



---

DuyHai DOAN, Technical Advocate

# Agenda

## Architecture

- 
- 

## Data model

- ( ),
- ( , , , )
-

# Who Am I ?

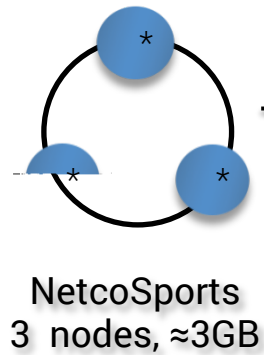
## Duy Hai DOAN

- , ,
- - (Achilles, )
- - 👉 [duy\\_hai.doan@datastax.com](mailto:duy_hai.doan@datastax.com)
  - 👉 [@doanduyhai](#)

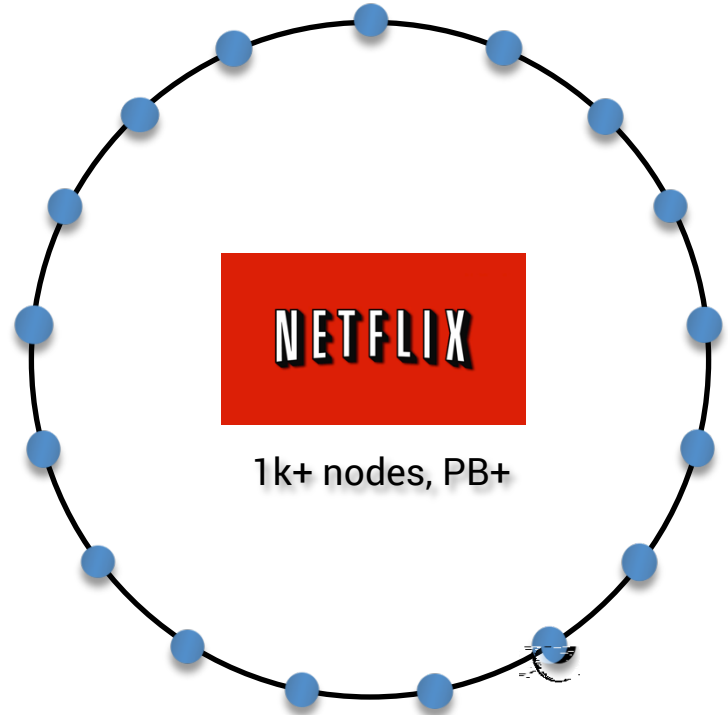
- April 2010
- a lot
- 400+ (25 100), 200+
- 
- London, France Germany
- Datastax Enterprise = + extra features

# Cassandra 5 key facts

1:



YOU



# Cassandra 5 key facts

- 2: ( ) (≈100% - )



# Cassandra 5 key facts

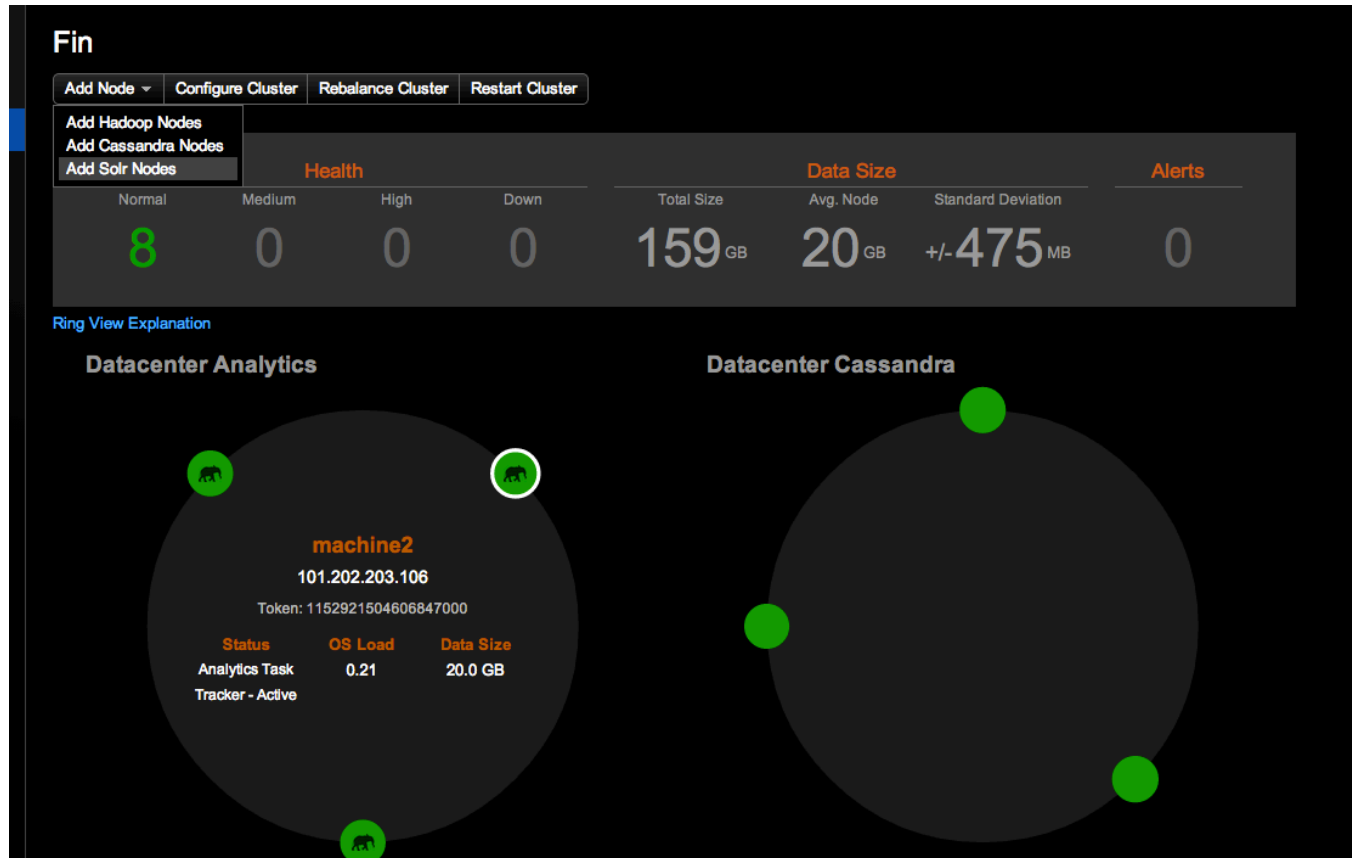
- 3: - - - ( )
- -
- /
-

# Cassandra 5 key facts

- $\frac{4}{1} = 1 + 2 + 1$  ( + )
- 
-



# Cassandra 5 key facts



# Cassandra 5 key facts

5:

