



2110625 - Data Science Architecture

## Wide Column Stores

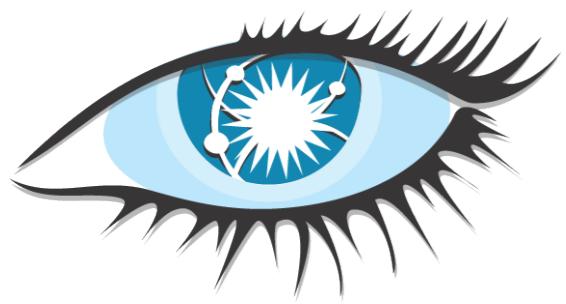
**Asst.Prof. Natawut Nupairoj, Ph.D.**

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

# Wide Column Store

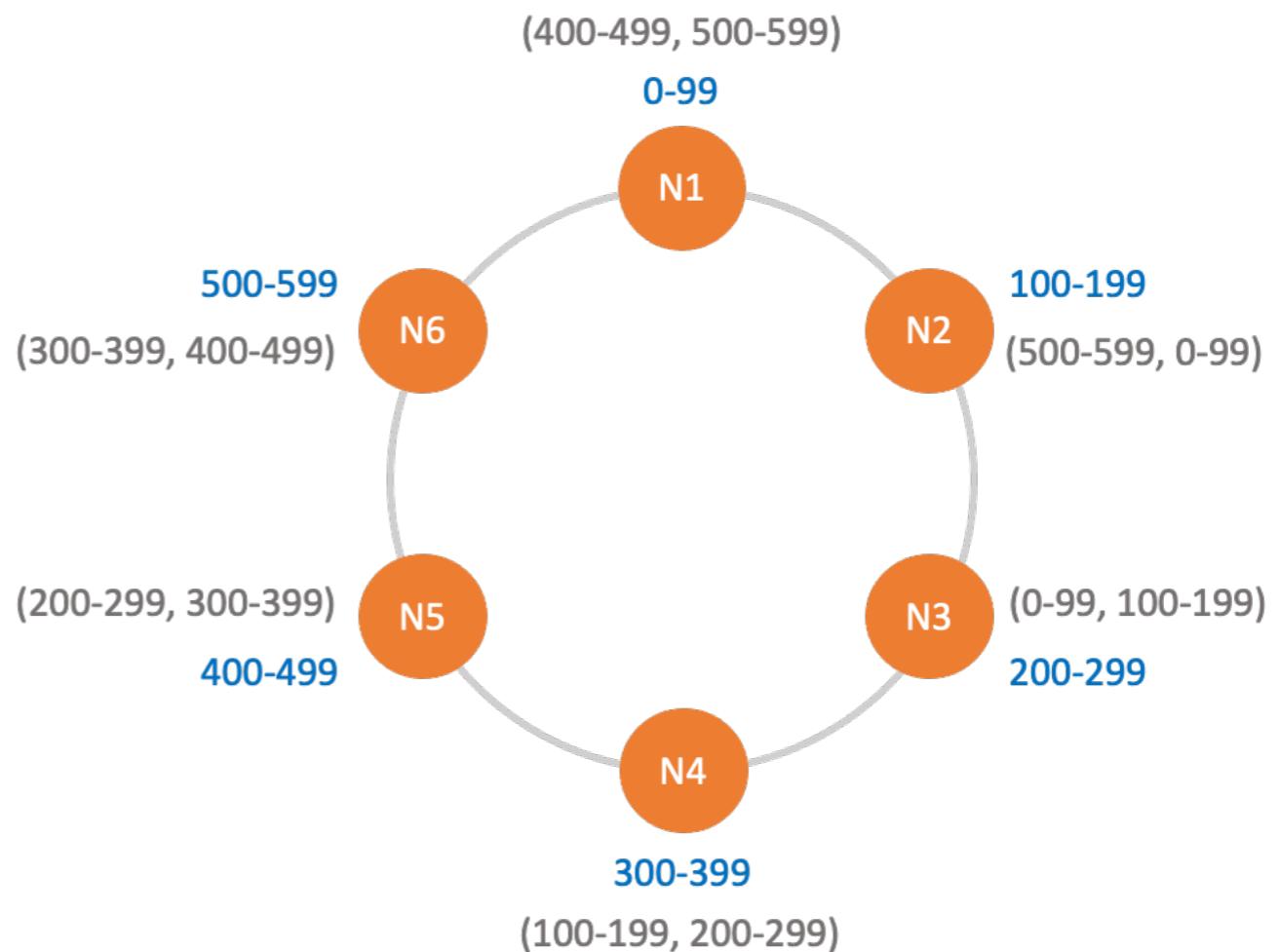
ផ្លាស់ប្តូរ column

- A column-oriented DBMS with tables, rows, and columns
- Names and format of the columns can vary from row to row in the same table
- A wide-column store can be interpreted as a two-dimensional key-value store with group of columns are stored in the same server
  - Often support the notion of column families
    - Each column family typically contains multiple columns that are used together
    - Within a given column family, all data is stored in a row-by-row fashion, such that the columns for a given row are stored together



