



2110625 - Data Science Architecture

## **Big Data File Systems**

**Asst.Prof. Natawut Nupairoj, Ph.D.**

Department of Computer Engineering  
Chulalongkorn University  
[natawut.n@chula.ac.th](mailto:natawut.n@chula.ac.th)

# Typical Big Data File Systems

---

- Based on scale-out architecture *បើស្ថាប្រាក់ scale out*
- Must support structured and unstructured data  
*(អង់, រាង, ស៊ីម)*
- Processing historical data
  - Hardly changed *ស្ម័គ្រប់នៅលើវា នូវការវិនិច្ឆ័យដែលត្រួតពិនិត្យ*
  - Whole-file scanning accessing pattern *ចាប់ពីផ្លូវ ដល់ byte ទុកការឱ្យលើវា*
- Fault tolerance is required *កំណត់ស្រីរការពិនិត្យ នូវវិធានា និង Reprocess*

# HDFS : Hadoop Distributed File System

---

- A distributed file system in Apache Hadoop Core project designed to run on commodity hardware *( សាស្ត្រិយកម្មវិធាន នូវ តាម )*  
*( 3 copy )*
- Highly fault-tolerant and high throughput access to application data  
*bandwidth នៃ CPU, I/O នូវ  
HDD*
- Supports a traditional hierarchical file organization (directories and files) with user quotas and access permissions (does not support hard links or soft links)

*⇒ directory នៃឯណុខ file នៅលើ hard / soft link (shortcut)*

# Data Lake

---

- A system or repository of data stored in its natural/raw format, usually object blobs or files
  - All data is loaded from source systems, no data is turned away
  - Data is stored at the leaf level in an untransformed or nearly untransformed state
  - Data is transformed and schema is applied to fulfill the needs of analysis
- Can be on premise (HDFS is quite popular), or on cloud (Amazon S3, Google Cloudstore, etc.)

# Assumptions and Goals

---

- Hardware failure is the norm, focus on quick failure detection and automatic recovery *masking* តួនាទី វិវាទ node សំណង់រាយ  
*batch*
- Applications on HDFS are assumed to need streaming data access, batch (not interactive), high throughput (not low latency)  
ក្រុមហ៊ុនកំសែងបាន client ស្តី មានការទេរស្ថិត និងវឌ្ឍនភាព throughput ស្តី និង latency ស្តី
- HDFS is **for large data sets** (GBs to TBs) and provides high aggregated bandwidths from hundreds of nodes in a single cluster
- Simple Coherency Model with write-once-read-many access model, a file once created and closed can only append or truncate *single update* (មិន 2 គ្រឹះទៅលើ) *(strong consistency)*
- Moving Computation is Cheaper than Moving Data *HDFS ឈរតារាង move ការគាំទ្រូចចែង*  
*move ធំម្អាត*
- Portability Across Heterogeneous Hardware and Software Platforms  
*HDFS ទទួលឱ្យការងារ*

# Play with HDFS and Hadoop in Containers

---

- Note: this still does not work on Apple Silicon
- Use [big-data-europe/docker-hive](https://github.com/big-data-europe/docker-hive) on GitHub

```
git clone https://github.com/big-data-europe/docker-hive.git
```

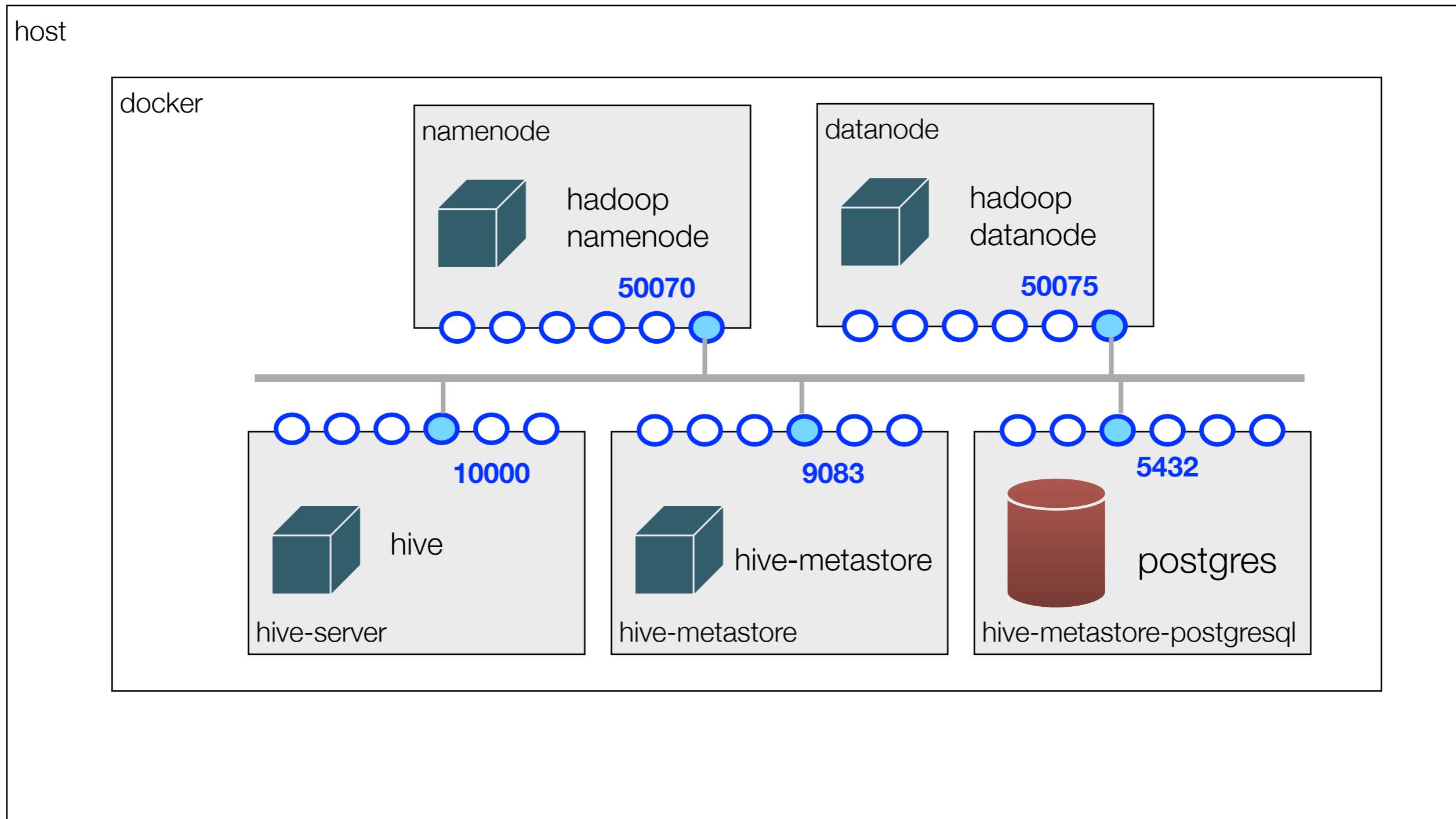
- Use docker-compose to create all containers

```
sudo docker-compose up
```

- Use docker-compose to create all containers

- Namenode: http://localhost:50070
- Datanode: http://localhost:50075
- Hive-server: http://localhost:10000
- Hive-metastore: http://localhost:9083

# Sample Hadoop Cluster Architecture



# Hadoop FS Commands behave similar to Unix Commands

---

```
sudo docker-compose exec namenode bash
```

```
hadoop version
```

```
hdfs dfs -mkdir /path/directory_name
```

```
hdfs dis -ls /path
```

```
hdfs dfs -copyFromLocal <localsrc> <hdfs destination>
```

```
hdfs dfs -copyToLocal <hdfs source> <localdst>
```

```
hdfs dfs -cat /path_to_file_in_hdfs
```

```
hdfs dfs -mv <src> <dest>
```

```
hdfs dfs -cp <src> <dest>
```

# HDFS Sample Commands

---

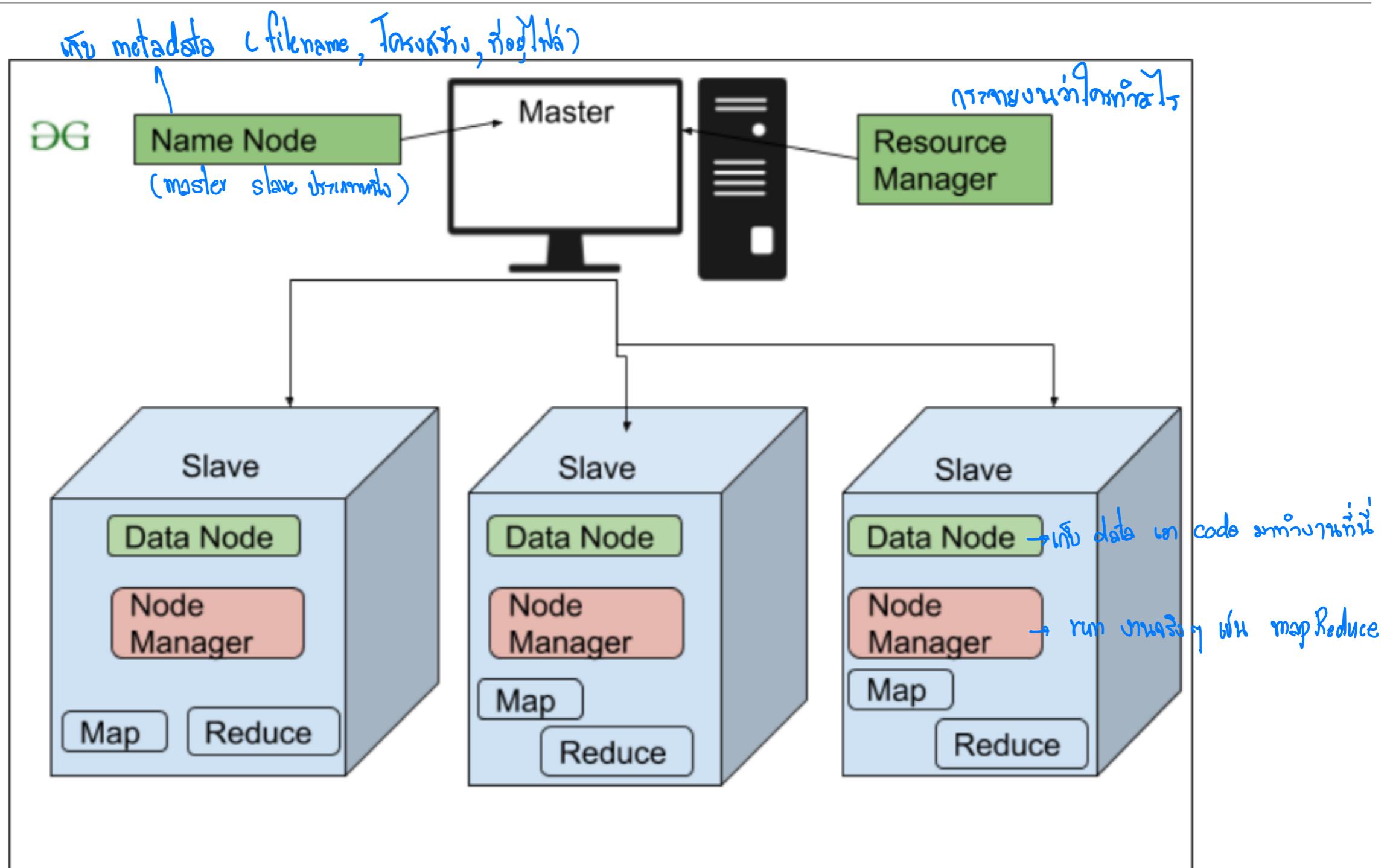
```
root@2602201f31b9:~# hdfs dfs -ls /
Found 3 items
drwxr-xr-x  - root supergroup          0 2021-08-17 11:45 /bigdata
drwxrwxr-x  - root supergroup          0 2021-03-20 06:31 /tmp
drwxr-xr-x  - root supergroup          0 2021-03-16 04:05 /user
root@2602201f31b9:~# hdfs dfs -mkdir /myhome
root@2602201f31b9:~# hdfs dfs -copyFromLocal NetflixOriginals.csv /myhome
root@2602201f31b9:~# hdfs dfs -ls /myhome
Found 1 items
-rw-r--r--  3 root supergroup      38678 2021-08-18 02:37 /myhome/NetflixOrigin
als.csv
root@2602201f31b9:~#
```

