



2110625 - Data Science Architecture

Coordination

Asst.Prof. Natawut Nupairoj, Ph.D.
Department of Computer Engineering
Chulalongkorn University
natawut.n@chula.ac.th

cluster បានចូលរួមរក្សាទុក នៅក្នុង config file យើង និងកំណត់ scale ទៅក្នុង config file
ក្នុងក្រុង service គឺជាបិធីរាយ share config នៃក្នុង member ដើម្បីដោកនឹងទី។ ដើម្បីក្នុងក្នុងក្រុង config

Needs of Cluster Coordination

- Distributed systems, especially large-scale clusters, need centralized coordination services *share config*
- Sharing data among nodes in the cluster e.g. membership, node status, configuration, metadata
- Coordinating between nodes e.g. leader election, agreement resolution
 - leader រឿងនៃ agreement และ solution
ឯកសារកំណត់សំវិធាន - ឬលី (sequencer) ឬនិង conductor
- Resource management e.g. compute nodes, data nodes
 - agreement នៅ 2 clients ទិន្នន័យ node នូវ update data ឱ្យត្រួតពេញ
 - total ordering !!!

យុលវិគ កីឡា peer-to-peer មាន leader កីឡាគិស្ស consistency ឌាក៉ីន (AP)
តាត 2 client នៅ update ត្រូវឱ្យក្នុង ឱ្យគម្រោន node Cassandra ឯកសារ

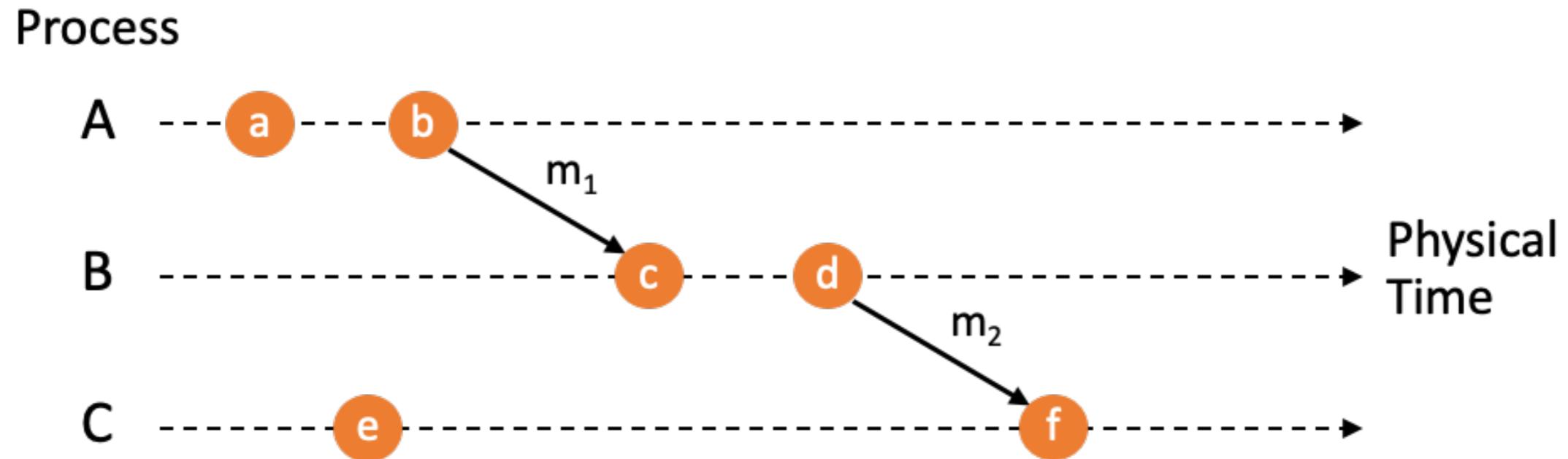
Understand Ordering

- In distributed system, there is no global clock and events can happen concurrently at different nodes
- Ordering is the result of execution based on certain order of events (e_z is executed/happened before e_y) that nodes agree
- Scope of agreements can be varied
 - FIFO ordering - orders are defined at a node that events are issued
 - Causal ordering - orders are defined based on happened-before relations
 - Total ordering - all nodes in the system agrees in one ordering - can be arbitrary-total (any order is OK, as long as everyone agrees), FIFO-total, and causal-total
- Ordering is very important to ensure certain system behaviors

Happened-Before Relation Revisited

- **happened-before** relation is
 - A partial order of events that reflects a flow of information
 - HB1: for any events in process A, if e is executed before e' then $e \rightarrow e'$ (e is happened before e')
 - HB2: for any message m , $send(m) \rightarrow recv(m)$
 - HB3: if $e \rightarrow e'$ and $e' \rightarrow e''$ then $e \rightarrow e''$