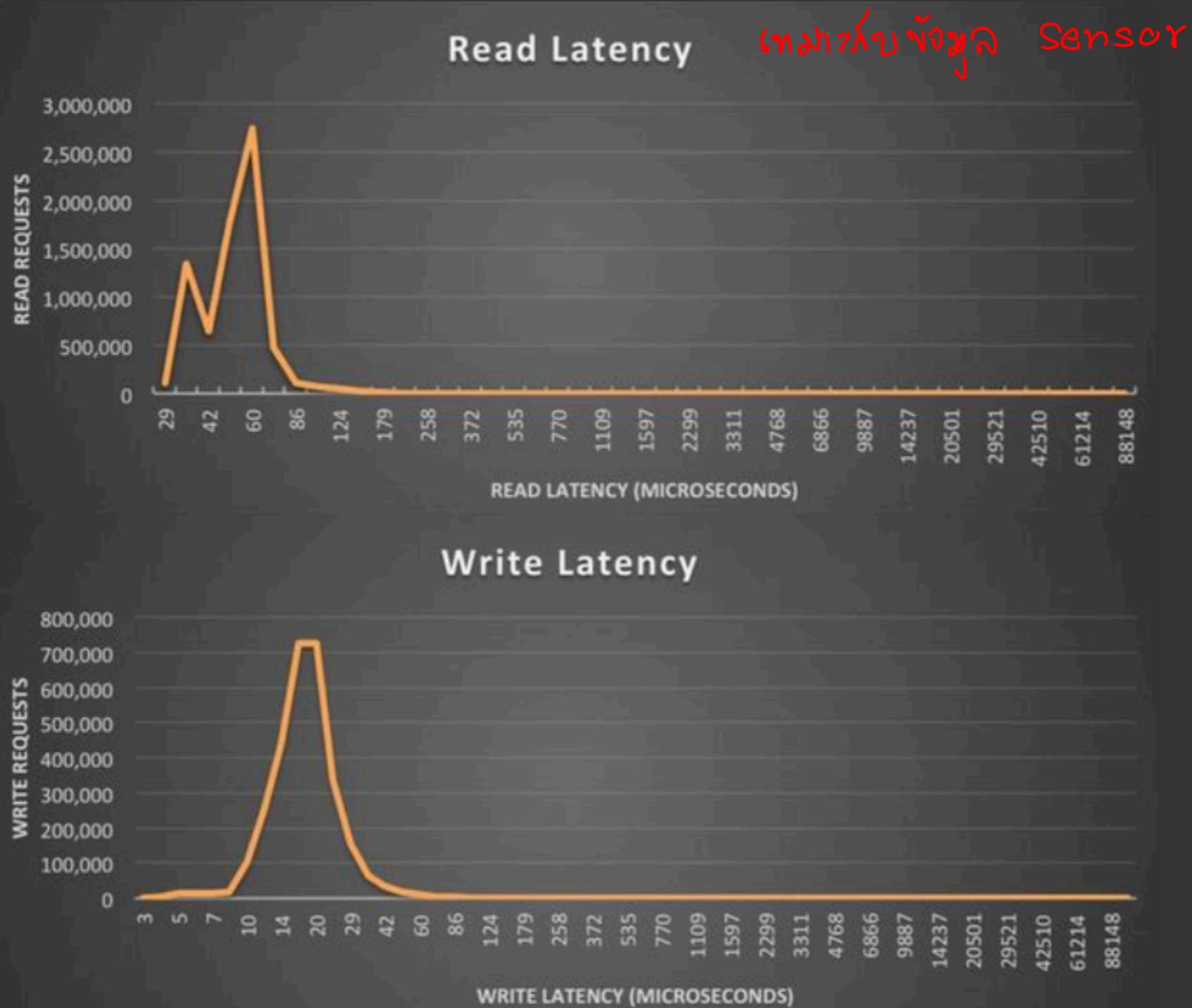
# Apache **CASSANDRA**

peer-to-peer  
ມະນຸຍາ  
ກົດຕະຫຼາດໄປ່ງ່າຍ໌



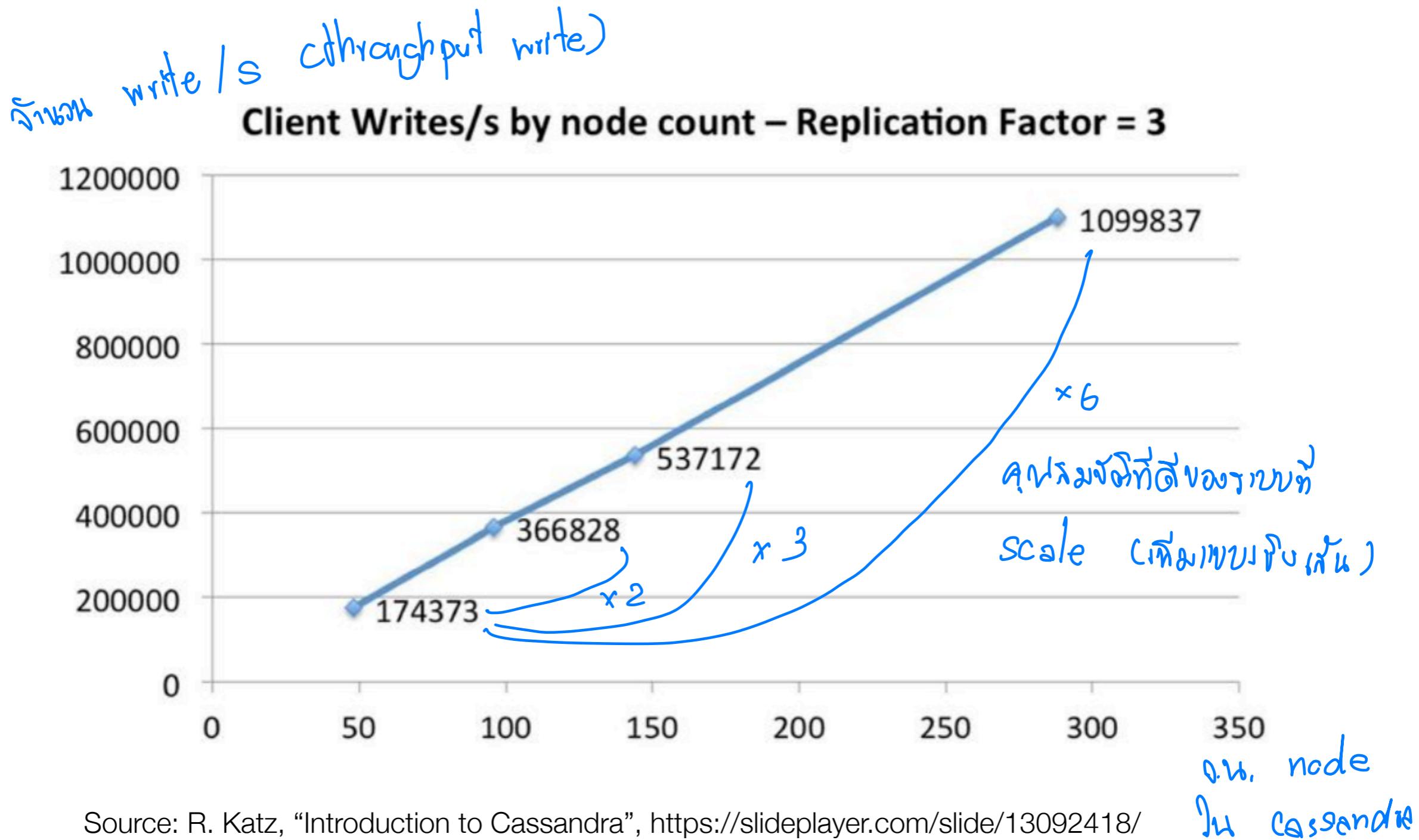
Apache Cassandra

# Cassandra



Source: R. Katz, "Introduction to Cassandra", <https://slideplayer.com/slide/13092418/>

# Cassandra is highly scalable



Source: R. Katz, "Introduction to Cassandra", <https://slideplayer.com/slide/13092418/>

# RDBMS vs. Cassandra

---

RDBMS

Database

Table

Row / Tuple

Column

Table Join

Primary Key

Cassandra

Keyspace

Table

Row

Column

Denormalization

Partition Key + Clustering Column

primary key    sort key    dynamoDB

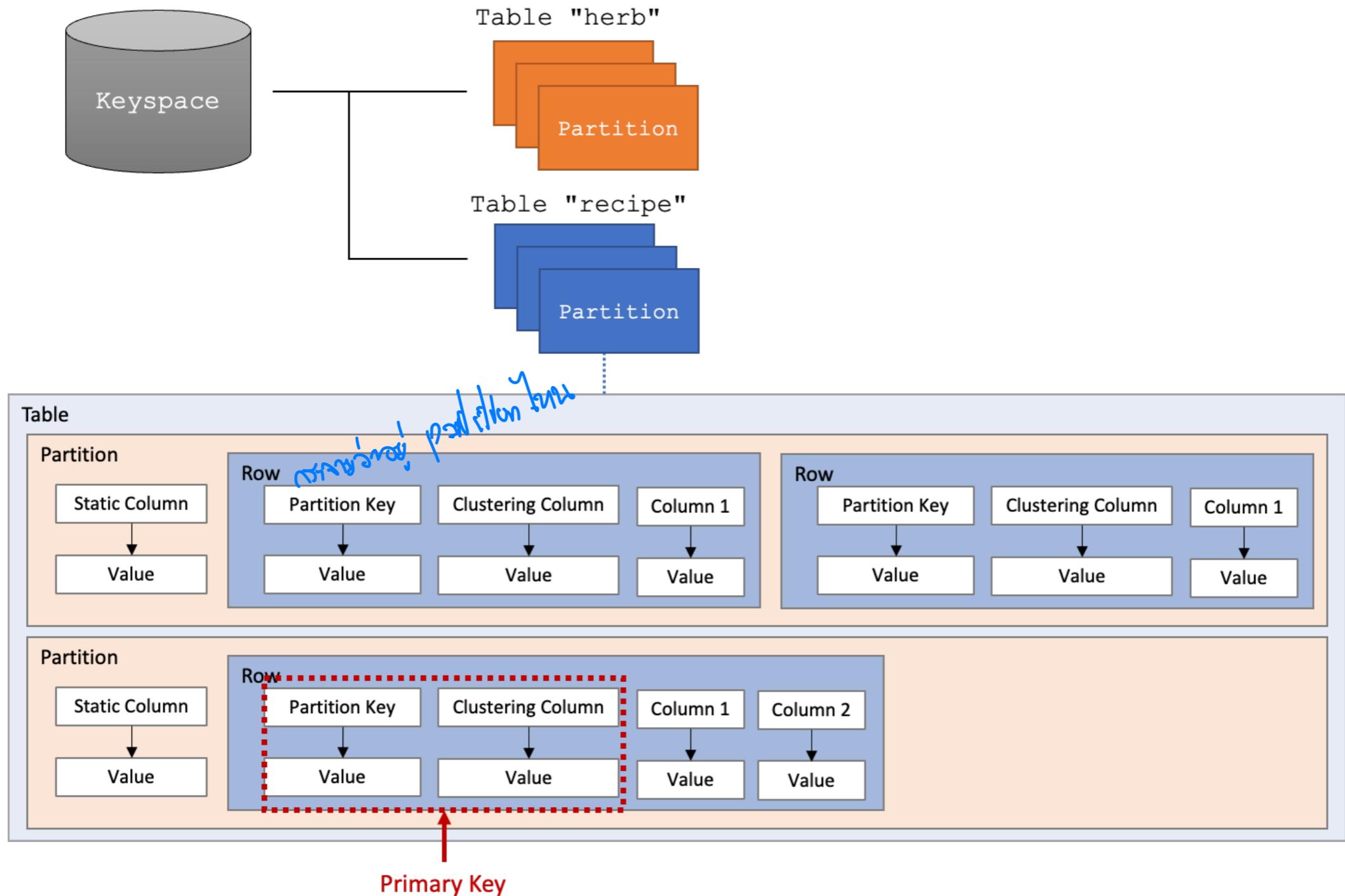
# Cassandra Data Model

ສູນໄຍ້ column ໃຊ້ 1 row ໂດຍໆໄວ້ row ອັນດຸ

column ຢາຍເກີງ

- Data model is query-driven, access patterns and application queries determine the structure of data
- Data in a table is divided into partitions using a column called ‘partition key’, which can be a composite key (multiple columns)
- Data in the same partition are sorted ascending or descending order based on the clustering columns