# All Hadoop FS Commands

---

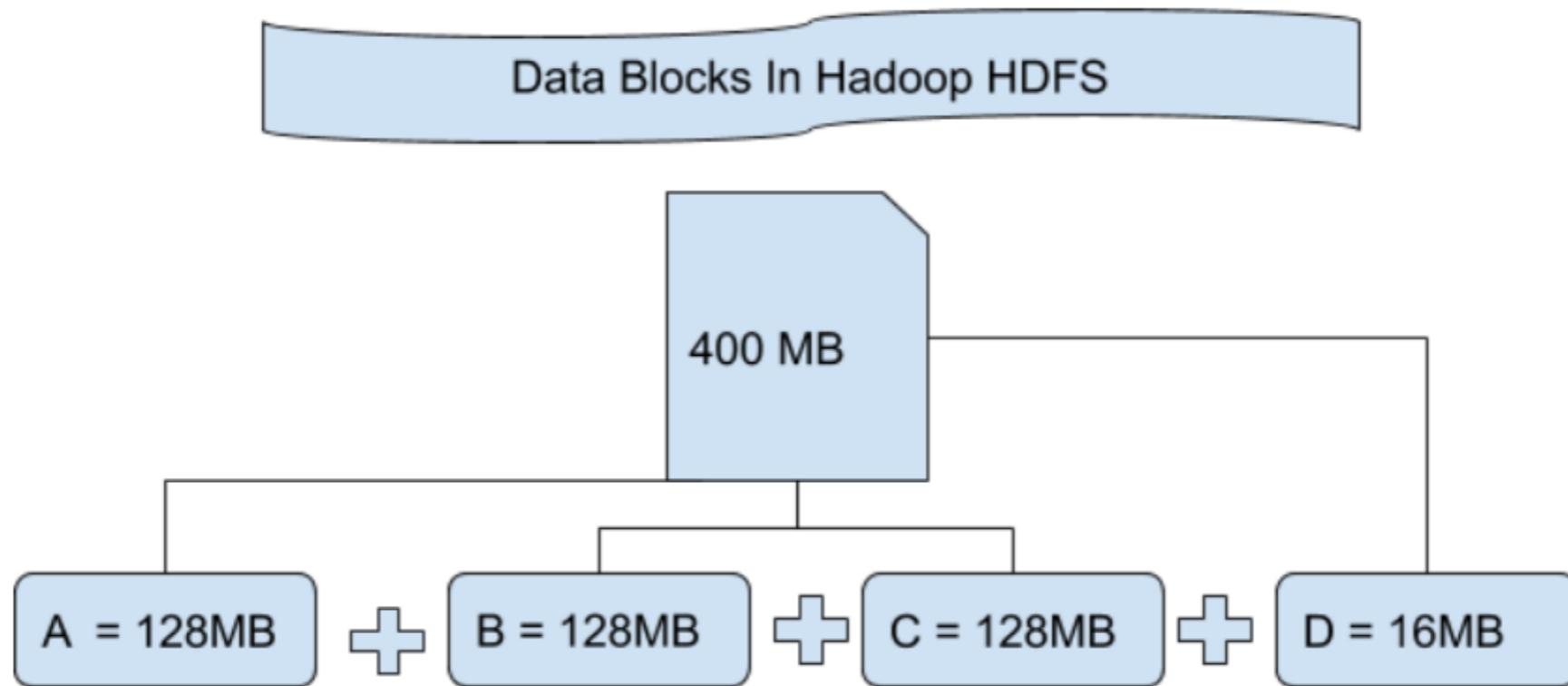
appendToFile	df	lsr	setrep
cat	du	mkdir	stat
checksum	dus	moveFromLocal	tail
chgrp	expunge	moveToLocal	test
chmod	find	mv	text
chown	get	put	touch
copyFromLocal	getfacl	renameSnapshot	touchz
copyToLocal	getfattr	rm	truncate
count	getmerge	rmdir	usage
cp	head	rmr	
createSnapshot	help	setfacl	
deleteSnapshot	ls	setfattr	

# High Level Architecture Of Hadoop



# File Block in HDFS

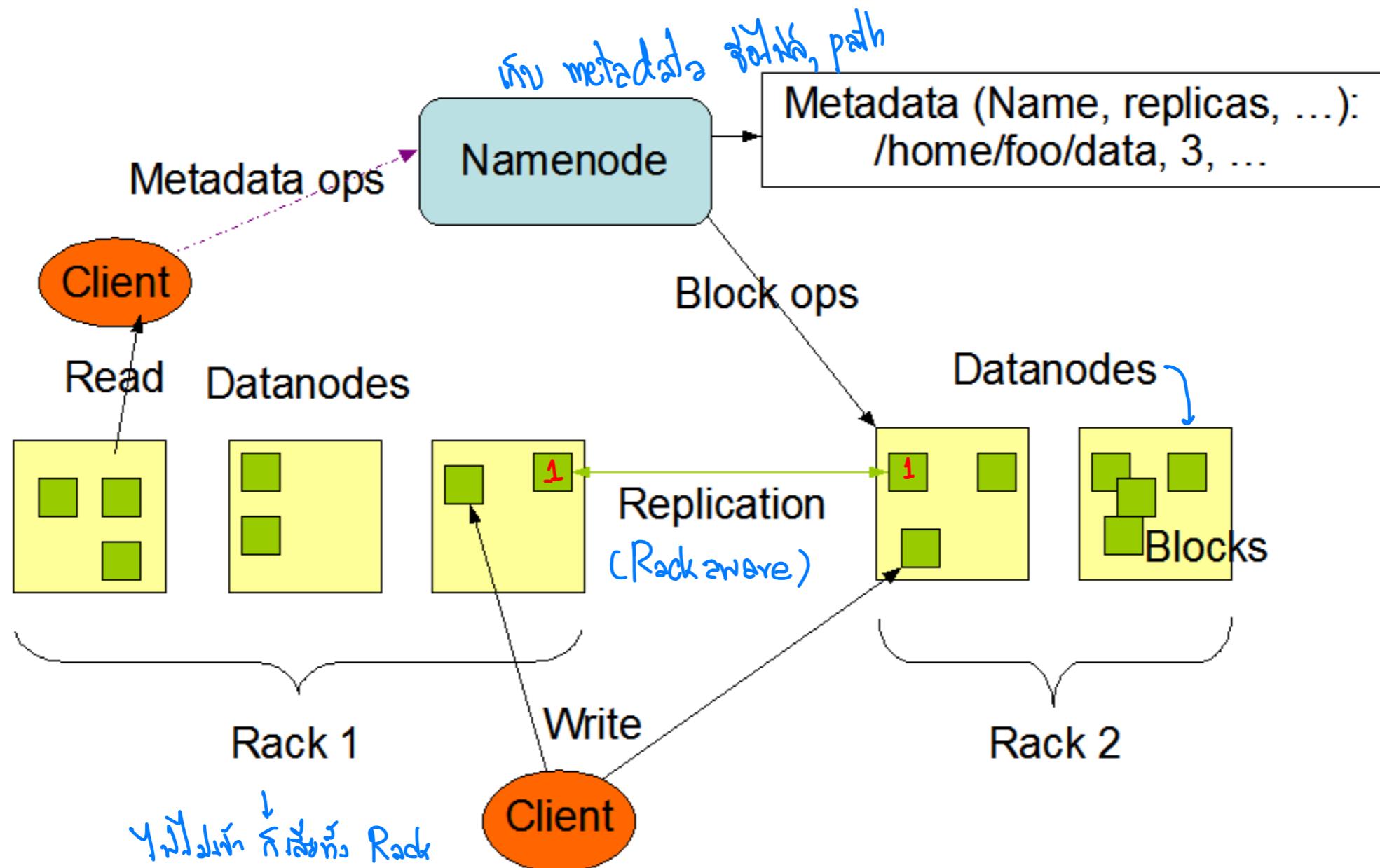
---



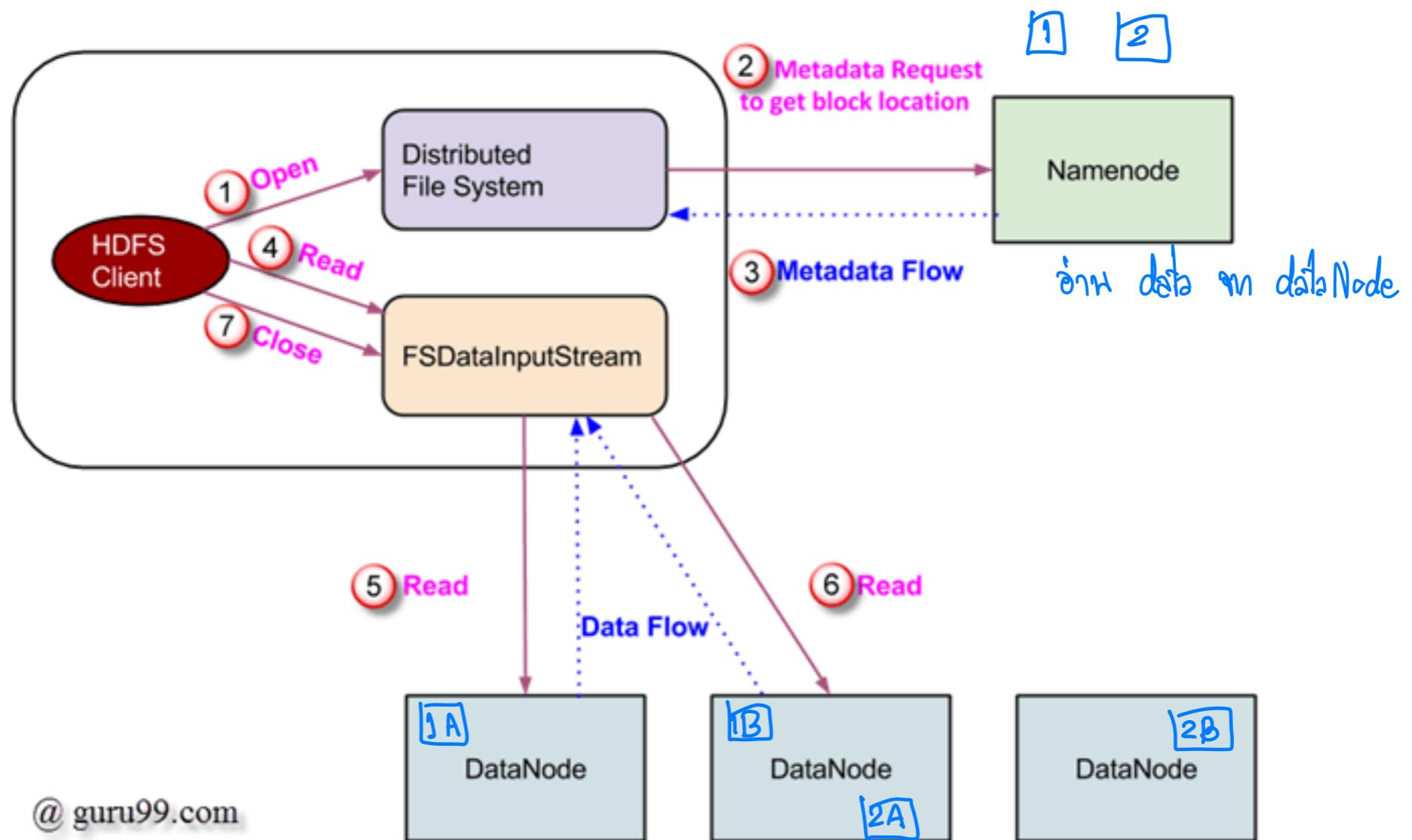
Source: <https://www.geeksforgeeks.org/hadoop-architecture/>

- A data file is splitter into blocks
- Each block is 128 MBs (can be configured), except the last block
- This is much larger comparing to 4 KBs of normal Linux file system

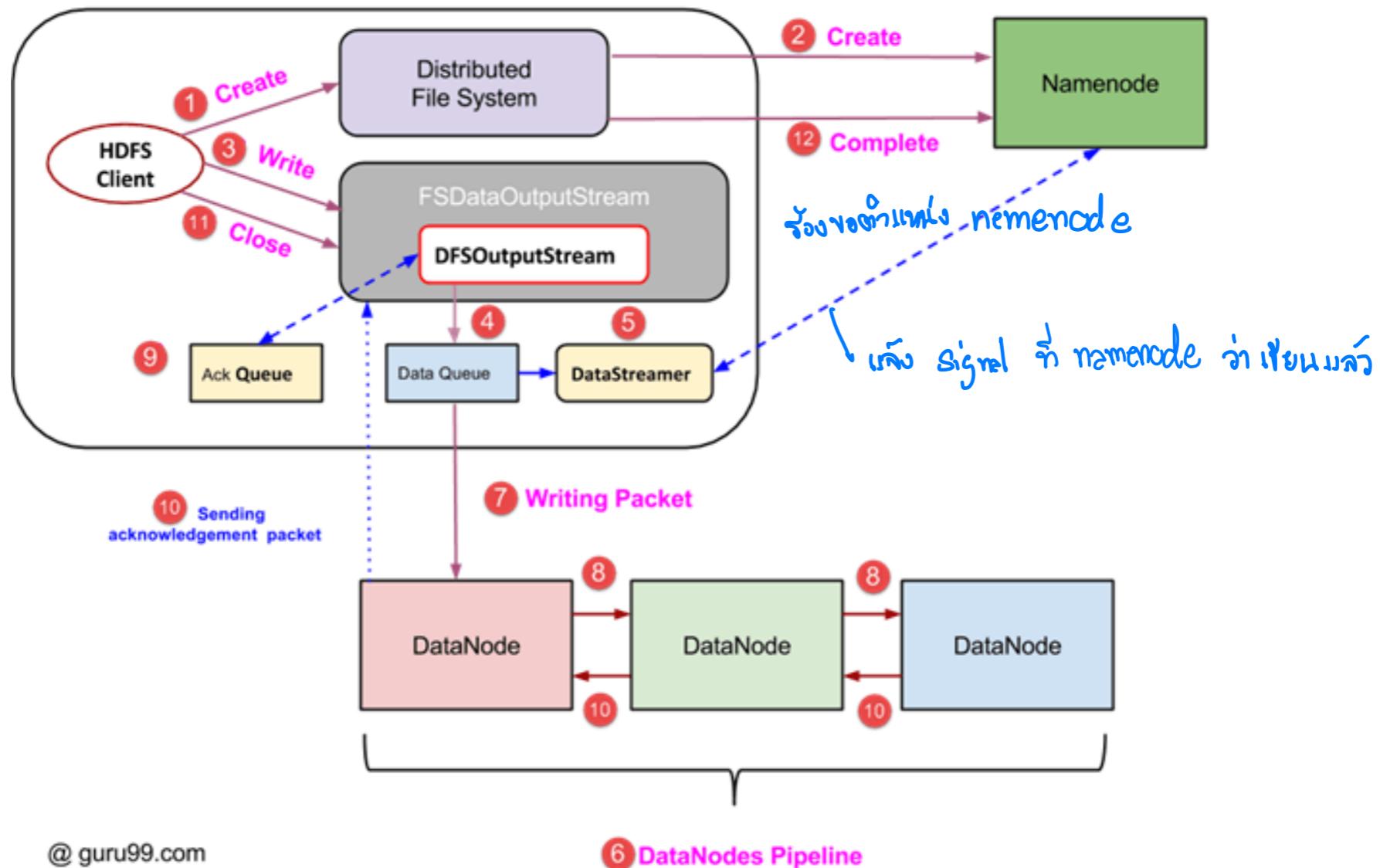
# HDFS Architecture



# HDFS File Read Operation



# HDFS File Write Operation



@ guru99.com

6 DataNodes Pipeline

file จะถูกจัดเก็บในรูปแบบพาร์ทิชัน 3 สำหรับ  
(copy คราว 3 datanode)

