




# TEACHING BIG DATA WITH LIMITED RESOURCES: PRACTICAL LESSONS FROM A SCALED-DOWN LAB

**Assoc. Prof. Dr. Azlan Ismail**

- Faculty of Computer and Mathematical Sciences, Universiti Teknologi MARA
  - Institute for Big Data Analytics and Artificial Intelligence, Universiti Teknologi MARA
- 

# WHAT THIS TALK IS ABOUT

- ❑ Sharing experience from teaching a data science technology course
  - Students are a mix of background, not all from computer science.
  - Course focuses on real big data tasks using Cloudera Hadoop tools.
- ❑ Sharing experience of using a scaled-down lab
  - Lab refers to: **a virtualized teaching environment that simulates real big data systems.**
  - The lab setup is modest, but designed to reflect real-world system behavior and performance patterns.



# WHY WE BUILT THIS MODEST INFRASTRUCTURE

1. Concept-heavy → students see theory, not systems
2. Abstract infra → Real infrastructure remains invisible
3. Simplified tools → Systems thinking & troubleshooting hidden
4. Cloud platforms → Great for production, but abstract away the details we want students to see



Taken from: <https://www.linkedin.com/pulse/understanding-types-scalability-storage-bq6zf/>



Taken from: <https://www.istockphoto.com/photos/distributed-computing>

# A MODEST LAB THAT MIRRORS REAL-WORLD SYSTEMS

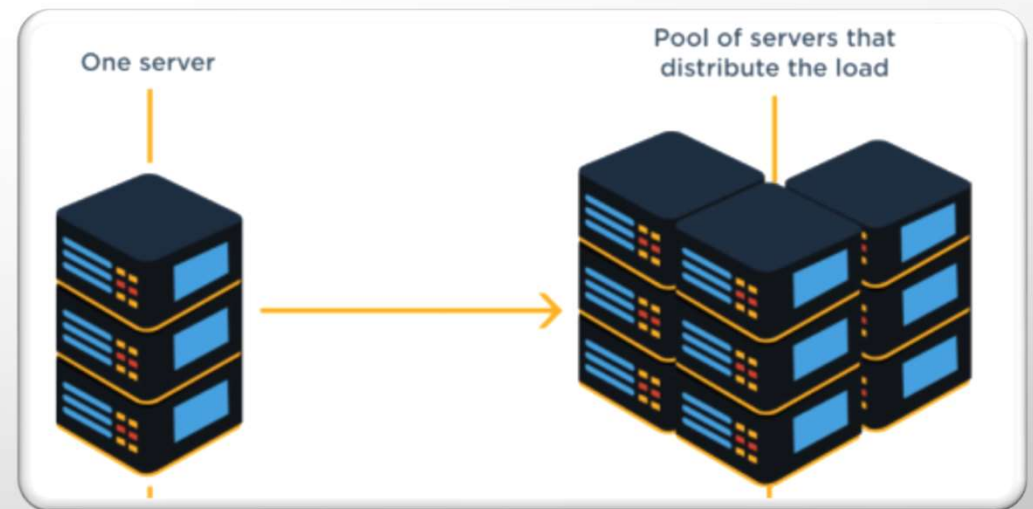
## What we set up?

- Designed a **teaching lab environment** to simulate real big data analytic workflows.
- Deployed a computing cluster with vendor support, which allowed us to focus on designing for student learning.
- **Platform:**
  - Cloudera hadoop distribution (CDH)
  - Tools: hive, spark, impala, sqoop, HDFS, YARN, etc
- **Infrastructure:**
  - 3 physical servers with 5 virtual machines
  - Master node + data nodes configured for cluster-based processing
- **Student access modes:**
  - Web UI (for hive, impala, sparksql, cloudera manager)
  - SSH terminal (for system-level exploration)
  - FTP client (for file exchange)

# A MODEST LAB: SOFTWARE PLATFORM - HADOOP

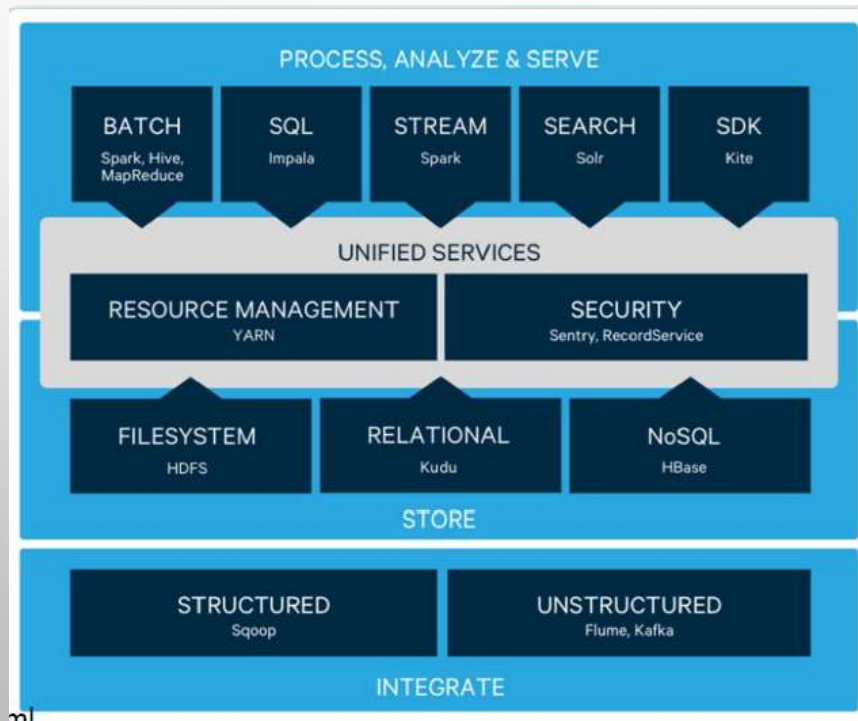
- **What is Hadoop?**

- One of the first widely adopted open-source frameworks for big data.
- Provides the foundation for modern big data platforms.
- Makes it possible to store and process massive datasets across many machines.
- Key idea: scale out by adding nodes, not scale up with one big server.



# A MODEST LAB: SOFTWARE PLATFORM

## Cloudera Hadoop distribution (CDH) Platform



### What ?

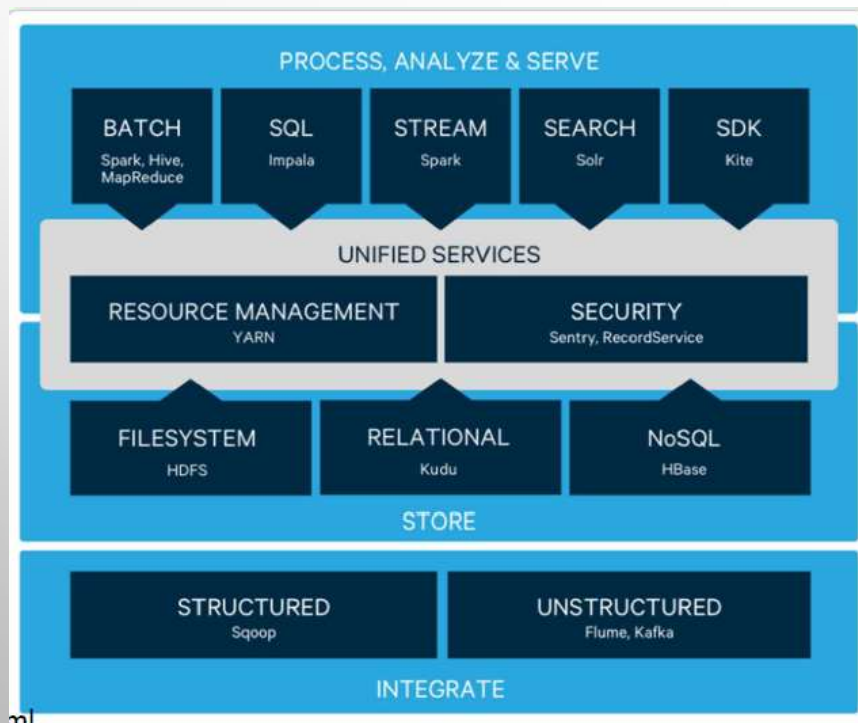
- ☐ An integrated Big Data ecosystem (Hive, Spark, Impala, Sqoop, HDFS, YARN).
- ☐ Supports the full workflow: ingest → store → process → analyze.
- ☐ Provides students with a realistic environment, not just single tools.

