# Data Engineering

204426

# Big Data

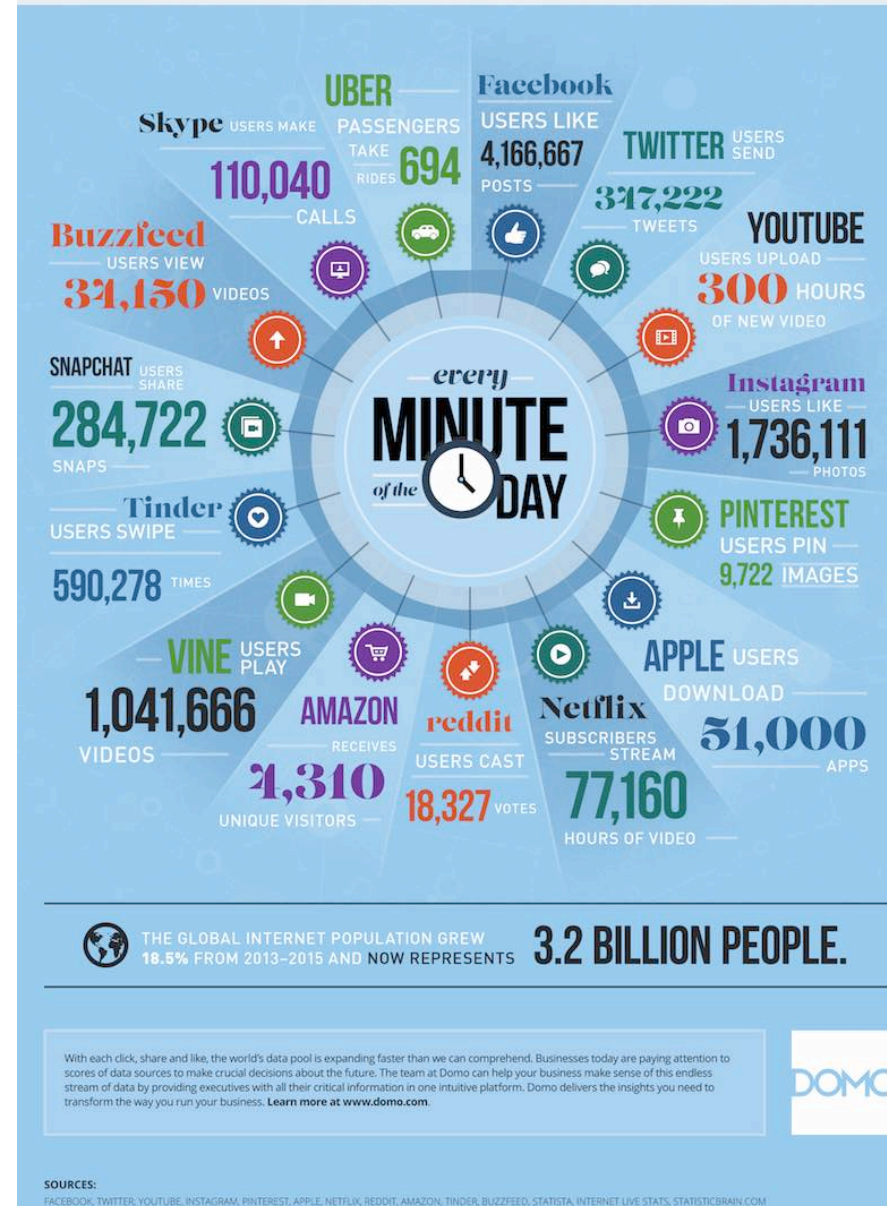# Big Data

- Too large + Complex
- Big data was originally associated with three key concepts: volume, variety, and velocity.

