# Combating COVID-19 Outbreaks

WQD7005 Data Mining
Master of Data Science | University of Malaya

**Part B:** Management of Data Assignment

Group Members:
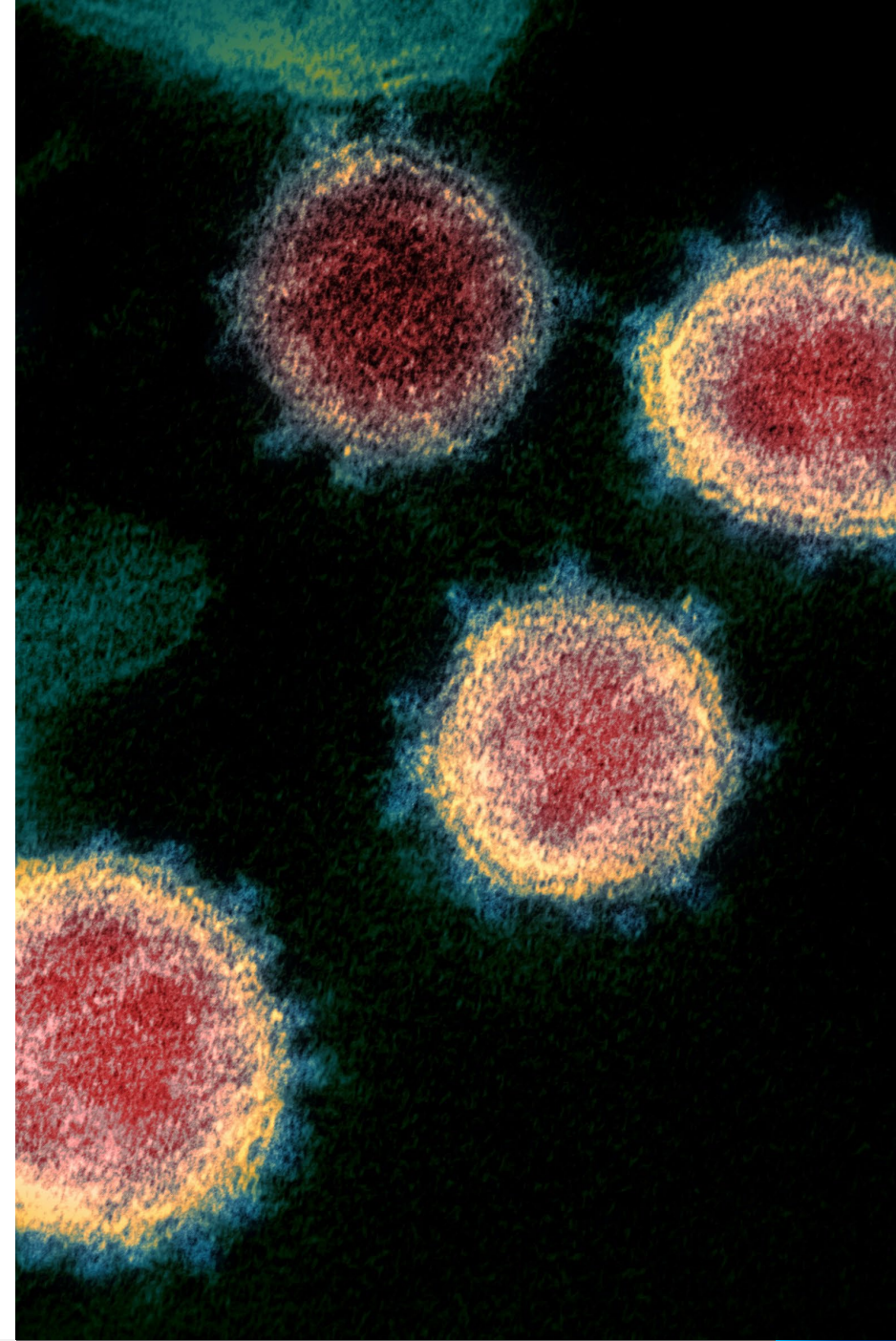Azwa bin Kamaruddin (WQD170089)
Kok Hon Loong (WQD170086)