- Cassandra + Spark = awesome !
- realtime streaming/

# Cassandra architecture

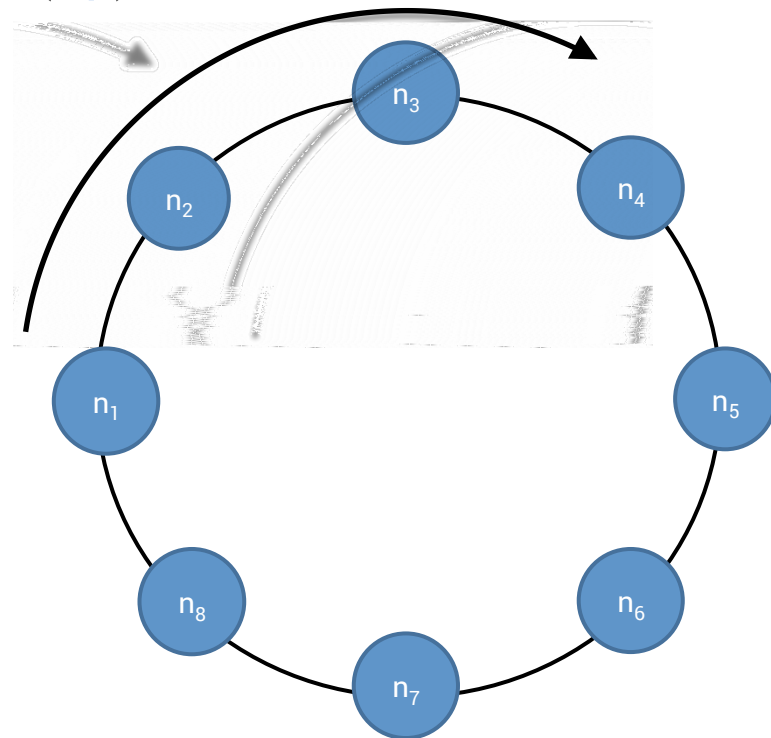
Cluster  
Replication

# Cassandra architecture

- DynamoDB
- masterless
- -
- Big Table
- /

# Data distribution

: #partition → token = (#p)  
:- ,  
=  $(2^{64}/2)$



# Token Ranges

A: 0, /8

B: /8, 2 /8

C: 2 /8, 3 /8