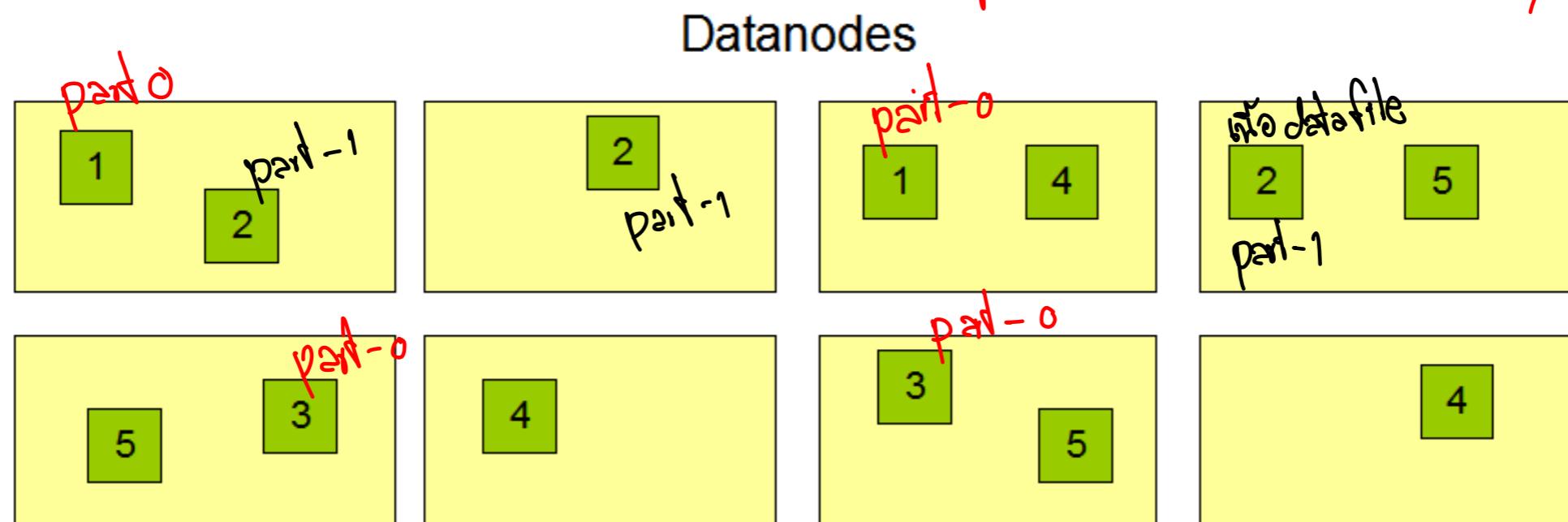
Source: <https://www.guru99.com/learn-hdfs-a-beginners-guide.html>

# Data Replication

## Block Replication

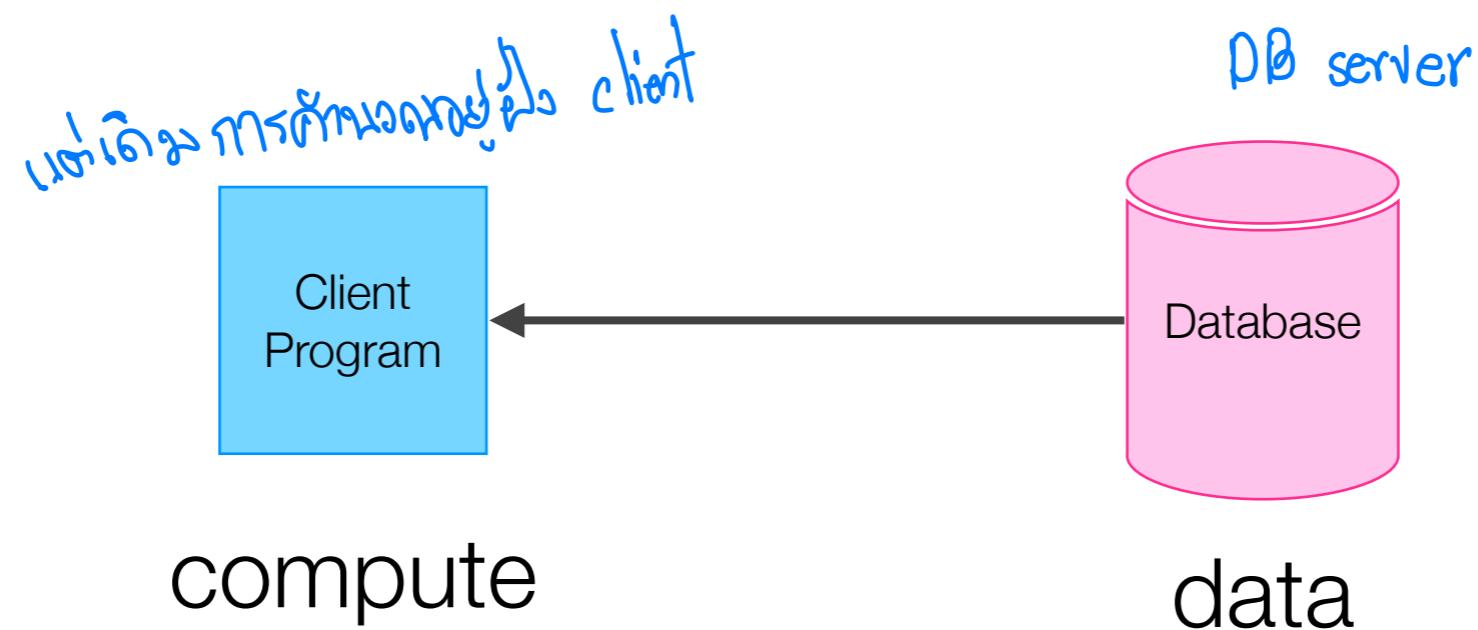
Namenode (Filename, numReplicas, block-ids, ...)  
/users/sameerp/data/part-0, r:2, {1,3}, ... } file 2 files  
/users/sameerp/data/part-1, r:3, {2,4,5}, ... }

replication factor = n.u, copy



# HDFS and Data Locality in Hadoop

- Data locality in Hadoop is the process of moving computation close to the actual data on the node
  - Faster execution and higher throughput
- Typical transaction processing system move data to the computation e.g. database and client program



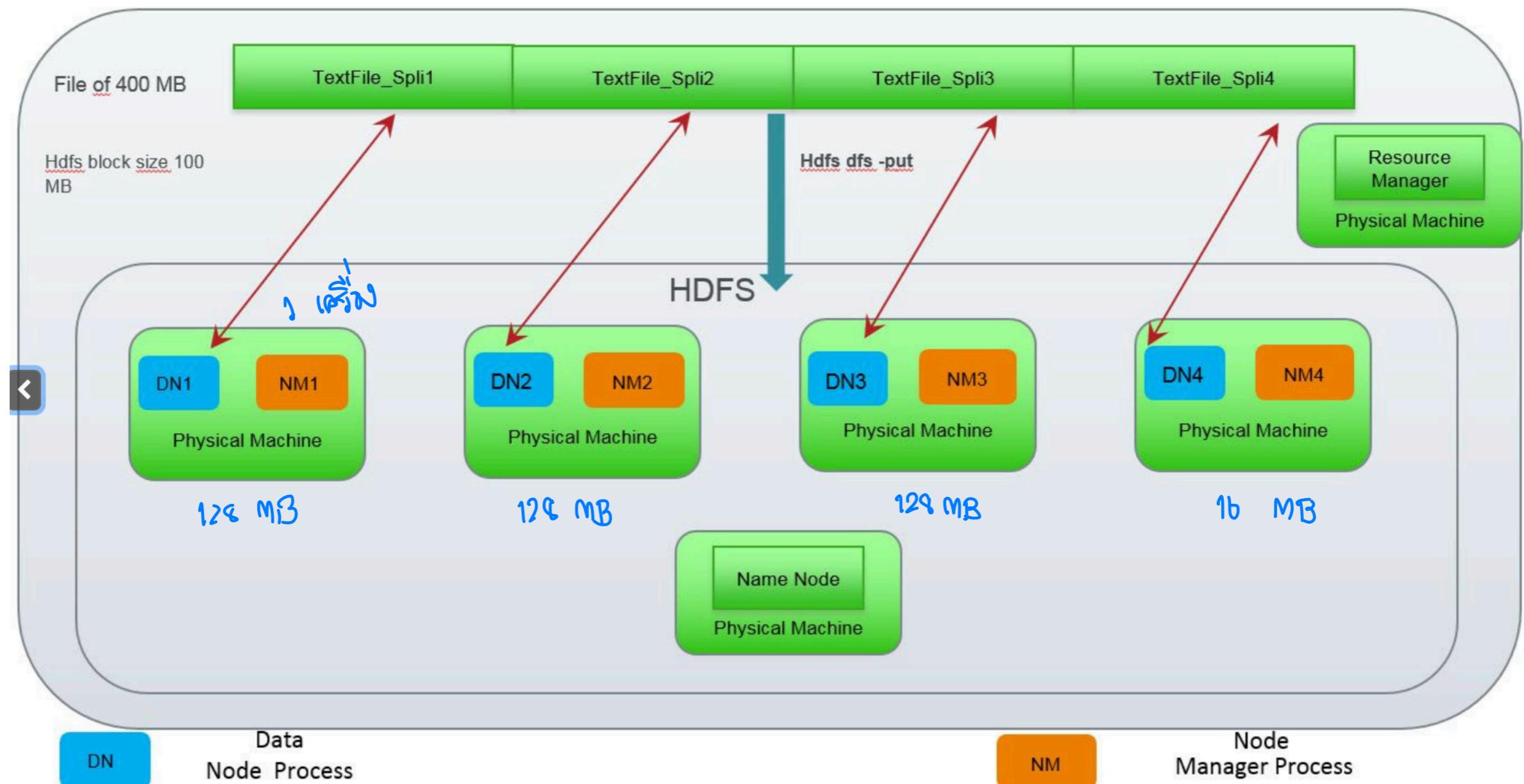
# HDFS and Data Locality in Hadoop

---

- In Hadoop, data is much bigger than code – moving code is much more efficient
  - Data is stored in blocks (with multiple copies) across DataNodes in HDFS
  - NameNode moves code of a MapReduce job to DataNodes to process data blocks (called data local data locality – preferred scenario)
  - If data local data locality is not possible, topology (intra-rack and inter-rack) comes into consideration for code execution placements