Causal Ordering Revisited

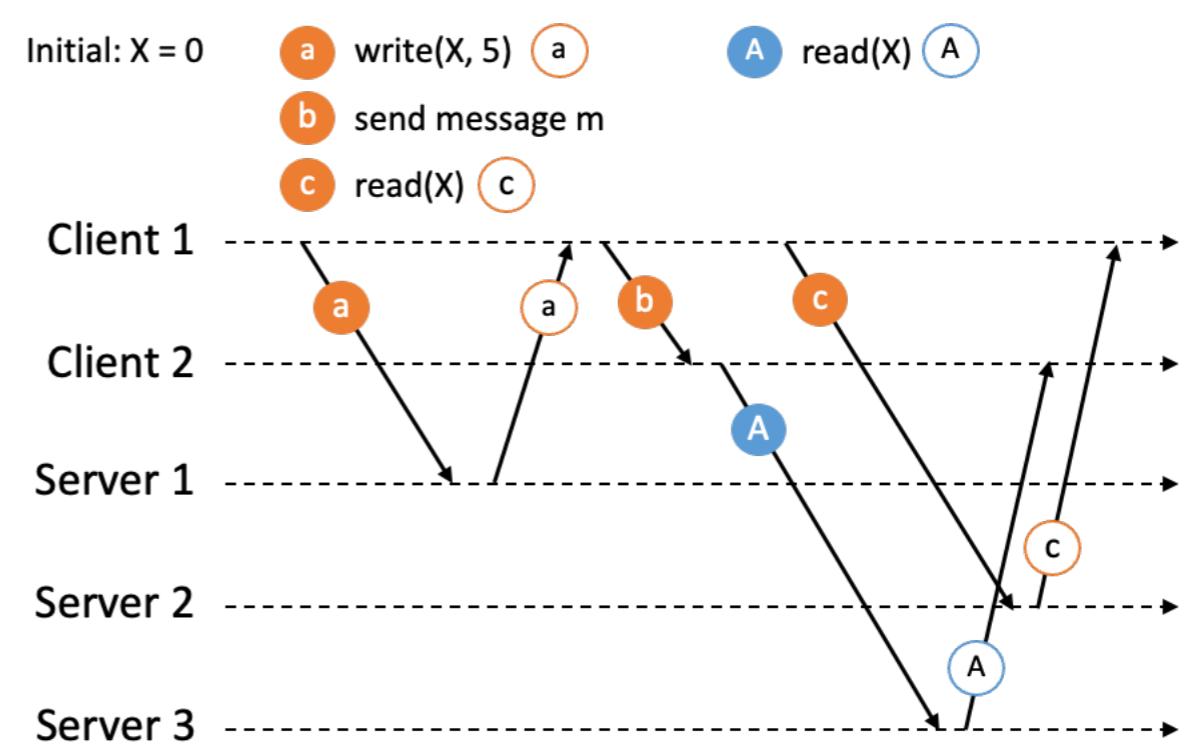
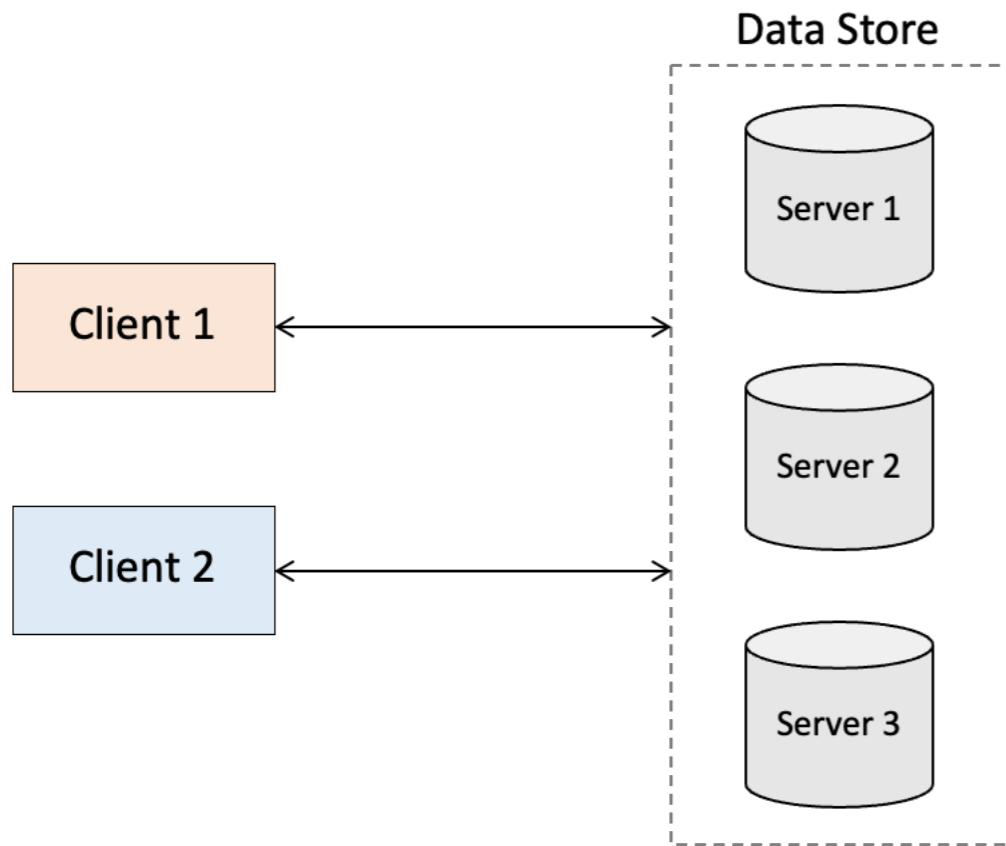


From HB1: $a \rightarrow b$ (at A) $c \rightarrow d$ (at B) $e \rightarrow f$ (at C)

From HB2: $b \rightarrow c$ (at m_1) $d \rightarrow f$ (at m_2)

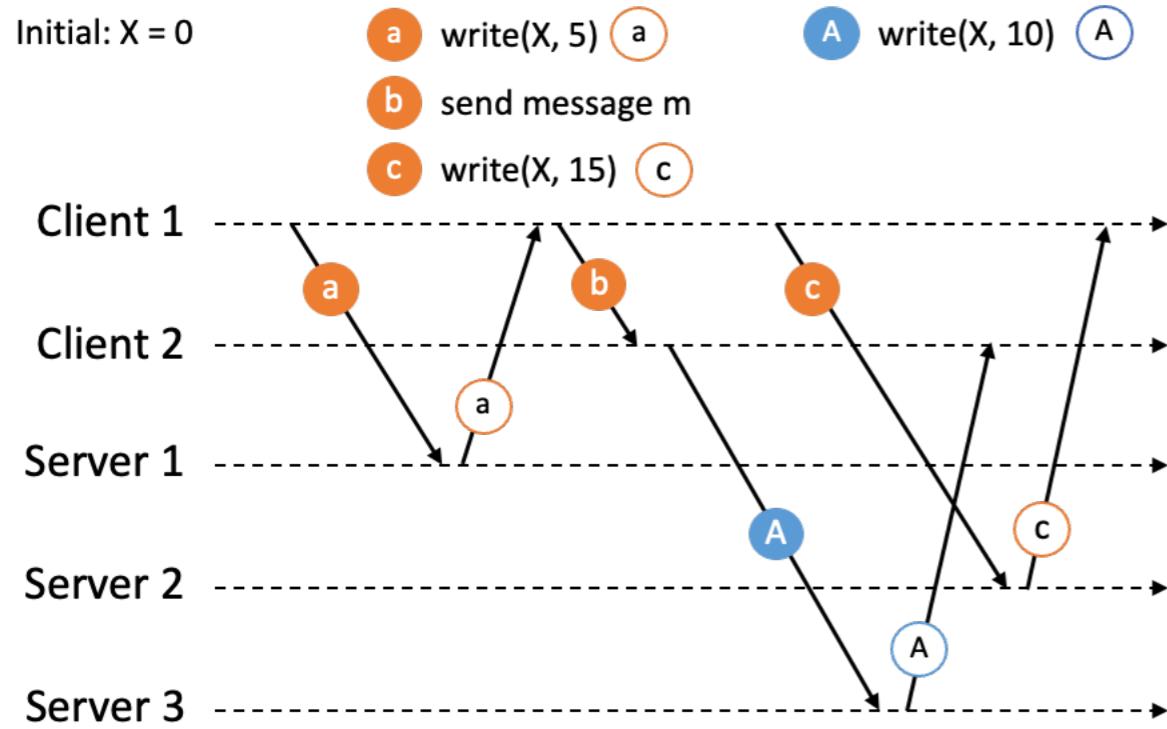
From HB3: $a \rightarrow c$ $b \rightarrow d$... $a \rightarrow f$

Ordering and Data Store



- Suppose a client can send a request to any server and the system guarantees ordering
- For FIFO ordering: as long as execute (c) after (a)
- Causal ordering - as long as execute (A) after (a)

Total Ordering กำหนดลำดับการ update ก่อนหลัง



- All servers will execute requests in the same order
- Arbitrary-total - in any order
- FIFO-total - as long as execute (c) after (a)
- Causal-total - as long as execute (A) after (a)

ກຳນົດ Goggle Chubby
Apache Zookeeper

Coordination



Apache ZooKeeper

ສະ folder ແລະ file ມີບັນ tree

- A distributed hierarchical key-value store with services for configuration management and group services
strong consistency ຮັງການ cluster ແລະ fully replicated
- Provide services similar to Google's Chubby for coordinating processes of distributed applications
 - group messaging ສ້າສະລະ node
 - wait-free shared registers (key-value datastore) ສ້າສະ node ໃນ tree
 - distributed lock services (distributed mutex) locking ດ້ວຍຈຸດເຊີ້ມກັນໄປ tree ຈຳລັມດັດນິ້ນ
ນີ້ແມ່ນມີການ
- Strong consistency and ordering guaranteed
- Perform replication among ZooKeeper nodes to provide reliability and performance improvement
ທີ່ມີ date ເພື່ອຈຳນາດສົມ ແລະ ຕຽບຄົນຮັມຂອງ