# Cassandra Data Model



# Cassandra Data Types

---

- Cassandra supports several data types including
  - Native types (see a table of examples in the next page)
  - Collection types (lists, sets, maps)
  - User defined types (a set of fields)
  - Tuple types (a list of anonymous fields of different types)
  - Custom types (Java class loaded by Cassandra)

Type	Constants supported	Description
bigint	integer	64-bit signed long
blob	blob	Arbitrary bytes (no validation)
boolean	boolean	Either <b>true</b> or <b>false</b>
counter	integer	Counter column (64-bit signed value)
date	integer, string	A date (with no corresponding time value)
decimal	integer, float	Variable-precision decimal
double	integer float	64-bit IEEE-754 floating point
duration	duration,	A duration with nanosecond precision
inet	string	An IP address: IPv4 (4 bytes long) or IPv6 (16 bytes long)
int	integer	32-bit signed int
text	string	UTF8 encoded string
timestamp	integer, string	A timestamp (date and time) with millisecond precision
uuid	uuid	A <b>UUID</b> (of any version)
varchar	string	UTF8 encoded string
varint	integer	Arbitrary-precision integer

There are other data types, refer to: <https://cassandra.apache.org/doc/latest/cassandra/cql/types.html>

# Playing with Cassandra

---

- Server-side: use docker compose file from cassandra folder in datasci\_architecture repo (this will use official cassandra image from docker hub)

```
cd cassandra
```

```
sudo docker compose up
```

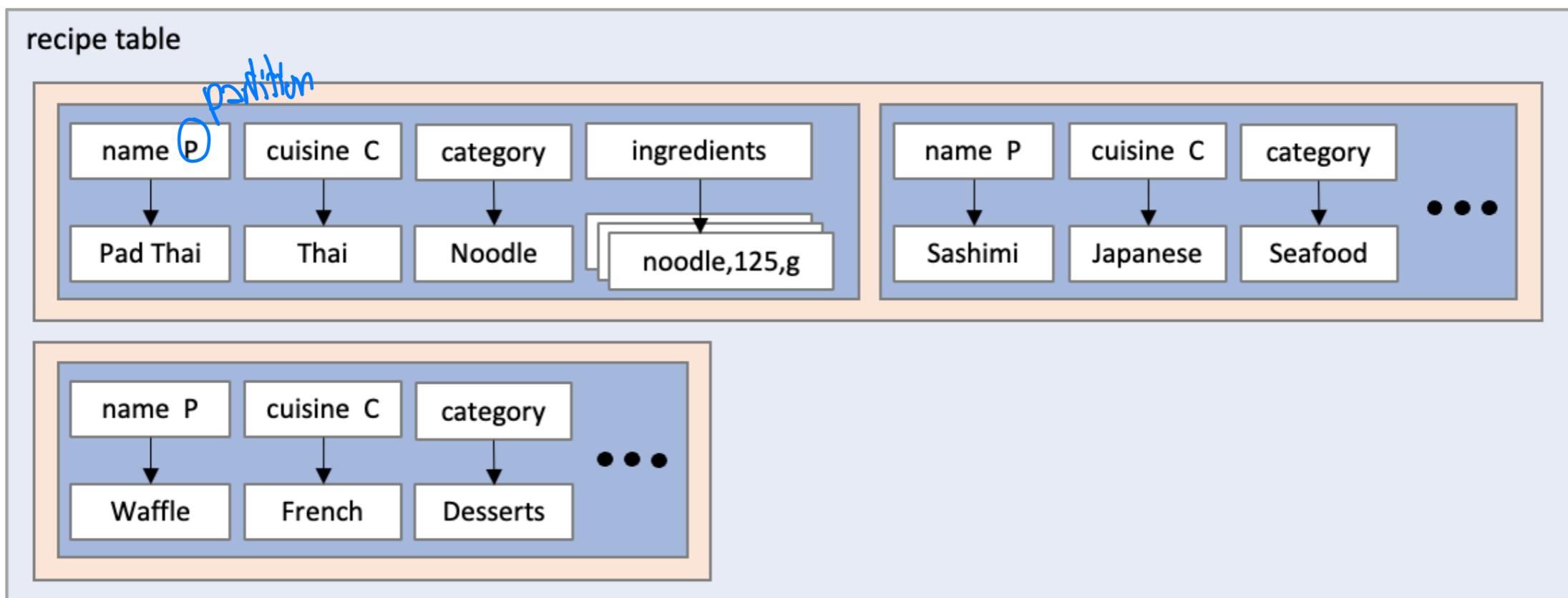
- From client-side, you use cqlsh in one of the cassandra nodes

```
sudo docker exec -it cs_0 cqlsh
```

- Note

- You will need to set your docker desktop to have at least 4GBs
- You can use nodetool to examine the cluster, see nodetool.sh for more details

# Example



# Example CQL – See init.cql

---

```
CREATE KEYSPACE restaurant WITH replication = {'class': 'SimpleStrategy',
'replication_factor' : 2};

USE restaurant;

CREATE TYPE ingredient (item text, quantity int, unit text);

CREATE TABLE recipe (name text, cuisine text, category text, ingredients
list<frozen<ingredient>>, primary key ((name), cuisine));

CREATE INDEX ON recipe (category);

INSERT INTO recipe (name, cuisine, category, ingredients) VALUES ('Pad_Thai',
'Thai', 'Noodle', [{item: 'noodle', quantity: 125, unit: 'g'}, {item: 'shrimp',
quantity: 150, unit: 'g'}, {item: 'brown sugar', quantity: 3, unit: 'tbsp'}]);

INSERT INTO recipe (name, cuisine, category, ingredients) VALUES ('Sashimi',
'Japanese', 'Seafood', [{item: 'rice', quantity: 100, unit: 'g'}, {item:
'salmon', quantity: 50, unit: 'g'}]);

INSERT INTO recipe (name, cuisine, category, ingredients) VALUES ('Chicken_Rice',
'Thai', 'Rice', [{item: 'rice', quantity: 300, unit: 'g'}, {item: 'chicken',
quantity: 100, unit: 'g'}]);

INSERT INTO recipe (name, cuisine, category, ingredients) VALUES ('Tom_Yum',
'Thai', 'Soup', [{item: 'shrimp', quantity: 200, unit: 'g'}, {item: 'water',
quantity: 500, unit: 'g'}, {item: 'lemon grass', quantity: 10, unit: 'leaves'}]);
```

# Example CQL – See init.cql

---

```
SELECT * FROM recipe;
```

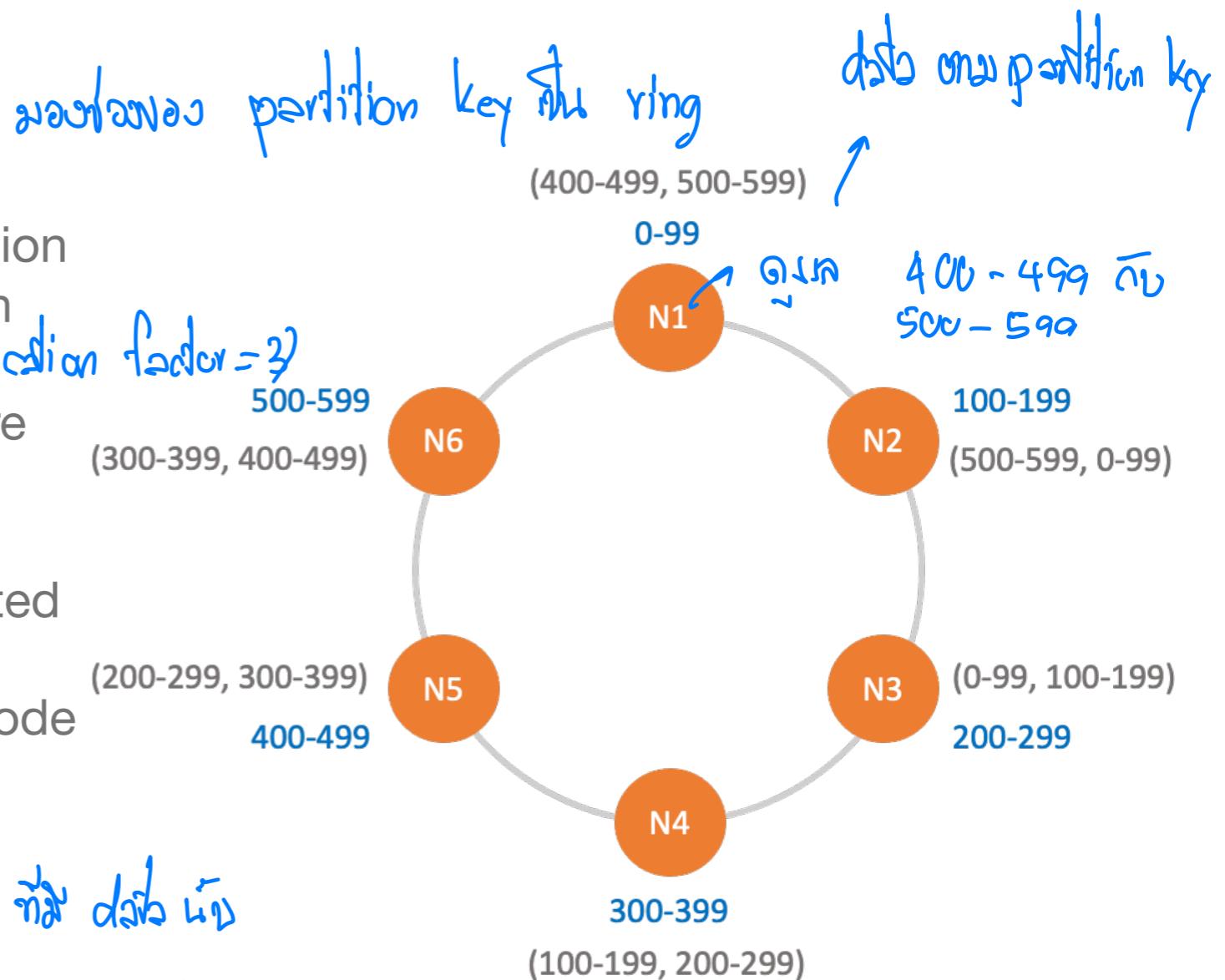
```
SELECT * FROM recipe WHERE name='Pad_Thai';
```

```
SELECT * FROM recipe WHERE category='Noodle';
```

```
SELECT token(name), name, cuisine FROM recipe;
```