# 3Vs Properties

**Volume**
- Scale of data
- A huge amount of data
- If the volume of data is very large then it is actually considered as a 'Big Data'
- Terabyte/Petabyte/Exabyte

**Variety**
- Different form of data
- Different function of data
- Different data sources

**Velocity**
- High speed of accumulation of data
- A massive and continuous flow of data.
- The potential of data that how fast the data is generated and processed to meet the demands.

# The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
**4.4 MILLION IT JOBS**
will be created globally to support big data, with 1.9 million in the United States

## Volume
### SCALE OF DATA

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

**6 BILLION PEOPLE**
have cell phones

WORLD POPULATION: 7 BILLION

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

## Velocity
### ANALYSIS OF STREAMING DATA

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

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

## Variety
### DIFFERENT FORMS OF DATA

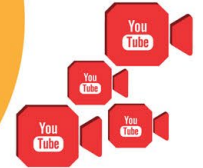As of 2011, the global size of data in healthcare was estimated to be
**150 EXABYTES**
[ 161 BILLION GIGABYTES ]

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

**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

## Veracity
### 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

Poor data quality costs the US economy around
**$3.1 TRILLION A YEAR**

IBM

The Five V's of Big Data

Source: https://morioh.com/p/ca19c6b8c0fe

# Components of the Big Data Ecosystem

**Data Sources**

- Social media
- IoT
- Web
- Business transaction
- Etc.

**Big Data Ecosystem**

Data Ingest

Stream Processing

Batch Processing

Data Organization

Distributed File System

Computer, Storage, Network Infrastructure

**Data Consumption**

- Data analysis
- Data mining
- Visualization
- Report
- Application
- Etc.

# Components of the Big Data Ecosystem

**Data Ingestion**
- Transportation of data from assorted sources to a storage medium.
- **Batch processing**
  - Ingestion layer periodically collects and groups source data and sends it to the destination system.
  - Groups may be processed based on any logical ordering, the activation of certain conditions, or a simple schedule.
- **Stream processing**
  - Data is sourced, manipulated, and loaded as soon as it's created or recognized by the data ingestion layer.

**Data Organization**
- Database - the method of classifying and organizing data sets to make them more useful.

# NoSQL

- Non-tabular databases and store data differently than relational tables.

- Types of NoSQL databases:
  - Key-Value Store Databases
  - Document Store Databases
  - Graph Databases
  - Column-Oriented Databases

# NoSQL

**Key-Value Store Databases**

- Data is represented as a collection of key–value pairs
- The key–value model can be extended to a discretely ordered model that maintains keys in lexicographic order.
- Example:
  - DynamoDB
  - Voldemort
  - Redis

| Key | Value |
|-----|-------|
| K1 | AAA,BBB,CCC |
| K2 | AAA,BBB |
| K3 | AAA,DDD |
| K4 | AAA,2,01/01/2015 |
| K5 | 3,ZZZ,5623 |

# NoSQL

**Document Store Databases**

- Documents encapsulate and encode data (or information) in some standard format or encoding.
- Encodings in use include XML, YAML, JSON
- Use indexing to read/write data in form of document object.
- Example
  - MongoDB
  - CouchDB

**Students Collection**

```
{
Student_ID: 1
Name: Jean Grey
D
I
p
C
A
}
```

```
{
Student_ID: 2
Name: Scott Summers
DateOfBirth: 12-10-1968
IDCard: 765414A
Supervisor: {
 Name: Emma Frost
 DateOfBirth: 1-1-1936
 IDCard: 222222
 }
}
```

**Professors Collection**

```
{
Professor_ID: 1
Name: Charles Xavier
D
I
p
C
}
```

```
{
Professor_ID: 2
Name: Emma Frost
DateOfBirth: 1-1-1936
IDCard: 222222
}
```

# NoSQL

## Graph Databases

- Use graph structures for semantic queries with nodes, edges, and properties to represent and store data.

- The graph relates the data items in the store to a collection of nodes and edges, the edges representing the relationships between the nodes.

