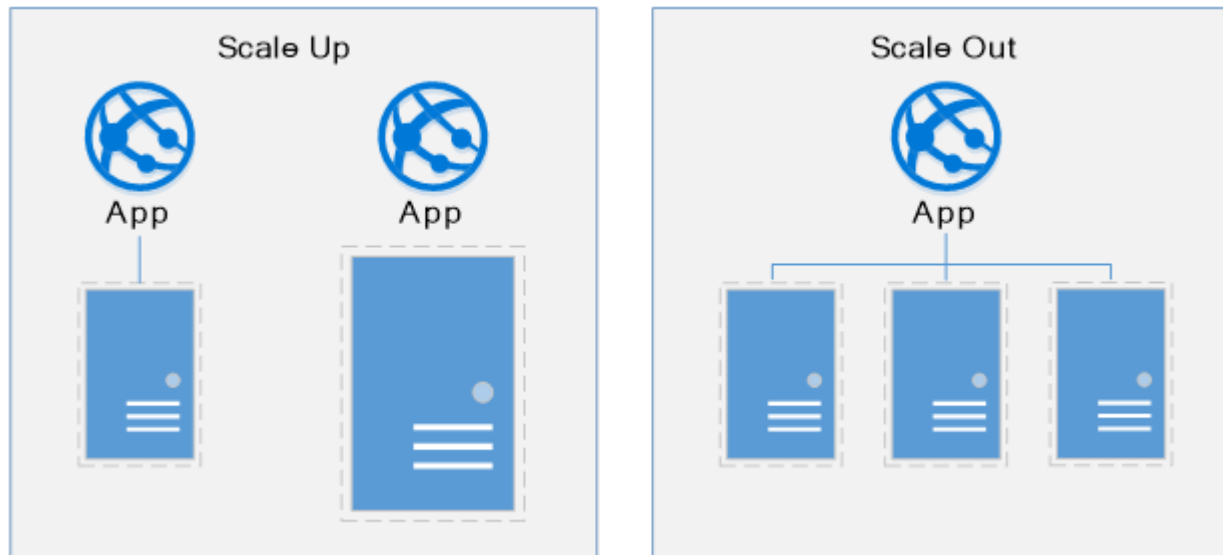
❖ Fault tolerance system

- System that performs tasks even when faults occur
- Faults: Hardware malfunction, Software errors, Computing resources, Data corruption, etc
- Problem



How to design Big Data System?

- ❖ Cost effective system
 - Low cost system



How to design Big Data System?

❖ Cost effective system

- Big data systems are not necessarily expensive
- Reduce cost by configuring the sufficient system for tasks
- Considerations: data type, processing time, data throughput, and other system requirements
- Advantage
 - Low initial cost → Competitive price

How to design Big Data System?

- ❖ Compatibility with existing systems
 - Compatibility with existing systems for big data collection
 - E.g., existing RDBMS, Hadoop system, applications, etc

How to design Big Data System?

❖ Opensource

- Available to anyone (without any restrictions)
- This is compatible with most existing systems
- Open source license: GPL, LGPL, MPL, BSD, MIT, Apache

	GPL	LGPL	MPL	BSD	MIT	Apache
Download	O	O	O	O	O	O
Deploy	O	O	O	O	O	O
Modify	X	X	X	O	O	O
Noncommercial	X	X	X	X	X	X
No Derives	X	X	X	X	X	X

Summary and Discussions

❖ Intro to Big Data

- Data most commonly refers to information that is transmitted or stored electronically
- 4Vs of Big Data

❖ What is Big Data Analytics

- Big Data analytics is a process used to extract meaningful insights, such as hidden patterns, unknown correlations, market trends, customer preferences

Questions?

SEE YOU NEXT TIME!