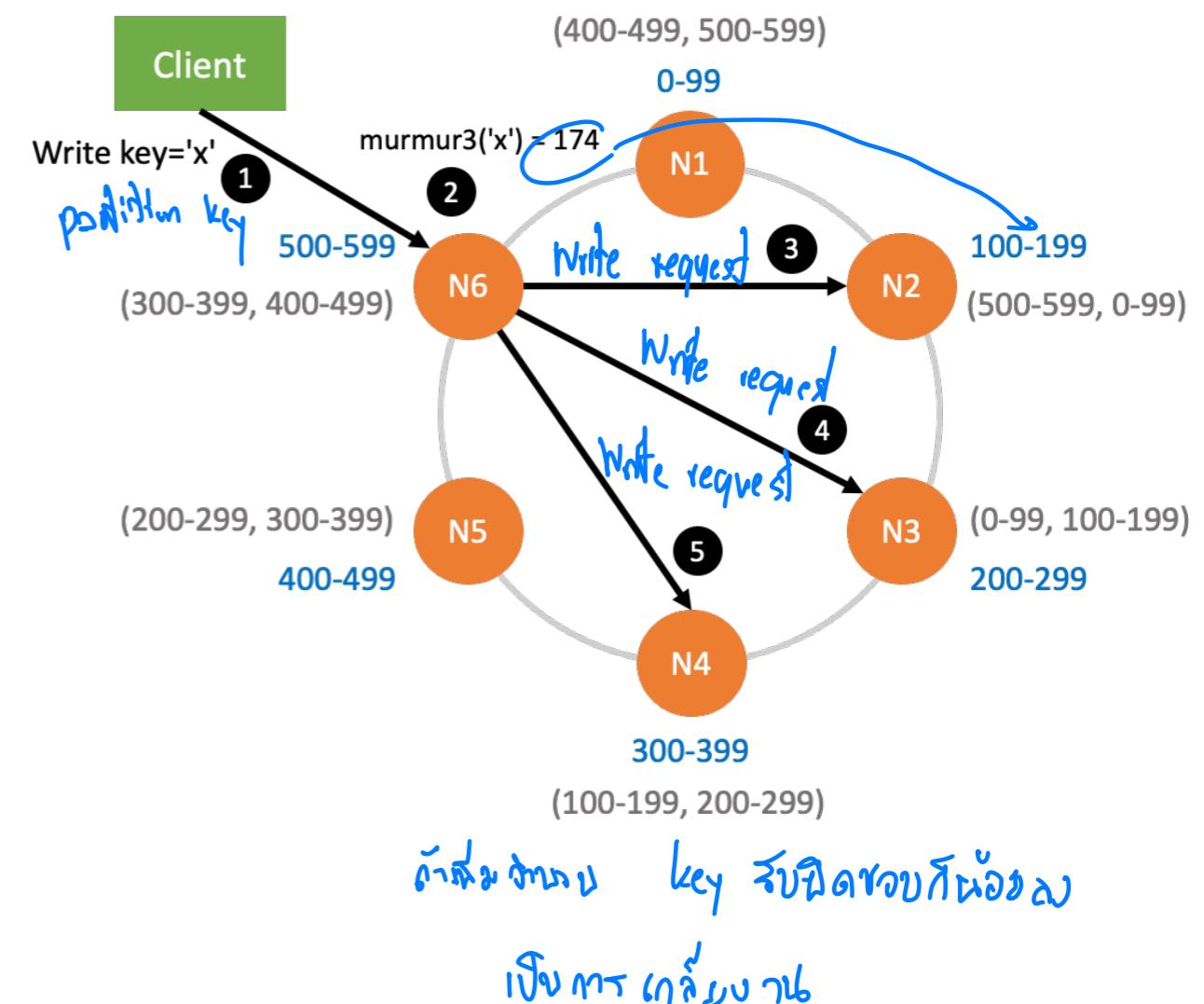
# Cassandra Cluster

- Nodes in cluster form a ring topology
- A token is generated by hashing partition key with Murmur3 partitioning function  
map 400 500 N5 N6 N1 (replication factor=3)
- Token value indicates the node to store the data
- Key space may not be evenly distributed
- Client can request to any node (that node becomes an “coordinator”) and the request will be routed to the node responsible for storing data
- Ring can be rebalanced
  - expand to accommodate more nodes
  - collapse to remove broken nodes



# Replication

- Cassandra replicates data to provide availability
- Client can write to any node in the cluster (that node becomes a “coordinator”)
- The coordinator node creates a token by hashing a partition key
- Using the token, the coordinator can find all replicas



# Read Operation

ជាឯកសារ និង Consistency

(ខាងក្រោម)

- In Read Operation there are three types of read requests that a coordinator can send to a replica. The node that accepts the write requests called coordinator for that particular operation

គរបៀប

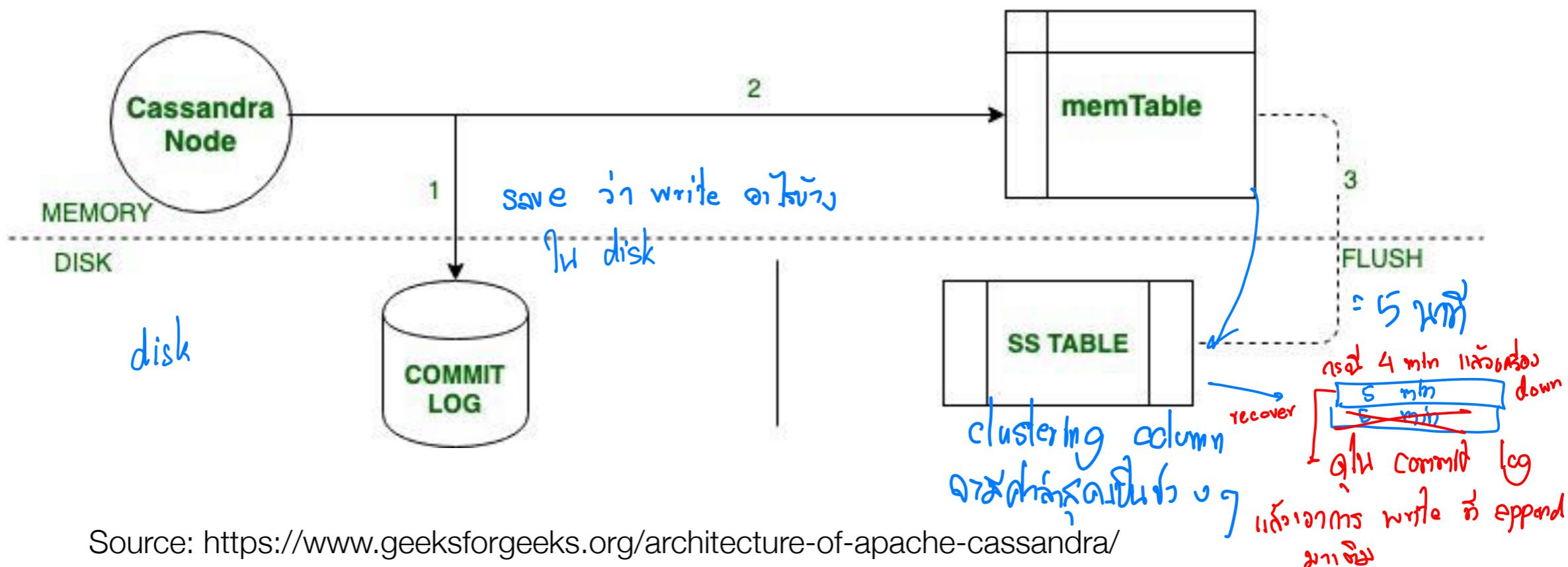
- Direct Request: In this operation coordinator node sends the read request to one of the replicas
- Digest Request: In this operation coordinator will contact to replicas specified by the consistency level. For Example: CONSISTENCY TWO; It simply means that Any two nodes in data center will acknowledge  $R+W > N$
- Read Repair Request: If there is any case in which data is not consistent across the node then background Read Repair Request initiated that makes sure that the most recent data is available across the nodes