### Three key software layers:

- ☐ **Integrate** – software to transfer data into ecosystem
- ☐ **Store** – software to store and manage data
- ☐ **Process, Analyze & Serve** – software to process and analyze data

# A MODEST LAB: WHY CLOUDERA HADOOP

## Cloudera Hadoop distribution (CDH) Platform



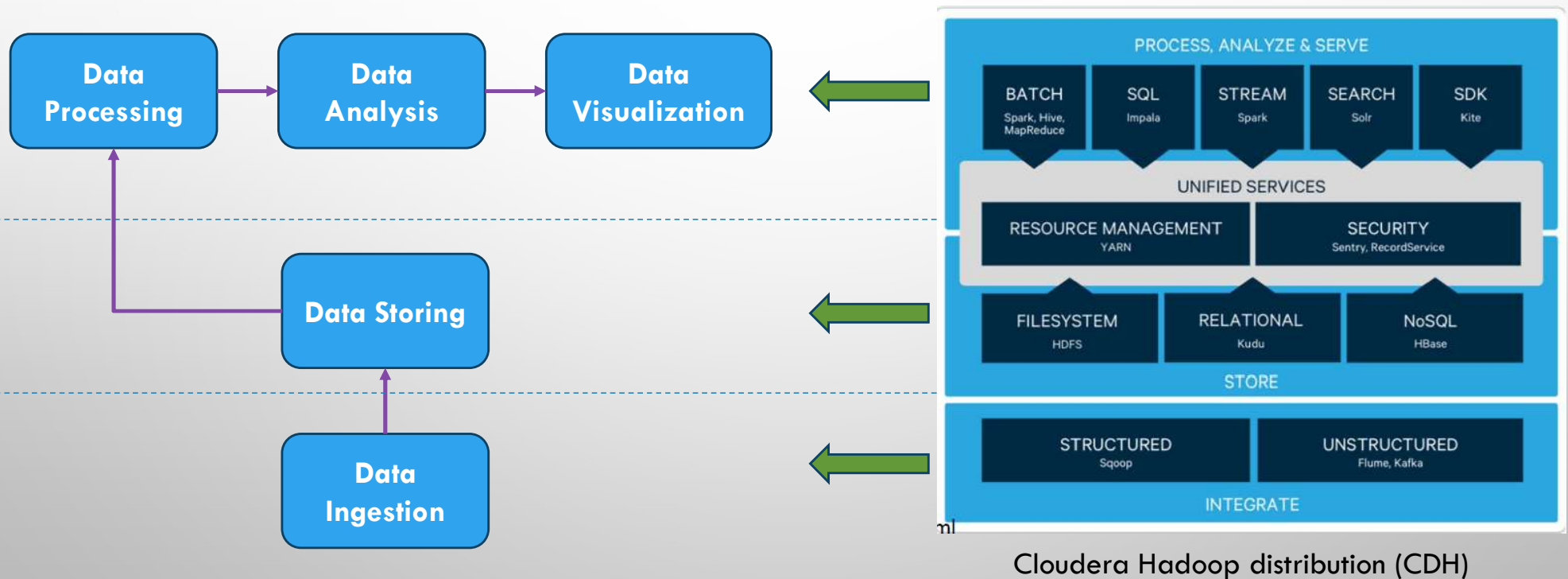
### Why ?

- ☐ Shows ecosystem integration, not isolated skills.
- ☐ Reflects industry practice but scaled down for education.
- ☐ Cloudera Manager made system health and job tracking visible.
- ☐ Reusable and repeatable across courses and student cohorts.

### Other reasons:

1. manageable for teaching.
2. no messy installations on student laptops.

# A MODEST LAB: CLOUDERA HADOOP SUPPORTS BIG DATA ANALYTICS PIPELINE





# A MODEST LAB: MAPPING CLOUDER HADOOP WITH BIG DATA ANALYTICS PIPELINE

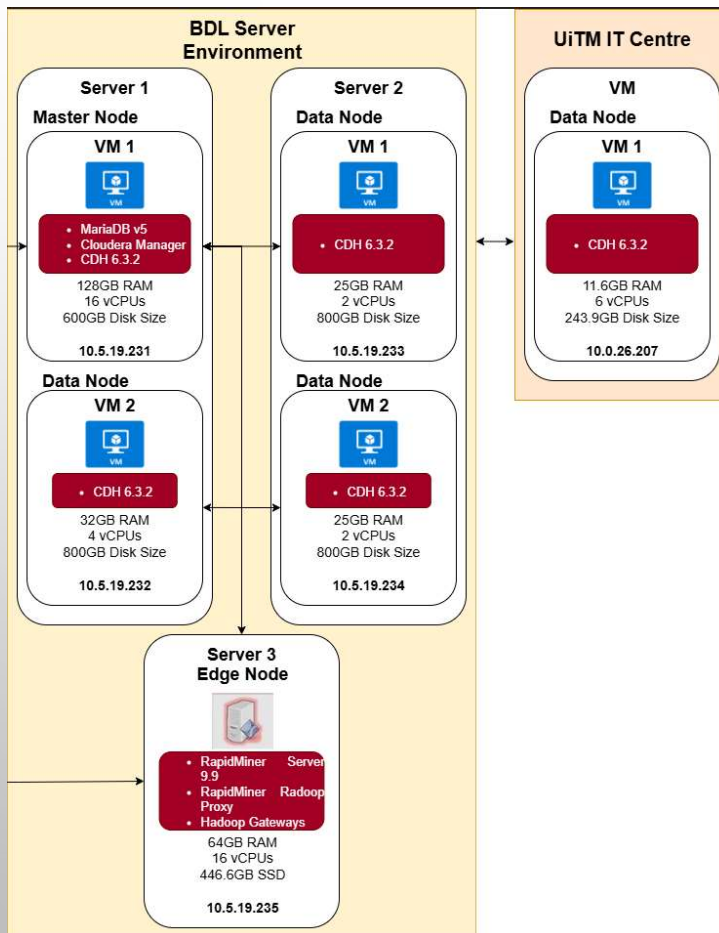
Pipeline Stage	Purpose	Cloudera Tool(s)
1. Data Ingestion	Ingest data from databases, logs, APIs	Sqoop (RDBMS → HDFS) Flume (logs → HDFS) Kafka (real-time streams) Hue Upload
2. Data Storing	Store raw or semi-structured data	HDFS (distributed file system) HBase (NoSQL, columnar)
3. Data Processing	Clean, transform, join, filter data	Hive (HiveQL) Spark (batch or structured streaming)
4. Data Analysis	Run SQL, ML, or statistical analysis	Impala (interactive SQL) Spark MLlib Hive (batch analytics)
5. Data Visualization	Present insights to users/admins	Hue (query interface + visualizations)

# A MODEST LAB THAT MIRRORS REAL-WORLD SYSTEMS

## What we set up?

- Designed a **teaching lab environment** to simulate real big data analytic workflows.
- Deployed a computing cluster with vendor support, which allowed us to focus on designing for student learning.
- **Platform:**
  - Cloudera hadoop distribution (CDH)
  - Tools: hive, spark, impala, sqoop, HDFS, YARN, etc
- **Infrastructure:**
  - 3 physical servers with 5 virtual machines
  - Master node + data nodes configured for cluster-based processing
- **Student access modes:**
  - Web UI (for hive, impala, sparksql, cloudera manager)
  - SSH terminal (for system-level exploration)
  - FTP client (for file exchange)

# A MODEST LAB: THE INFRASTRUCTURE



3 physical servers with 5 virtual machines

- ❑ **Master Node** – the *coordinator*, keeps track of what job goes where.
- ❑ **Data Nodes** – the *workers*, actually store and process pieces of the data.
- ❑ **Edge Node** – the *entry point*, where students accessed tools and submitted jobs.

# A MODEST LAB: THE INFRASTRUCTURE (WHY)

## Why Set Up This Infrastructure (3 servers, 5 VMs)?

### 1. Pedagogical Need

1. Concept-heavy courses and notebooks hide the system layer.
2. Students rarely see how distributed systems *actually* work.
3. The infra made the *invisible visible*: nodes, job scheduling, monitoring.

Actions for Selected ▾								
<input type="checkbox"/>	Status ▾	Name	IP	Roles	Commission State	Last Heartbeat	Load Average	Disk Usage
<input type="checkbox"/>	✓	bigdatalab-cdh-dn1.uitm.edu.my	10.5.19.232	5 Role(s)	Commissioned	10.71s ago	0.00 0.01 0.05	320.7 GiB / 791.6 GiB
<input type="checkbox"/>	✓	bigdatalab-cdh-dn2.uitm.edu.my	10.5.19.233	5 Role(s)	Commissioned	4.19s ago	0.04 0.07 0.10	174.4 GiB / 791.6 GiB
<input type="checkbox"/>	✓	bigdatalab-cdh-dn3.uitm.edu.my	10.5.19.234	5 Role(s)	Commissioned	4.84s ago	0.25 0.11 0.09	270.9 GiB / 791.6 GiB
<input type="checkbox"/>	✓	bigdatalab-cdh-mn1.uitm.edu.my	10.5.19.231	22 Role(s)	Commissioned	7.57s ago	0.27 0.32 0.41	125.8 GiB / 591.7 GiB
<input type="checkbox"/>	✓	bigdatalab-rm-en1.uitm.edu.my	10.5.19.235	11 Role(s)	Commissioned	8.21s ago	0.01 0.03 0.05	156.1 GiB / 413 GiB
<input type="checkbox"/>	✓	huefskm.uitm.edu.my	10.0.26.207	1 Role(s)	Commissioned	7.78s ago	0.00 0.04 0.05	47.7 GiB / 243.9 GiB