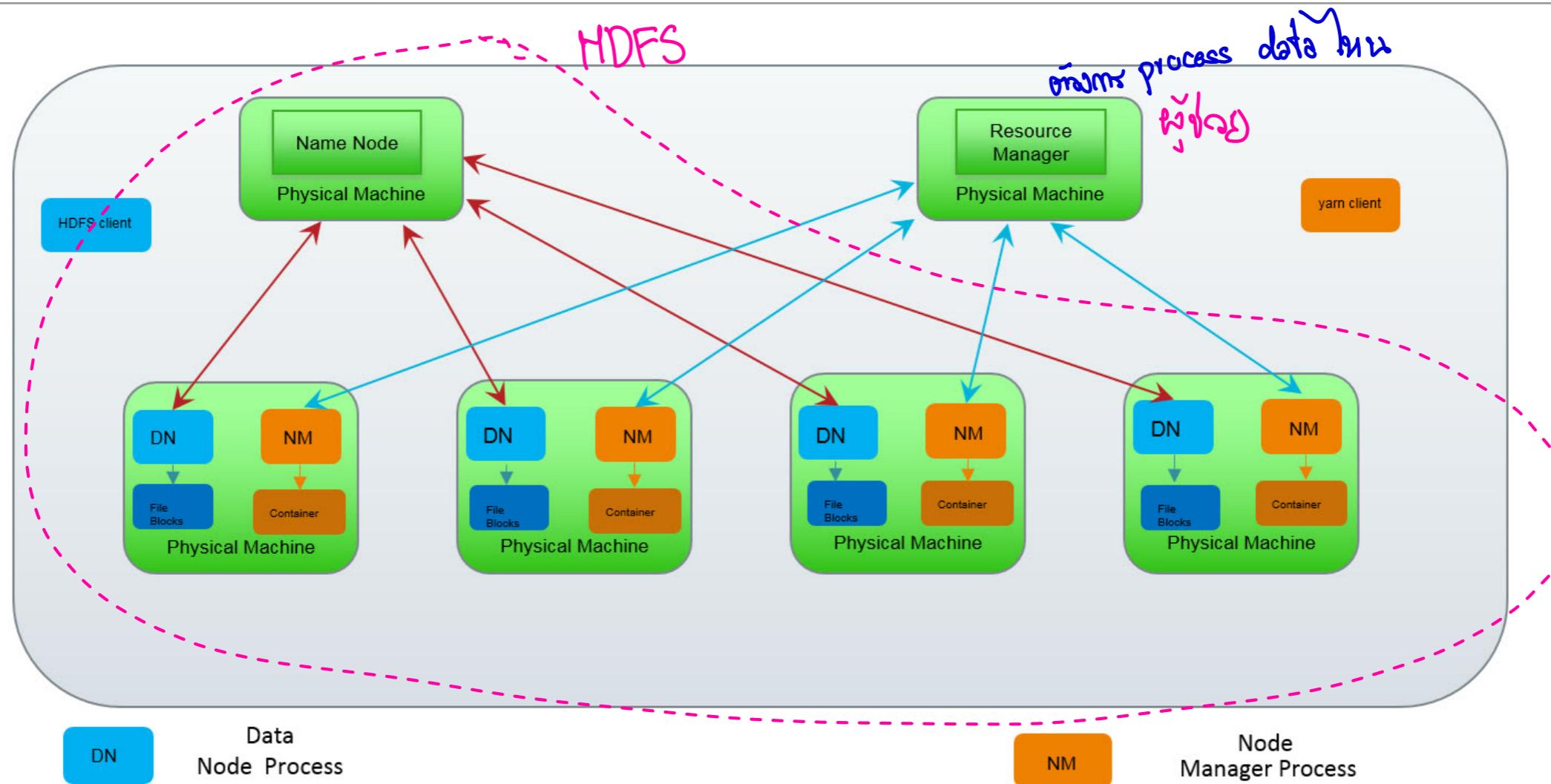
# HDFS and Data Locality in Hadoop

Source: <https://community.cloudera.com/t5/Community-Articles/Understanding-basics-of-HDFS-and-YARN/ta-p/248860>



# HDFS and Data Locality in Hadoop

Source: <https://community.cloudera.com/t5/Community-Articles/Understanding-basics-of-HDFS-and-YARN/ta-p/248860>



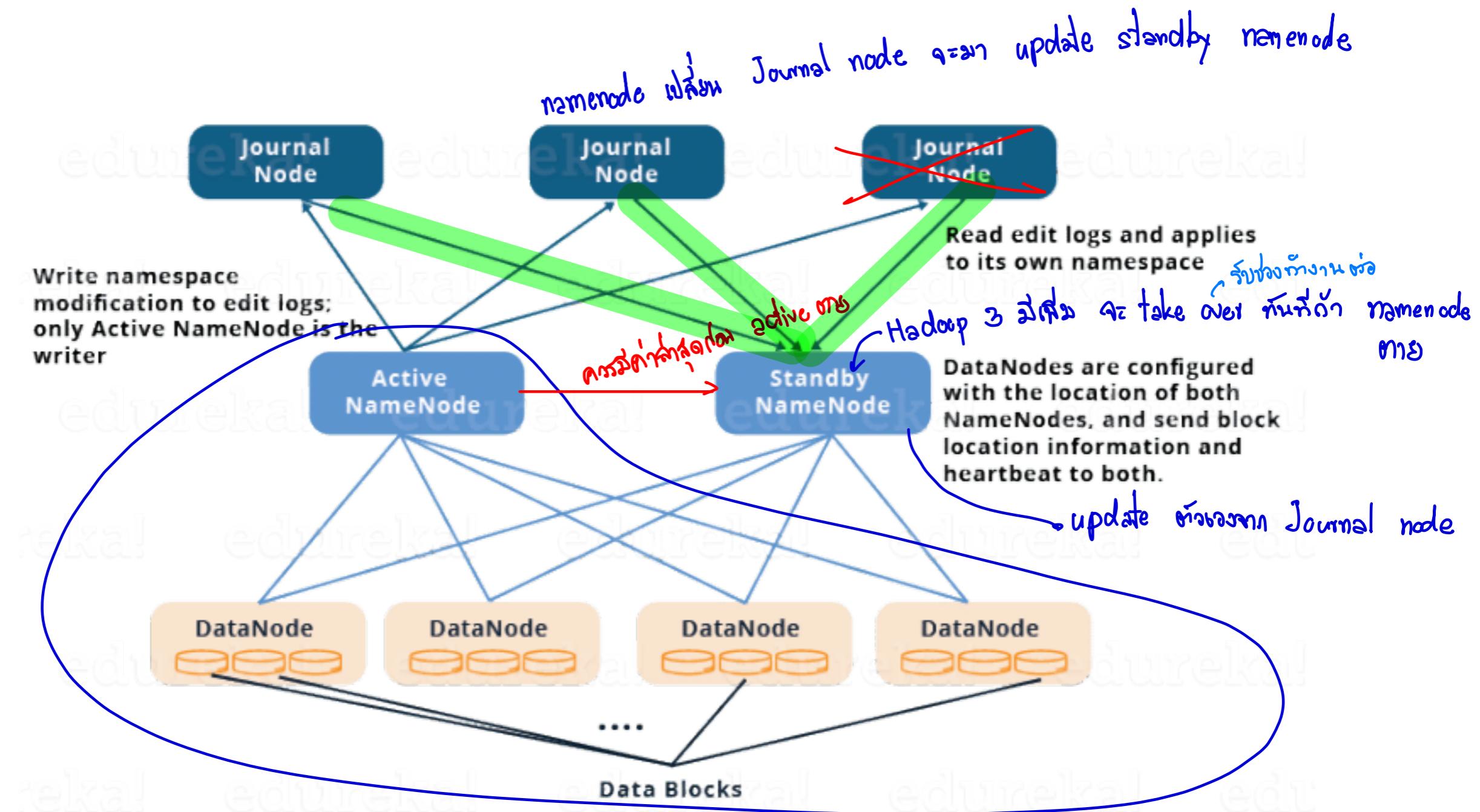
- MapReduce Job Manager is a yarn client that submit codes to resource manager
- Resource manager distributes codes to selected node managers to run in yarn containers
- A Yarn container consists of codes and HDFS client to access data blocks in the same (preferred) or other DataNode

# HDFS Pros and Cons

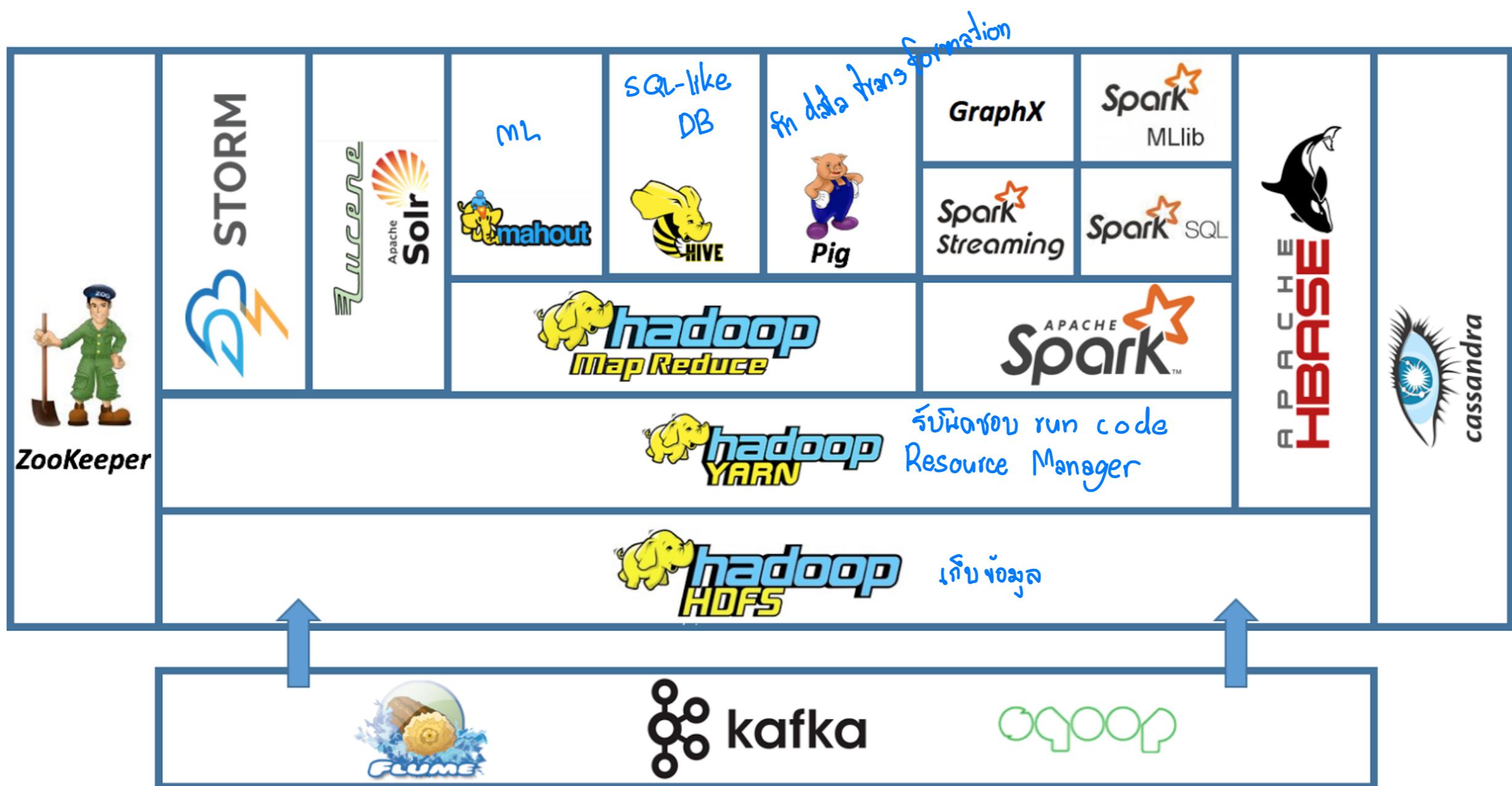
---

- Pros
  - Cost effective កំណត់ស្ថាបន និង បញ្ចូន
  - Highly scalable and high throughput
  - Fault tolerance និង 3 replication (default)
- Cons
  - Problem with small files - NameNode bottleneck, low performance NameNode ជួលចំនួនឯកសារក្នុងមានសម្រាប់
  - NameNode is a single point of failure NameNode មិនមែន backup
  - Dependency on disks - more overhead for read/write
  - Suitable for only batch processing មិនអាចប្រើបាយក្នុងការទាញយក  
interactive

# Hadoop 3 : HDFS HA Architecture



# Apache Hadoop Ecosystem



Hive

ស៊ីន DW ឬន HDFS ផ្ទែកអ៊ី SQL លើវា

---

- An information platform built by a team from Facebook to create a data warehouse framework on top of Hadoop
- A mechanism to impose structure on a variety of data formats stored directly in Apache HDFS or in other data storage systems such as Apache HBase
- Enable analysts with strong SQL skills to run HPL-SQL queries via Apache Tez, Apache Spark, or MapReduce