Zookeeper Primitive: Znodes

- ZooKeeper provides in-memory data objects (znodes) with 1MB storage and meta-data

- Znodes are organized in hierarchical name space

- Node lifetime

- Permanent (manually delete) อยู่ต่อไปจนกว่าจะลบออกมือ

- Ephemeral (creator's session lifetime) ถ้า node เกิดขึ้นในตอนที่ล็อกอิน zk

- Node's name (sequential flag)

- Flag not set: normal name ชื่อปกติ

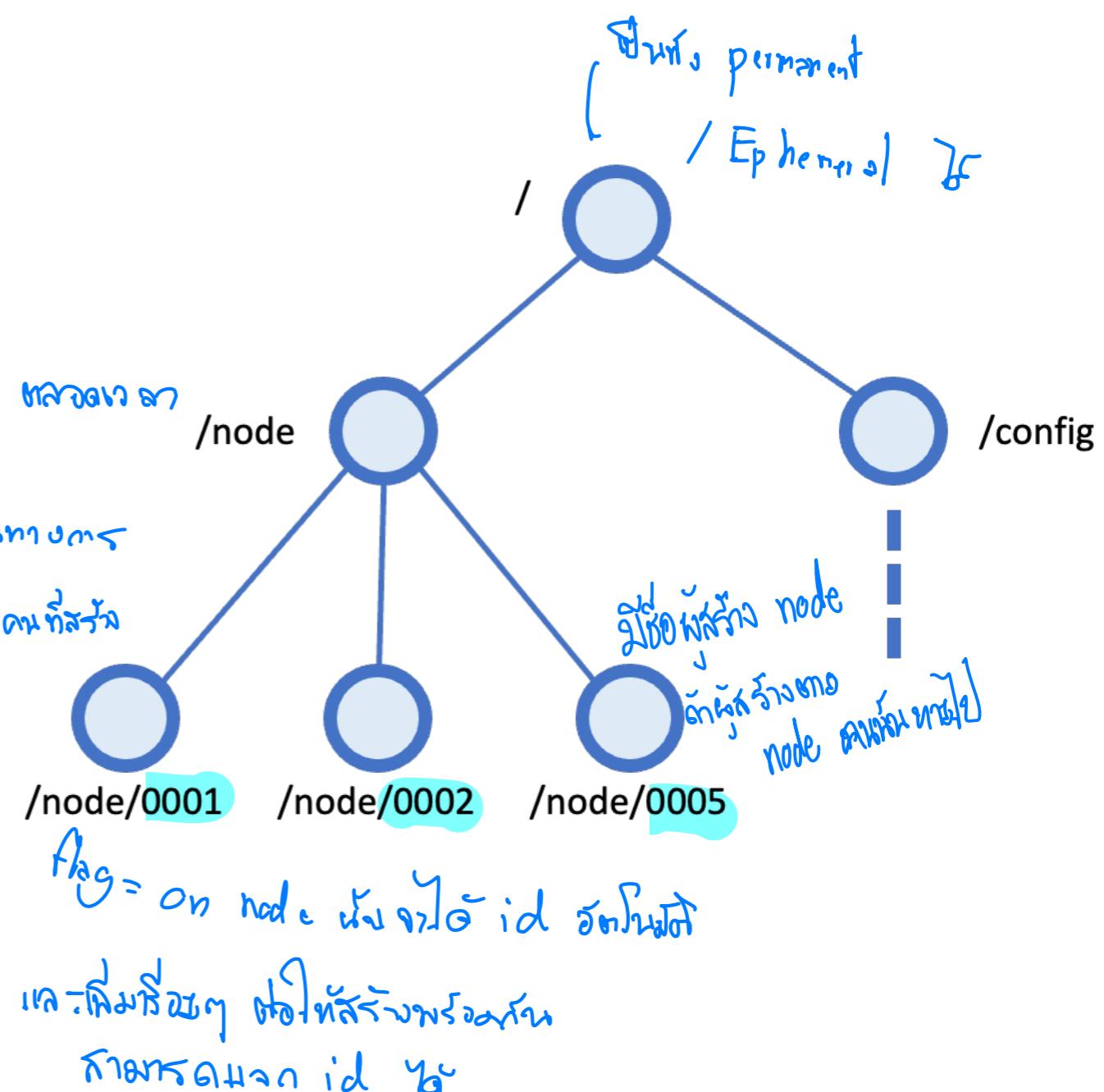
- Flag set: a new node will have a monotonically increasing counter appended to its name ตัวลับซึ่งเพิ่มขึ้นเรื่อยๆ

- Client can ask zookeeper to notify if there is any change to a specific znode (watch) ขอเฝ้า node นี้ว่ามีอะไรเปลี่ยนแปลง

zookeeper & heartbeat წິ້ນ node ແລະ ຂອງຂາຍ

ໃນ zk
(node status)

ກັບ node ສະບັບ



ZooKeeper Session

- Client is assigned a session when connecting to ZooKeeper
- Client becomes faulty if fails to communicate with ZooKeeper within a period of time
Client មិនចិត្តទៅ zk វិញ session នេះ
or faulty (បាន)
- Session ends when client closes the session or ZooKeeper detects that a client is faulty
Client សម្រេចបើកត្រូវ, បន្ថែម session, ឬមានការពិន័យ
terminate session
- Session persists across ZooKeeper servers allowing client to connect to any server using the same session
ការរួចរាល់ session នៅលើកណ្តាល

ZooKeeper API

create(path, data, flags)

getData(path, watch)

setData(path, data, version)

delete(path, version)

exists(path, watch)

getChildren(path, watch)

sync(path)

Running Zookeeper

- The simplest way to run a zookeeper instance is to use docker

```
docker pull wurstmeister/zookeeper
```

```
docker run -d --rm --name zookeeper -p 2181:2181 wurstmeister/zookeeper
```

- This will start zookeeper in your docker at port 2181 and map the port to your localhost
- You can also use docker-compose.yml and other example files in redis folder in datasci_architecture repo
- Kazoo - standard python package for Zookeeper client

```
pip install kazoo
```

ZooKeeper API - Python Kazoo

create(path, data, flags)

get(path, watch)

set(path, data, version)

delete(path, version)

exists(path, watch)

get_children(path, watch)

sync(path)

