# Big Data Analytics

Dr Sirintra Vaiwsri | Email: sirintra.v@itm.kmutnb.ac.th

# Databases

# Relational Database Management System
**(Holmes, 2017; Staragile, 2023)**

- The major problem with using RDBMS for Big Data is scalability.

- RDBMS was designed to be used on one server which can be improved but still has its limits, called vertical scaling.

- RDBMS has **ACID** properties which are :

  - **A**tomicity - the incomplete transaction cannot update the database.

  - **C**onsistency - data is consistent before and after a transaction.

  - **I**solation - each transaction does not interfere with others.

  - **D**urability - a database must be updated to ensure data will not be lost if the system fails.
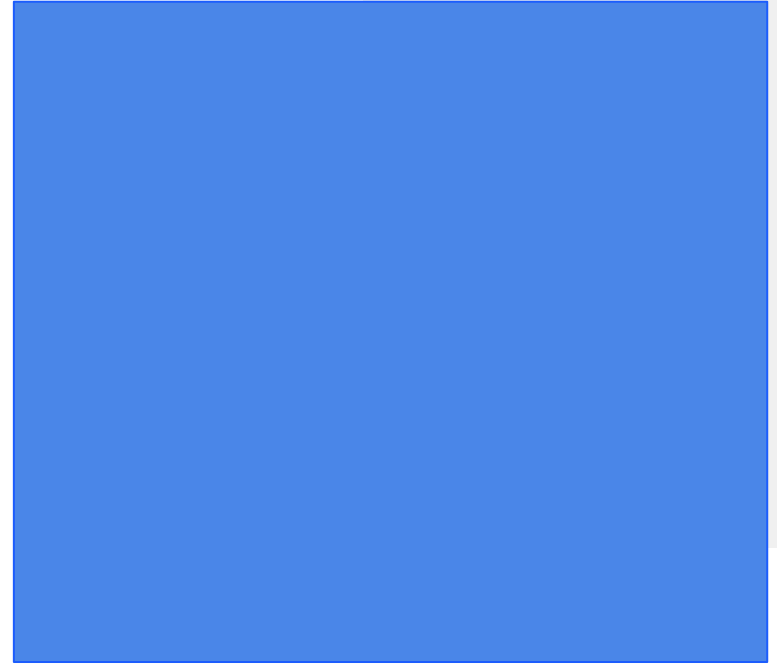
Dr Sirintra Vaiwsri

# Not only SQL (NoSQL) (Holmes, 2017)

- Non-relational database

- NoSQL has **BASE** properties which are :

  - **B**asically **A**vailable - the system is always available even when a network failure occurs.

  - **S**oft state - means it has flexibility with consistent requirements.

  - **E**ventually consistent - the system eventually becomes consistent.

- NoSQL is horizontal scaling which allows more data to be inserted into the database.

- Working with NoSQL must concern the CAP theorem.

Dr Sirintra Vaiwsri

# RDBMS VS NoSQL (Edward and Sabharwal, 2015)

|  | RDBMS | NoSQL |
|---|---|---|
| Schema flexibility | Inflexible, often ends up creating new tables | Column-oriented which allows adding more columns. It also supports semi-structured data. |
| Complex query | Often uses complex JOIN queries which are difficult to implement and maintain | It does not support relationships and foreign keys, thus, no complex query |
| Data update | If the system does not allow for updating multiple nodes at the same time, there is a risk of node failure | Synchronisation across nodes is challenging. However, NoSQL solutions offer synchronisation options. |
| Scalability | Low speed for large amounts of data | Provide great scalability |

Dr Sirintra Vaiwsri

# NoSQL

# Key-Value Databases
**(Bahga and Madisetti, 2019; DEV, 2023; KDnuggets, 2023; Janev et al., 2020)**

- It is the simplest NoSQL database.

- It stores data in the form of key-value pairs.

- A key is unique for each data.

- Key is usually a string or an integer.

- A value contains data which can be in the form of attributes or collections.

- Value can be any type of data.

Dr Sirintra Vaiwsri

# Key-Value Databases (DEV, 2023; KDnuggets, 2023)

- Advantages:
  - Scalability - it is horizontal scaling through partitioning and replication. It also has low overhead.
  - Mobility - it is easy to move from one to another system without changing in code/architecture required.
- Disadvantages:
  - All joins must be done in code.
  - No complex query filters.

# **Document Databases** (Bahga and Madisetti, 2019)

- It is similar to key-value databases in that each document has a unique key (ID).

- Each document can store any type of data.

- Its query is JSON-like documents.

- Therefore, it requires a data format that a database can understand.