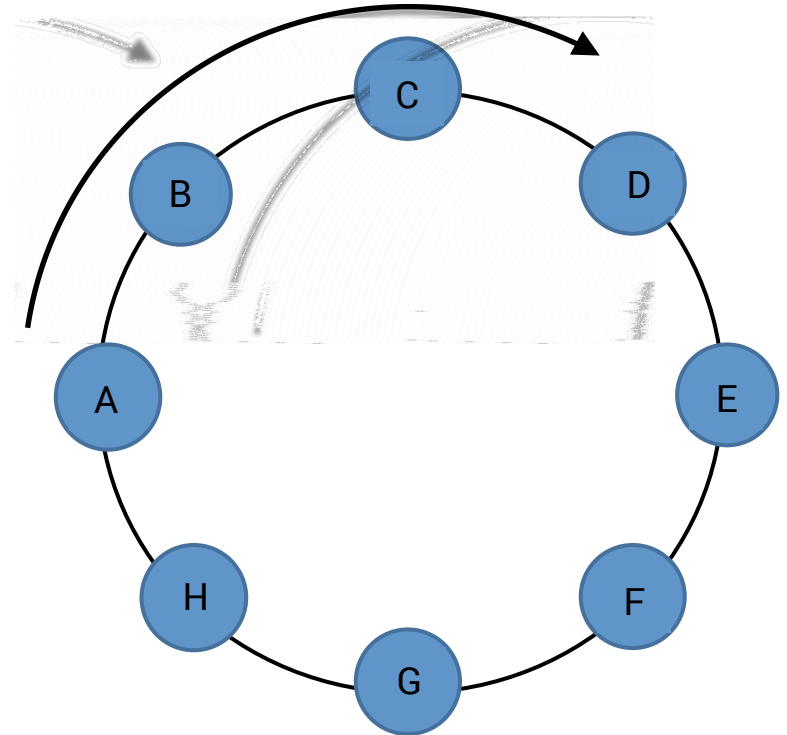
D: 3 /8, 4 /8

E: 4 /8, 5 /8

F: 5 /8, 6 /8

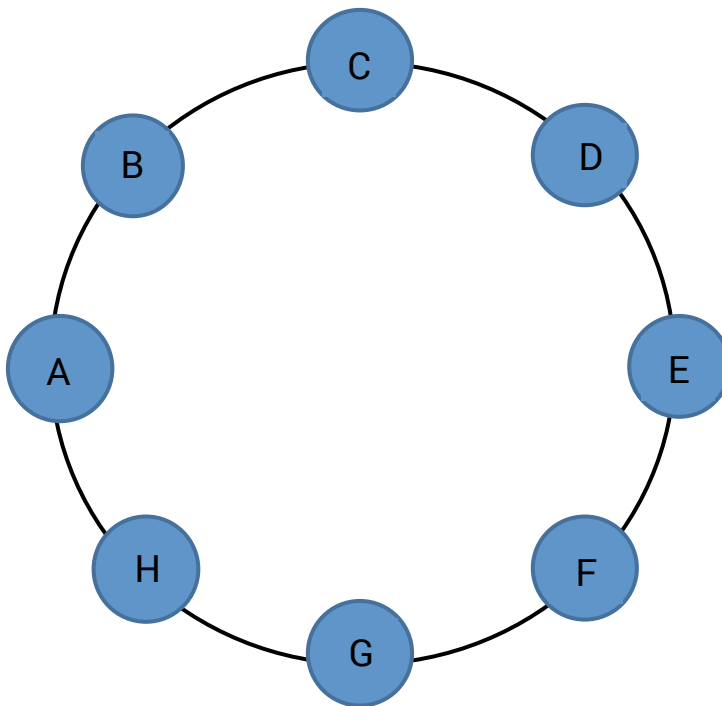
G: 6 /8, 7 /8

H: 7 /8,

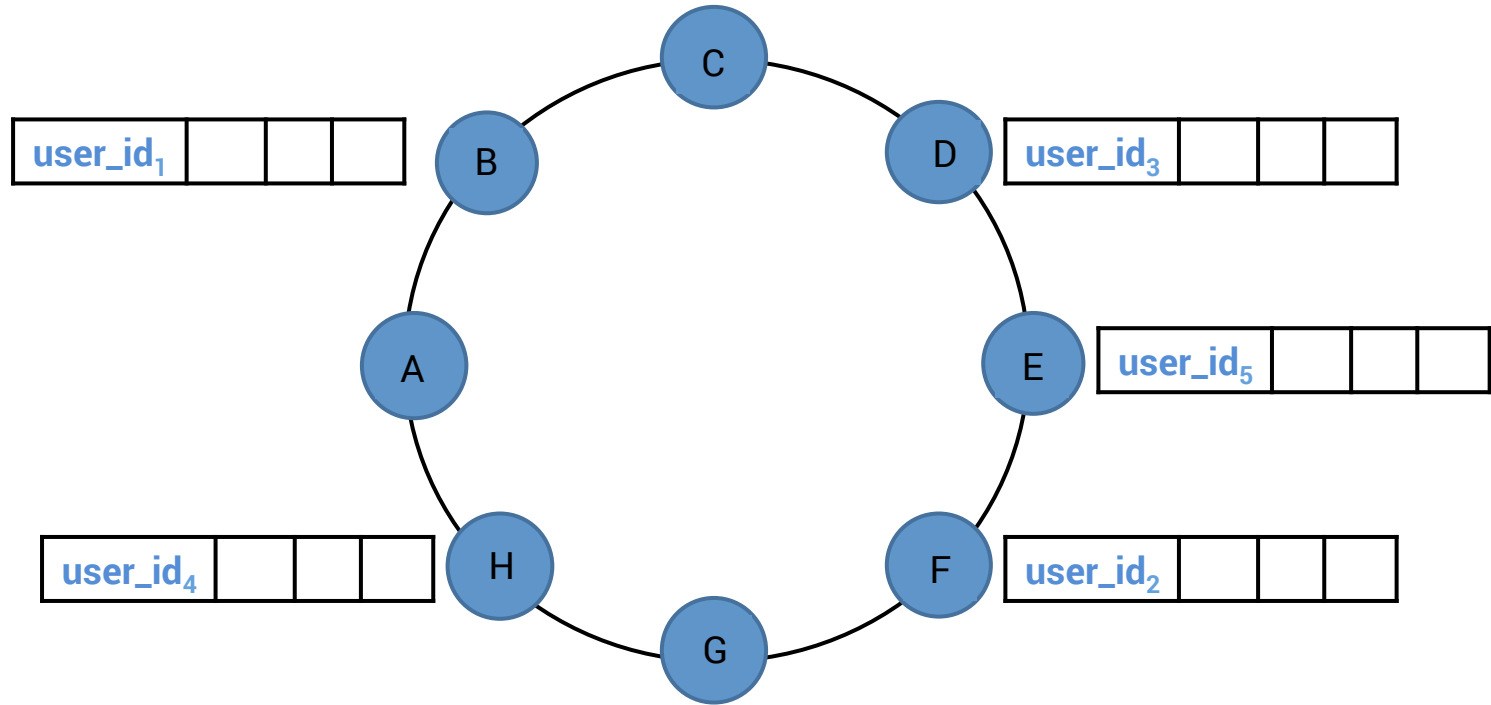


# Distributed Table

user_id <sub>1</sub>			
user_id <sub>2</sub>			
user_id <sub>3</sub>			
user_id <sub>4</sub>			
user_id <sub>5</sub>			