Example: Simple Kazoo Client

```
from kazoo.client import KazooClient

def zk_tree(zk, path, ident=''):
    print(ident, path)
    nextident = ident + ' '
    for c in zk.get_children(path):
        nextpath = path + ('' if path.endswith('/') else '/') + c
        zk_tree(zk, nextpath, nextident)

zk = KazooClient(hosts='localhost')
zk.start()
zk.ensure_path('/mycluster/nodes') creates path
zn_node = zk.create('/mycluster/nodes/', ephemeral=True, sequence=True, makepath=True)
print('zn_node =', zn_node)
zk_tree(zk, '/mycluster')

zk2 = KazooClient(hosts='localhost')
zk2.start()
zn_node2 = zk2.create('/mycluster/nodes/', ephemeral=True, sequence=True, makepath=True)
print('zn_node2 =', zn_node2)
zk_tree(zk, '/mycluster')
zk2.stop()
print('after zk2 stops')
zk_tree(zk, '/mycluster')
zk.stop()

zk.start()
print('after both zk and zk2 stop')
zk_tree(zk, '/mycluster')
zk.delete('/mycluster', recursive=True)
zk.stop()
```

*sequence flag
c =on*

Results

```
...
zk.start()
zk.ensure_path('/mycluster/nodes')
zn_node = zk.create('...')
print('zn_node =', zn_node)
zk_tree(zk, '/mycluster')

zk2 = KazooClient(...)
zk2.start()
zn_node2 = zk2.create('...')
print('zn_node2 =', zn_node2)
zk_tree(zk, '/mycluster')
zk2.stop()
print('after zk2 stops')
zk_tree(zk, '/mycluster')
zk.stop()

zk.start()
print('after both zk and zk2 stop')
zk_tree(zk, '/mycluster')
...

```

Session 2

zn_node = /mycluster/nodes/0000000000
/mycluster
/mycluster/nodes
/mycluster/nodes/0000000000
zn_node2 = /mycluster/nodes/0000000001
/mycluster
/mycluster/nodes
/mycluster/nodes/0000000000
/mycluster/nodes/0000000001
after zk2 stops
/mycluster
/mycluster/nodes
/mycluster/nodes/0000000000
after both zk and zk2 stop
/mycluster
/mycluster/nodes

โน้ต node
co.a._o

บัญชีนี้จะลบหลังจาก session 2 ถูกปิด

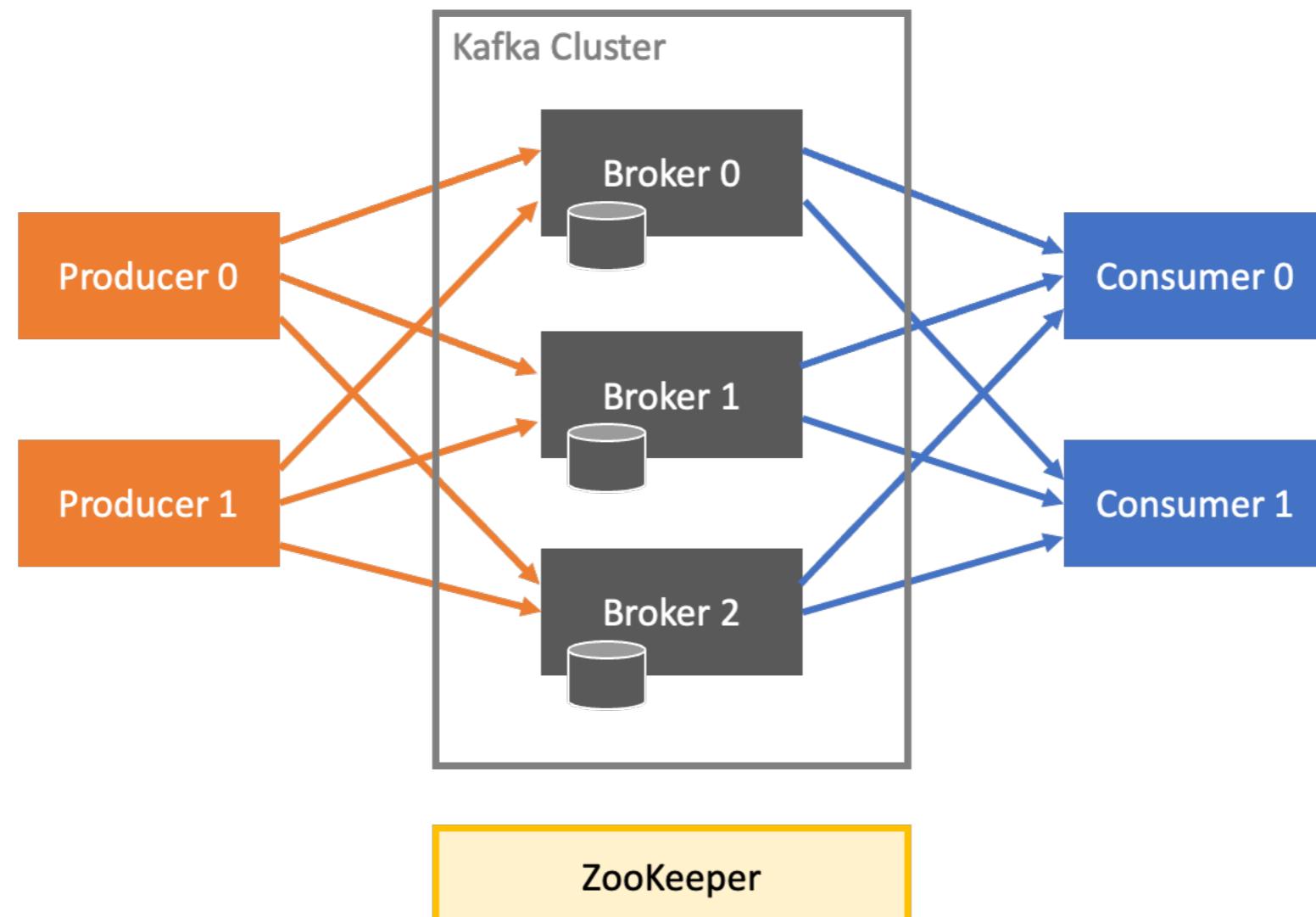
ก็ต้าบๆ ที่ terminal นี้ มี 2 รายการรันร่วมกัน id = 0 กับ id = 2

Example: Leader Election with Zookeeper

- Zookeeper can simplify several group services e.g. electing a leader of the group
បានកែត្រូវ node នៃក្រុង group ដែលមែន leader ត្រូវគេចូលរួមជាលើកដែលត្រូវការពារ (ephemeral) ដើម្បីតាមរយៈការចូលរួមដែលត្រូវការពារ (sequential) ដើម្បីតាមរយៈការចូលរួមដែលត្រូវការពារ (sequential)
- Elect a leader:
 - All members join a group and get unique ids by creating their own ephemeral znode in /members/ with sequential flag set and using the znode name as id
 - All members enter election by creating an ephemeral znode /leader; only one member will be successful and that member becomes the leader
 - Other non-leader nodes watch /leader znode, if the leader dies, /leader znode will be deleted and everyone re-enter election process again
- See more details in `leader_election.py`