2 node ឬចំនួន ពីរចាប់ពីមុន

# Write Operation

ກົດຍໍາ memory ອອນນາກ ຕ້າ crash data ອຸທະນະເວັບໄສ 3 copy

- Key concept: Data is kept in memory and lazily written to the disk Cassandra ອິນເນັດ data persistent 2 step ອິນເປັນ
- Step-1: In Write Operation as soon as we receives request then it is first dumped into commit log to make sure that data is saved ແຈ້ງໃຫຍ່ write request ອີ່ save request ໃນ commit log ສໍາເລັດ
- Step-2: Insertion of data into table that is also written in MemTable that holds the data till it's get full update table ໃນ memory
- Step-3: If MemTable reaches its threshold then data is flushed to SS Table



Cassandra & data center-aware

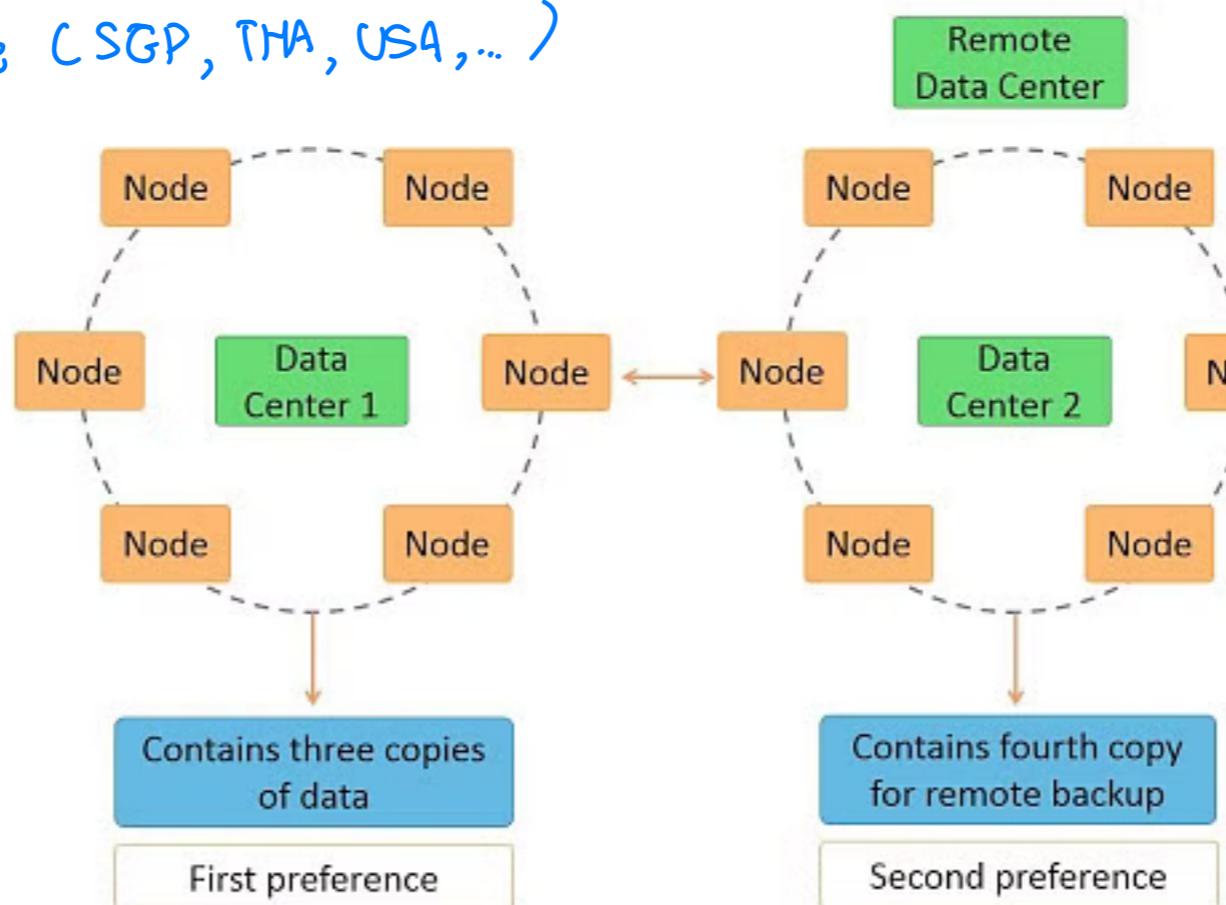
## Additional Features

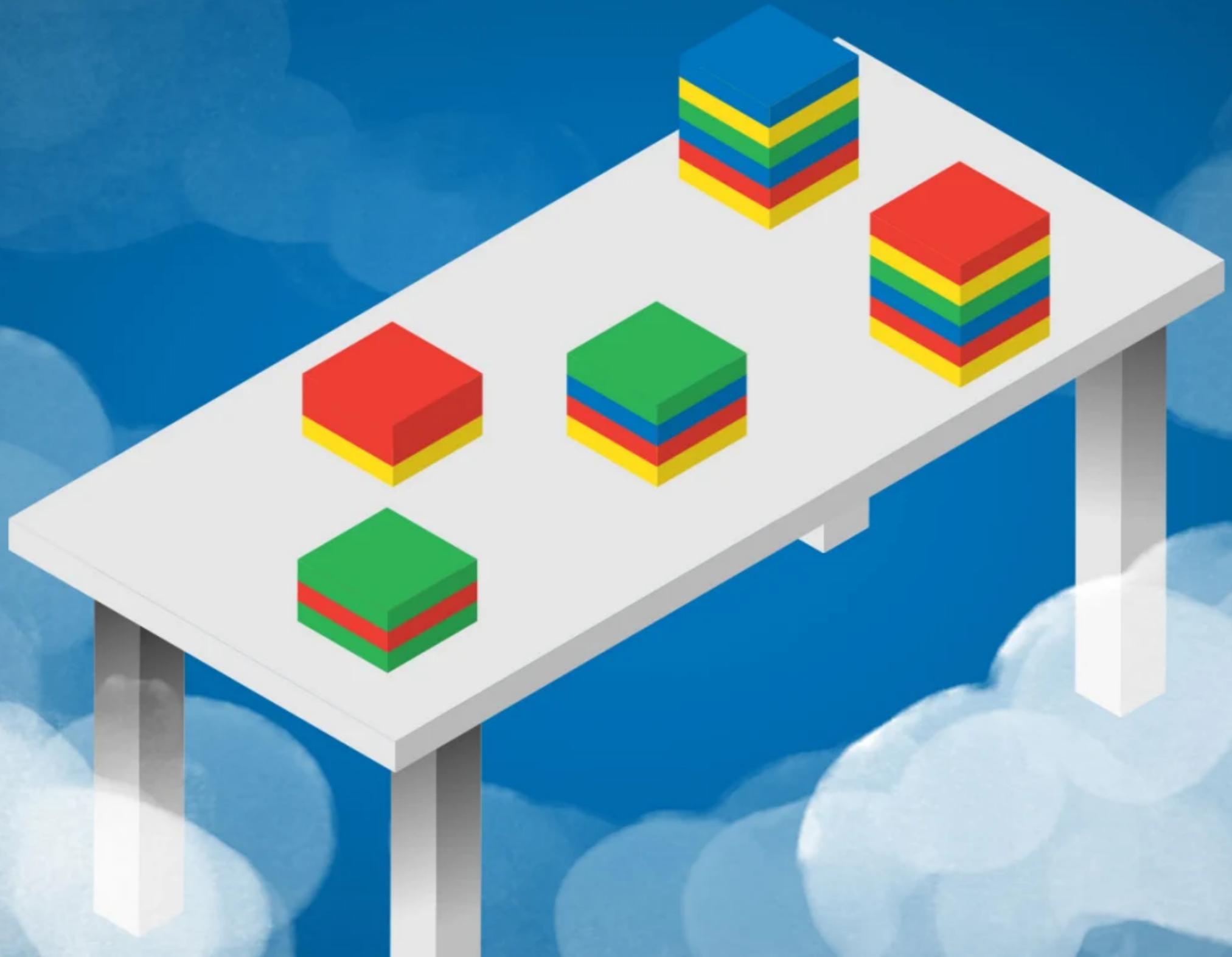
& 1 copy within rack

- Cassandra is rack-aware for availability (rack failure) and read performance (locality)
- Cassandra supports multiple data centers which data can be replicated across data centers
- See <https://www.simplilearn.com/tutorials/big-data-tutorial/cassandra-architecture>

Within copy within zone (SGP, THA, USA, ...)