# Hive Architecture

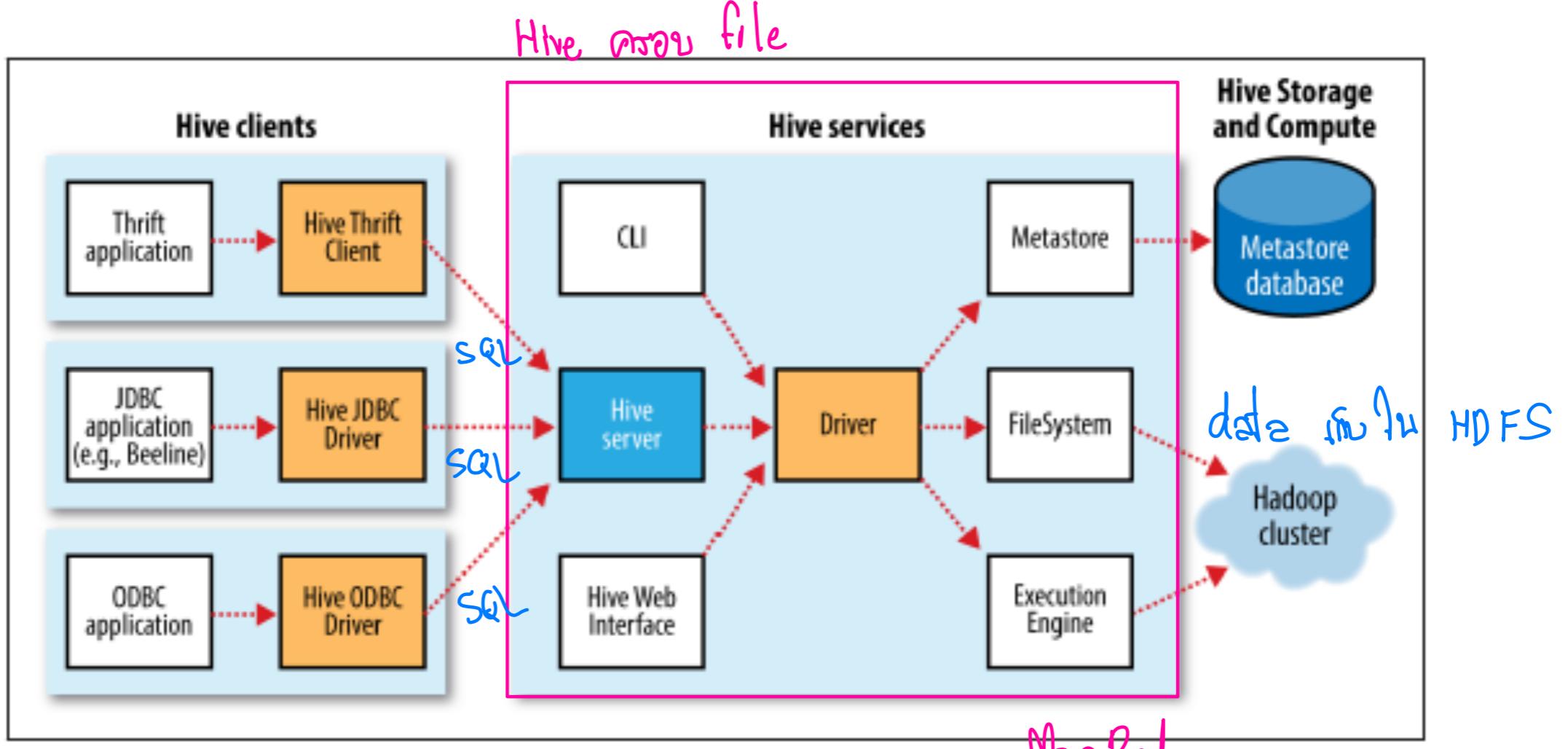


Figure 17-1. Hive architecture

Source: Hadoop: The Definitive Guide, 4th Edition

မျှော်းပြန်သူ mapReduce ဂိုဏ်ပြု  
HDFS ပြာ

# Play with Hive

```
sudo docker-compose exec hive-server bash  
more /opt/hive/examples/files/kv1.txt  
hive
```

```
469val_469  
145val_145  
495val_495  
37val_37  
327val_327  
281val_281  
277val_277  
209val_209  
15val_15  
82val_82  
403val_403  
166val_166  
417val_417  
430val_430  
root@b1c82e30e1c3:/opt#
```

key value

```
SHOW TABLES;
```

```
CREATE TABLE pokes (foo INT, bar STRING);
```

```
DESCRIBE pokes;
```

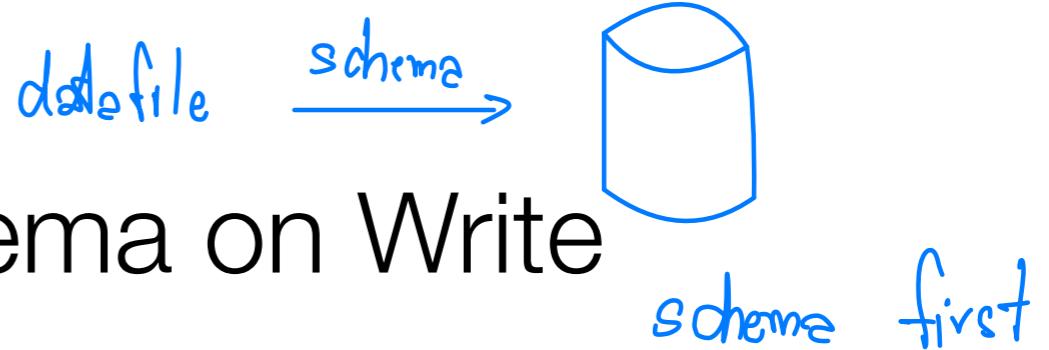
```
hive> desc pokes  
>;  
OK  
key          int  
value        string
```

```
LOAD DATA LOCAL INPATH '/opt/hive/examples/files/kv1.txt'  
OVERWRITE INTO TABLE pokes; in data file នឹងកែចរណ៍
```

```
SELECT * FROM pokes;
```

```
hadoop fs -ls /user/hive/warehouse/pokes
```

```
root@b1c82e30e1c3:/opt# hdfs dfs -ls /user/hive/warehouse/pokes  
Found 1 items  
-rwxrwxr-x  3 root supergroup      5812 2024-04-06 07:04 /user/hive/warehouse/  
pokes/kv1.txt
```



## Traditional Database = Schema on Write

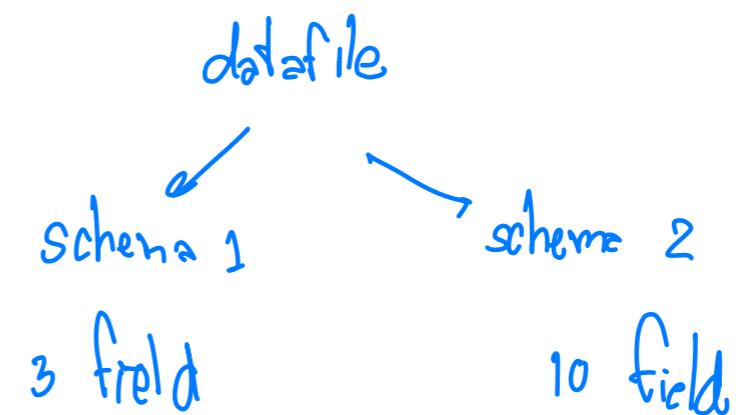
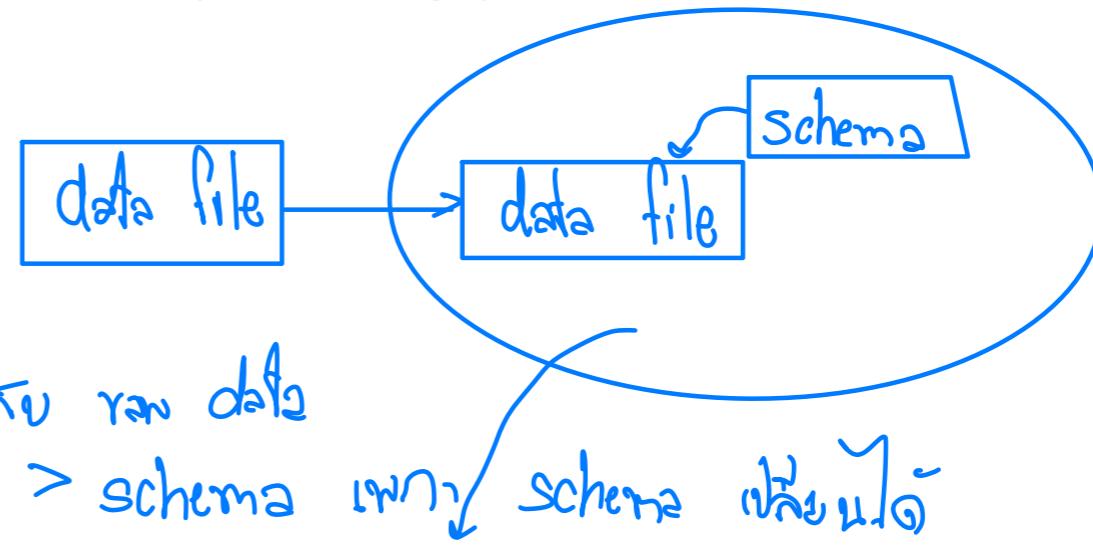
---

- Table's schema is enforced at data load time ពេលដាក់ការងារនៅក្នុងក្រុងការការពារ
- Stop if data being loaded does not conform to the schema data មិនត្រូវបានរាយទៅត្រូវ
- Pros
  - Faster query performance
  - Known level of data quality before executing any queries
- Cons
  - Long time to load ពេលការអាងដែលមិនមែន schema

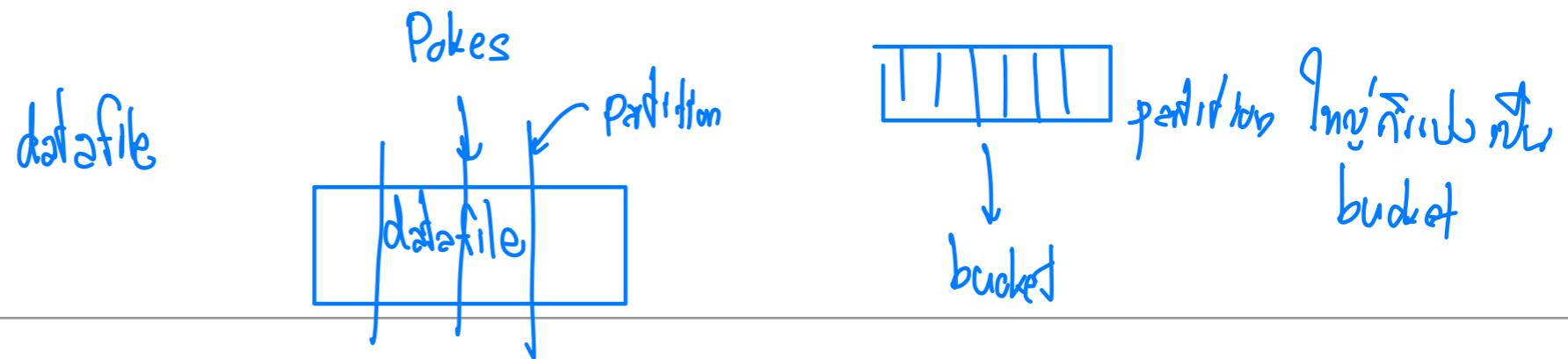
# Hive = Schema on Read

ແຈລງ schema ໂອນນິ້ມມີ ສູງກະເມີນ

- Do not verify the data at load time, just copy the data file into the data store
- Enforce schema when read
- Pros
  - Faster initial load
  - More flexible for having multiple schemas on the same underlying data
  - Good for scenarios where the schema or index is not known at load time
- Cons
  - Slower query performance