# Assignment Milestones

WQD7005 Data Mining | Semester 2 Session 2019/2020



- Part A: (Group)
  - Web Crawling of Real-time Data

- **Part B: (Group)**
  - Management of Data using Hadoop Data Warehouse or Data Lake

- Part C: (Group)
  - Accessing and Processing of Data from Hadoop Data Warehouse or Data Lake

- Part D: (Individual)
  - Interpretation and Communication of Data Insights

- Part E: (Group)
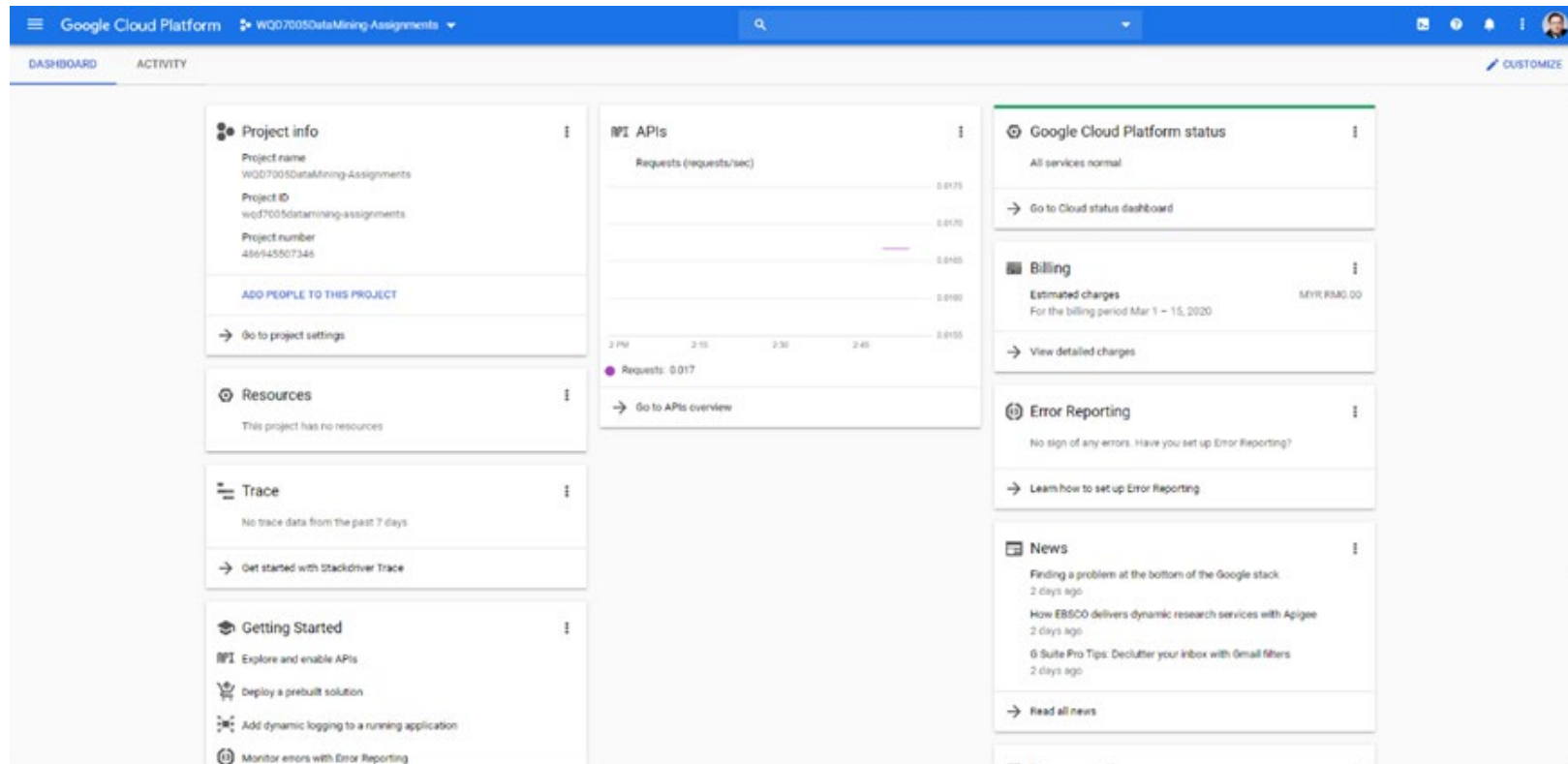  - Deployment of the Data Mining Results on Web (Flask) and Mobile Application (Kivy)

# Assignment Background

After we have acquired the data from the websites using the web crawling method as explained in Part A of the assignment, our group will be setting up a storage repository so that we can proceed in processing the data over and over to extract more data as we learn more about the contents.