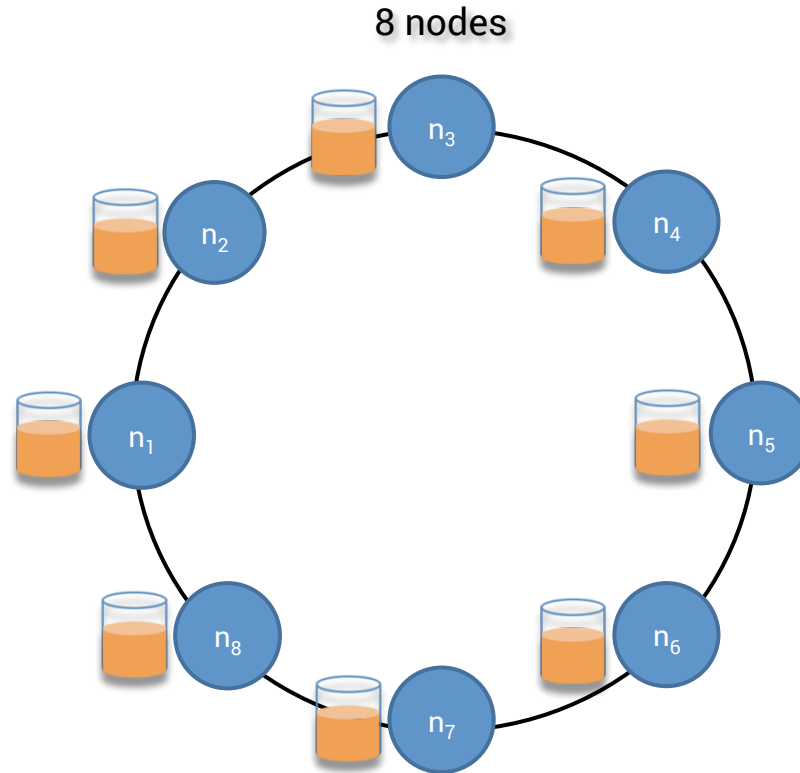


# Distributed Table



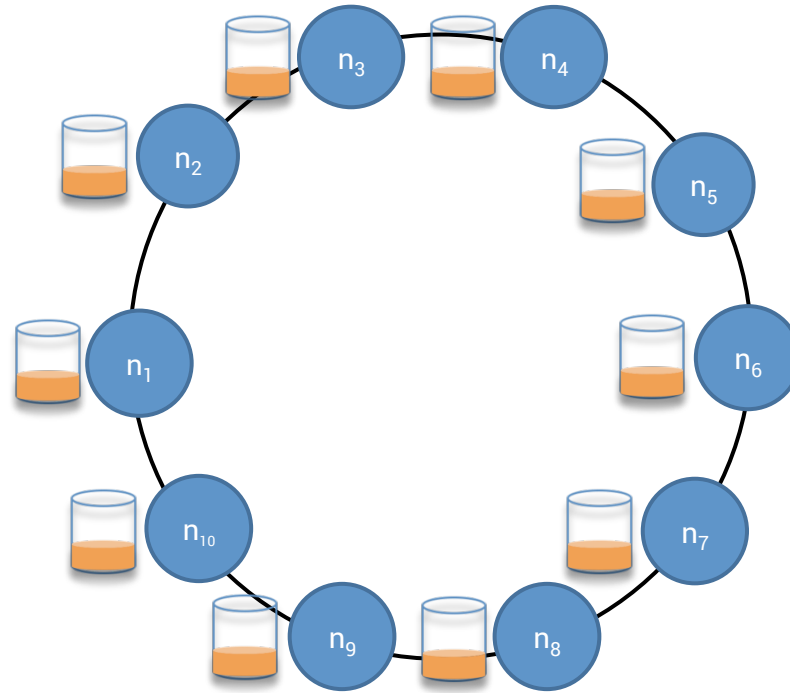


# Linear scalability



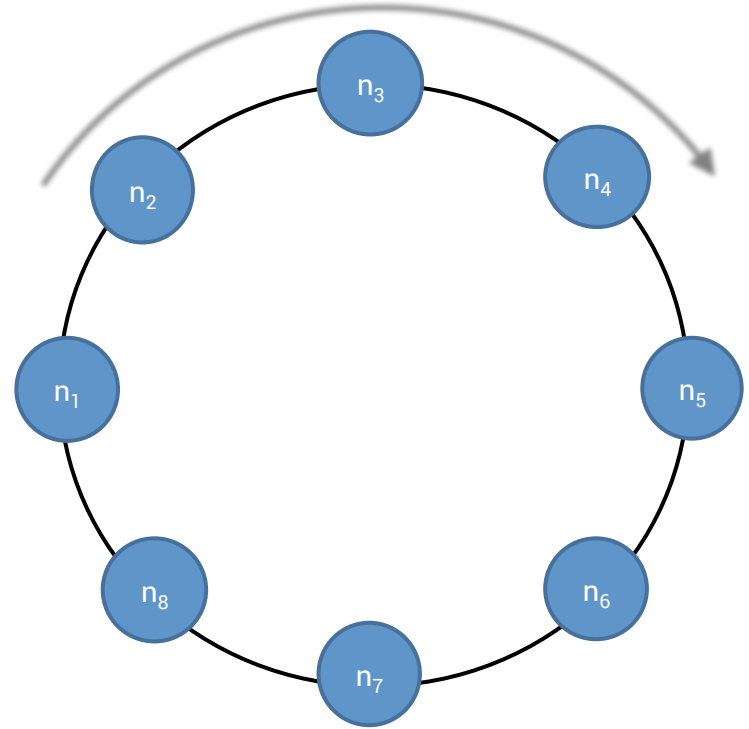
# Linear scalability

10 nodes



# Failure tolerance

$$(RF) = 3$$

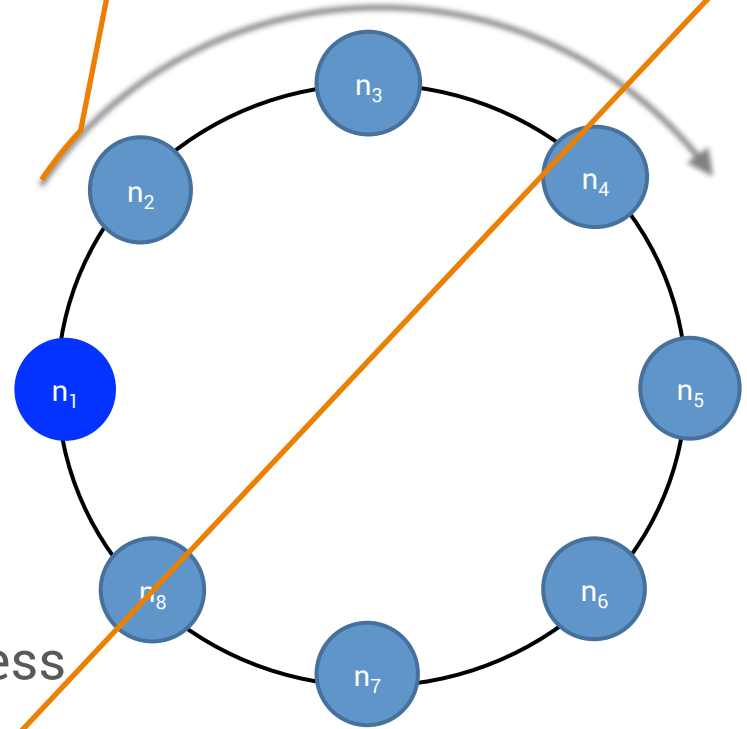


# Coordinator node

( / )

Coordinator

coordinator → masterless



# Consistency

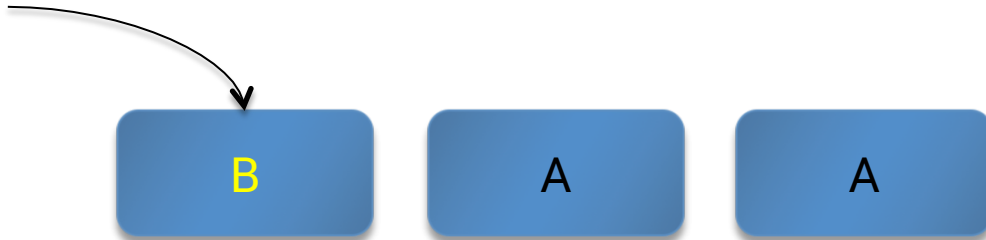
- ONE
- QUORUM (strict majority . . . RF)
- ALL