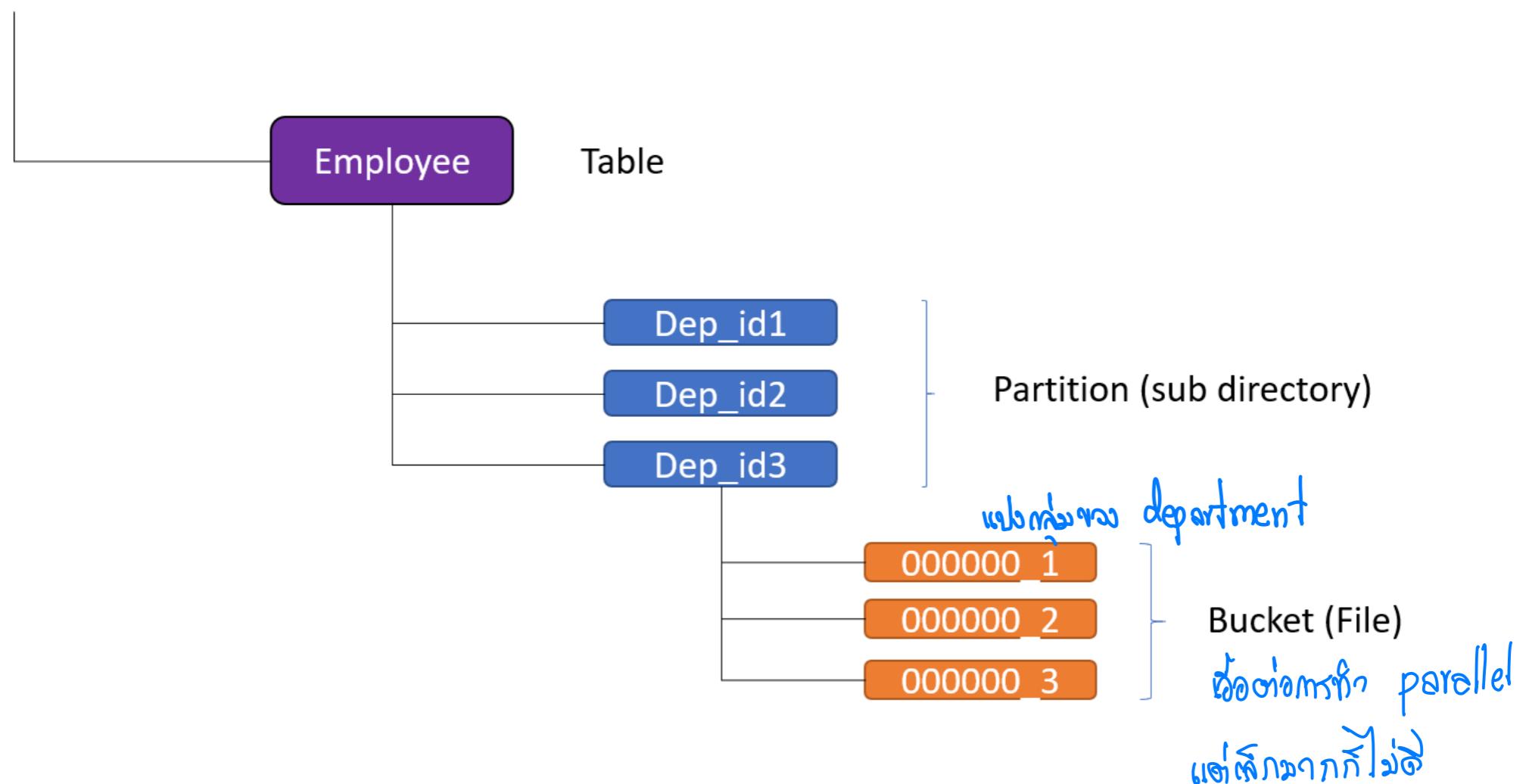
# Data Units



- **Databases**: Namespaces function to avoid naming conflicts for tables, views, partitions, columns, and so on
- **Tables**: Homogeneous units of data which have the same schema
- **Partitions**: Each Table can have one or more partition Keys which determines how the data is stored, query only on the relevant partition of the table, thereby speeding up the analysis significantly
- **Buckets (or Clusters)**: Data in each partition may in turn be divided into Buckets based on the value of a hash function of some column of the Table
- Note that it is not necessary for tables to be partitioned or bucketed

# Data Units and HDFS Storages

ເທົ່ານັ້ນ DB ອຸ່ນ



# Hive - Primitive Types

ມີເອົ້າໃຈ DB ກໍາລຖງ

- Integers (TINYINT, SMALLINT, INT, BIGINT)
- Boolean type (BOOLEAN—TRUE/FALSE)
- Floating point numbers (FLOAT, DOUBLE)
- Fixed point numbers (DECIMAL)
- String types (STRING, VARCHAR, CHAR)
- Date and time types (TIMESTAMP, TIMESTAMP WITH LOCAL TIME ZONE, DATE)
- Binary types (BINARY)

# Hive - Complex Types

---

- **Structs**: the elements within the type can be accessed using the DOT (.) notation
- **Maps** (key-value tuples): The elements are accessed using ['element name'] notation
- **Arrays** (indexable lists): The elements in the array can be accessed using the [n] notation where n is an index (zero-based) into the array

# Built In Operators and Functions

---

- Relational Operators    *where, JOIN*
- Arithmetic Operators
- Logical Operators
- Operators on Complex Types
- Built In Basic Functions
- Built In Aggregate Functions
- WHERE, SELECT, JOIN, SORT BY, GROUP BY
- User-Defined Functions

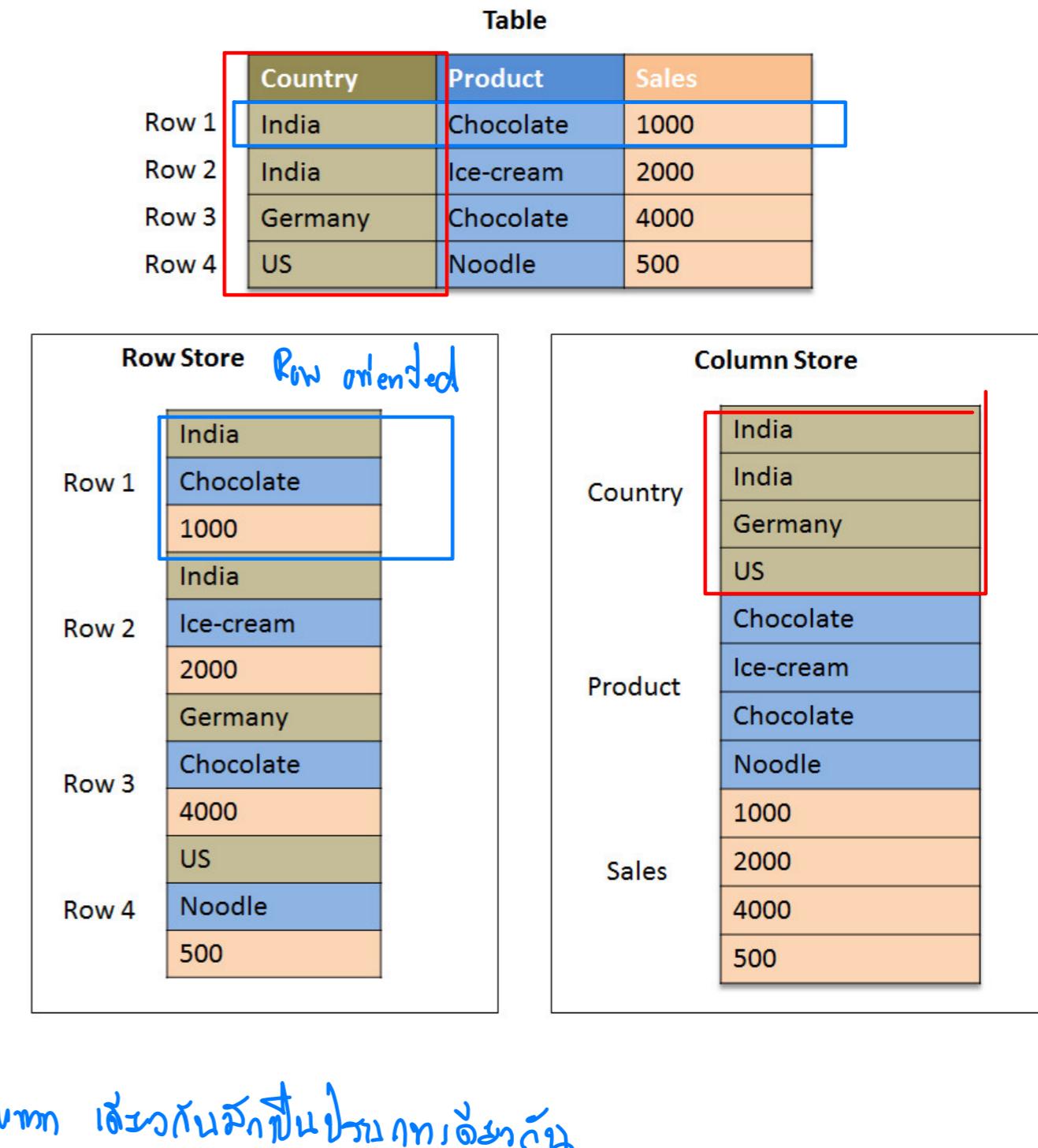
# Big Data File Formats

---

- HDFS treats files as objects, can be any file formats
- Popular file formats
  - Text-based: txt, csv, json, XML
  - Binary-based: avro, parquet, orc *ที่ random access ได้ด้วย binary base*
- Benefits of binary-based file formats
  - Random-access, support parallel processing
  - On-the-fly compression, smaller size but not at the cost of reading speed  
*ดู data รูปแบบ compress แล้ว unzip สะดวกทางด้าน ไม่ต้อง unzip กันแล้ว*
  - Embedded schema
  - Support table-oriented data store

# Table-Oriented Data Store

- Typical data analytics are done on “table” of data
- For big data, a table may not be fit in a single file
- Multiple files can be processed in parallel effectively
- Row-store file format is more efficient for add/update/delete operations
- Column-store file format is more efficient for data analytic operations



# AVRO

---

- Row-based schema-oriented data serialization framework
  - Use JSON to define schema កំណត់ schema ទៅ JSON នៅ
  - Serialize data in a compact binary format
  - Similar to Java Object Serialization Pack row ឲ្យជា data ការណើន
  - Support both primitive and complex data types (allows hierarchical data structure)
    - integer
    - list, hierarchy
- Primary usage in Big Data
  - Apache Hadoop for persistent data (input/output/checkpoint)
    - A wire format for communication (e.g. Kafka, Spark Streaming)
      - network
      - socket
  - Have libraries for many programming languages

Create

Home

Competitions

Datasets

Code

Discussions

Courses

More

Your Work

## RECENTLY VIEWED

Squid Game Netflix Tw...

Netflix Movie Rating D...

Netflix Movies and TV ...

HPA 2020 16-Bit Traini...

Netflix Original Films &amp;...

Dataset

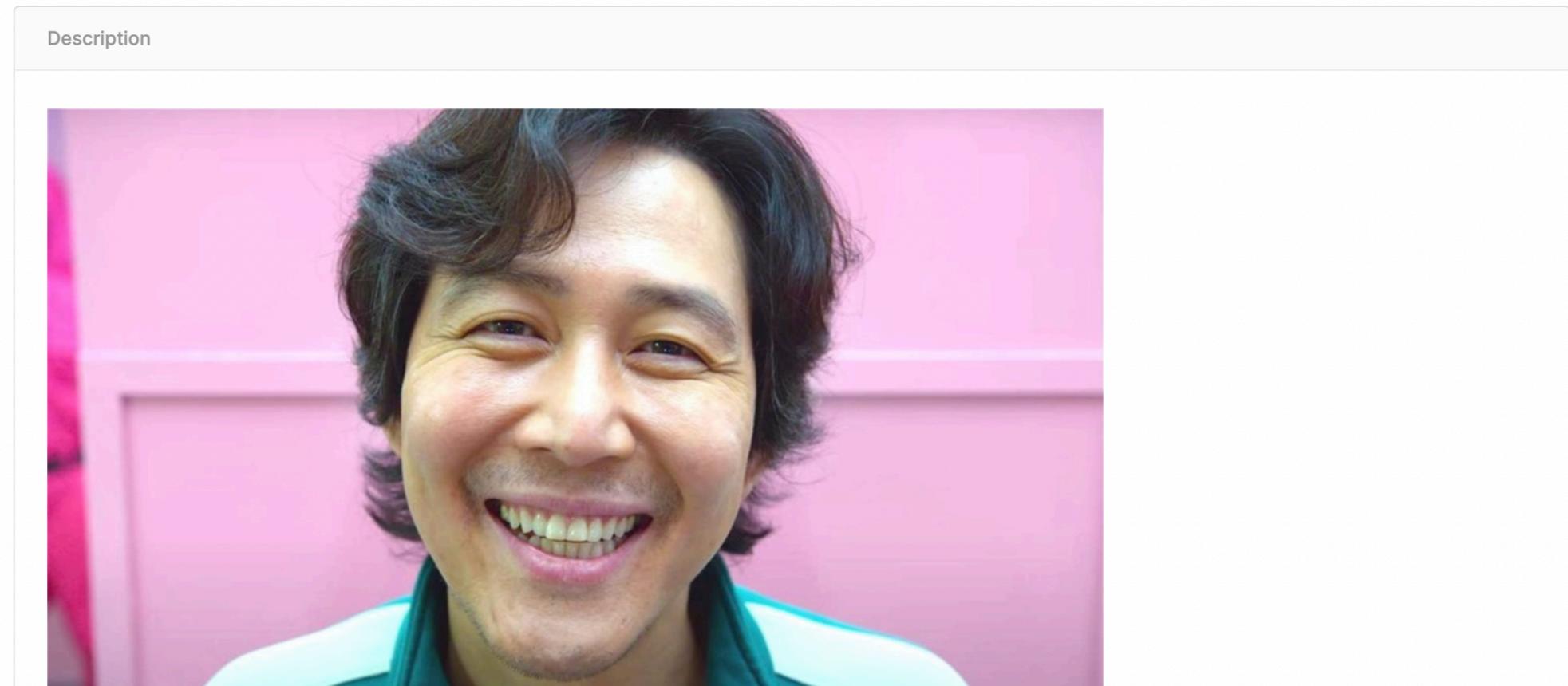
## Squid Game Netflix Twitter Data

This data set contains twitter dump for the hashtag #squidgame.

Deep Contractor • updated 18 days ago (Version 11)

Data    Code (2)    Discussion    Activity    Metadata    Download (26 MB)    New Notebook    :

Usability 10.0    License CC0: Public Domain    Tags arts and entertainment, online communities, text data, beginner, exploratory data analysis



- The dataset contains the recent tweets about the record-breaking Netflix show "Squid Game"
- The data is collected using tweepy Python package to access Twitter API.