<https://bigdatalab-cdh-mn1.uitm.edu.my:7183/cmf/login>

hadoop										All Applications									
Cluster										Cluster Metrics									
About Nodes Node Labels Applications										Apps Submitted Apps Pending Apps Running Apps Completed Containers Running Memory Used Memory To									
NEW NEW SAVING SUBMITTED ACCEPTED RUNNING FINISHED FAILED KILLED Scheduler										2 0 0 2 0 0 B 121 GB									
Tools										Cluster Nodes Metrics									
										Active Nodes Decommissioning Nodes Decommissioned Nodes Lost Nodes									
										4 0 0 0 0 0									
										User Metrics for dr.who									
										Apps Submitted Apps Pending Apps Running Apps Completed Containers Running Containers Pending Containers Reserved Memory User									
										0 0 0 0 0 0 0 0 0 B									
										Scheduler Metrics									
										Scheduler Type Scheduling Resource Type Minimum Allocation M									
										Fair Scheduler [memory-mb (unit=M), vcores] <memory 512, vCores 1> <memory 8192, vCores 1>									
										Show 20 ▾ entries									
										ID User Name Application Type Queue Application Priority StartTime LaunchTime FinishTime State FinalStatus									
										application_1755574613124_0002 hive SELECT actor1code, actor1name, ...10 (Stage-1) MAPREDUCE root users.student4 0 Sun Aug 31 01:40:34 +0800 2025 Sun Aug 31 01:40:35 +0800 2025 Sun Aug 31 01:42:12 +0800 2025 FINISHED SUCCEEDED									
										application_1755574613124_0001 hive SELECT actor1code, actor1name, ...10 (Stage-1) MAPREDUCE root users.student4 0 Sun Aug 31 01:38:21 +0800 2025 Sun Aug 31 01:38:22 +0800 2025 Sun Aug 31 01:39:50 +0800 2025 FINISHED SUCCEEDED									
										Showing 1 to 2 of 2 entries									

<https://bigdatalab-cdh-mn1.uitm.edu.my:8090/cluster>

# A MODEST LAB: THE INFRASTRUCTURE (WHY)

Why Set Up This Infrastructure (3 servers, 5 VMs)?

## 2. Alignment with Hadoop's Core Features

1. **Distributed:** Even a small cluster showed data split across nodes.
2. **Reliable:** Students observed fault tolerance *and* the reality of service restarts.
3. **Commodity Hardware:** Demonstrated that Big Data principles can be learned without enterprise infrastructure.



# A MODEST LAB: THE INFRASTRUCTURE (WHY)

**Why Set Up This Infrastructure (3 servers, 5 VMs)?**

## **3. Teaching Practicality**

1. Small enough to manage with limited budget and support.
2. Big enough to mimic industry workflows (HDFS + YARN + Hive/Spark).

## **4. Scalable for Education**

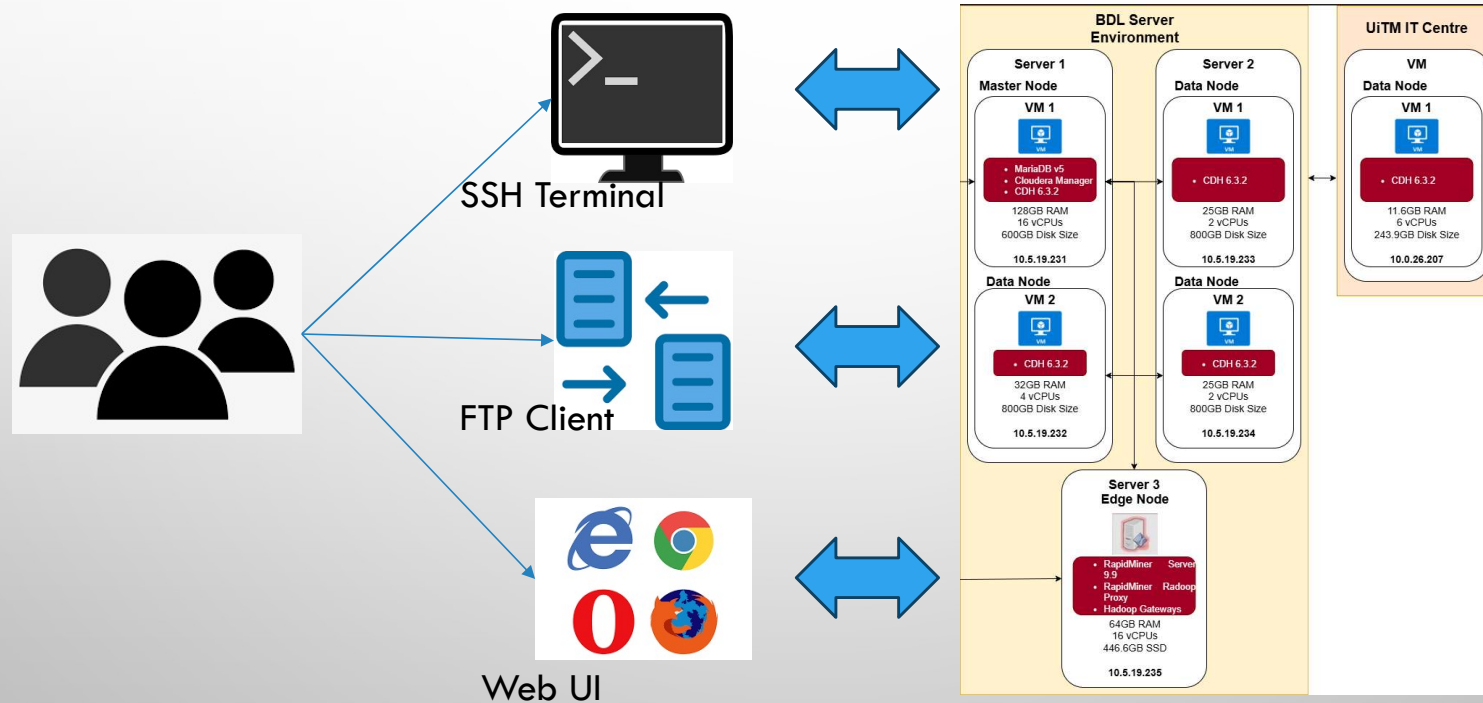
1. Reusable across courses, levels, and cohorts.
2. Flexible for both in-class and remote teaching.
3. Supports not just one-off labs, but assignments, projects, and FYPs.

# A MODEST LAB THAT MIRRORS REAL-WORLD SYSTEMS

## What we set up?

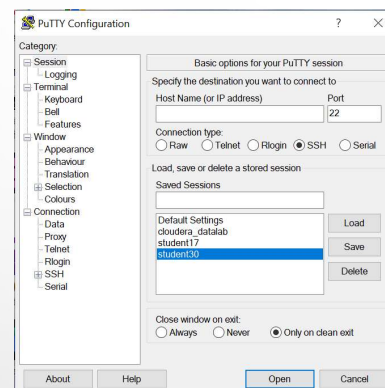
- Designed a **teaching lab environment** to simulate real big data analytic workflows.
- Deployed a computing cluster with vendor support, which allowed us to focus on designing for student learning.
- **Platform:**
  - Cloudera hadoop distribution (CDH)
  - Tools: hive, spark, impala, sqoop, HDFS, YARN, etc
- **Infrastructure:**
  - 3 physical servers with 5 virtual machines
  - Master node + data nodes configured for cluster-based processing
- **Student access modes:**
  - Web UI (for hive, impala, sparksql, cloudera manager)
  - SSH terminal (for system-level exploration)
  - FTP client (for file exchange)

# A MODEST LAB: STUDENT ACCESS MODES





# A MODEST LAB: STUDENT ACCESS MODES



Using putty

```
student30@bigdatalab-rm-en1:~$ ls
categories.java
check_hbase.py
checkpoints
check_pyspark.py
codegen_customers.java
consumer_task_events.py
consumer_to_hive.py
consume_tweets_hdfs.py
consume_tweets_hive.py
```

Accessing linux home directory

Accessing HDFS directory

```
[student30@bigdatalab-rm-en1 ~]$ hdfs dfs -ls
Found 21 items
drwxr-xr-x - student30 student30 0 2025-06-15 19:00 .Trash
drwxr-xr-x - student30 student30 0 2025-08-17 16:56 .sparkStaging
drwxr-xr-x - student30 student30 0 2025-05-09 22:50 .staging
drwxrwx--- - student30 student30 0 2025-06-21 12:15 checkpoint
drwxr-xr-x - student30 student30 0 2025-06-21 11:50 checkpoint
ory
drwxrwx--- - student30 student30 0 2025-06-21 11:53 checkpoint
drwxr-xr-x - student30 student30 0 2025-04-19 12:30 countfrom
drwxr-xr-x - student30 student30 0 2025-04-14 11:30 counts
```

```
Spark session available as 'spark'.
Welcome to

Spark version 2.4.0-cdh6.3.2

Using Scala version 2.11.12 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_181)
Type in expressions to have them evaluated.
Type :help for more information.

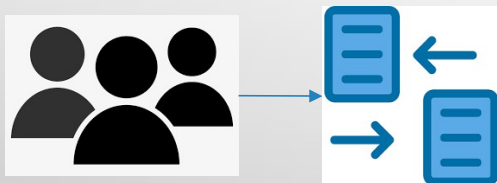
scala>
```