read & write

# Consistency in action

= 3, ONE, ONE

Write ONE: B



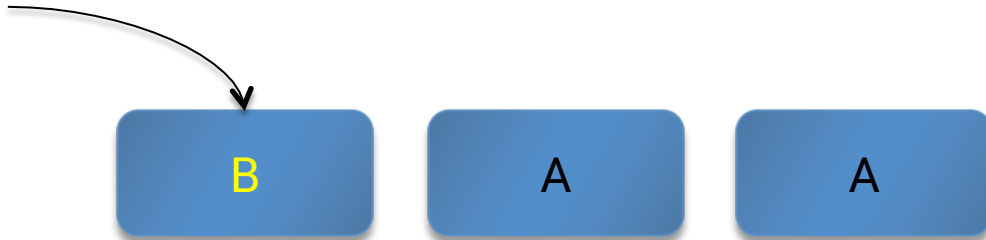
Read ONE: A



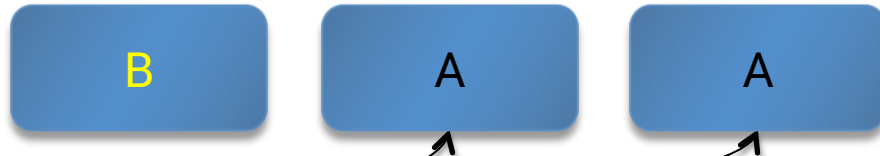
# Consistency in action

= 3, ONE, QUORUM

Write ONE: B



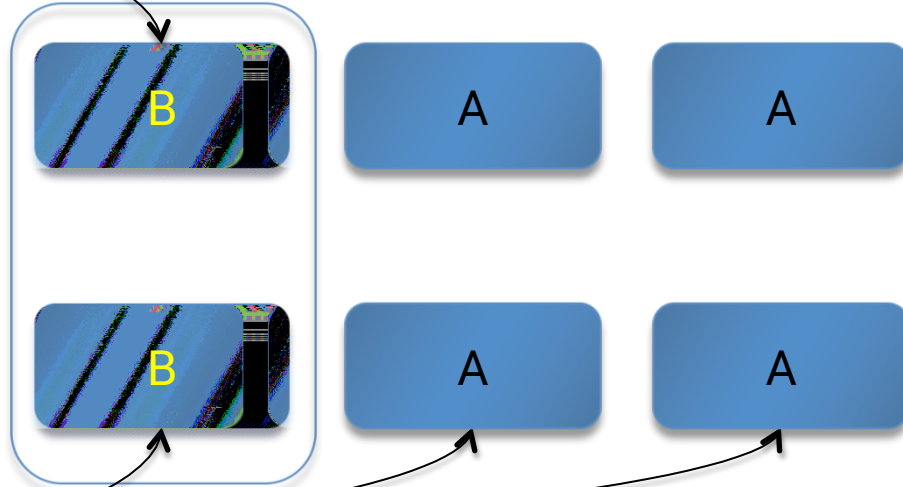
Read QUORUM: A



# Consistency in action

= 3, ONE, ALL

Write ONE: B



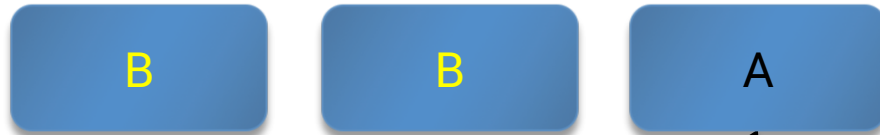
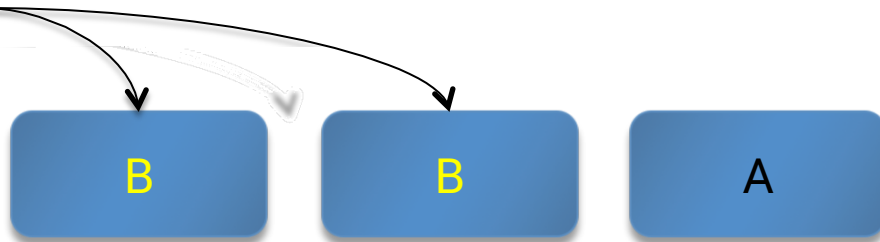
Read ALL: B



# Consistency in action

= 3, QUORUM, ONE

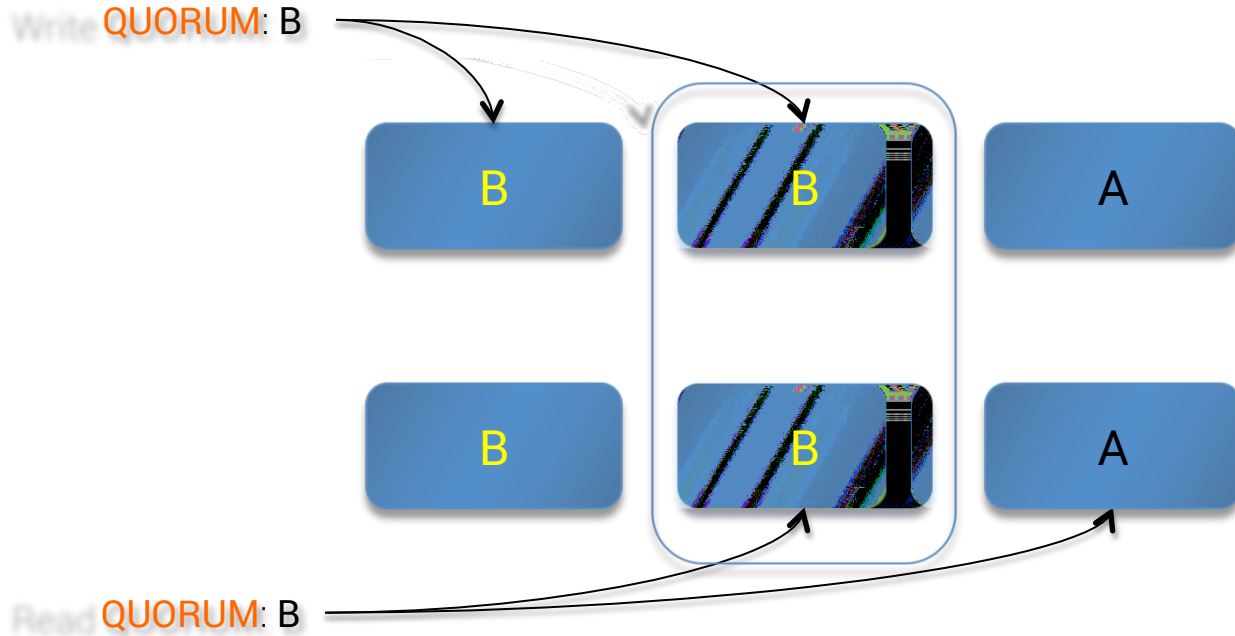
Write QUORUM: B



Read ONE: A

# Consistency in action

= 3, QUORUM, QUORUM



# Consistency trade-off

**Latency**

**Consistency**



# ONE

Fast, may not read latest written value

## QUORUM

Strict majority w.r.t. Replication Factor

Good balance

**ALL**

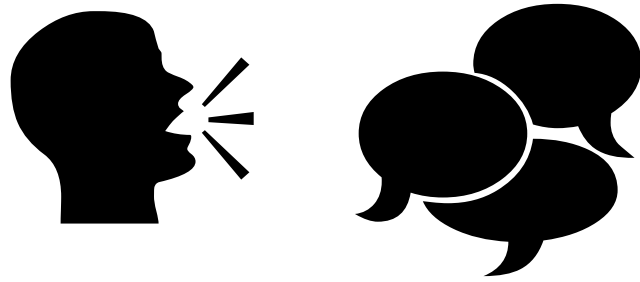
Paranoid

Slow, no high availability

# Consistency summary

**ONE**<sub>Read</sub> + **ONE**<sub>Write</sub>  
available / (N-1)

**QUORUM**<sub>Read</sub> + **QUORUM**<sub>Write</sub>  
available / 1+



Q & R



# Data model

Last Write Win

CQL basics

Clustered tables

Lightweight transactions

# Last Write Win (LWW)

```
(login, , ) (jdoe, John DOE, 33);
```

#partition

jdoe		
	33	John DOE

# Last Write Win (LWW)

(login, , ) (jdoe, +

jdoe	(t <sub>1</sub> )	(t <sub>1</sub> )
	33	John DOE

# Last Write Win (LWW)

= 34      login = jdoe;

SSTable<sub>1</sub>

jdoe	(t <sub>1</sub> )	(t <sub>1</sub> )
	33	John DOE

SSTable<sub>2</sub>

jdoe	(t <sub>2</sub> )
	34

# Last Write Win (LWW)

```
login = jdoe;
```

SSTable<sub>1</sub>

jdoe	(t <sub>1</sub> )	(t <sub>1</sub> )
	33	John DOE

SSTable<sub>2</sub>

jdoe	(t <sub>2</sub> )	
	34	

tombstone

SSTable<sub>3</sub>

jdoe	(t <sub>3</sub> )	
	⊗	

# Last Write Win (LWW)

```
login = jdoe;
```

?