Dr Sirintra Vaiwsri

# Document Databases (DEV, 2023)

- Advantages:
  - It collects data from RAM which is fast to access.
  - It is horizontal scaling.
- Disadvantages:
  - Selecting data from multiple collections requires multiple queries.
  - Data duplication can occur which makes it difficult to handle.

Dr Sirintra Vaiwsri

# Column-Oriented Databases
**(Bahga and Madisetti, 2019; DEV, 2023; KDnuggets, 2023; Janev et al., 2020)**

- It stores data in the form of columns.

- Column - it contains name, value, and timestamp.

- Row - it contains one or more columns where different rows are not necessary to contain the same number of columns. It has row-key as a unique key (ID).

- Column family - it contains multiple rows where each row can contain multiple column families.

- Keyspace - it contains multiple column families.

Dr Sirintra Vaiwsri

# Column-Oriented Databases (DEV, 2023)

- Advantages:
  - It is scalable and flexible.
  - Load and aggregation times are very fast.
- Disadvantages:
  - It is slow when deleting rows.
  - It can be slow when querying data using a join query.

Dr Sirintra Vaiwsri

# Graph Databases (Bahga and Madisetti, 2019; DEV, 2023)

- It stores data that has a graph structure.

- It shows a relationship between data.

- Nodes have a set of attributes.

- Edges (links) also have a set of attributes.

Dr Sirintra Vaiwsri

# Graph Databases (DEV, 2023)

- Advantages:
  - It is easy to understand data and has descriptive queries.
  - It is flexible.
- Disadvantages:
  - It is difficult to scale.
  - It does not have a standard language.

Dr Sirintra Vaiwsri

# MongoDB

# Install MongoDB

Download MongoDB from here.

Dr Sirintra Vaiwsri

# Install MongoDB

Dr Sirintra Vaiwsri

# Install MongoDB

# Install MongoDB

Dr Sirintra Vaiwsri

# Install MongoDB

# Install MongoDB



MongoDB Compass is being installed.

It will launch once it is done.

# New Connection

# Create Database

Dr Sirintra Vaiwsri

# Create Database

Dr Sirintra Vaiwsri

# Import Data

# Import Data



Dr Sirintra Vaiwsri

# Select All

SQL: select * from <table>          MQL: { }

Dr Sirintra Vaiwsri

# Select ..... Where

SQL: select * from <table> where <column> = <value>      MQL: {<column>:<value>}

Dr Sirintra Vaiwsri

# Select ….. Where ….. And

SQL: select * from <table> where <column1> = <value1> and <column2> = <value2>

MQL: {<column1>:<value1>, <column2>:<value2>}



Dr Sirintra Vaiwsri

# Select ..... Where ..... Or

SQL: select * from <table> where <column1> = <value1> or <column2> = <value2>

MQL: {$or: [ {<column1>:<value1>}, {<column2>:<value2>} ] }

Dr Sirintra Vaiwsri

# Select ..... Where ..... And ..... Or

SQL: select * from <table> where (<column1> = <value1>) and (<column2> = <value2> or <column3> = <value3>)



MQL: {{<column1>:<value1>}, $or: [ {<column2>:<value2>},{<column3>:<value3>} ] }

Dr Sirintra Vaiwsri

# Assignment (1 point)

- Please select all records that have a number of likes equal to 500 and also have a number of reactions greater than 3,000 or a number of comments greater than 10

```
_id: ObjectId('668e4c01f5b615d3c9216f5b')
status_id : "246675545449582_515902548526879"
status_type : "video"
status_published : "3/6/2014 5:29"
num_reactions : 500
num_comments : 16
num_shares : 0
num_likes : 500
num_loves : 0
num_wows : 0
num_hahas : 0
num_sads : 0
num_angrys : 0
```

Dr Sirintra Vaiwsri

# References

- Holmes, D. E. (2017). Big data: a very short introduction. Oxford University Press.
- Staragile. A Brief comparison between ACID vs BASE database model. https://staragile.com. Accessed: 2023-08-05.
- Edward, S. G., & Sabharwal, N. (2015). *Practical MongoDB: Architecting, Developing, and Administering MongoDB*. Apress.
- Bahga, A., & Madisetti, V. (2019). Big Data analytics: A hands-on approach.
- DEV. Intro to 4 types of NoSQL databases. https://dev.to. Accessed: 2023-08-14.
- Janev, V., Graux, D., Jabeen, H., & Sallinger, E. (2020). Knowledge graphs and big data processing (p. 209). Springer Nature.
- KDNuggets. https://www.kdnuggets.com/. Accessed: 2023-08-14.

Dr Sirintra Vaiwsri