- Example
  - Neo4J
  - FlockDB

Id: 2
Name: Bob
Age: 22

Id: 100
Label: knows
Since: 2001/10/03

Id: 101
Label: knows
Since: 2001/10/04

Id: 105
Label: is_member
Since: 2011/02/14

Id: 104
Label: Members

Id: 1
Name: Alice
Age: 18

Id: 103
Label: Members

Id: 102
Label: is_member
Since: 2005/7/01

Id: 3
Type: Group
Name: Chess

# NoSQL

## Column-Oriented Databases

- Stores data tables by column rather than by row.

- By storing data in columns rather than rows, the database can more precisely access the data it needs to answer a query rather than scanning and discarding unwanted data in rows.

- Stores each column continuously. i.e. on disk or in-memory each column on the left will be stored in sequential blocks.

| RowId | EmpId | Lastname | Firstname | Salary |
|-------|-------|----------|-----------|--------|
| 001 | 10 | Smith | Joe | 60000 |
| 002 | 12 | Jones | Mary | 80000 |
| 003 | 11 | Johnson | Cathy | 94000 |
| 004 | 22 | Jones | Bob | 55000 |

```
10:001,12:002,11:003,22:004;
Smith:001,Jones:002,Johnson:003,Jones:004;
Joe:001,Mary:002,Cathy:003,Bob:004;
60000:001,80000:002,94000:003,55000:004;
```

# Distributed File System

- File system that is distributed on multiple file servers or multiple locations.

- Allows programs to access or store isolated files as they do with the local ones, allowing programmers to access files from any network or computer.
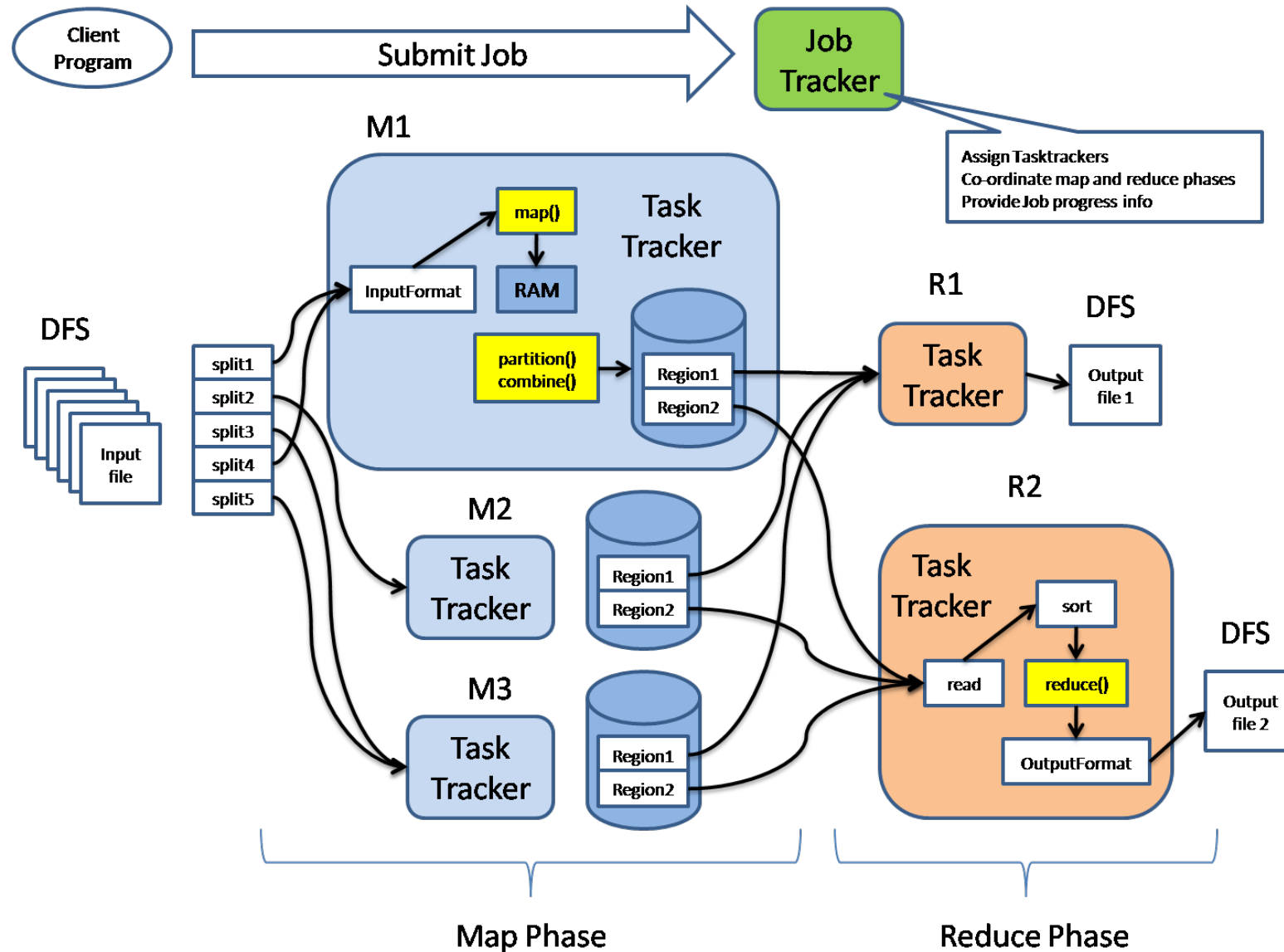
# Hadoop: A bigdata framework