, message broker

Example: Kafka and ZooKeeper



Example: Kafka and ZooKeeper

- Data in ZooKeeper for kafka cluster
 - 2 brokers
 - nodes contain metadata
 - ephemeral nodes
 - 3 topics
 - contain responsible info
 - regular nodes (ephemeralOwner)

```
/cluster
/cluster/id
/controller_epoch
/controller
/brokers
/brokers/ids
/brokers/ids/1002 } broker id in cluster
/brokers/ids/1001
/brokers/topics
/brokers/topics/twitter topic 1
/brokers/topics/twitter/partitions
/brokers/topics/twitter/partitions/0
/brokers/topics/twitter/partitions/0/state
/brokers/topics/sample topic 2
/brokers/topics/sample/partitions
/brokers/topics/sample/partitions/0
/brokers/topics/sample/partitions/0/state
/brokers/topics/sample/partitions/1
/brokers/topics/sample/partitions/1/state
/brokers/topics/avro topic 3
/brokers/topics/avro/partitions
/brokers/topics/avro/partitions/0
/brokers/topics/avro/partitions/0/state
/brokers/topics/avro/partitions/1
/brokers/topics/avro/partitions/1/state
/brokers/seqid
```

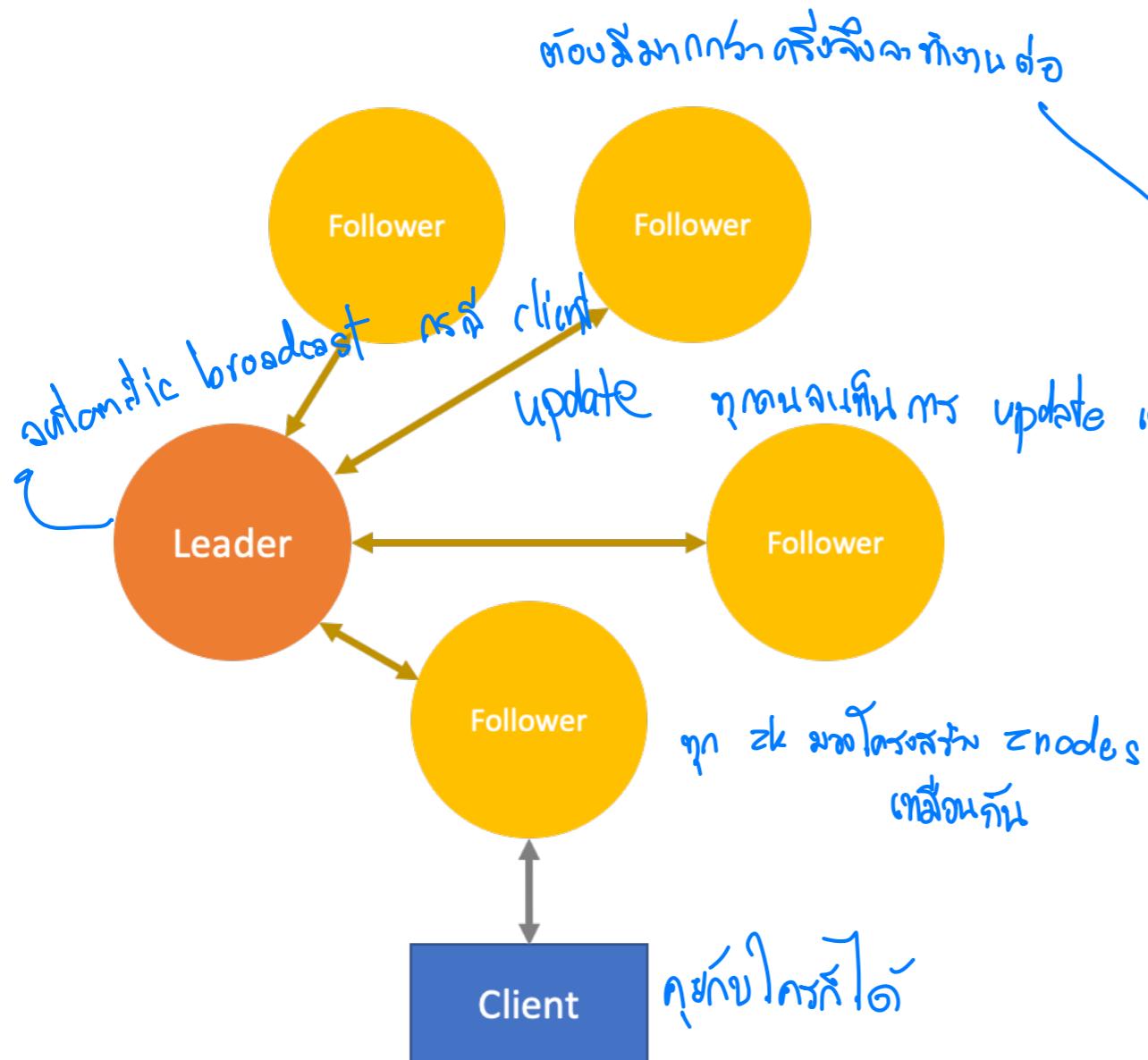
save ລາ ຊກ ໄນດັບສິ້ນເຊີ່ງ

client may zk send kafka Topic

```
zk.get('/brokers/ids/1001')
(b'{"features":{}, "listener_security_protocol_map":{"PLAINTEXT":"PLAINTEXT"}, "endpoints":["PLAINTEXT://192.168.68.1
37:9093"], "jmx_port":-1, "port":9093, "host":"192.168.68.137", "version":5, "timestamp":"1630329937692"},'
ZnodeStat(czxid=41, mzxid=41, ctime=1630329937720, mtime=1630329937720, version=1, cversion=0, aversion=0, ephemeralOwner=72065036800950272, dataLength=212, numChildren=0, pzxid=41))

topic : twitter
zk.get("/brokers/topics/twitter")
(b'{"version":2, "partitions":{"0": [1002]}, "adding_replicas":{}, "removing_replicas":{}',
ZnodeStat(czxid=72, mzxid=72, ctime=1630329946816, mtime=1630329946816, version=0, cversion=1, aversion=0, ephemeralOwner=0, dataLength=83, numChildren=1, pzxid=74))
```

ZooKeeper Architecture



- ZooKeeper cluster is called **ensemble**
- Ensemble is still working as long as half of the members are still functioned
- There is one leader who handles all information updating
- Once updated, leader sends an atomic broadcast to update all other nodes (followers)
- Both leader and followers provide services to clients
- If the leader fails, one of the followers will become a leader
- ZooKeeper guarantees
 - Linearizable writes
 - FIFO client order