Accessing Spark

```
[azlanismail@bigdatalab-rm-en1 ~]$ beeline -u jdbc:hive2://bigdatalab-cdh-mni.uitm.edu.my:10000 -n cloudera -p cloudera
WARNING: Use "yarn jar" to launch YARN applications.
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jarfile:/opt/cloudera/parcels/CDH-6.3.2-1-cdh6.3.2-p6-1605554/jars/log4j-slf4j-impl-2.8.2-jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jarfile:/opt/cloudera/parcels/CDH-6.3.2-1-cdh6.3.2-p6-1605554/jars/slf4j-log4j12-1.7.25-jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Connecting to jdbc:hive2://bigdatalab-cdh-mni.uitm.edu.my:10000
Connected to: Apache Hive (version 2.1.1-cdh6.3.2)
Driver: Hive JDBC (version 2.1.1-cdh6.3.2)
Transaction isolation: TRANSACTION_SERIALIZABLE READ
Beeline version 2.1.1-cdh6.3.2 by Apache Hive
0: jdbc:hive2://bigdatalab-cdh-mni.uitm.edu.my>
```

Accessing Hive Beeline

# A MODEST LAB: STUDENT ACCESS MODES



PySpark – student30@10.5.19.235 – WinSCP

Local Mark Files Commands Tabs Options Remote Help

student30@10.5.19.235 X New Tab

Local directory: C:\Users\azlan\Documents\DataScienceProject\PySpark\

Name	Size	Type
spark_app_csv1.py	1 KB	Python File
spark_app_housing.py	2 KB	Python File
spark_correlation_housingdata.py	1 KB	Python File
spark_datahandling_housingdata.py	2 KB	Python File
spark_explore_housingdata.py	1 KB	Python File
spark_impute_housingdata.py	1 KB	Python File
spark_mllib.py	2 KB	Python File
spark_mllib_chart_save.py	2 KB	Python File
spark_mllib_csv.py	2 KB	Python File
spark_mllib_csv_chart_save.py	2 KB	Python File
spark_statistic_housingdata.py	1 KB	Python File
spark_streaming_app.py	2 KB	Python File
spark_structured_streaming.py	1 KB	Python File
spark_structured_streaming_check...	2 KB	Python File
spark_structured_streaming_check...	2 KB	Python File

0 B of 17.9 MB in 0 of 42

Remote directory: /home/student30/

Name	Size	Changed	Rights
check_hbase.py	1 KB	30/11/2024 2:33:26 PM	rw-rw-r-
check_pyspark.py	1 KB	6/12/2024 11:05:32 AM	rw-rw-r-
codegen_customers.java	29 KB	23/4/2022 6:54:27 AM	rw-rw-r-
consume_tweets_hdfs.py	3 KB	20/5/2023 9:23:32 PM	rw-rw-r-
consume_tweets_hive.py	3 KB	27/5/2023 5:39:57 PM	rw-rw-r-
consumer_task_events.py	4 KB	15/5/2023 10:32:20 AM	rw-rw-r-
consumer_to_hive.py	2 KB	18/5/2023 9:58:50 AM	rw-rw-r-
covid_19_data.csv	22,008 KB	24/6/2021 4:27:52 AM	rw-rw-r-
crypto_data.csv	10 KB	11/4/2025 9:13:10 PM	rw-rw-r-
customer.java	25 KB	11/4/2023 9:30:14 PM	rw-rw-r-
customers.avsc	2 KB	22/4/2022 10:30:20 PM	rw-rw-r-
customers.java	29 KB	7/11/2021 12:21:16 PM	rw-rw-r-
departments.java	12 KB	28/4/2022 3:53:12 PM	rw-rw-r-
derby.log	1 KB	3/11/2023 2:51:08 PM	rw-rw-r-
describe_output.txt	2 KB	12/11/2024 2:18:44 PM	rw-rw-r-
extract_data.py	1 KB	11/4/2025 8:27:25 PM	rw-rw-r-

0 B of 401 MB in 0 of 91

SFTP-3 0:00:40

17 hidden

Local directory

Remote directory

# A MODEST LAB: STUDENT ACCESS MODES



**HUE**  
Query. Explore. Repeat.

Username  
Password

Sign In

**data engineer/analyst  
perspective**



Hue interface showing a list of tables and a file browser view.

Tables (12) +

- customers
- customers\_buck\_orc
- customers\_bucketed
- customers\_new
- customers\_part\_buck
- movie\_rating
- movie\_rating\_av
- movie\_rating\_dynpart
- movie\_rating\_ex
- orders
- orders\_buck\_orc
- orders\_bucketed

File browser view showing a directory structure with files like .Trash, .sparkStaging, .staging, and checkpoint-csv.

<https://bigdatalab-rm-en1.uitm.edu.my:8889/hue/accounts/login?next=/>

**Cloudera Manager**

Username  
Password

☐ Remember me

Sign In

**System administrator  
perspective**



Cloudera Manager interface showing cluster status and charts.

Cluster 1 Status

CDH 6.3.2 (Parcels)

- 0 Hosts
- HBase
- HDFS
- Hive
- Hue
- Impala
- Kafka
- Oozie
- Sentry
- Spark
- Sqoop 1 Client

Charts: Cluster CPU, Cluster Disk IO

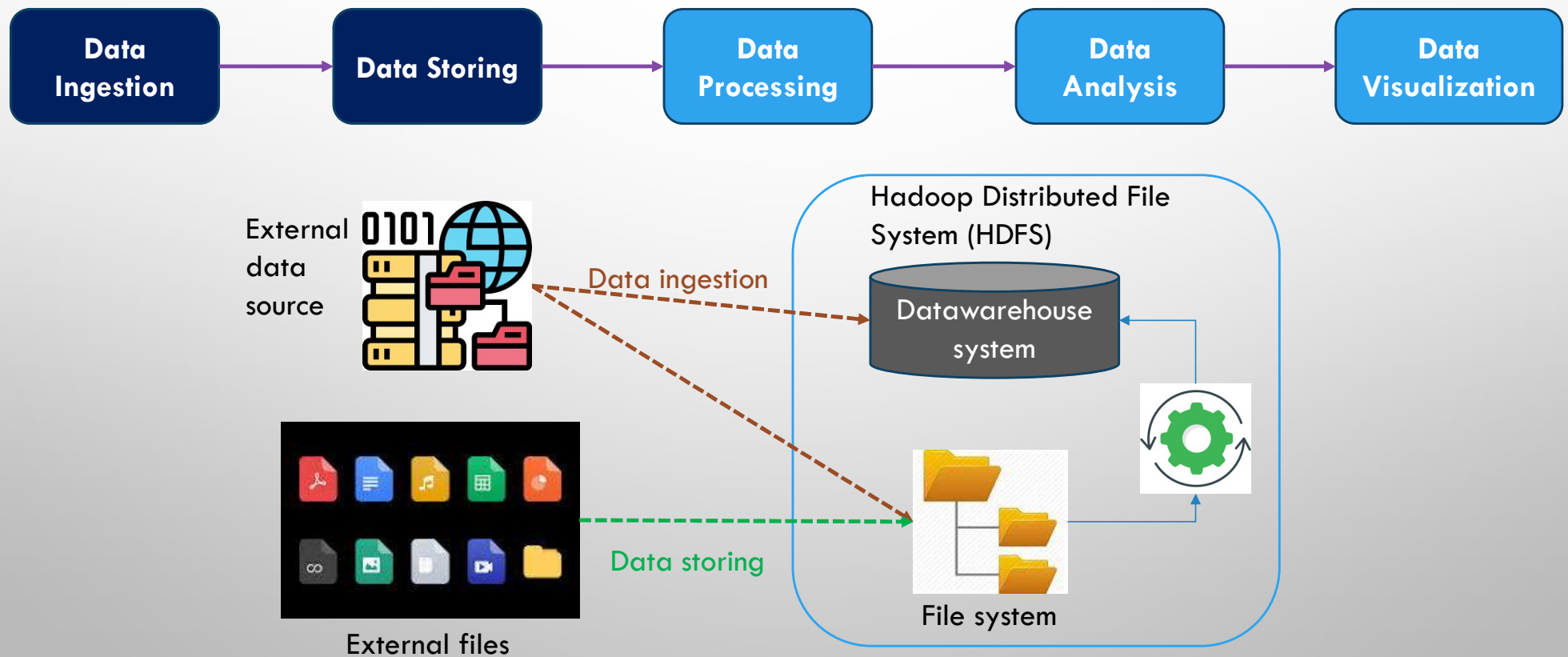
<https://bigdatalab-cdh-mn1.uitm.edu.my:7183/cm/#!/login>

# FROM DOING TO UNDERSTANDING TO ASSIGNMENTS



- ❑ **From Theory to Action** – students performed end-to-end Big Data tasks (ingest, store, process, analyze).
- ❑ **From Execution to Insight** – they reflected on results and connected theory with system behavior.
- ❑ **Structured Assignments** – structured tasks captured these insights through benchmarking and reporting.