View Active Events

**Data Explorer**

201 MB

/ tweets.v0.csv (201 MB)

1 / 57

Detail Compact Column

សារិក (motto)

Tweet

A user_name	A user_locati...	A user_descr...	D user_creat...	# user_follo...	# user_friends	# user_favou...	✓ user_verified	D date	A text
the _ündér-ratèd niggáh		@ManUtd die hard❤️❤️👉👉 YOLO J'ai besoin de quelqu'un qui peut m'aimer au pire😊 Non, je ne...	2019-09-06 19:24:57+00:00	581	1035	8922	False	2021-10-06 12:05:38+00:00	When life hits and the same time poverty strikes you Gong Yoo : Lets play a game #SquidGame #Netfli...
Best uncle on planet earth			2013-05-08 19:35:26+00:00	741	730	8432	False	2021-10-06 12:05:22+00:00	That marble episode of #SquidGame ruined me. 😭😭
marcie		animal crossing. chicken nuggets. baby yoda. smol animals. tv shows. 🏳 pronouns: any	2009-02-21 10:31:30+00:00	562	1197	62732	False	2021-10-06 12:05:22+00:00	#Squidgame time
YoMo.Mdp	Any pronouns	Where the heck is the karma I'm going on my school grave brb #Technosupport	2021-02-14 13:21:22+00:00	3	277	1341	False	2021-10-06 12:05:04+00:00	//Blood on 1st slide I'm joining the squidgame thing, I'm already dead by sugar honeycomb ofc #Squi...
Laura Reactions	France	I talk and I make reactions videos about shows I love #theexpans #peakyblinders #thelastkingdom #la...	2018-12-19 20:38:28+00:00	330	152	2278	False	2021-10-06 12:05:00+00:00	The two first games, players were killed by the mask guys ; the bloody night and the third game, the...
Peyman KAI	United Kingdom	Official @KardiaChain SKAI Ambassador Marketing Advisor @kephigallery Graphics and Film Artist https...	2018-01-27 12:07:31+00:00	546	318	6265	False	2021-10-06 12:04:54+00:00	\$THG Going to explode to 4B Marketcap very soon. The world first MOBA This game is on another level!...
Aeriaaaa♡		Fujoshi 🐶/ Thai BL- obsessed/Always distracted by poetry 🌸/ CAP 🐴	2021-06-01 14:08:10+00:00	14	110	518	False	2021-10-06 12:04:45+00:00	@B_hundred_Hyun pls use that gun on me. 😞 #BAEKHYUN #EXO #weareoneEXO #SquidGame https://t.co/ksk...

# AVRO Schema File

JSON schema

complex types (record, enum, array, map, union, and fixed)

```
1 {  
2     "namespace": "example.avro.chula",  
3     "type": "record",  
4     "name": "squid_tweets",  
5     "fields": [  
6         {"name": "user_name", "type": "string"},  
7         {"name": "user_location", "type": ["null", "string"]},  
8         {"name": "user_description", "type": "string"},  
9         {"name": "user_created", "type": {"type": "long", "logicalType": "local-timestamp-millis"}},  
10        {"name": "user_followers", "type": "int"},  
11        {"name": "user_friends", "type": "int"},  
12        {"name": "user_favourites", "type": "int"},  
13        {"name": "user_verified", "type": "boolean"},  
14        {"name": "date", "type": {"type": "long", "logicalType": "local-timestamp-millis"}},  
15        {"name": "text", "type": "string"},  
16        {"name": "source", "type": "string"},  
17        {"name": "is_retweet", "type": "boolean"}  
18    ]  
19}
```

optional

primitive types (null, boolean, int, long, float, double, bytes, and string)

# AVRO Example in Python

```
In [1]: import pandas as pd  
import fastavro as favro  
import json, os
```

```
In [2]: schema = json.load(open('tweets_v8.avsc')) 加载 AVRO schema JSON file  
parsed_schema = favro.parse_schema(schema)
```

```
In [3]: parse_dates = ['user_created', 'date']  
df = pd.read_csv('tweets_v8.csv', parse_dates=parse_dates, keep_default_na=False)  
df.shape
```

```
Out[3]: (80019, 12)
```

```
In [4]: records = df.to_dict('records')
```

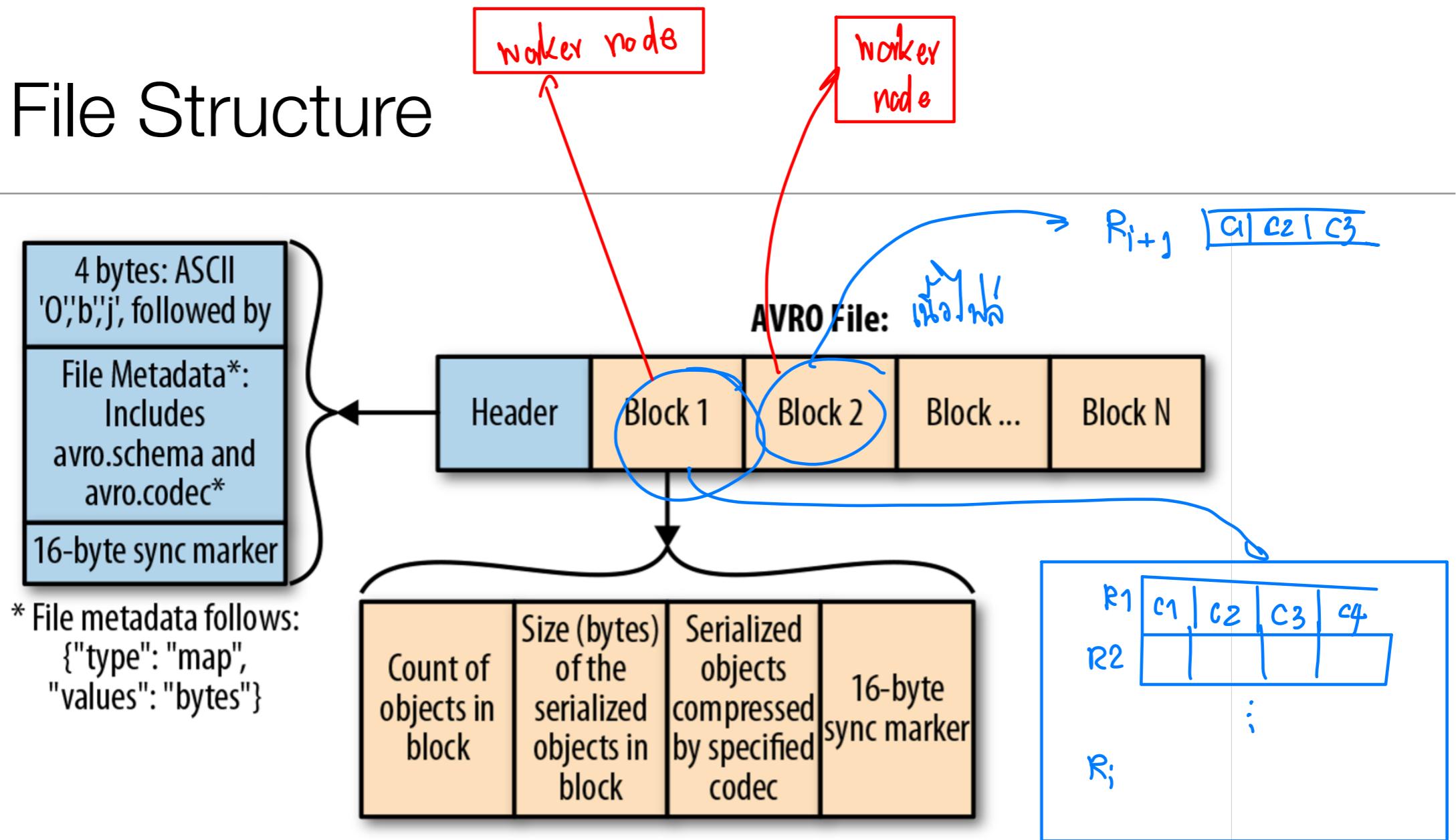
```
In [5]: with open('tweets_v8.avro', 'wb') as out:  
    favro.writer(out, parsed_schema, records) } avro file
```

```
In [6]: with open('tweets_v8_deflate.avro', 'wb') as out:  
    favro.writer(out, parsed_schema, records, codec='deflate') } compressed
```

```
In [7]: print('tweets_v8.csv = {:.5.2f}MB'.format(os.path.getsize('tweets_v8.csv')/(1024*1024)))  
print('tweets_v8.avro = {:.5.2f}MB'.format(os.path.getsize('tweets_v8.avro')/(1024*1024)))  
print('tweets_v8_deflate.avro = {:.5.2f}MB'.format(os.path.getsize('tweets_v8_deflate.avro')/(1024*1024)))
```

tweets\_v8.csv = 24.89MB  
tweets\_v8.avro = 20.49MB *avro file*  
tweets\_v8\_deflate.avro = 11.51MB *compressed*

# AVRO File Structure



- Each row is “serialized”, multiple rows store in each block
- Data compressed per block
- Block's sync marker is for MapReduce file splitting

00000000	4F 62 6A 01 04 14 61 76 72 6F 2E 63 6F 64 65 63	Obj...avro.codec
00000010	08 6E 75 6C 6C 16 61 76 72 6F 2E 73 63 68 65 6D	.null.avro.schem
00000020	61 DC 0A 7B 22 74 79 70 65 22 3A 20 22 72 65 63	a.{"type": "rec
00000030	6F 72 64 22 2C 20 22 6E 61 6D 65 22 3A 20 22 65	ord", "name": "e
00000040	78 61 6D 70 6C 65 2E 61 76 72 6F 2E 63 68 75 6C	xample.avro.chul
00000050	61 2E 73 71 75 69 64 5F 74 77 65 65 74 73 22 2C	a.squid_tweets",
00000060	20 22 66 69 65 6C 64 73 22 3A 20 5B 7B 22 6E 61	"fields": [{"na
00000070	6D 65 22 3A 20 22 75 73 65 72 5F 6E 61 6D 65 22	me": "user_name"
00000080	2C 20 22 74 79 70 65 22 3A 20 22 73 74 72 69 6E	, "type": "strin
00000090	67 22 7D 2C 20 7B 22 6E 61 6D 65 22 3A 20 22 75	g"}, {"name": "u
000000A0	73 65 72 5F 6C 6F 63 61 74 69 6F 6E 22 2C 20 22	ser_location", "
000000B0	74 79 70 65 22 3A 20 5B 22 6E 75 6C 6C 22 2C 20	type": ["null",
000000C0	22 73 74 72 69 6E 67 22 5D 7D 2C 20 7B 22 6E 61	"string"]}, {"na
000000D0	6D 65 22 3A 20 22 75 73 65 72 5F 64 65 73 63 72	me": "user_descr
000000E0	69 70 74 69 6F 6E 22 2C 20 22 74 79 70 65 22 3A	iption", "type":
000000F0	20 22 73 74 72 69 6E 67 22 7D 2C 20 7B 22 6E 61	"string"}, {"na
00000100	6D 65 22 3A 20 22 75 73 65 72 5F 63 72 65 61 74	me": "user_creat
00000110	65 64 22 2C 20 22 74 79 70 65 22 3A 20 7B 22 6C	ed", "type": {"l
00000120	6F 67 69 63 61 6C 54 79 70 65 22 3A 20 22 6C 6F	ogicalType": "lo

19 #SquidGame... https://t.co/N4UGv9hxx8",Twitter Web App,False  
 20 Laura Reactions,France,I talk and I make reactions videos about shows I love #theexpansé #peakyblinders #thelastkingdom  
 #lacasadepapel #atla #theboys #thewitcher,2018-12-19 20:38:28+00:00,330,152,2278,False,2021-10-06 12:05:00+00:00,"The two  
 first games, players were killed by the mask guys ; the bloody night and the third game, they killed each o...  
 https://t.co/Qf057XDJ7C",Twitter Web App,False  
 21 Peyman KAI,United Kingdom,"Official @KardiaChain \$KAI Ambassador  
 22 Marketing Advisor @kephigallery  
 23 Graphics and Film Artist https://t.co/0sT0SEKrHt",2018-01-27 12:07:31+00:00,546,318,6265,False,2021-10-06  
 12:04:54+00:00,"\$THG