- In this assignment, our team have decided to proceed in setting up the Hadoop Data Warehouse leveraging on the Google Cloud Platform (GCP).

- We will be using Hive as the SQL-like scripting language interface to query and analyze the data acquired from the websites in Hadoop HDFS, and the scripting code is also uploaded in our group assignment GitHub at the link below:
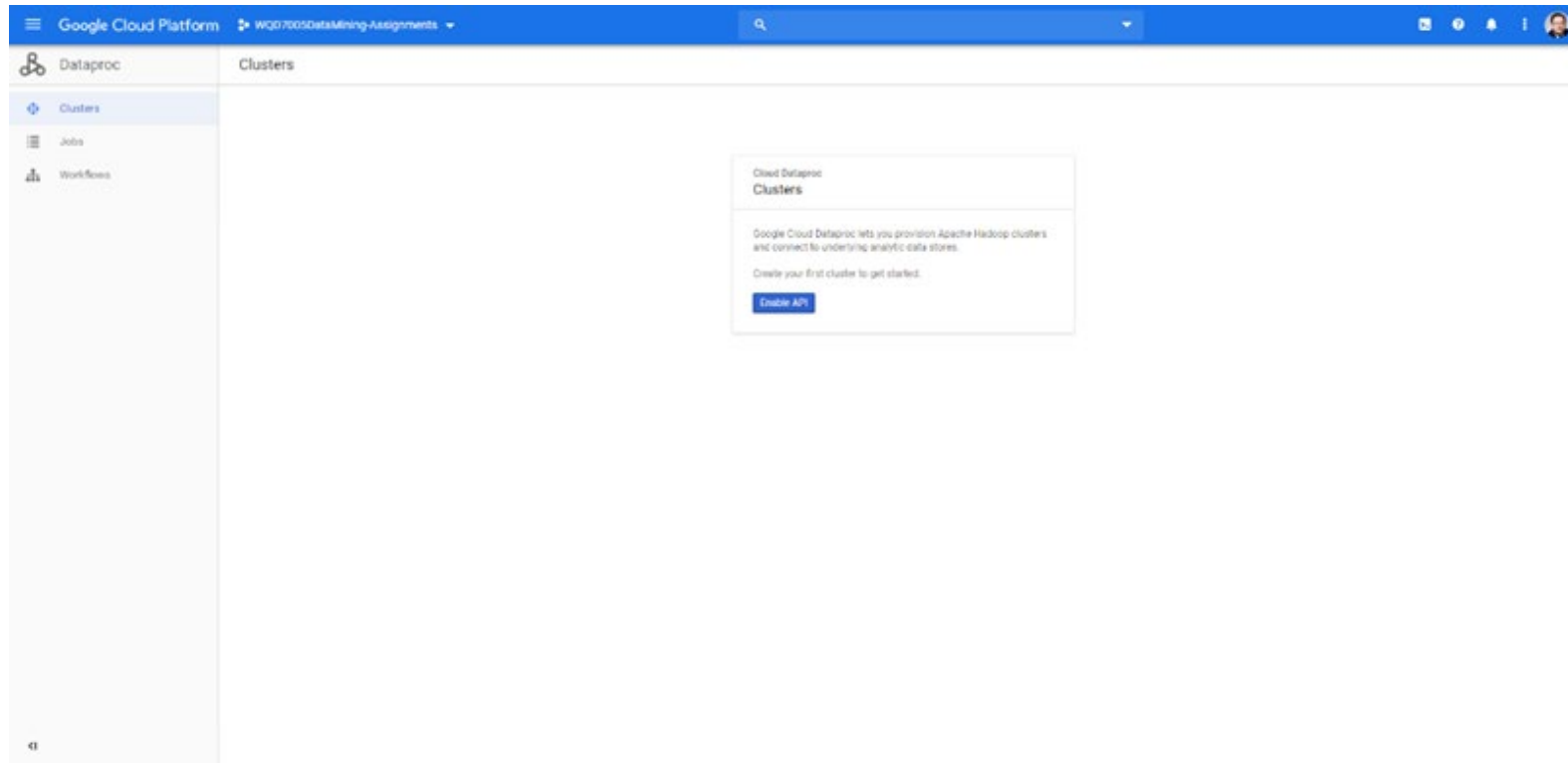
  o https://github.com/scholarazwa/wqd7005-assignment