Within update node





Google Bigtable

# Google Bigtable

---

- A fully managed wide-column and key-value NoSQL database service by Google since 2015  
*សំវិជ្ជកម្មនៃការបង្កើតគេហទ័រនៃក្រុមហ៊ុនកូរែង*
- Main focuses
  - **Scalability:** stores petabytes of data on thousands of nodes (machines)
  - Wide applicability: serves throughput oriented batch processing and latency sensitive services
  - High performance and high availability

# Bigtable Data Model

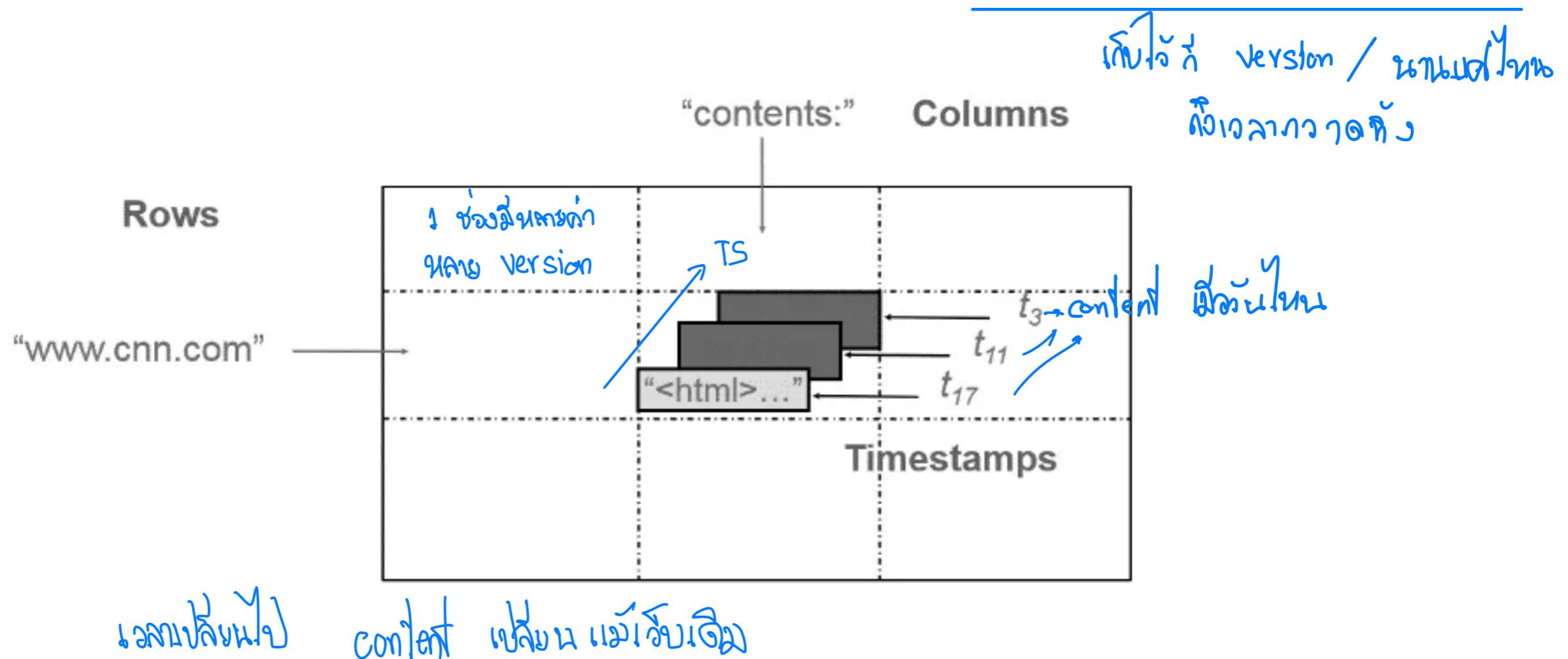
เก็บ content website

เก็บ 2 มิติ

- Data is an uninterpreted string stored in a cell indexed using row, column, and timestamp

$(row : string, column : string, \boxed{timestamp : int64}) \rightarrow string$

- Thus, a single cell can store multiple values (limited by Garbage Collection (GC) policy)



# Bigtable Rows and Tablets

tablet គឺជាមុនក្រាស handle data នៃ node  
ទីផ្សារ Bigtable (រាយការណ៍ row key ដូចការណ៍ set មួយ)

- Each table has only one index, the row key; each row key must be unique
- Rows are sorted lexicographically by row keys and rows with consecutive keys are grouped into a tablet for better system performance
- A tablet can hold up to 100-200 MB of data and can be handled by one node at any time
- If a tablet becomes too big or has significant more loads than other tablets, it can be further partitioned into multiple tablets

សមារណការណ៍ tablet អំពី

Group នៃ row = tablet

# Bigtable Columns and Column Families

---

- A column key is a two level naming structure, consisted of column family and optional qualifiers

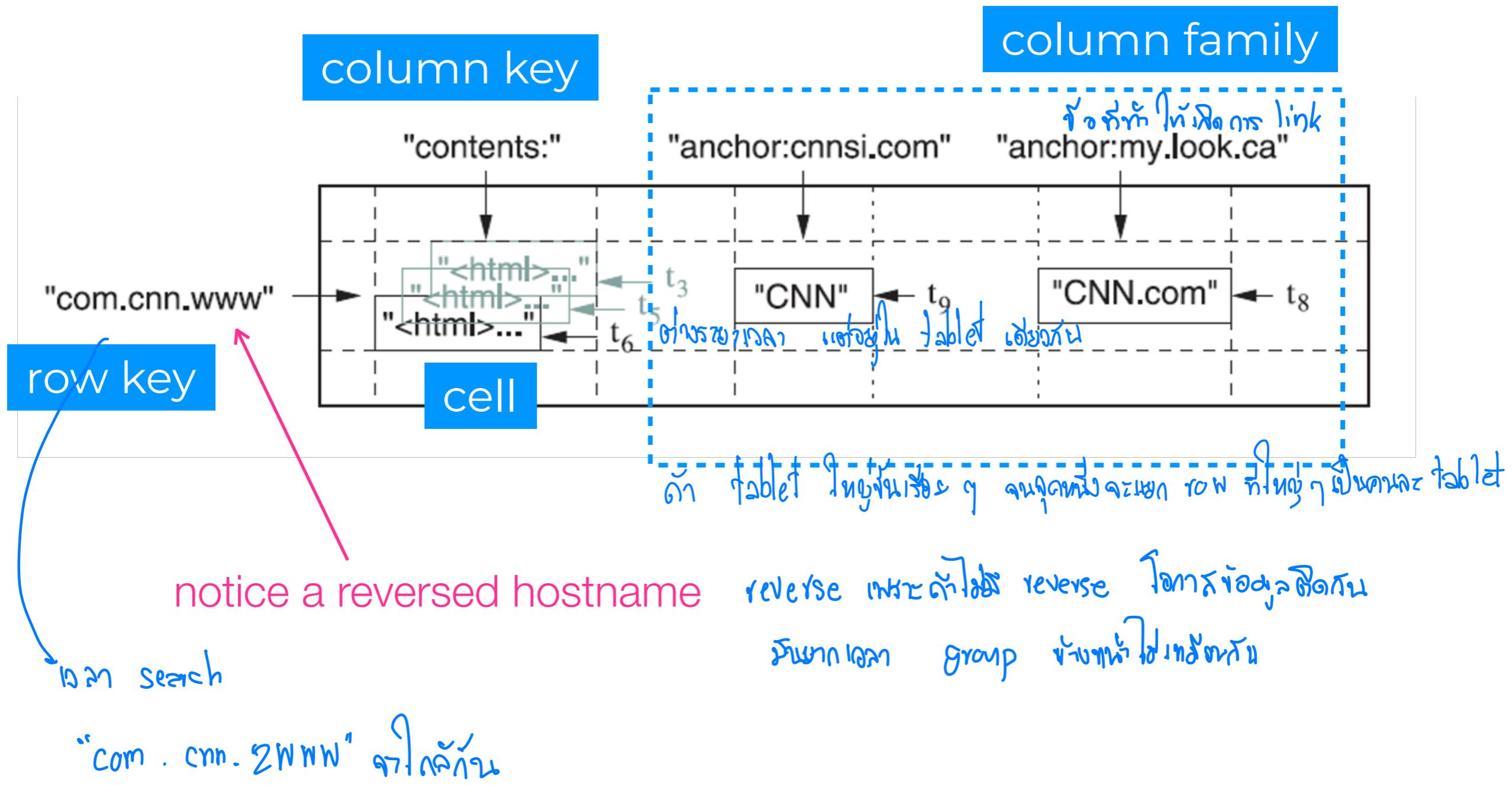
*Column name → Column family : optional qualifier*