Leader handles all update requests
Leader handles all update requests

Apache Yarn

继 Hadoop v.2

Coordination



HDFS ໄວ້ data ໂປ່ນ dataNode

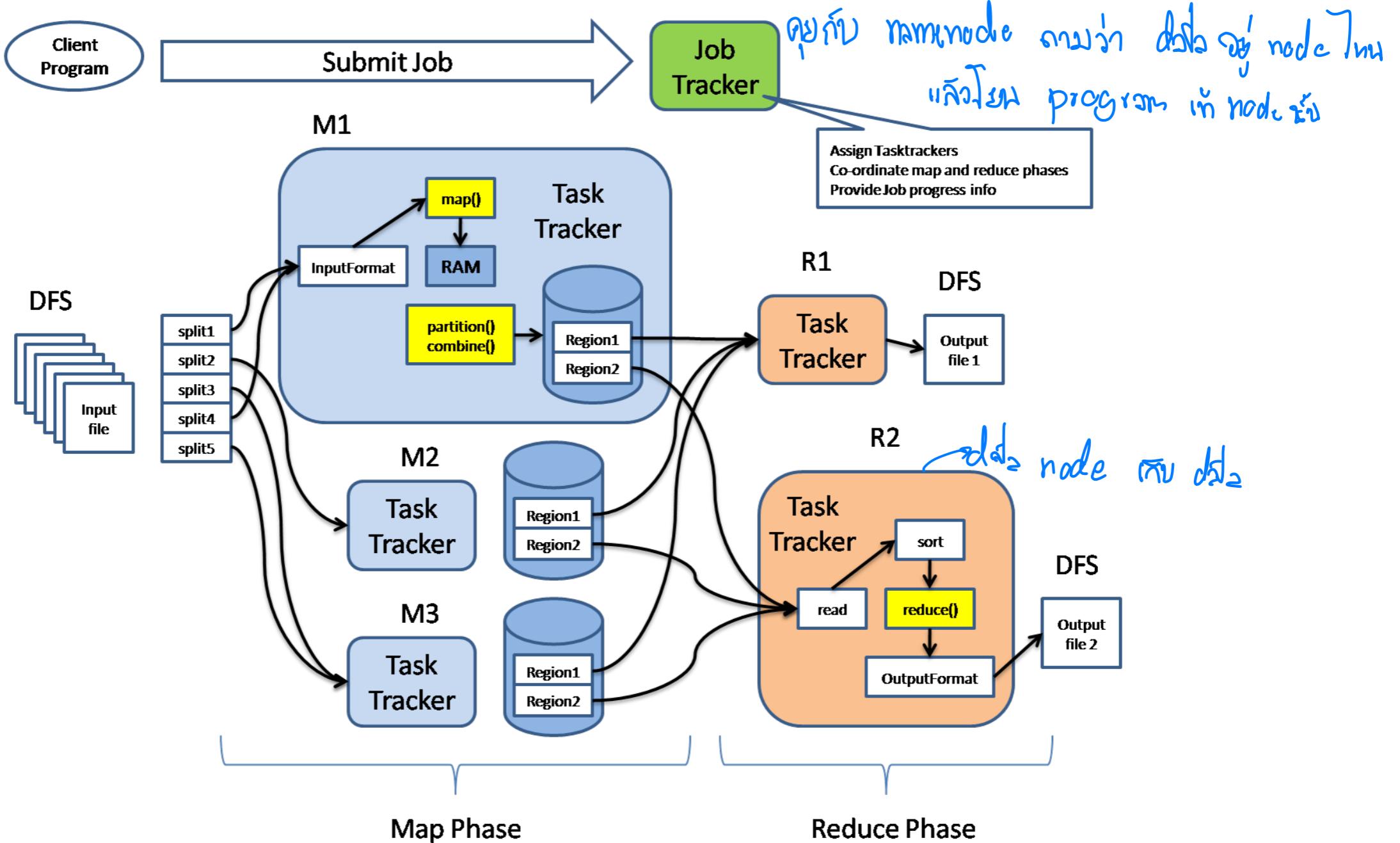
ເຊື່ອຮັບ mapReduce Hadoop ຈະໂປ່ນ code ອາດໄວ້ (ເລືອງ node ກໍ່ມີ data
ກໍ່ຕໍ່ທີ່ພາຍຫຼຸດ run ເງິນາກໂປ່ງໂດຍຢູ່)

Apache Yarn

- Hadoop version 2 introduced a new resource manager called Yarn (Yet Another Resource Negotiator)
ເພື່ອເຄີຍ JobTracker ທີ່ມີ mapReduce ໂພນທີ່ມີ YARN ເສົ່າງຄວາມທຳການທາງ model
- Replace JobTracker, which was inflexible and limited scalability
- Provide lots of functionality: scalability, multi-tenancy, reliability, etc.
ສໍາຄັນການກັບຮົງຂອງ YARN ເທົ່າລົງເກມອາຄອງຢູ່ດັນໃຈ programming model
- Provide REST API
- Flexible for many programming models e.g. MapReduce, Spark, Flint, etc.