Warehouse vs. Data Lakes

# **Configuring of the Hadoop & Hive Data Warehouse Cluster**
## using Google Cloud Platform (GCP)



- For this course assignment, we have created a new project for our team to work on using GCP.

- The screen on the left illustrate a brand new project dashboard that our team have created.

# Configuring of the Hadoop & Hive Data Warehouse Cluster
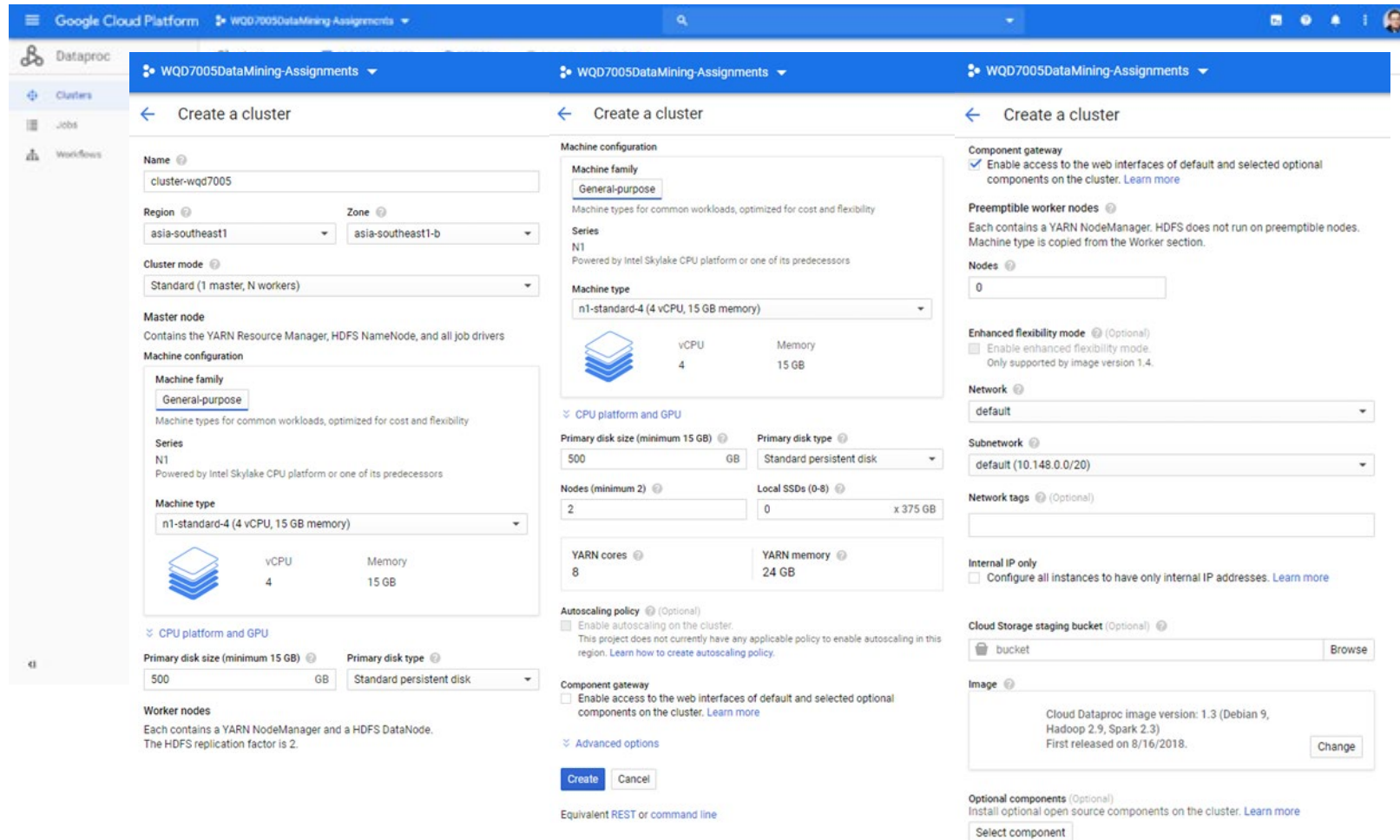using Google Cloud Platform (GCP)



- To configure Hadoop and Hive Data Warehouse in GCP, click on the *Navigation Menu* at the top left corner.