# MapReduce

- Enhances the processing of massive data using dispersed and parallel algorithms in the Hadoop ecosystem.

- Process large datasets across <u>computer clusters.</u>

- Two primary tasks in MapReduce: <u>Map</u> and <u>Reduce</u>

- In the map job
  - Split the input dataset into chunks.
  - Task processes these chunks in parallel.
  - Use outputs as inputs for the reduce tasks.

- For reducer
  - Process the intermediate data from the maps into smaller tuples, that reduces the tasks, leading to the final output of the framework.

# MapReduce



Source:
http://a4academics.com/images/hadoop/Hadoop-Mapreduce-Architecture.png

# MapReduce

**Simplified flow diagram for the MapReduce program**



Input → Map → Shuffling & Sorting → Reduce → Output

- A dataset is split into equal units called chunks (input splits) in the splitting step.
- Hadoop consists of a RecordReader that uses TextInputFormat to transform input splits into key-value pairs.
- The key-value pairs are then used as inputs in the mapping step.
- The mapping step contains a coding logic that is applied to these data blocks.
- The mapper processes the key-value pairs and produces an output of the same form (key-value pairs).

# MapReduce

**Simplified flow diagram for the MapReduce program**



- It consists of two main steps: sorting and merging.
- In the sorting step, the key-value pairs are sorted using the keys. Merging ensures that key-value pairs are combined.
- The shuffling phase facilitates the removal of duplicate values and the grouping of values.
- Different values with similar keys are grouped. The output of this phase will be keys and values, just like in the Mapping phase.

# MapReduce

**Simplified flow diagram for the MapReduce program**



- The output of the shuffling phase is used as the input.
- The reducer processes this input further to reduce the intermediate values into smaller values.
- It provides a summary of the entire dataset. The output from this phase is stored in the HDFS.