Family នៃបានក្នុង access control នៅលើវា មានការបញ្ជាក់ស្ថិតិយវិធី

- Column family is the basic unit of access control and must be created before storing data under column key in that family
- Column families are not sorted, but columns are sorted in lexicographic order within the family
- Bigtable tables are sparse; a column does not take up space in a row that does not use the column

# Example: Web Table

Bigtable នឹងរាយកើត គិតថវិកសំខាន់ដ៏ខ្ពស់



# Important Notes

---

- All operations are atomic at the row level; an operation affects either an entire row or none of the row
  - Design row key based on the queries you will use
    - Row key can consist of multiple delimited values e.g. thai#pad\_thai or seafood:sashimi or thai/soup/tom\_yum
    - Well-designed row key prefix can improve query performance as they can be used to retrieve data without a full table scan e.g. retrieve all Thai dishes in the table
    - Put related columns with similar data retention needs in the same column family
- ເສົ່ານີ້ກັບທີ່ໄວ້ ໂພງເຄົ້າຫຼັກຂຶ້ນກັນ ລາຍງ່າ tablet ເຕືອນ ວຳກັນ  
ວ່ານກັບລົງຈູດ  
dat = ສະບຸ

# Playing with Bigtable Emulator

---

- Server-side: use docker compose file from bigtable folder in datasci\_architecture repo (this will use bigtable emulator image from docker hub)

```
cd bigtable
```

```
sudo docker compose up
```

- From client-side, install google bigtable python package
- Run a python example script, simple\_client.py to create an example recipe table and demonstrate how to retrieve data from bigtable emulator

```
pip install google-cloud-bigtable
```

# Google Bigtable Architecture

Cassandra មាន service ណាយគេ បានកំពុងរួច

API forward request ទៅលើ big table node / tablet server

ត្រូវកំណត់លើ row key ទៅ consult សែនី chubby service នៃ row key ដើម្បី tablet server ដោយបាន

- Front-end servers (API)

- Tablet servers (Bigtable nodes)

- Manage a set of tablets
- សំណង R/w tablet  
ទៅ Bigtable node នៅអីណា !!

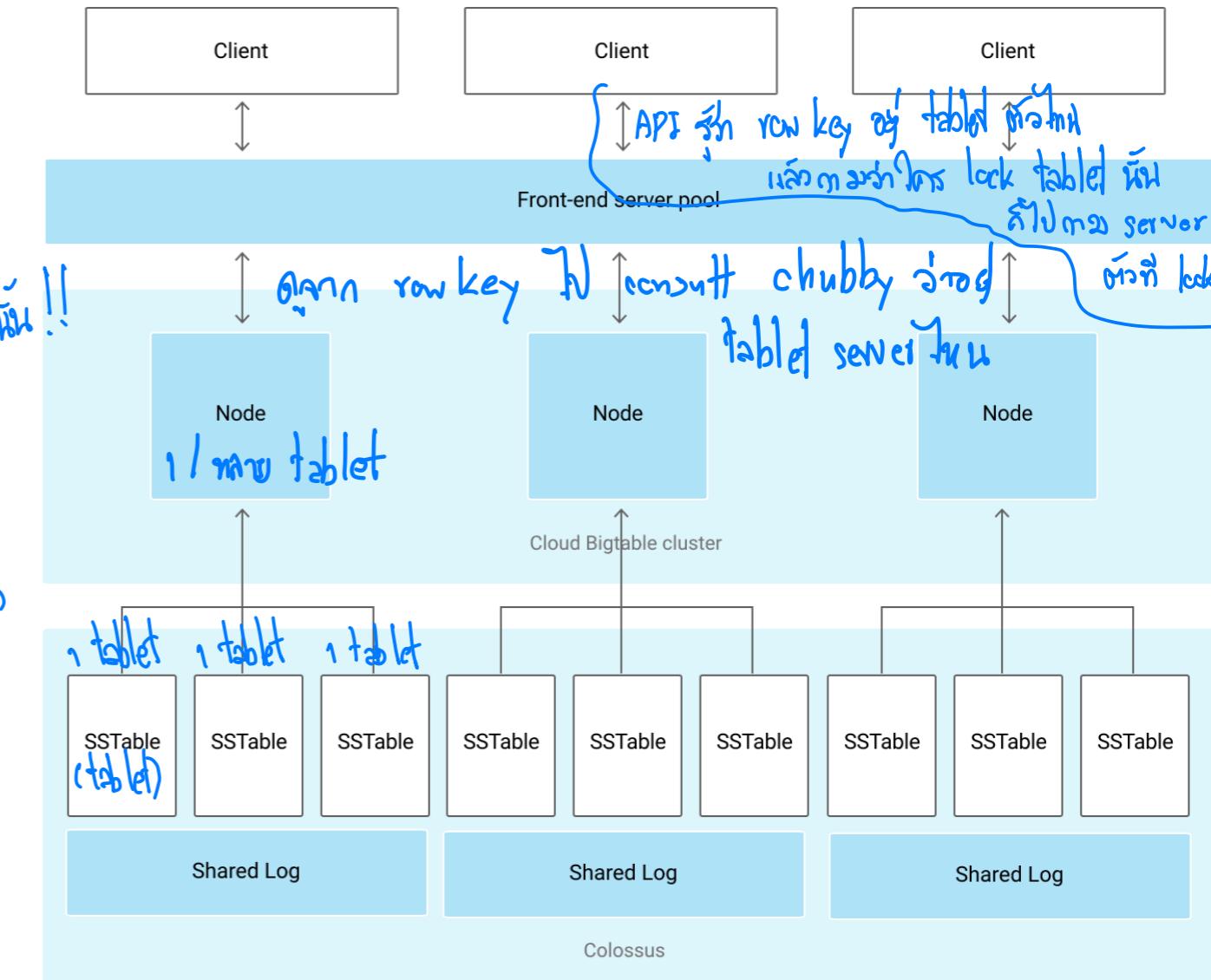
- Handle read and write requests to the tablets
- Split tablets that have grown too large

- File Systems (Colossus) SS Table

- Store logs and data files
- Tablet is stored in SSTable (Sorted String Table)
- Immutable map from keys to values

- Chubby

data เហេងន៍ Colossus



Colossus ក្រោចប្រើប្រាស់

1 tablet ត្រូវបានខូចខាងក្រោម tablet server ដោយការងារនៃក្រុងក្រោមរាយការណ៍

1 tablet = 1 SSTable

ឬ 1 tablet server មិនមែន  $> 1$  tablet នៅក្នុងក្រោមរាយការណ៍

tablet server

cloud

cloud

Chubby = Coordinator Service

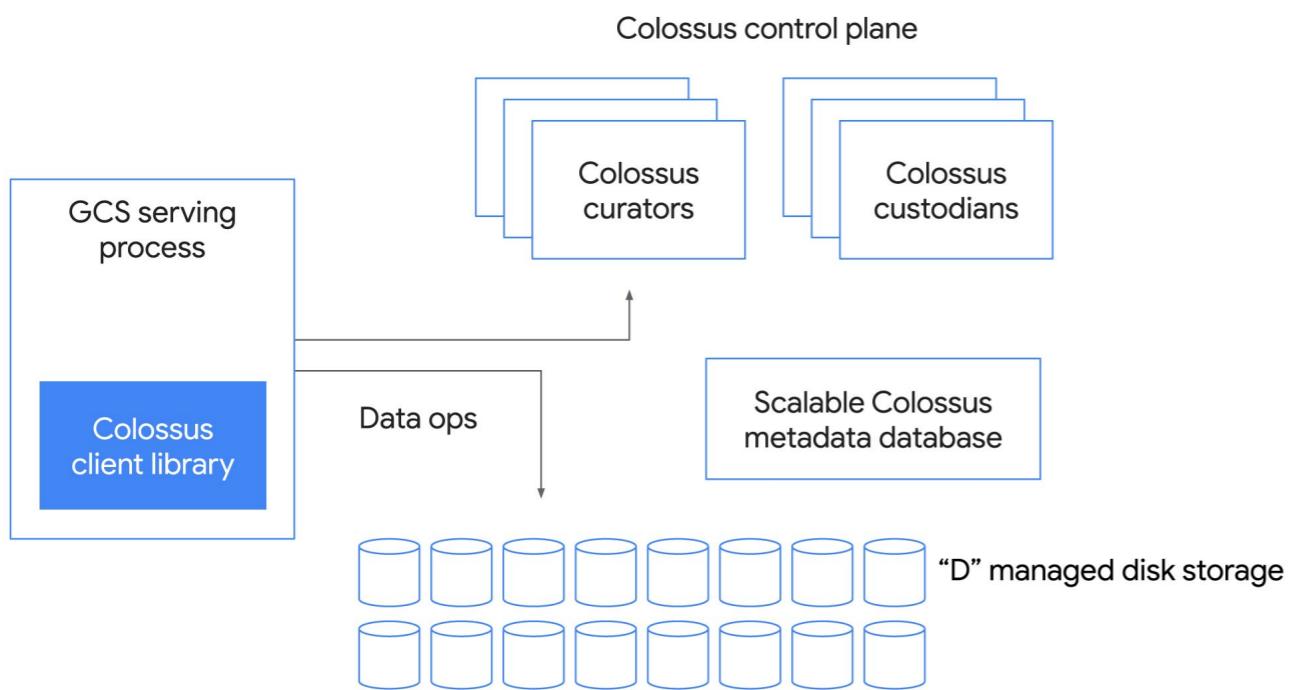
Chubby<sup>②</sup> maintain mapping between tablet server and tablet