SSTable<sub>1</sub>

jdoe	(t <sub>1</sub> )	(t <sub>1</sub> )
	33	John DOE

?

SSTable<sub>2</sub>

jdoe	(t <sub>2</sub> )	
	34	

?

SSTable<sub>3</sub>

jdoe	(t <sub>3</sub> )	
	<del>34</del>	

# Last Write Win (LWW)

login = jdoe;

×

SSTable<sub>1</sub>

jdoe	(t <sub>1</sub> )	(t <sub>1</sub> )
	33	John DOE


×

SSTable<sub>2</sub>

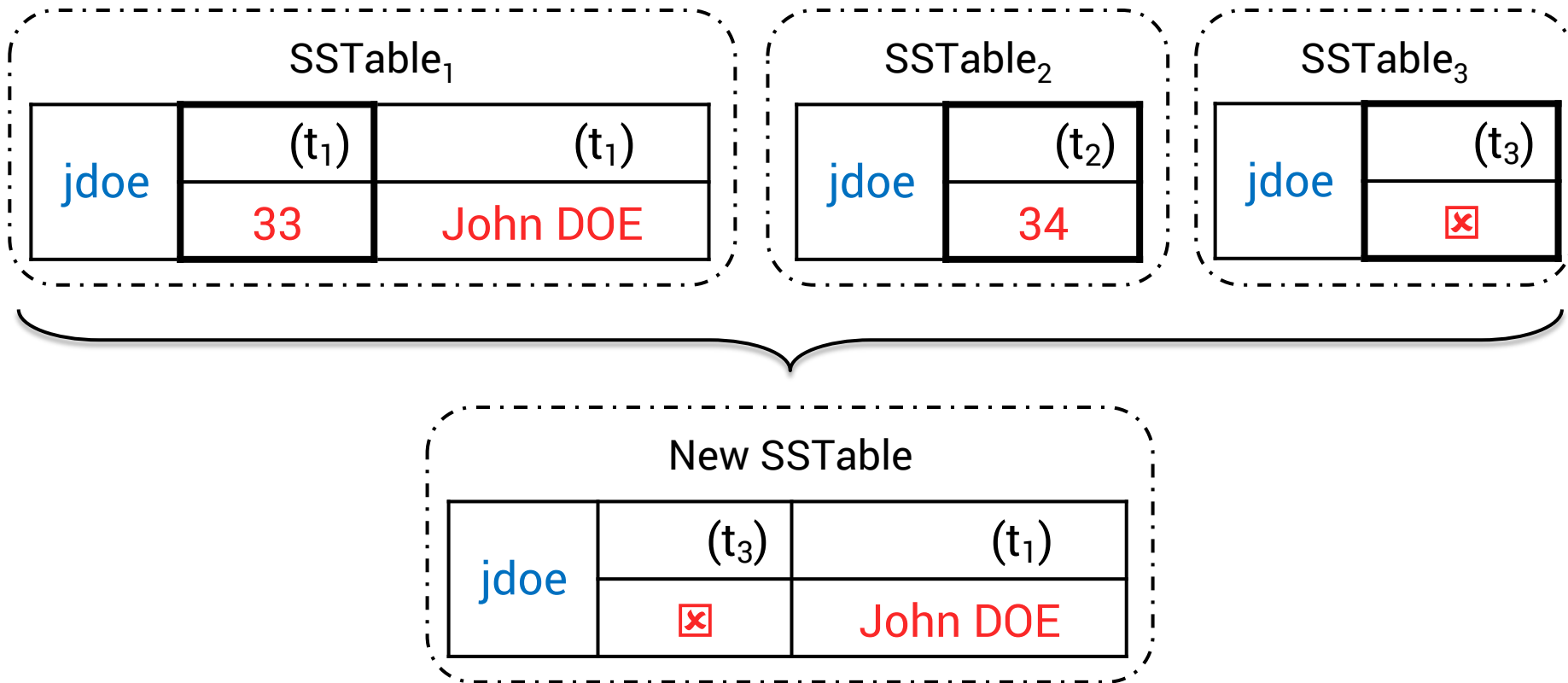
jdoe	(t <sub>2</sub> )	
	34	

✓

SSTable<sub>3</sub>

jdoe	(t <sub>3</sub> )	
		

# Compaction

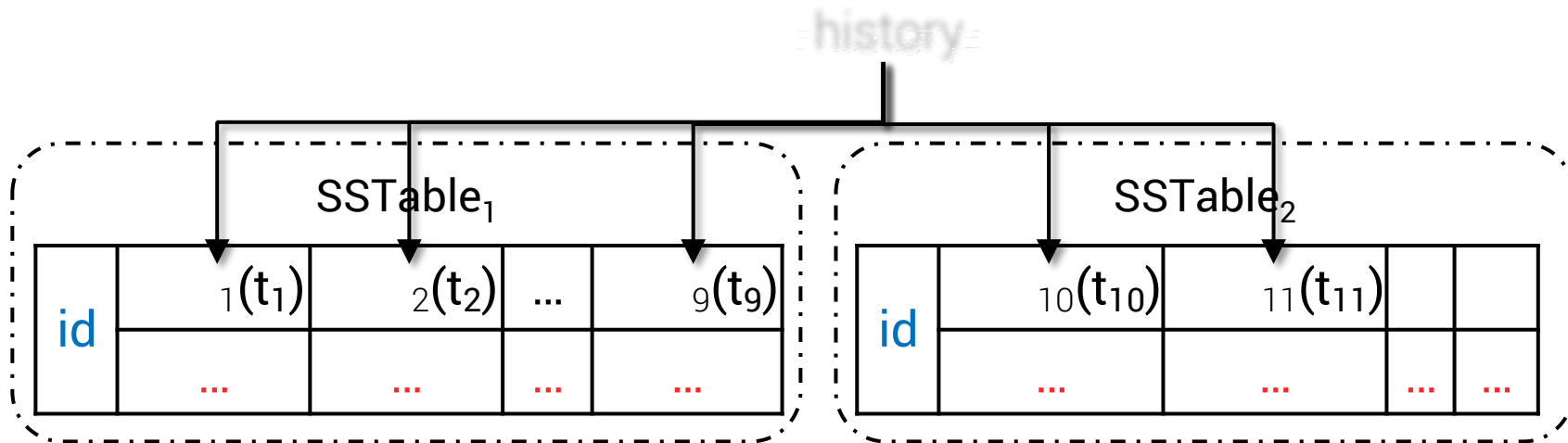




# Historical data

?