# MapReduce

**Simplified flow diagram for the MapReduce program**
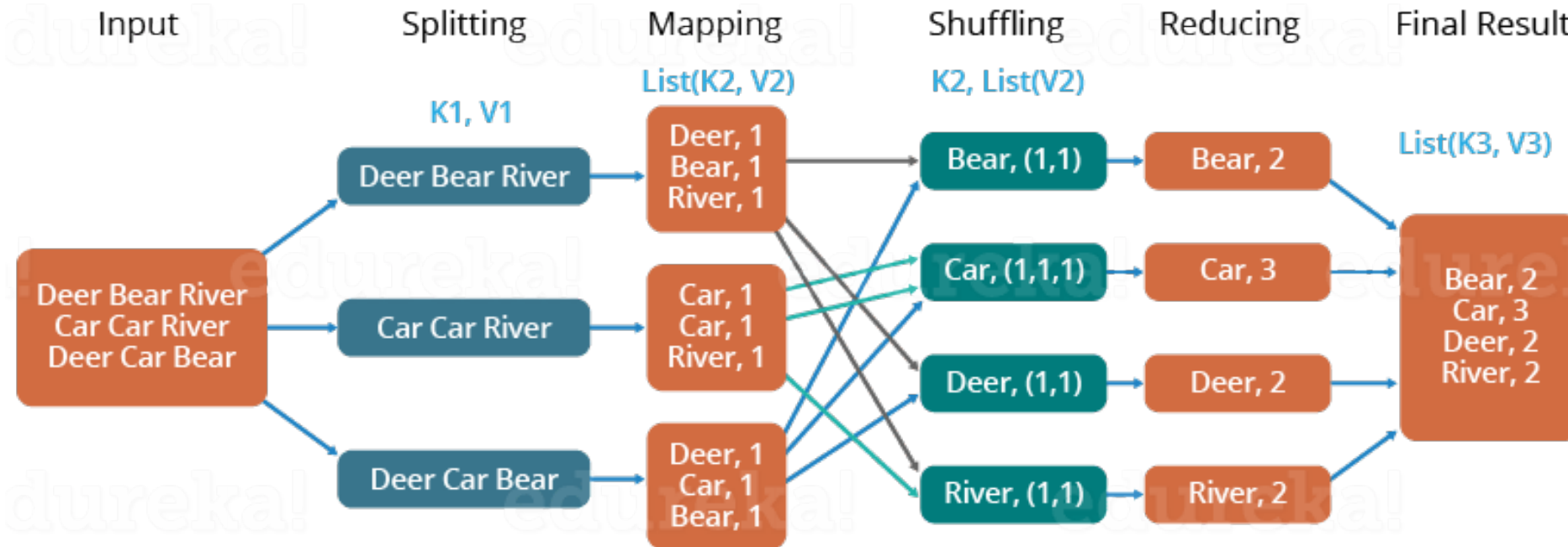


**Combiner phase**
- Optional phase that's used for optimizing the MapReduce process.
- It's used for reducing the pap outputs at the node level.
- In this phase, duplicate outputs from the map outputs can be combined into a single output.
- The combiner phase increases speed in the Shuffling phase by improving the performance of Jobs.

# MapReduce



The Overall MapReduce Word Count Process

edureka!

| Input | Splitting | Mapping | Shuffling | Reducing | Final Result |
|---|---|---|---|---|---|

K1, V1

List(K2, V2)

K2, List(V2)

List(K3, V3)

Deer Bear River
Car Car River
Deer Car Bear

Deer Bear River → Deer, 1 / Bear, 1 / River, 1

Car Car River → Car, 1 / Car, 1 / River, 1

Deer Car Bear → Deer, 1 / Car, 1 / Bear, 1

Bear, (1,1) → Bear, 2

Car, (1,1,1) → Car, 3

Deer, (1,1) → Deer, 2

River, (1,1) → River, 2

Bear, 2
Car, 3
Deer, 2
River, 2

# References

- วราภรณ์ พรหมวิอินทร์ (2562). **Big Data Analytics.** นนทบุรี: คอร์ฟังก์ชั่น
- [https://www.section.io/engineering-education/understanding-map-reduce-in-hadoop/](https://www.section.io/engineering-education/understanding-map-reduce-in-hadoop/)
- [https://www.stitchdata.com/resources/data-ingestion/](https://www.stitchdata.com/resources/data-ingestion/)