# FROM THEORY TO ACTION



# FROM THEORY TO ACTION

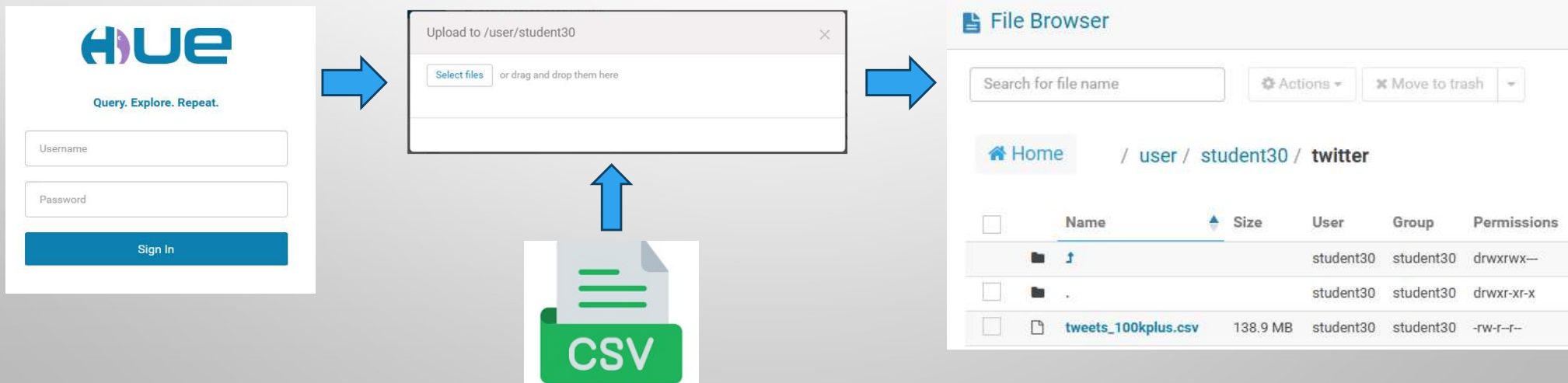


```
[student30@bigdatalab-rm-en1 ~]$ sqoop import --connect jdbc:mysql://bigdatalab-rm-en1:3306/retail_db --username student --password p@ssw0rd --table customers --target-dir /user/student30/customers_comp --compression-codec SnappyCodec --delete-target-dir
```

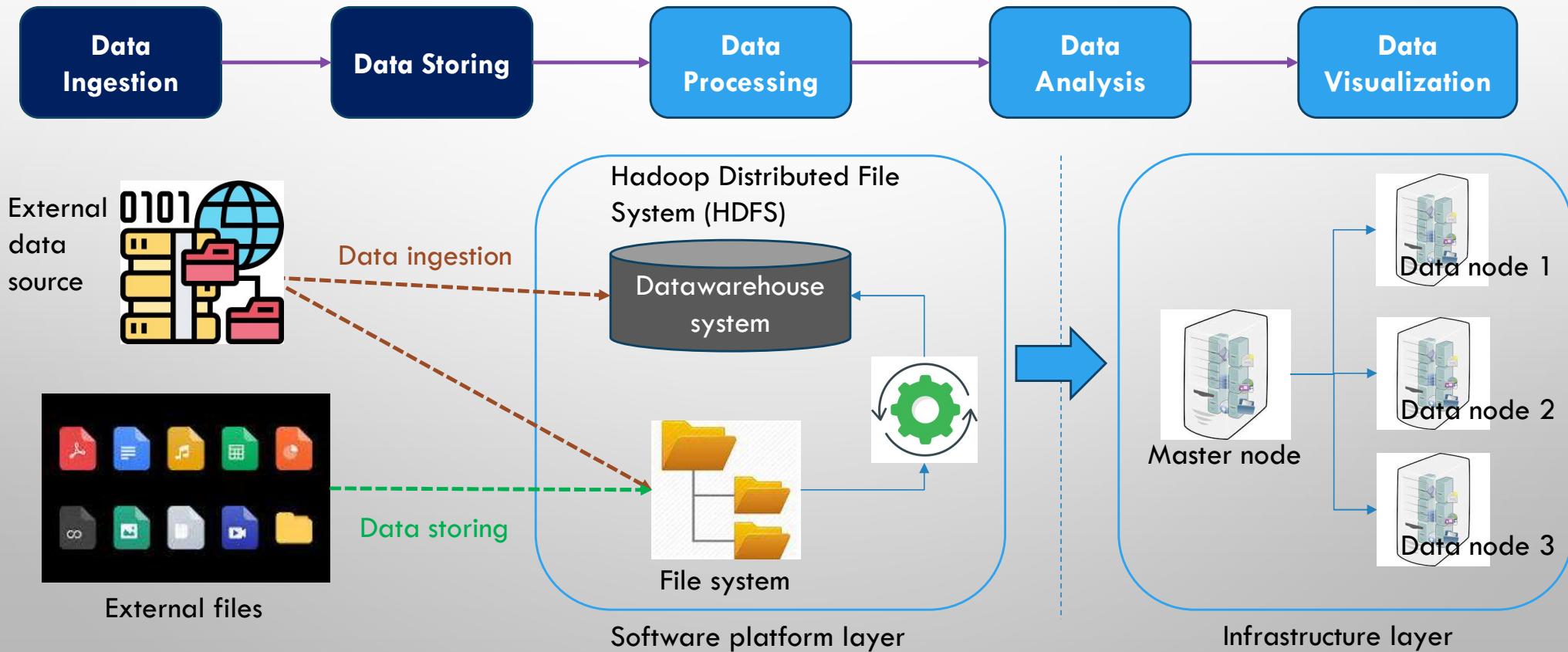


```
Map input records=12435
Map output records=12435
Input split bytes=479
Spilled Records=0
Failed Shuffles=0
Merged Map outputs=0
GC time elapsed (ms)=330
CPU time spent (ms)=4870
Physical memory (bytes) snapshot=1147584512
Virtual memory (bytes) snapshot=13925732352
Total committed heap usage (bytes)=1716518912
Peak Map Physical memory (bytes)=298045440
Peak Map Virtual memory (bytes)=3485007872
File Input Format Counters
  Bytes Read=0
File Output Format Counters
  Bytes Written=415825
22/04/22 22:19:58 INFO mapreduce.ImportJobBase: Transferred 406.0791 KB in
8 seconds (24.9925 KB/sec)
22/04/22 22:19:58 INFO mapreduce.ImportJobBase: Retrieved 12435 records.
```