- do not
-  time-series



# CRUD operations

```
(login, , ) (jdoe, John DOE, 33);
```

```
= 34 login = jdoe;
```

```
login = jdoe;
```

```
login = jdoe;
```

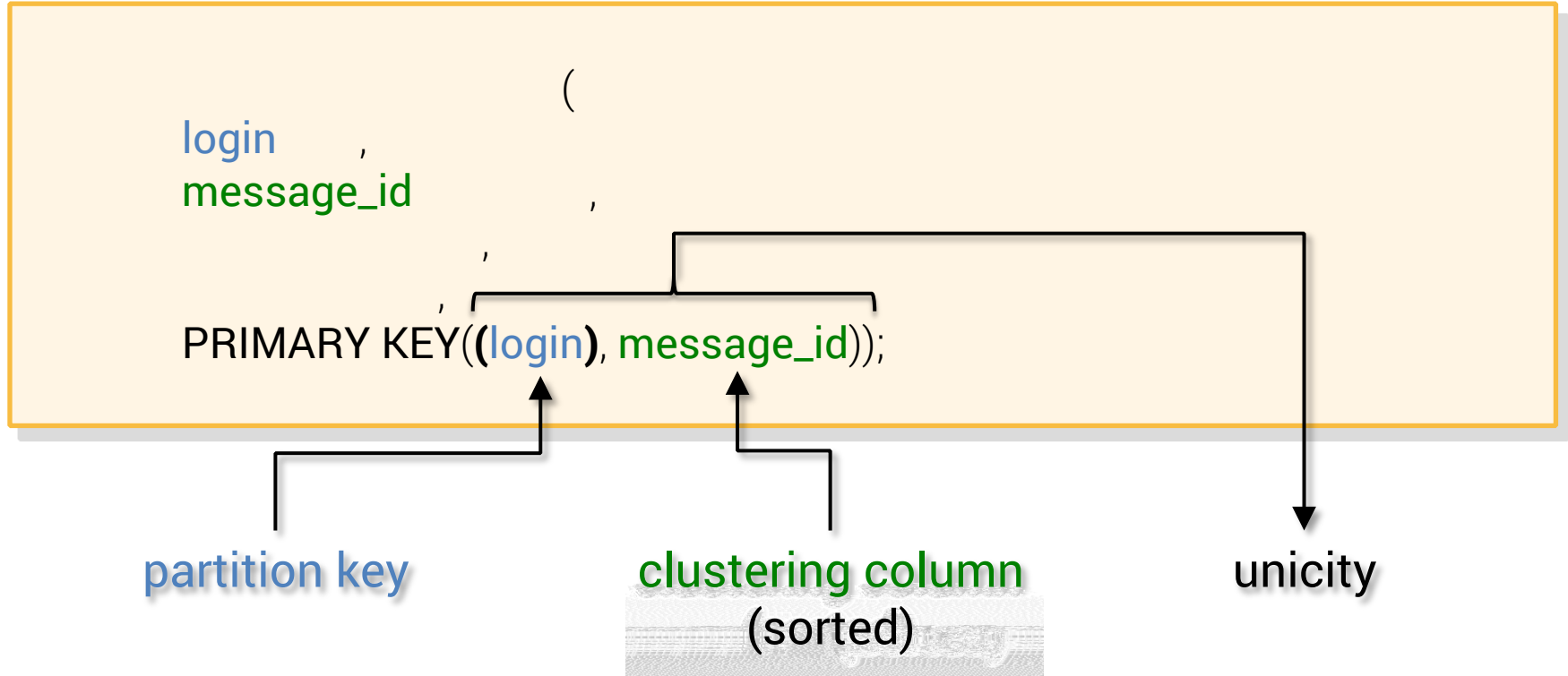
# Simple Table

```
login  
,  
,  
,  
PRIMARY KEY(login));
```

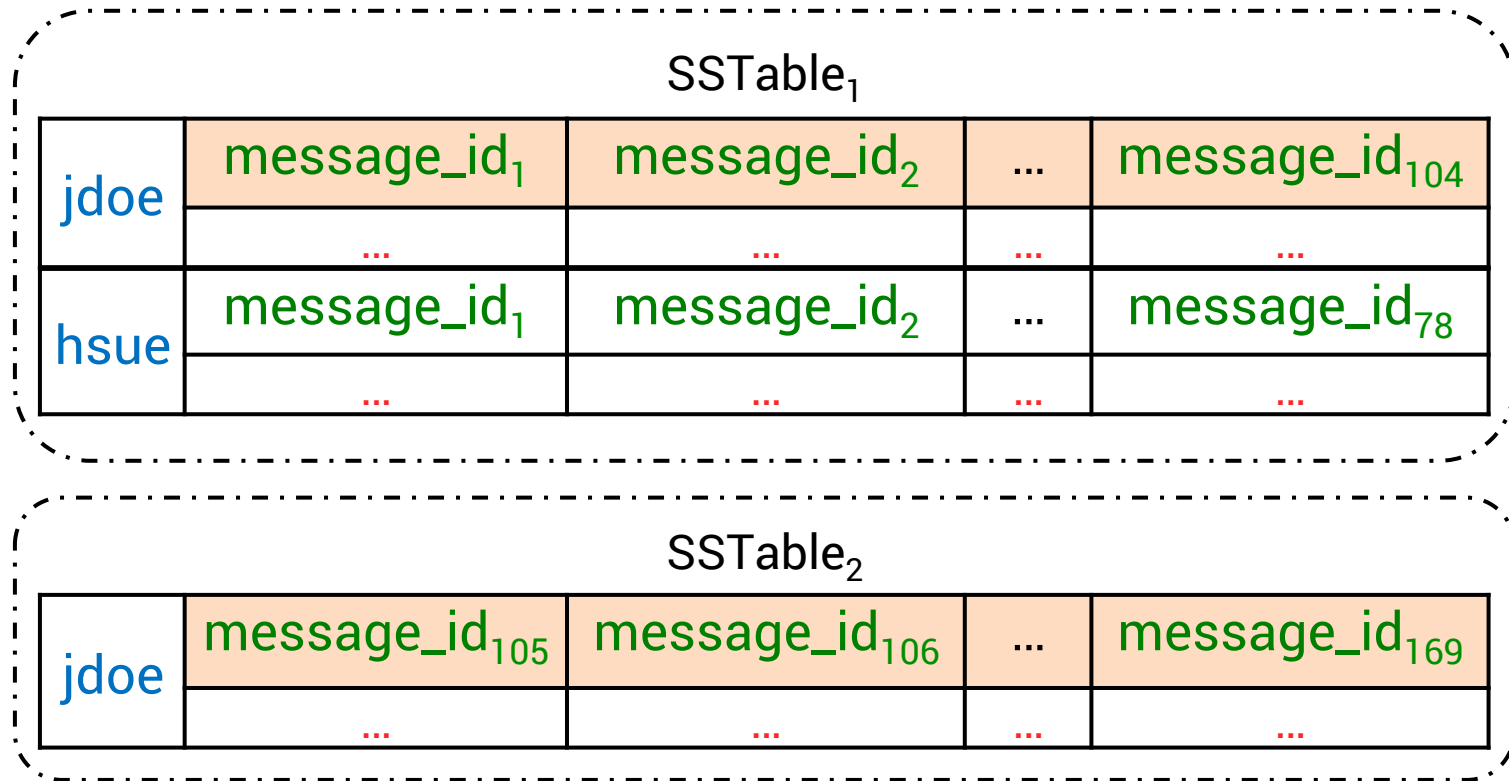
partition key (#partition)