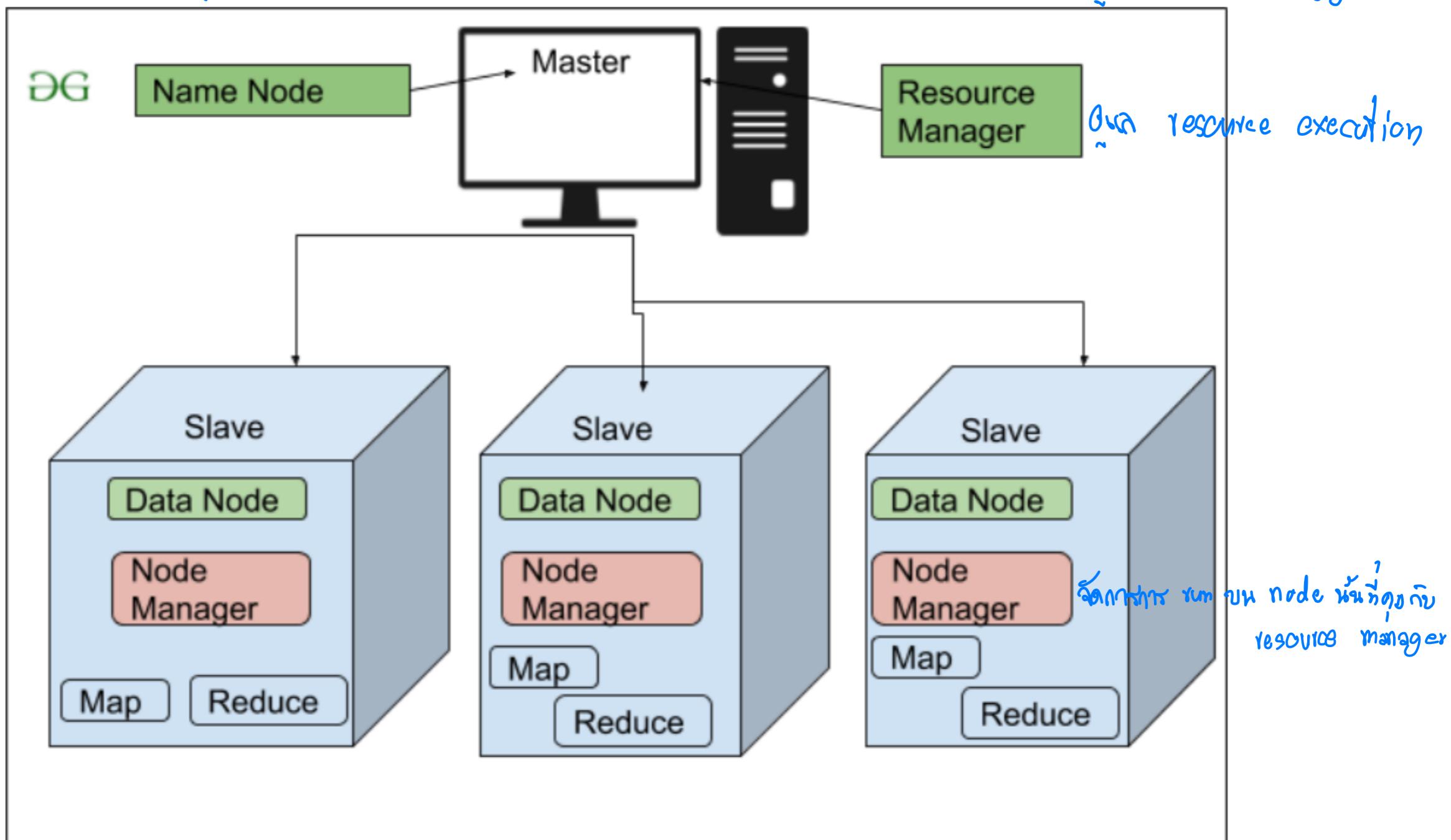
Hadoop Version 1 - JobTracker



Hadoop Version 2 and Later

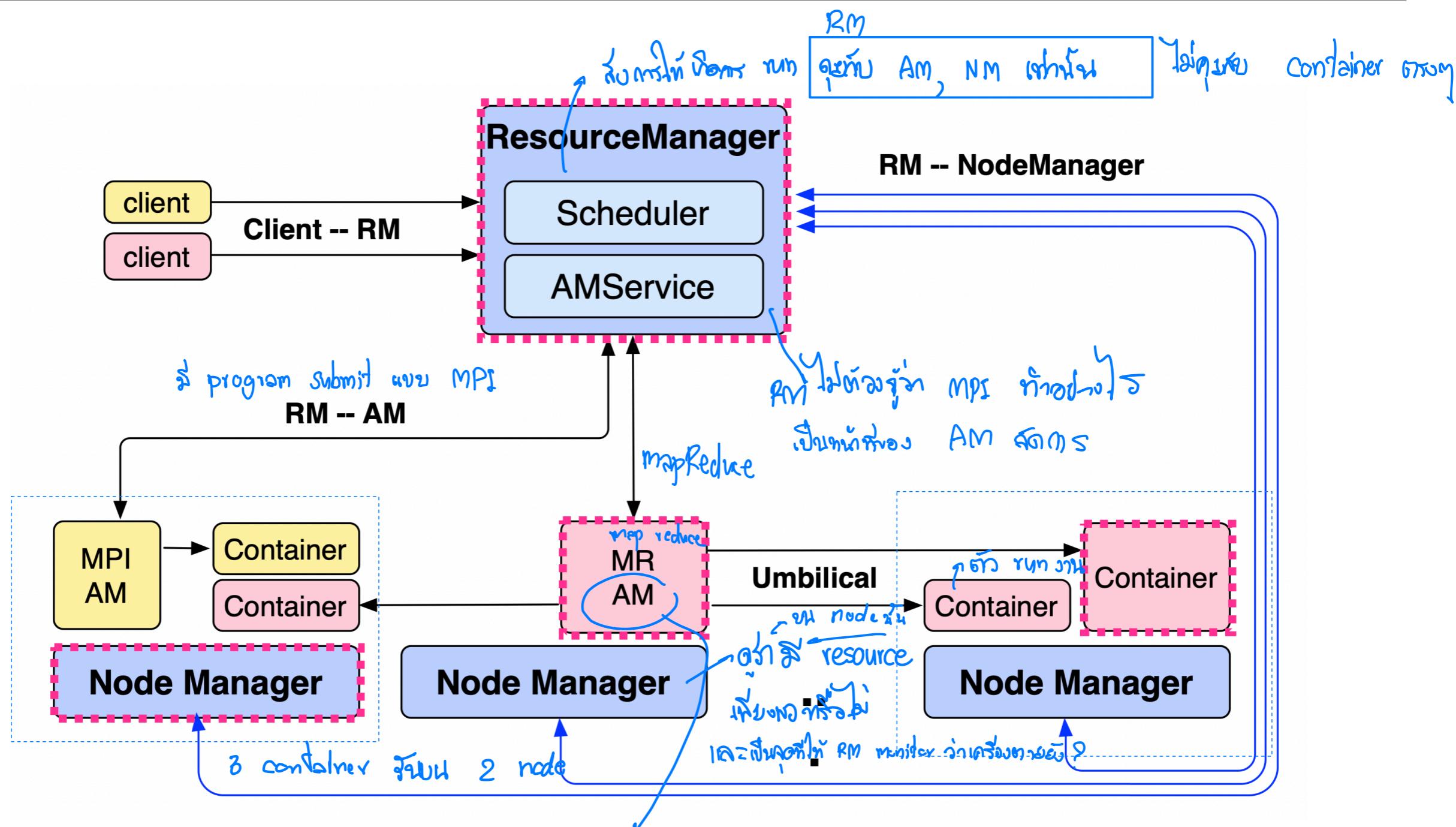
HADOOP

HDFS ດຸກ ດັບ ລາຍລະອຽດ



RM ส่งคำสั่งให้ AM run และดูแลการ run ของ AM ให้เสร็จ

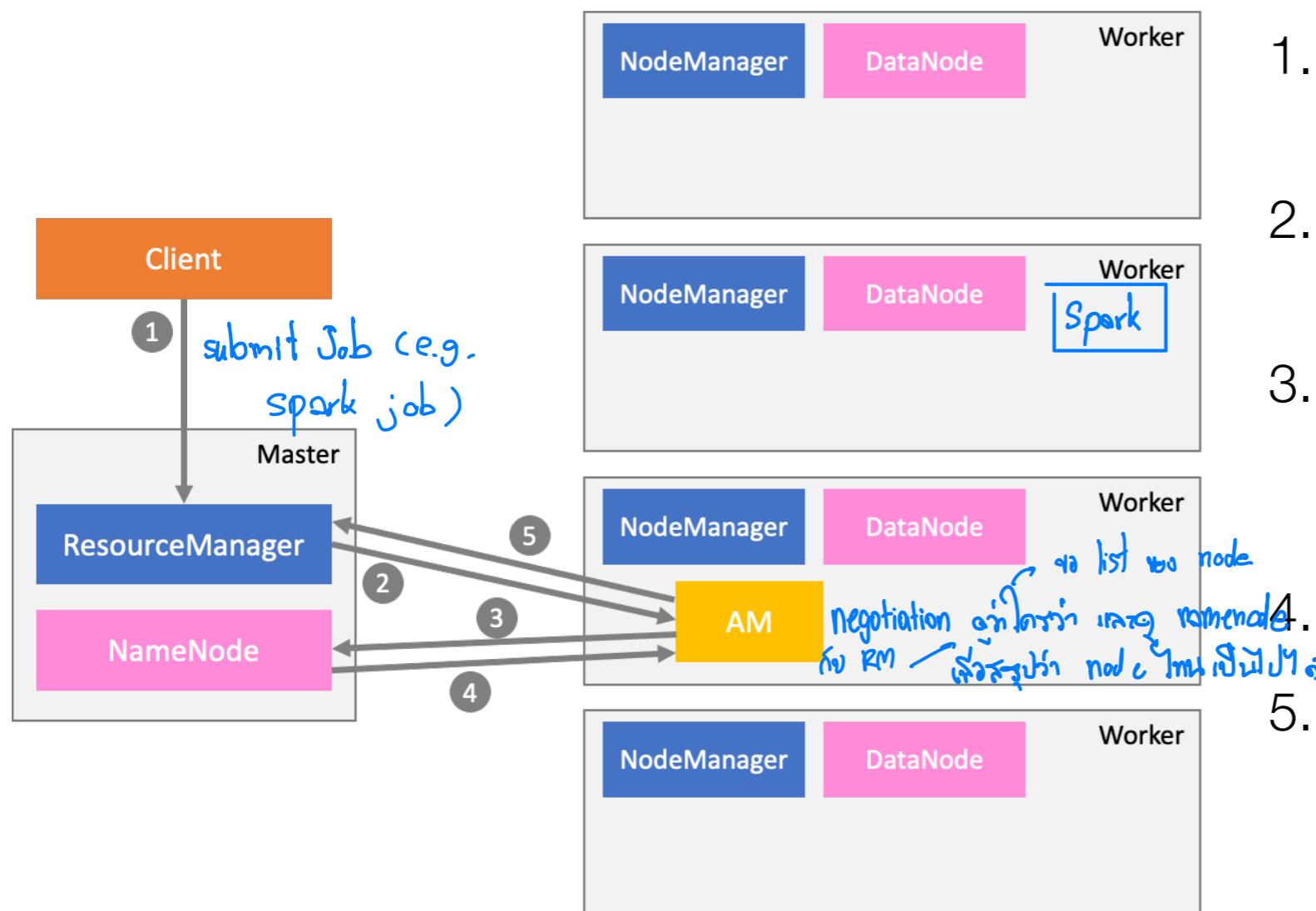
YARN Architecture



AM = ApplicationMaster

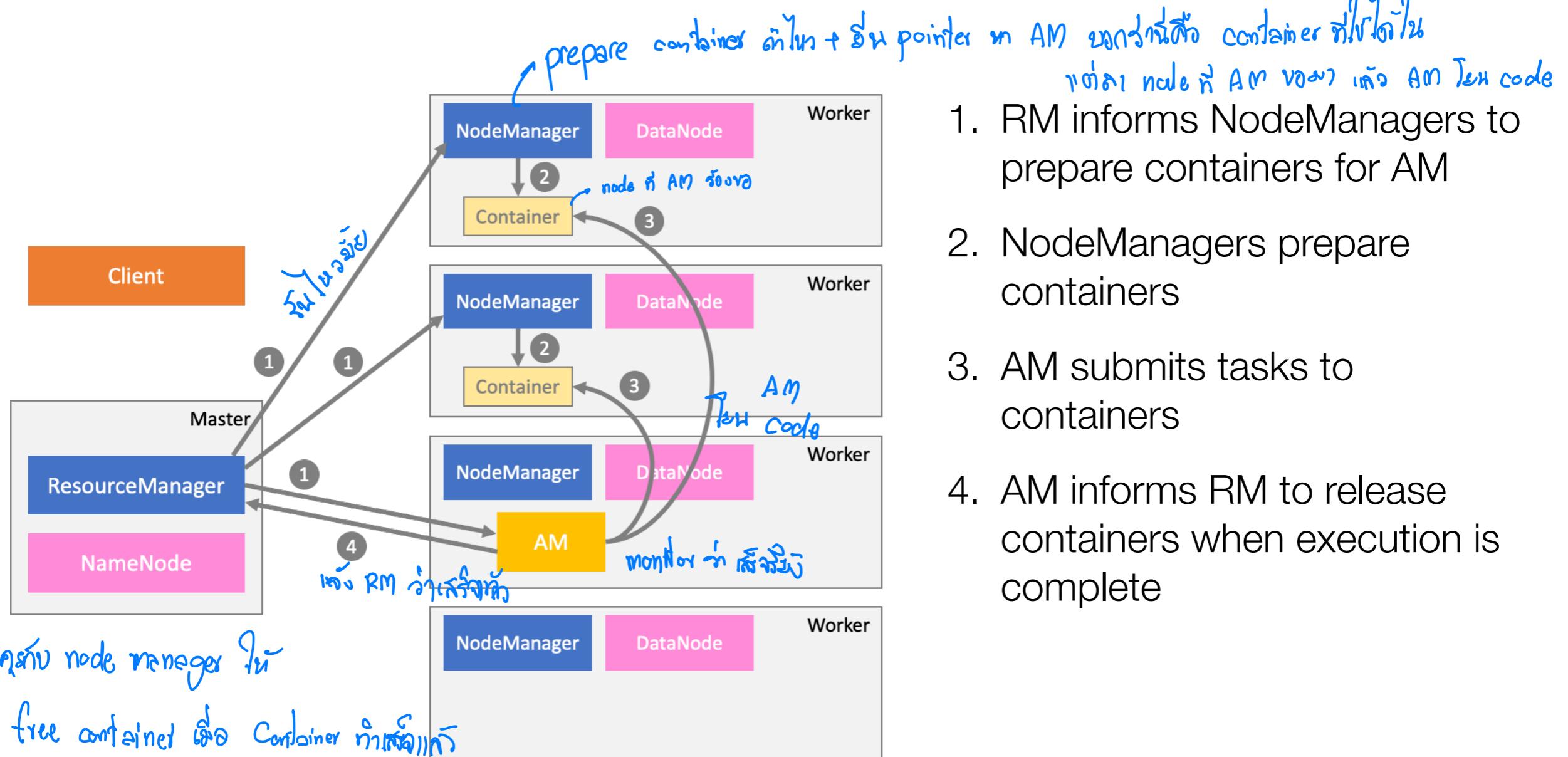
Application Manager ผันแปร Programming Paradigm
ที่ map reduce AM ต้องมีใน map Reduce คือ map/reduce

Resource Negotiation



1. Client requests resources from RM
2. RM reserves a container from a worker to run AM *in node to worker*
3. AM asks NameNode a list of DataNode with data to-be-processed in HDFS
4. NameNode provides the list
5. AM chooses DataNodes and sends request to RM

Resource Usage



Yarn Scalability and Reliability

- There are 4 possible failure points ចំណោម
 - Container - AM asks for a new container from RM AM និង container នេះ
 - AM - RM monitors AM status and restart AM in a new container (if necessary). It is AM responsibility to recover from the failures RM ចក្ខុង AM នៅ node នៃវេសាន AM នៅ AM ត្រូវរាយការណ៍ recovery នៅ
 - NodeManager - RM monitors NodeManager and removes failed NodeManager and its containers. RM informs AM to recovery from the failures RM ចក្ខុង NodeManager (ដែល RM គឺជាដែល container ត្រូវ) នៅ AM វា NM ត្រូវ ប្រព័ន្ធបន្ទាត់ AM ក្នុង RM នៅក្នុង
 - RM - used to be a **single point of failure**, however, has been developed with active-standby model in the later versions
 - Hadoop v.3 RM នឹងមកជាមួយ RM backup (Journal node)
 - RM នឹង Active, standby

Conclusion

- In large-scale cluster, coordinating between components become important, including:
 - group membership and leader election
 - resource management / placement
 - meta-data / configuration information storages
 - locking mechanism
- Apache ZooKeeper provides centralized services for coordination and synchronization
- Apache Yarn provides resource management services with flexibility
- All centralized services must have failure recovery processes as failure can happen

Reading Assignments

- P. Hunt et al, “ZooKeeper: Wait-free coordination for Internet-scale systems”, USENIX annual technical conference. Vol. 8. No. 9. 2010.
- V. Vavilapalli et al, “Apache Hadoop YARN: Yet Another Resource Negotiator”, Proceedings of the 4th annual Symposium on Cloud Computing. 2013.