# FROM THEORY TO ACTION



# FROM THEORY TO ACTION





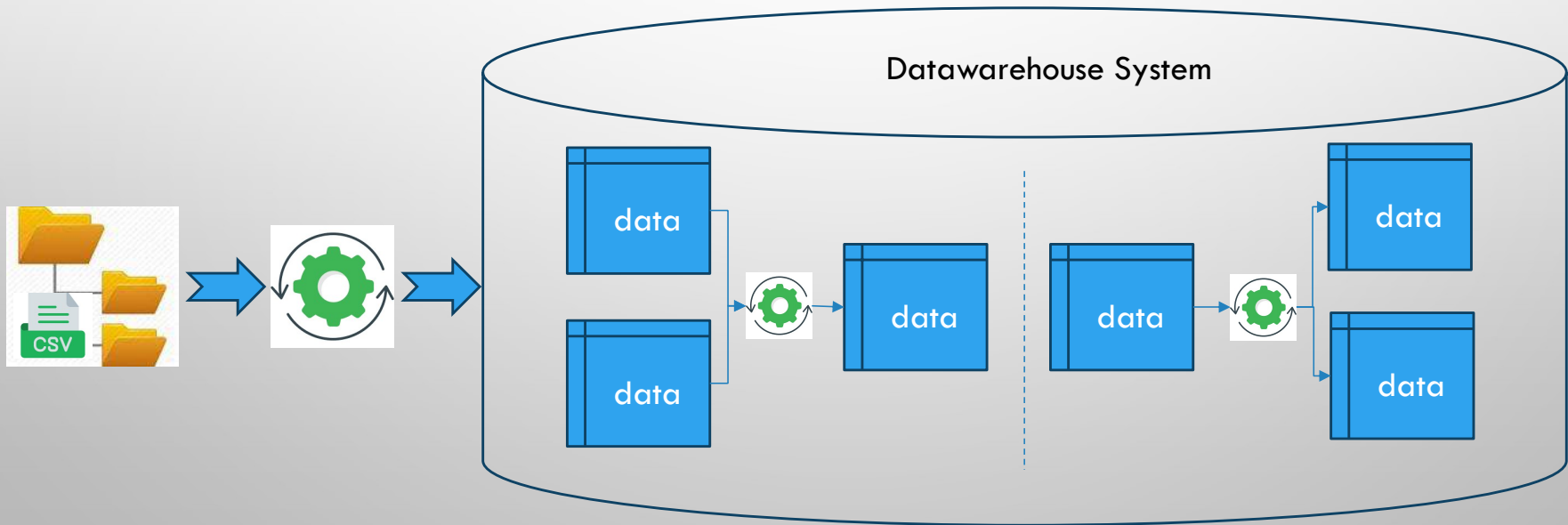
# FROM THEORY TO ACTION



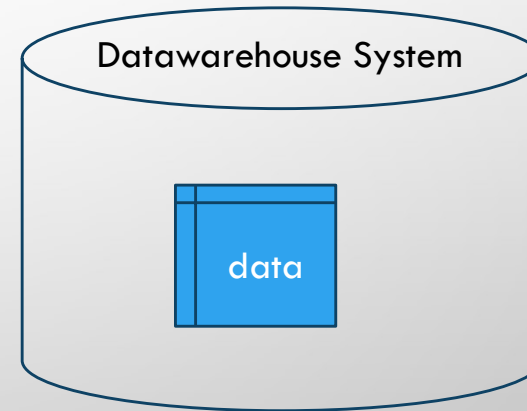
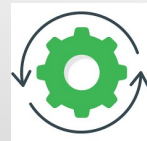
```
student30@bigdatalab-rm-en1:~  
ions=1&path=%2Fuser%2Fstudent30%2Ftwitter%2Ftweets  
FSCK started by student30 (auth:SIMPLE) from /10.5  
t Sun Sep 07 16:16:24 MYT 2025  
/user/student30/twitter/tweets_100kplus.csv 145628  
0 BP-57606532-10.5.19.204-1610901141600:blk_10833  
age[10.5.19.233:9866,DS-f3b2004f-6b99-4baf-8fc2-4f  
1 BP-57606532-10.5.19.204-1610901141600:blk_10833  
ge[10.5.19.233:9866,DS-f3b2004f-6b99-4baf-8fc2-4f
```

Status	Name	IP	R
✓	bigdatalab-cdh-dn1.uitm.edu.my	10.5.19.232	>
✓	bigdatalab-cdh-dn2.uitm.edu.my	10.5.19.233	>
✓	bigdatalab-cdh-dn3.uitm.edu.my	10.5.19.234	>
✓	bigdatalab-cdh-mn1.uitm.edu.my	10.5.19.231	>
✓	bigdatalab-rm-en1.uitm.edu.my	10.5.19.235	>
✓	huefskm.uitm.edu.my	10.0.26.207	>

# FROM THEORY TO ACTION



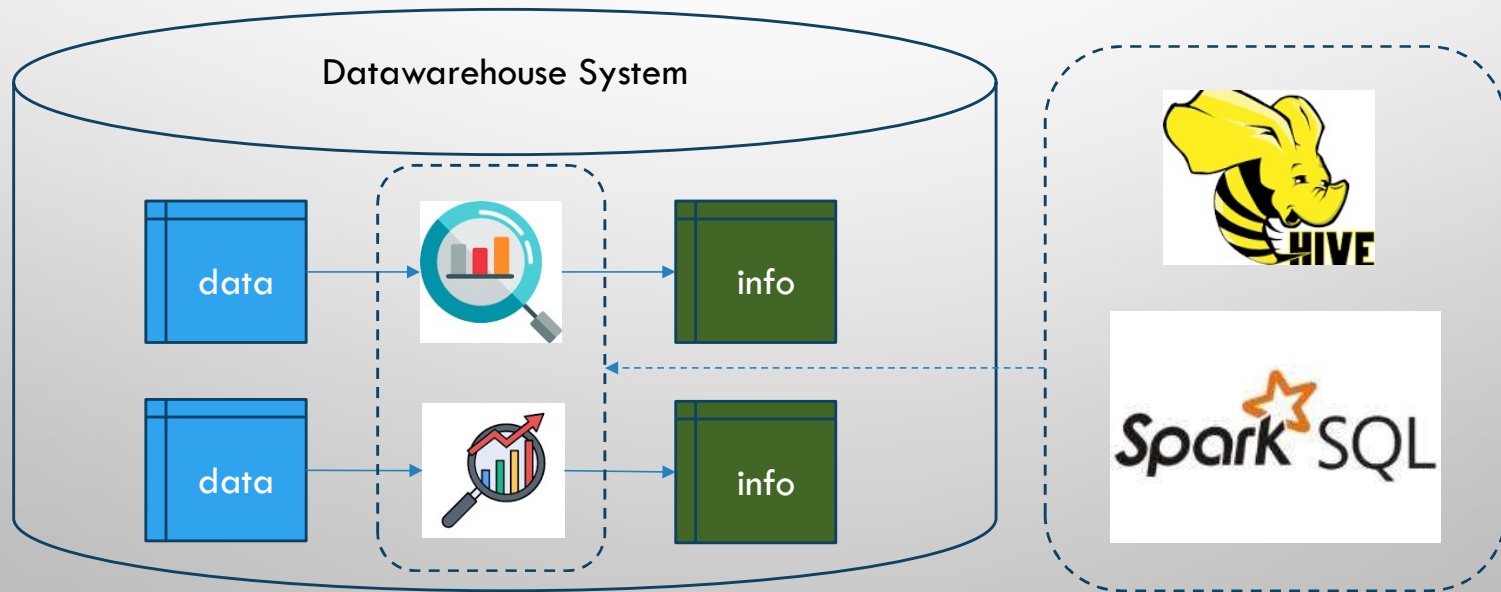
# FROM THEORY TO ACTION



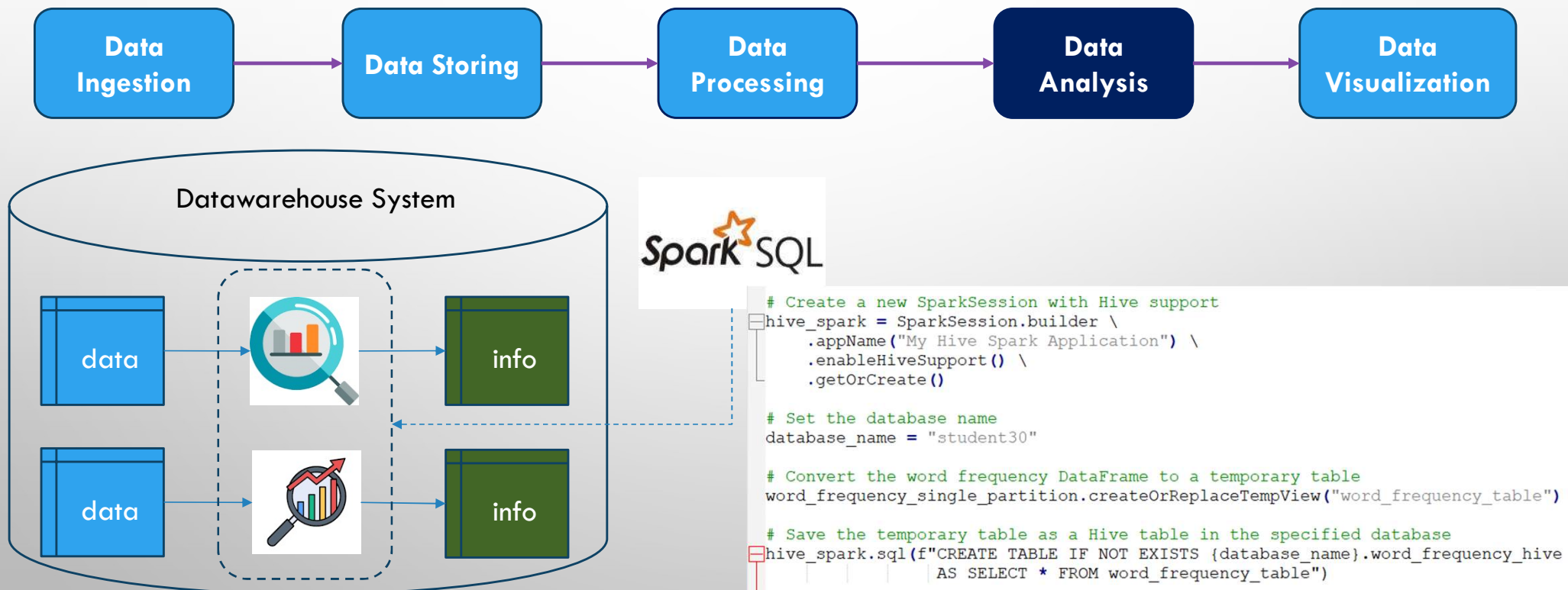
```
-mnl.uitm.edu.m> create external table movie_rating_ex (  
  .> userid int,  
  .> movieid int,  
  .> rating int,  
  .> unixtime string)  
  .> row format delimited  
  .> fields terminated by ','  
  .> location '/user/student30/movie_rating2';
```