- Scroll all the way down until you reach to the *Big Data* section, and then click on the *Dataproc* to create the cluster.

- Next, click on *Enable API* button to proceed.

# Configuring of the Hadoop & Hive Data Warehouse Cluster
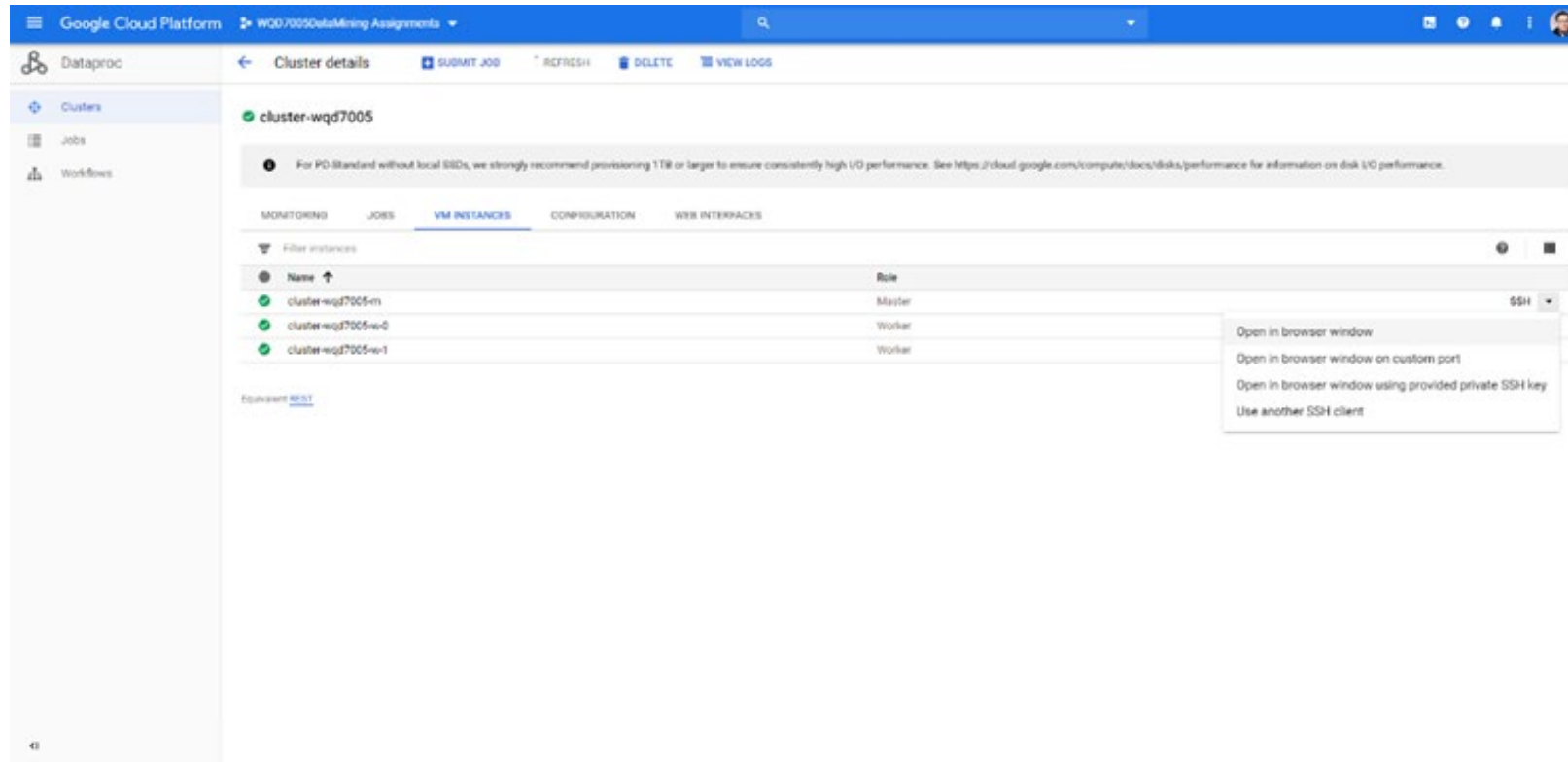## using Google Cloud Platform (GCP)



- Once provisioning the Apache Hadoop is initiated in GCP, click on the **Create cluster** button.