tablet នូវចាន់មិន register រវ Chubby ឬ heartbeat រវ Chubby នៅលាក់

- Highly-available and persistent distributed lock service with 5 replicas and one is elected as a master
  - Provide a namespace of nodes, similar to Linux file systems (directories and small files)
  - Node can be permanent or ephemeral; Ephemeral nodes go away as no client using the node go away
  - See <https://medium.com/coinmonks/chubby-a-centralized-lock-service-for-distributed-applications-390571273052> for more information
  - Bigtable uses Chubby for checking if tablet servers are active, ensuring there is only one tablet server master, keeping schema and ACL, etc.

# Colossus

---



- Developed by Google to replace Google File System (GFS) as its cluster-level file system
- Clients perform software RAID by directly reading and writing data from/to D file server to minimize latency
- Data typically written using Reed-Solomon (Error Correction Codes)
- Curators handle metadata layer similar to HDFS NameNode (but distributed)
- Custodians perform background tasks to maintain durability and availability of data as well as overall efficiency

Big table scale ជំនាញ ① ការអភិវឌ្ឍន៍នូវការសម្រេចក្នុង tablet ឬឯក-ឯក នៃទីតាំង ដើម្បី

② เลือกจักษุน์ tablet server ที่สามารถดู tablet ได้ (1-1, 1-n)

③ กรณี chubby และ master server ดูแลการ mapping ทำให้การจัดการมัน simple

# Tablet Servers

Tablet alarm lock tablet

- One master server - one server is elected to be a master
    - Assigning tablets to tablet servers ត្រូវការសេវាឌាមួយទៀតទៅ tablet
    - Detecting addition and expiration of tablet servers
    - Balancing tablet-server load
    - Garbage collection
  - Many tablet servers
    - Tablet servers handle read and write requests to its table
    - Splits tablets that have grown too large

Tablet servers ត្រូវបានសេវាដំឡើងទៅលើ tablet  
ឱ្យការងារនៃ tablet ត្រូវបានគ្រប់គ្រងដោយ tablet server  
មិនអាចបញ្ចប់ទីតាំង tablet បានប៉ុណ្ណោះទៀតទៅទេ

# Tablet Servers

---

- Communication
  - Clients communicate directly with tablet servers for reads and writes (not master)
  - Most clients never communicate with the master server, leaving it lightly loaded in practice

# Tablets

---

- Large tables broken into tablets at row boundaries
  - Tablet holds contiguous range of rows: choosing good row keys to achieve locality
  - Aim for ~100MB to 200MB of data per tablet
- Each tablet server responsible for ~100 tablets
- Fast recovery: 100 servers each pick up 1 tablet for 1 failed server
- Fine-grained load balancing:
  - Migrate tablets away from overloaded machine
  - Master makes load-balancing decisions

1 server 90 tablets 1cc tablet

# Tablet Assignment

ရုံးချုပ် master

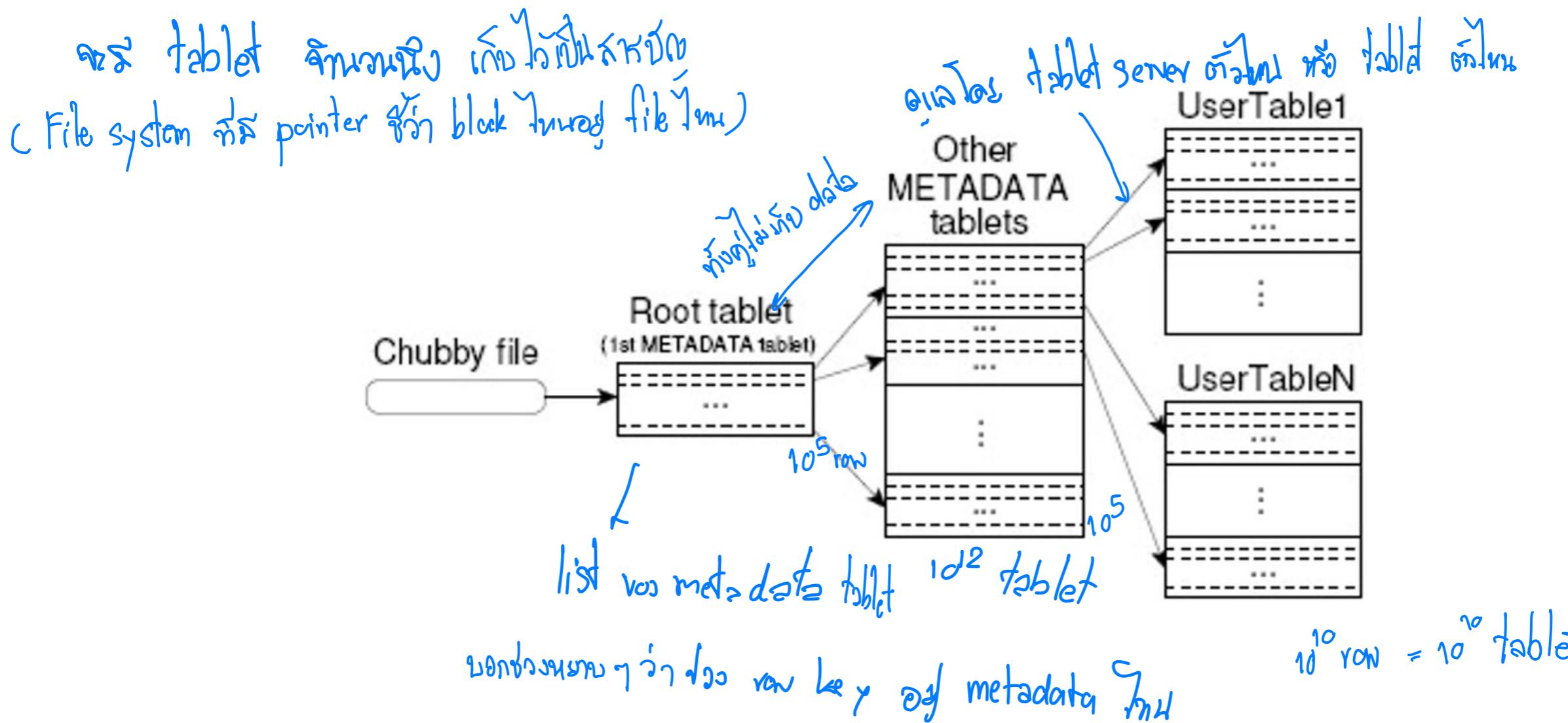
---

- Each tablet is assigned to one tablet server at a time
- Master server keeps track of the set of live tablet servers and current assignments of tablets to servers. Also keeps track of unassigned tablets
- When a tablet is unassigned, master assigns the tablet to an tablet server with sufficient room

API request ចាំងការ access row key ដែលក្នុងណា ?

## Tablet Location - Two-Level Index Scheme

- Since tablets move around from server to server, given a row, how do clients find the right machine?
- Need to find tablet whose row range covers the target row



# Performance Refinements

---

- Locality Groups - group multiple column families into a locality group for more efficient reads when accessing together
- Compression - reduce storage consumptions (10-to-1 ratio) with minimal speed losses
- Bloom Filter - reduce number of disk accesses by checking if an SSTable might contain data for a specified row/column pair
- Read assignment reading for more details

# Reading Assignment

---

- A. Lakshman, P. Malik, “Cassandra - A Decentralized Structured Storage System”, ACM SIGOPS operating systems review, 44.2 (2010), pp. 35-40.
- F. Chang et al, “Bigtable: A Distributed Storage System for Structured Data”, ACM Transactions on Computer Systems (TOCS), 26.2 (2008), pp.1-26.

# References

---

- Simplilearn, “Apache Cassandra Architecture From The Ground-Up”,  
<https://www.simplilearn.com/tutorials/big-data-tutorial/cassandra-architecture>
- J. Carpenter and E. Hewitt, Cassandra: The Definitive Guide: Distributed Data at Web Scale, 3rd Edition, 2022
- Google, “Bigtable overview”, <https://cloud.google.com/bigtable/docs/overview>
-