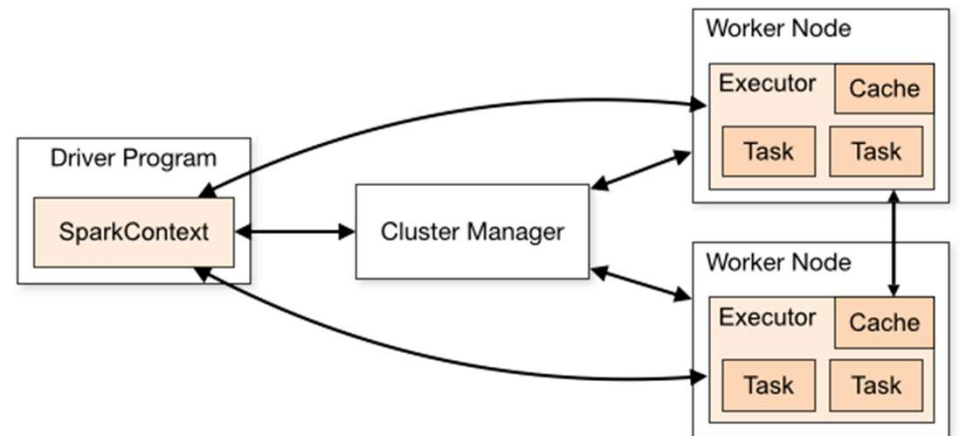
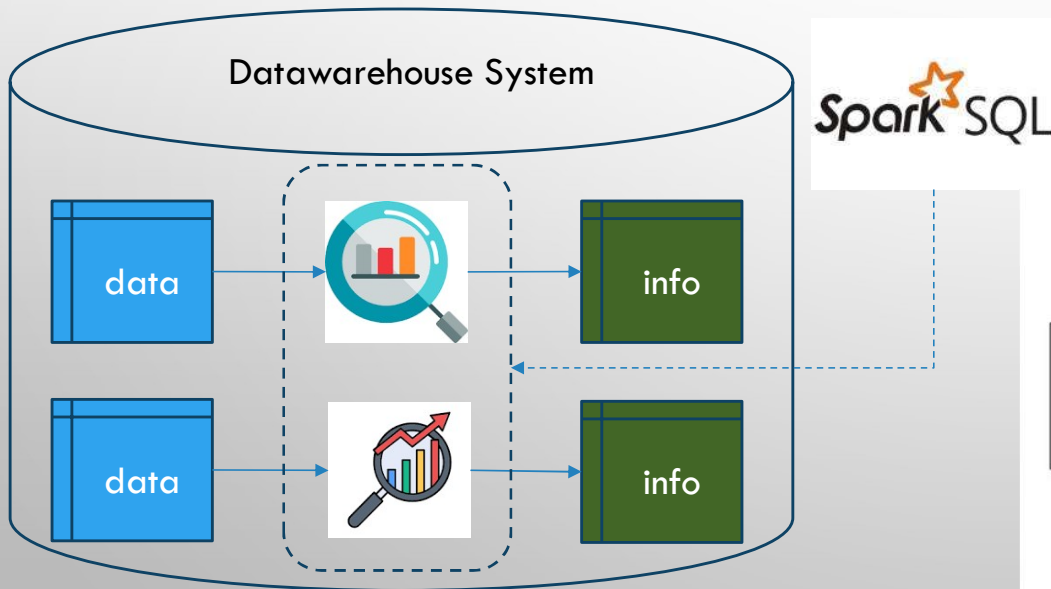
# FROM THEORY TO ACTION



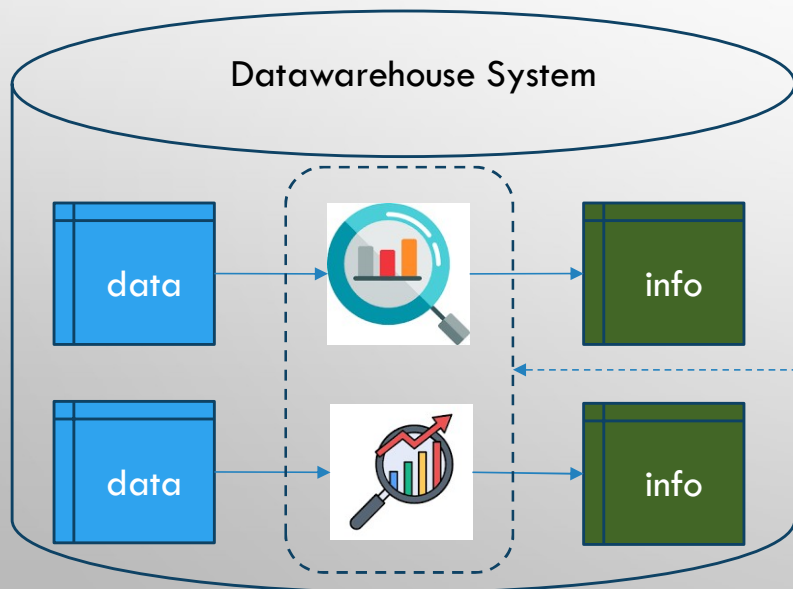
# FROM THEORY TO ACTION



# FROM THEORY TO ACTION



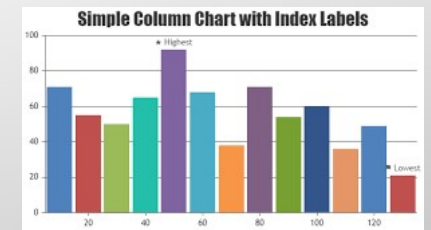
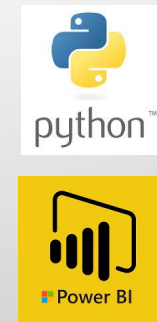
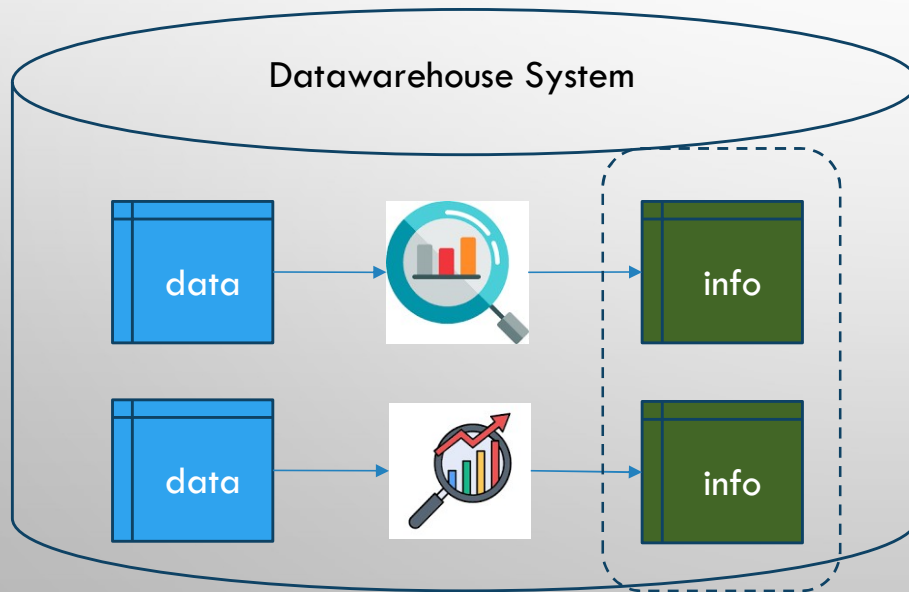
# FROM THEORY TO ACTION



Spark SQL

Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)
driver	bigdatalab-rm-en1.uitm.edu.my:40946	Active	0	0.0 B / 384.1 MB	0.0 B	0	0	0	0	0	0 ms (0 ms)
1	bigdatalab-rm-en1.uitm.edu.my:38982	Active	0	0.0 B / 413 MB	0.0 B	1	0	0	295	295	27 s (0.6 s)
2	bigdatalab-rm-en1.uitm.edu.my:40435	Active	0	0.0 B / 413 MB	0.0 B	1	0	0	283	283	26 s (0.5 s)
3	bigdatalab-rm-en1.uitm.edu.my:37471	Active	0	0.0 B / 413 MB	0.0 B	1	0	0	47	47	10 s (0.3 s)
4	bigdatalab-cdh-dn2.uitm.edu.my:41871	Active	0	0.0 B / 384.1 MB	0.0 B	1	0	0	0	0	0 ms (0 ms)
5	bigdatalab-cdh-dn1.uitm.edu.my:33286	Active	0	0.0 B / 384.1 MB	0.0 B	1	0	0	1	1	9 s (0.2 s)
6	bigdatalab-cdh-dn1.uitm.edu.my:44715	Active	0	0.0 B / 384.1 MB	0.0 B	1	0	0	0	0	0 ms (0 ms)

# FROM THEORY TO ACTION





# FROM EXECUTION TO INSIGHT

## What students learned from their actions

❑ **File formats matter** → observed clear performance differences (text vs ORC).

*“I didn’t expect file formats to matter so much, ORC was so much faster than Text.”*

❑ **Design choices impact outcomes** → partitioning & bucketing changed query speed and efficiency.

*“Partitioning and bucketing really changed the way queries performed.”*

❑ **Comparative mindset** → developed ability to analyze spark vs hive beyond syntax.