ROW ເຊິ່ງລາຍການ ເກີນເສັດກັນ

00000690	74 65 72 20 57 65 62 20   41 70 70 00   1E   4C 61 75
000006A0	72 61 20 52 65 61 63 74   69 6F 6E 73   02   0C 46 72
000006B0	61 6E 63 65 94 02 49 20   74 61 6C 6B 20   61 6E 64
000006C0	20 49 20 6D 61 6B 65 20   72 65 61 63 74   69 6F 6E
000006D0	73 20 76 69 64 65 6F 73   20 61 62 6F 75   74 20 73
000006E0	68 6F 77 73 20 49 20 6C   6F 76 65 20 23 74   68 65
000006F0	65 78 70 61 6E 73 65 20   23 70 65 61 6B 79   62 6C
00000700	69 6E 64 65 72 73 20 23   74 68 65 6C 61 73 74   6B
00000710	69 6E 67 64 6F 6D 20 23   6C 61 63 61 73 61 64 65
00000720	70 61 70 65 6C 20 23 61   74 6C 61 20 23 74   68 65
00000730	62 6F 79 73 20 23 74 68   65 77 69 74 63 68 65 72
00000740	C0 A6 87 83 F9 59   94 05   B0 02   CC 23   00   C0 E7 F0
00000750	D7 8A 5F   9C 02 54 68 65   20 74 77 6F 20 66 69 72
00000760	73 74 20 67 61 6D 65 73   2C 20 70 6C 61 79 65 72
00000770	73 20 77 65 72 65 20 6B   69 6C 6C 65 64 20 62 79
00000780	20 74 68 65 20 6D 61 73   6B 20 67 75 79 73 20 3B
00000790	20 74 68 65 20 62 6C 6F   6F 64 79 20 6E 69 67 68
000007A0	74 20 61 6E 64 20 74 68   65 20 74 68 69 72 64 20
000007B0	67 61 6D 65 2C 20 74 68   65 79 20 6B 69 6C 6C 65
000007C0	64 20 65 61 63 68 20 6F   E2 80 A6 20 68 74 74 70
000007D0	73 3A 2F 2F 74 2E 63 6F   2F 51 66 30 35 37 58 44
000007E0	4A 37 43 1E 54 77 69 74   74 65 72 20 57 65 62 20
000007F0	41 70 70 00 26 50 65 79   6D 61 6E 20 F0 9F 85 9A

ter Web App..Lau  
 ra Reactions..Fr  
 anceö.I talk and  
 I make reaction  
 s videos about s  
 hows I love #the  
 expansé #peakybl  
 inders #thelastk  
 ingdom #lacasad  
 papel #atla #the  
 boys #thewitcher  
 Laçâ·Yö..||#..L\_≡  
 ||è\_£|The two fir  
 st games, player  
 s were killed by  
 the mask guys ;  
 the bloody nigh  
 t and the third  
 game, they kille  
 d each oΓÇª http  
 s://t.co/Qf057XD  
 J7C.Twitter Web  
 App.&Peyman =fàÜ

# Parquet

---

- Columnar schema-oriented data store
- Analytic operations are usually column-based *pandas តារាំង per column  
ទម្លៃវិវាទនៅលើកបញ្ជីដែលបានរាយការណ៍ជាបន្ទាន់*
- Schema is embedded in the file
- Support nested schema using record shredding and assembly algorithm (Google Dremel) to map hierarchical data into column-based data store
- Support column-wise data compression *compress តារ៉ាំង column*
  - More efficient as each column shares the same datatypes
  - Run length encoding, dictionary encoding, delta encoding
- Support in multiple languages, as well as, pandas, MapReduced, and spark

# Parquet Example in Python

---

```
In [1]: import pandas as pd  
import os
```

```
In [2]: parse_dates = ['user_created', 'date']  
df = pd.read_csv('tweets_v8.csv', parse_dates=parse_dates, keep_default_na=False)  
df.shape
```

Out[2]: (80019, 12) *pandas parquet*

```
In [3]: df.to_parquet('tweets_v8.par', compression=None)
```

```
In [4]: df.to_parquet('tweets_v8_snappy.par')
```

```
In [5]: print('tweets_v8.csv = {:.2f}MB'.format(os.path.getsize('tweets_v8.csv')/(1024*1024)))  
print('tweets_v8.par = {:.2f}MB'.format(os.path.getsize('tweets_v8.par')/(1024*1024)))  
print('tweets_v8_snappy.par = {:.2f}MB'.format(os.path.getsize('tweets_v8_snappy.par')/(1024*1024)))
```

tweets\_v8.csv = 24.89MB  
tweets\_v8.par = 22.98MB  
tweets\_v8\_snappy.par = 14.19MB

17 I'm joining the squidgame thing, I'm already dead by sugar honeycomb ofc  
18  
19 #SquidGame... https://t.co/N4UGv9hxx8",Twitter Web App, False  
20 Laura Reactions, France, I talk and I make reactions videos about shows I love #theexpanses #peakyblinders #thelastkingdom  
#lacasadepapel #atla #theboys #thewitcher, 2018-12-19 20:38:28+00:00, 330, 152, 2278, False, 2021-10-06 12:05:00+00:00, "The two  
first games, players were killed by the mask guys ; the bloody night and the third game, they killed each o...  
https://t.co/Qf057xDJ7C",Twitter Web App, False  
21 Peyman KAI, United Kingdom, "Official @KardiaChain \$KAI Ambassador  
22 Marketing Advisor @kephigallery  
23 Graphics and Film Artist https://t.co/0sT0SEKrHt", 2018-01-27 12:07:31+00:00, 546, 318, 6265, False, 2021-10-06  
12:04:54+00:00, "\$THG

```

17 I'm joining the squidgame thing, I'm already dead by sugar honeycomb ofc
18
19 #SquidGame... https://t.co/N4UGv9hxx8",Twitter Web App, False
20 Laura Reactions, France, I talk and I make reactions videos about shows I love #theexpans... #peakyblinders #thelastkingdom
#lacasadepapel #atla #theboys #thewitcher, 2018-12-19 20:38:28+00:00, 330, 152, 2278, False, 2021-10-06 12:05:00+00:00, "The two
first games, players were killed by the mask guys ; the bloody night and the third game, they killed each o...
https://t.co/Qf057XDJ7C",Twitter Web App, False
21 Peyman KAI, United Kingdom, "Official @KardiaChain $KAI Ambassador
22 Marketing Advisor @kephigallery
23 Graphics and Film Artist https://t.co/0sT0SEKrHt", 2018-01-27 12:07:31+00:00, 546, 318, 6265, False, 2021-10-06
12:04:54+00:00, "$THG

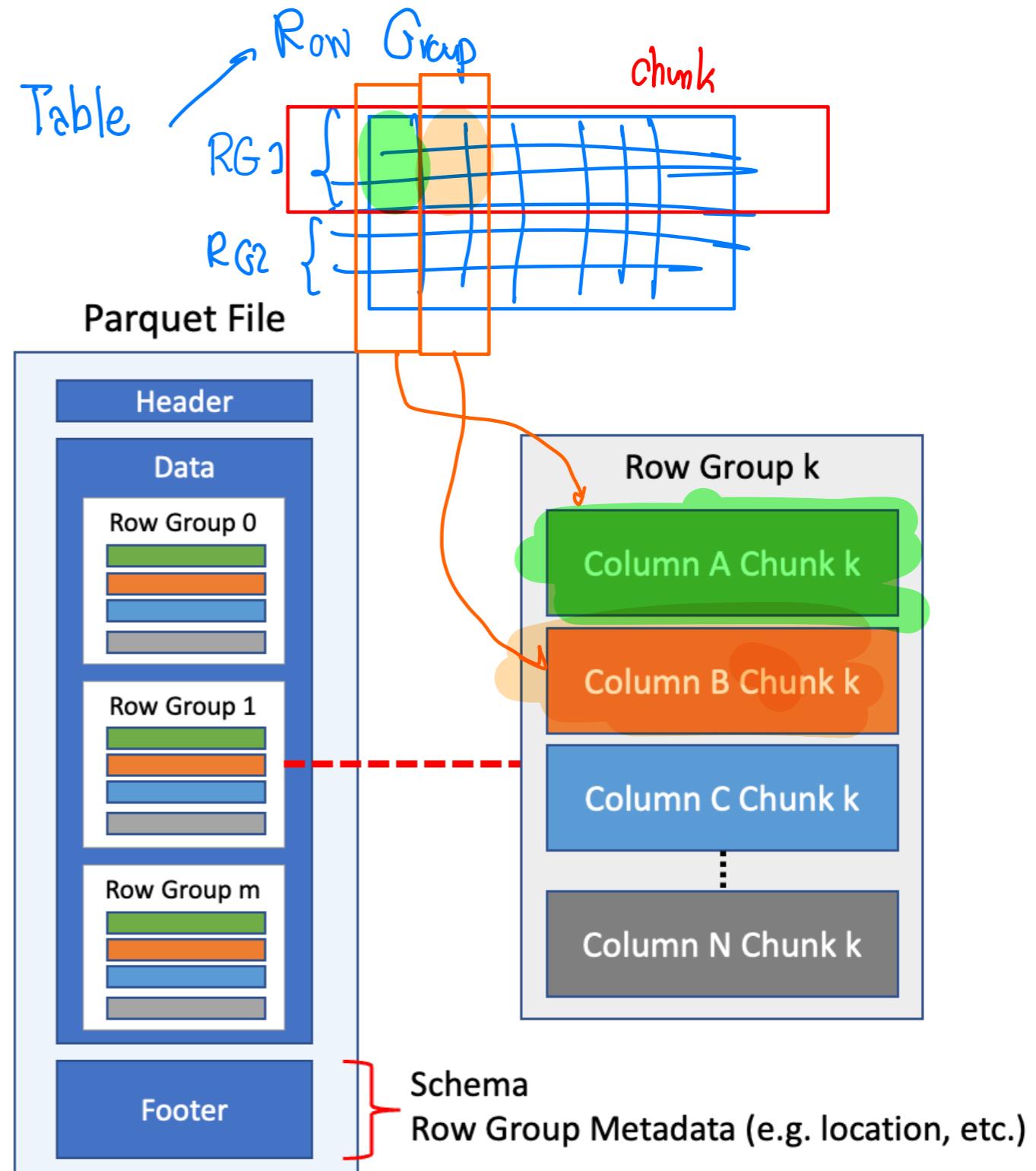
```

00292600	6E 6F 73 75 70 70 6F 72   74 8A 00 00 00 49 20 74	nosupporté... I t
00292610	61 6C 6B 20 61 6E 64 20   49 20 6D 61 6B 65 20 72	alk and I make r
00292620	65 61 63 74 69 6F 6E 73   20 76 69 64 65 6F 73 20	eactions videos
00292630	61 62 6F 75 74 20 73 68   6F 77 73 20 49 20 6C 6F	about shows I lo
00292640	76 65 20 23 74 68 65 65   78 70 61 6E 73 65 20 23	ve #theexpans...
00292650	70 65 61 6B 79 62 6C 69   6E 64 65 72 73 20 23 74	#peakyblinders #t
00292660	68 65 6C 61 73 74 6B 69   6E 67 64 6F 6D 20 23 6C	helastkingdom #l
00292670	61 63 61 73 61 64 65 70   61 70 65 6C 20 23 61 74	acasadepapel #at
00292680	6C 61 20 23 74 68 65 62   6F 79 73 20 23 74 68 65	la #theboys #the
00292690	77 69 74 63 68 65 72 76   00 00 00 4F 66 66 69 63	witcher... Official
002926A0	69 61 6C 20 40 4B 61 72   64 69 61 43 68 61 69 6E	@KardiaChain
002926B0	20 24 4B 41 49 20 41 6D   62 61 73 73 61 64 6F 72	\$KAI Ambassador
002926C0	0A 4D 61 72 6B 65 74 69   6E 67 20 41 64 76 69 73	.Marketing Advis
002926D0	6F 72 20 40 6B 65 70 68   69 67 61 6C 6C 65 72 79	or @kephigallery
002926E0	0A 47 72 61 70 68 69 63   73 20 61 6E 64 20 46 69	.Graphics and Fi
002926F0	6C 6D 20 41 72 74 69 73   74 20 68 74 74 70 73 3A	lm Artist https:
00292700	2F 2F 74 2E 63 6F 2F 4F   73 54 4F 53 45 4B 72 48	//t.co/0sT0SEKrH
00292710	74 4C 00 00 00 46 75 6A   6F 73 68 69 20 F0 9F 99	tL... Fujoshi ≡fÖ
00292720	88 2F 20 54 68 61 69 20   42 4C 2D 6F 62 73 65 73	ê/ Thai BL-obses
00292730	73 65 64 2F 41 6C 77 61   79 73 20 64 69 73 74 72	sed/Always distr
00292740	61 63 74 65 64 20 62 79   20 70 6F 65 74 72 79 20	acted by poetry
00292750	F0 9F 8C B8 2F 20 43 41   50 20 E2 99 91 F0 9F 90	≡fîç/ CAP ΓÖæ≡fÉ

માર્ગદર્શિકા કોમ્પ્લેક્સ

# Parquet File Structure

- Header (4 bytes) 'PAR1' — indicate Parquet file format
- Data contains set of row groups
  - Partition data into **row groups** (m row groups in the example) — for MapReduce parallelization
  - Each row group contains all columns (from A to N in the example) and metadata of the group
  - Column data in the same row group is stored in **a column chunk** with metadata of the chunk — for I/O parallelization
- Footer contains schema and metadata of the file (e.g. location of metadata of each row group and each column in row group)



# Assignment Readings

---

- K. Shvachko, et al. "The hadoop distributed file system." 2010 IEEE 26th symposium on mass storage systems and technologies (MSST), 2010.
- A. Thusoo, et al. "Hive-a petabyte scale data warehouse using hadoop", 2010 IEEE 26th international conference on data engineering (ICDE 2010), 2010.

# References

---

- Hadoop: The Definitive Guide, 4th Edition
- <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsDesign.html>
- <https://www.geeksforgeeks.org/anatomy-of-file-read-and-write-in-hdfs/>
- <https://data-flair.training/blogs/top-hadoop-hdfs-commands-tutorial/>
- <https://data-flair.training/blogs/data-locality-in-hadoop-mapreduce/>
- <https://community.cloudera.com/t5/Community-Articles/Understanding-basics-of-HDFS-and-YARN/ta-p/248860>
- <https://cwiki.apache.org/confluence/display/Hive/Tutorial>
- <https://www.javatpoint.com/dynamic-partitioning-in-hive>