- Our team has named our assignment cluster as **cluster-wqd7005** and change the **region** and **zone** align to our current location.

- Expand the **Advanced options** and click on the **Enable access to web interface** checkbox. Leave the rest of the setting as default.

- Next, click on **Create** to proceed.

# Configuring of the Hadoop & Hive Data Warehouse Cluster
using Google Cloud Platform (GCP)



- It will takes a while for clusters to be created and once it's completed, the *green checkbox* as shown will appear.

- To verify on the cluster provisioned, click on the *cluster name* and it will show the Hadoop services are started.

- Click on the *Cluster* on the left, the Master and Workers nodes should be active as well.

- To test out the service, click on the *VM Instances* tab and from the *SSH* pull-down menu on the far right, select *Open in browser window*.

# Configuring of the Hadoop & Hive Data Warehouse Cluster
using Google Cloud Platform (GCP)



- The *connection to the VM instances* processes kicks in.

- Once the connection is done, the *SSH command line* screen will appear.

- To test if the Hive Data Warehouse is created, type *Hive* from the command line.

- From here, we can start the *Hive SQL-like script* to load the data we have acquired from web crawling.

# Using Hive to Query Hadoop Data Warehouse Cluster
using Google Cloud Platform (GCP)



- After the crawl datasets are stored in Hadoop Data Warehouse, we can start querying the dataset by creating a table using the Hive SQL-like script in Hive shell.

  CREATE TABLE asean_confirmed_cases

  (CaseDate STRING,

  Thailand INT, Singapore INT, Malaysia INT, Cambodia INT,

  Philippines INT, Indonesia INT, Brunei INT, Vietnam INT)

  ROW FORMAT DELIMITED FIELDS TERMINATED BY ','

  STORED AS TEXTFILE

  TBLPROPERTIES("skip.header.line.count"="1");

- Next, we load the dataset from the Hadoop Data Warehouse to the Hive Table

  LOAD DATA INPATH '/user/honloong/t_confirmed_cases_asean.csv'

  OVERWRITE INTO TABLE asean_confirmed_cases;

- To test out the query from the Hive Table, we can use the following SQL-like scripts:

  SELECT * FROM asean_confirmed_cases;

  Or

  SELECT CaseDate, Malaysia FROM asean_confirmed_cases;

# Summary

- To acquire content from a large number of data sources, we need to also prepare the data acquisition and ingestion tools.

- While data scraping can happen in any data array and can be done manually, web scraping or crawling takes place only on the web pages and is performed by special robots i.e. crawlers/scrapers.

- It's important to start with an agile, flexible and adaptable data repository and can rapidly adapt to changes with the application stacks of choices.

- However, there are different characteristics to consider when choosing the repository to manage the data acquired.

| Characteristics | Data Warehouse | Data Lake |
|---|---|---|
| Data | Relational from transactional systems, operational databases, and line of business applications | Non-relational and relational from IoT devices, web sites, mobile apps, social media, and corporate applications |
| Schema | Designed prior to the DW implementation (schema-on-write) | Written at the time of analysis (schema-on-read) |
| Price/Performance | Fastest query results using higher cost storage | Query results getting faster using low-cost storage |
| Data Quality | Highly curated data that serves as the central version of the truth | Any data that may or may not be curated (ie. raw data) |
| Users | Business analysts | Data scientists, Data developers, and Business analysts (using curated data) |
| Analytics | Batch reporting, BI and visualizations | Machine Learning, Predictive analytics, data discovery and profiling |

- Depending on the requirements, typically it will require both data warehouse and data lake to serve different needs and use cases as illustrated in the table above.

# Thank You

Azwa Kamaruddin & Hon-Loong Kok (HL)