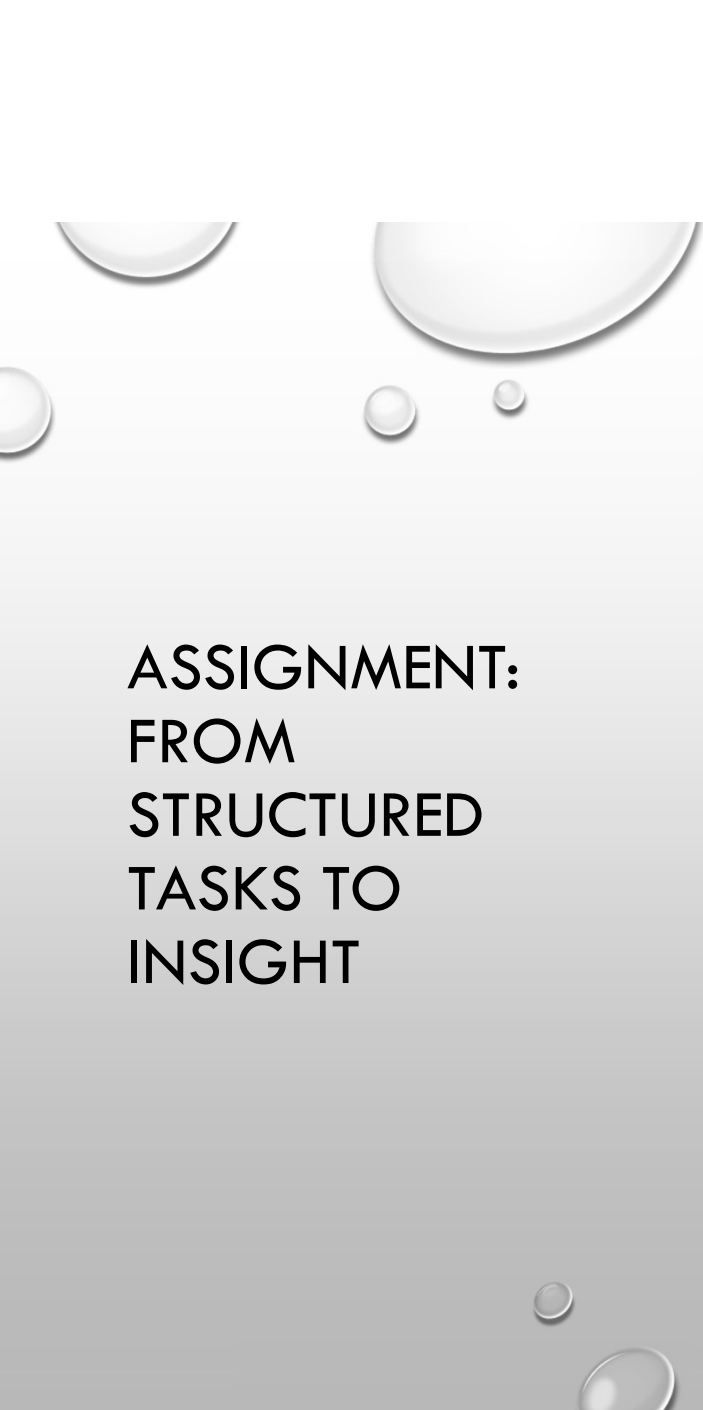
*“I could finally see why Spark behaves differently from Hive, even on the same data.”*

❑ **System awareness** → monitoring with cloudera manager showed how jobs consumed resources.

*“Watching Cloudera Manager showed me how jobs actually use memory and CPU.”*

❑ **Reflective learning** → moved from *running queries* to *explaining why results differ*.

*“It wasn’t just about running queries, I learned to explain why the results differed.”*



## ASSIGNMENT: FROM STRUCTURED TASKS TO INSIGHT

### ASSIGNMENT: SPARK VS HIVE – COMPARATIVE BENCHMARKING

❑ **Objective:** evaluate Sparksql vs Hive across the same table setups to understand architectural and performance differences.

❑ **Student tasks:**

- Get a dataset and upload it to HDFS.
- Formulate the comparative benchmarking framework.
- Create hive tables with different configurations (such as Internal/external, Textfile/parquet, Partitioned, bucketed).
- Run complex queries (with join & group-by) on each setup.
- Collect measurements: such as execution time, Storage/file size.
- Analyze trade-offs: query time vs file size vs table configuration.
- Interpret findings using charts + discussion.

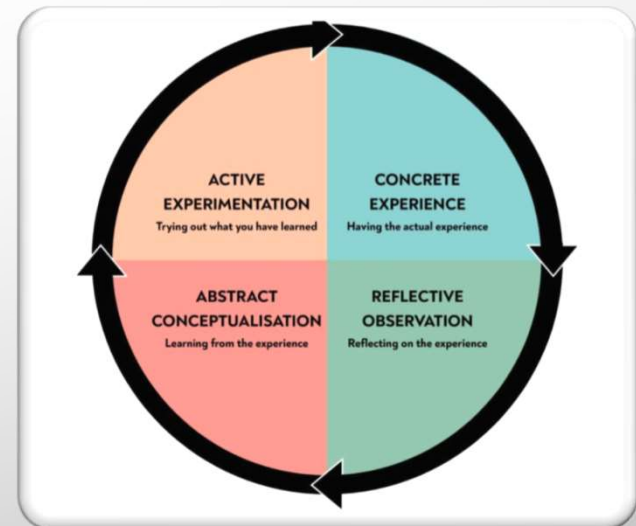
❑ **Learning focus:**

- Storage design, hands-on benchmarking, memory-aware computing, interpreting system performance and understanding distributed processing behavior.

# WHY KOLB'S EXPERIENTIAL CYCLE ANCHORED OUR COURSE DESIGN

## Why kolb?

- ❑ Kolb's experiential learning theory supports **learning by doing**, followed by **reflection and experimentation** which is ideal for hands-on, technical courses.
- ❑ It provides a complete cycle.
- ❑ It is widely applied in computing education.
- ❑ "Many of our students didn't just run code, they reflected on why certain configurations worked better, formed new hypotheses, and tested again. That's kolb in action."



# MAPPING KOLB'S CYCLE TO OUR ASSIGNMENTS

Kolb Stage	Student Activity
Concrete Experience	Students ran queries, loaded data, benchmarked Hive vs Spark
Reflective Observation	They reflected on why Parquet / ORC was faster, why bucketing / partitioning mattered
Abstract Conceptualization	They linked performance outcomes to data structure configurations / system design principles
Active Experimentation	They changed configs, reran jobs, tested new hypotheses



# WHY THIS APPROACH IS SCALABLE AND REUSABLE

## ❑ Modest infrastructure, real impact

- Deploying cluster environment: 3 servers, 5 VMs
- Stable enough for multi-user access big data operations

## ❑ Designed for teaching, not production

- No need for complex cloud services or autoscaling
- Students accessed real systems, not emulators or simplified GUIs

## ❑ Reusable across courses and levels

- Can be used for postgraduate and undergraduate related courses,
- Python analytics and benchmarking projects, and FYPs

## ❑ Remote and in-class friendly

- Easy access for students regardless of learning mode (Intranet, Internet)

# WHY NOT USE AWS OR GOOGLE CLOUD?



## **Cloud platforms are powerful, but not always practical for teaching**

Free tiers are limited or time-limited

Budget and billing control are difficult for classroom settings



## **Limited control and visibility**

Hard to expose students to system-level details (YARN memory, spark logs, file system behavior)



## **Repeatability and reset issues**

Labs are harder to reset or reuse consistently across semesters



## **Our setup prioritizes learning over cloud elasticity**

It's scalable for teaching, not for enterprise compute loads

Students interact directly with tools *and* the environment

The background of the slide features a light gray gradient with several realistic water droplets of varying sizes. Some droplets are at the top, some are in the middle, and a few are at the bottom right, creating a clean, modern aesthetic.

# FUTURE DIRECTIONS

## **Scaling up the lab**

- ☐ Expanding the current cluster to support larger datasets and multi-course usage.

## **Cloud integration**

- ☐ Introducing AWS academy / AWS educate modules to complement hadoop fundamentals.

## **Student pathways**

- ☐ Linking course outcomes to certifications (e.G. AWS data engineer certification).

## **Collaboration**

- ☐ Inviting colleagues to utilize the scaled-down lab for teaching and research testbed (e.G. For pilot experiments, applied domain research).

# CONCLUSION

1

**You don't need a cloud budget or enterprise cluster** to teach big data meaningfully.

2

A well-designed, scaled-down lab can expose students to:

- Distributed storage
- Query optimization
- Performance tuning
- System-level thinking

3

The learning was **not just technical**, it was **reflective, iterative, and empowering**





Thank  
you



[Home](#) > [Education and Information Technologies](#) > [Article](#)

## Data science technology course: The design, assessment and computing environment perspectives

Published: 24 January 2023

Volume 28, pages 10209–10234, (2023) [Cite this article](#)

[Download PDF](#) 

 Access provided by Universiti Teknologi MARA

[Azlan Ismail](#) , [Sofianita Mutalib](#) & [Haryani Haron](#)

Published article: <https://link.springer.com/article/10.1007/s10639-022-11558-8>