# Clustered table (1 – N)



# On disk layout



# Queries

( )

(#partition )

\*

message\_id = 2014-09-25 16:00:00 ; + 

(#partition )

\*

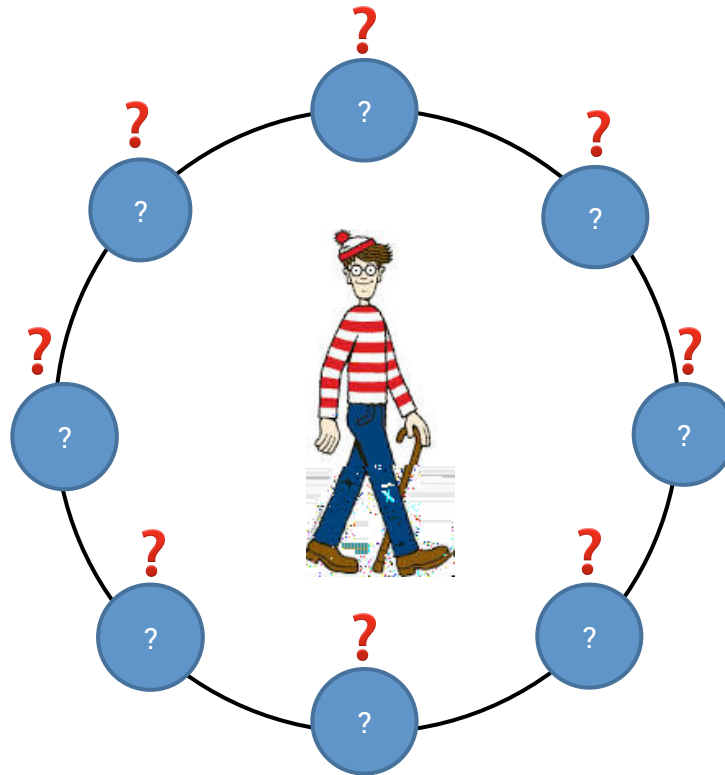
message\_id <= 2014-09-25 16:00:00 +   
message\_id >= 2014-09-20 16:00:00 ; +

# Without #partition

#



?





# The importance of #partition

In RDBMS, no primary key

 **full table scan**



With Cassandra, no **partition key**

 **full CLUSTER scan**



( #partition)

\*

login >= hsue

login <= jdoe;



( #partition)

\*

login

%doe%;



# Clustering order

```
login  
message_id  
  
((login), message_id))  
CLUSTERING ORDER BY (message_id );
```

# Reverse on disk layout

SSTable<sub>1</sub>

jdoe	message_id <sub>169</sub>	message_id <sub>168</sub>	...	message_id <sub>105</sub>
	...	...	...	...

SSTable<sub>2</sub>

jdoe	message_id <sub>104</sub>	message_id <sub>103</sub>	...	message_id <sub>1</sub>
	...	...	...	...

# WHERE clause restrictions

( / / / )

#partition

exact match (=) #partition,

(<, ≤, >, ≥)

-  full cluster scan

clustering columns,

(<, ≤, >, ≥)

exact match

WHERE

-

# WHERE clause restrictions

?

• ,

# WHERE clause restrictions

?

- ,   
 -   
 \_\_\_\_\_

👉 Apache Solr ( ) (Datastax Enterprise)

👉 , 1- -2- ( & )



# WHERE clause restrictions

?

•  
\_\_\_\_\_  
-

👉 Apache Solr ( ) (Datastax Enterprise)

👉 , 1- -2- ( & )

```
* solr_query = age:[33 TO *] AND gender:male ; +
```

```
* solr_query = lastname:*schwei?er ; +
```

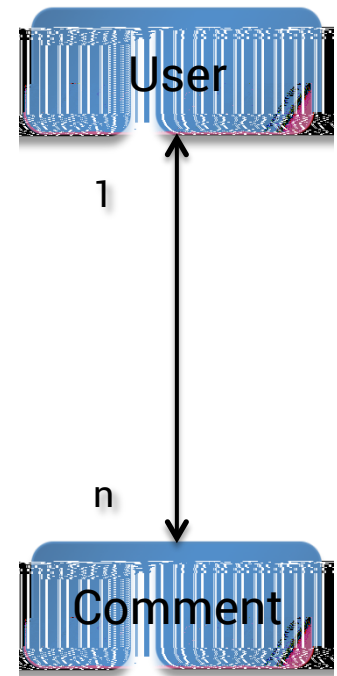
# Collections & maps

```
login      (  
    ,  
    ,  
    ,  
    set<text>,  
    list<text>,  
    map<int, text>,  
    PRIMARY KEY(login));
```

(≈ 1000)

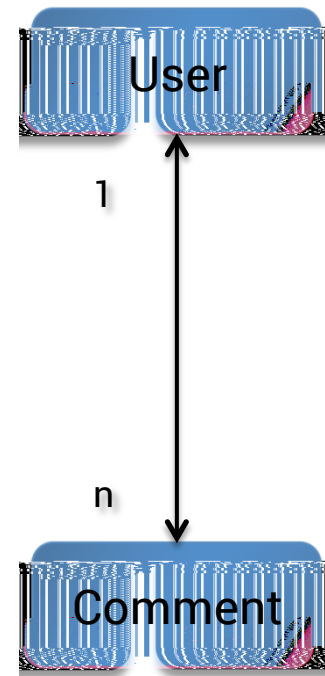
# From SQL to CQL

```
        (
    article_id      ,
    comment_id     ,
    author_id text, // typical join id
        ,
        ((article_id), comment_id));
```



# From SQL to CQL

```
(  
  article_id  
  comment_id  
  author_json text, // de-normalize  
  ((article_id), comment_id));
```



# Data modeling best practices

- 
- $1 \approx 1$

# Data modeling best practices

- 
- 1  $\approx$  1
- , necessary & immutable data
- / trade-off

# Data modeling best practices

## Article title



## Person JSON

- firstname/lastname
- date of birth
- gender
- mood
- location



**John DOE**

Male 33

At 21/03/2011: 10:23



**Helen SUE**

Female 27

At 21/03/2011: 10:12







# Lightweight Transaction (LWT)

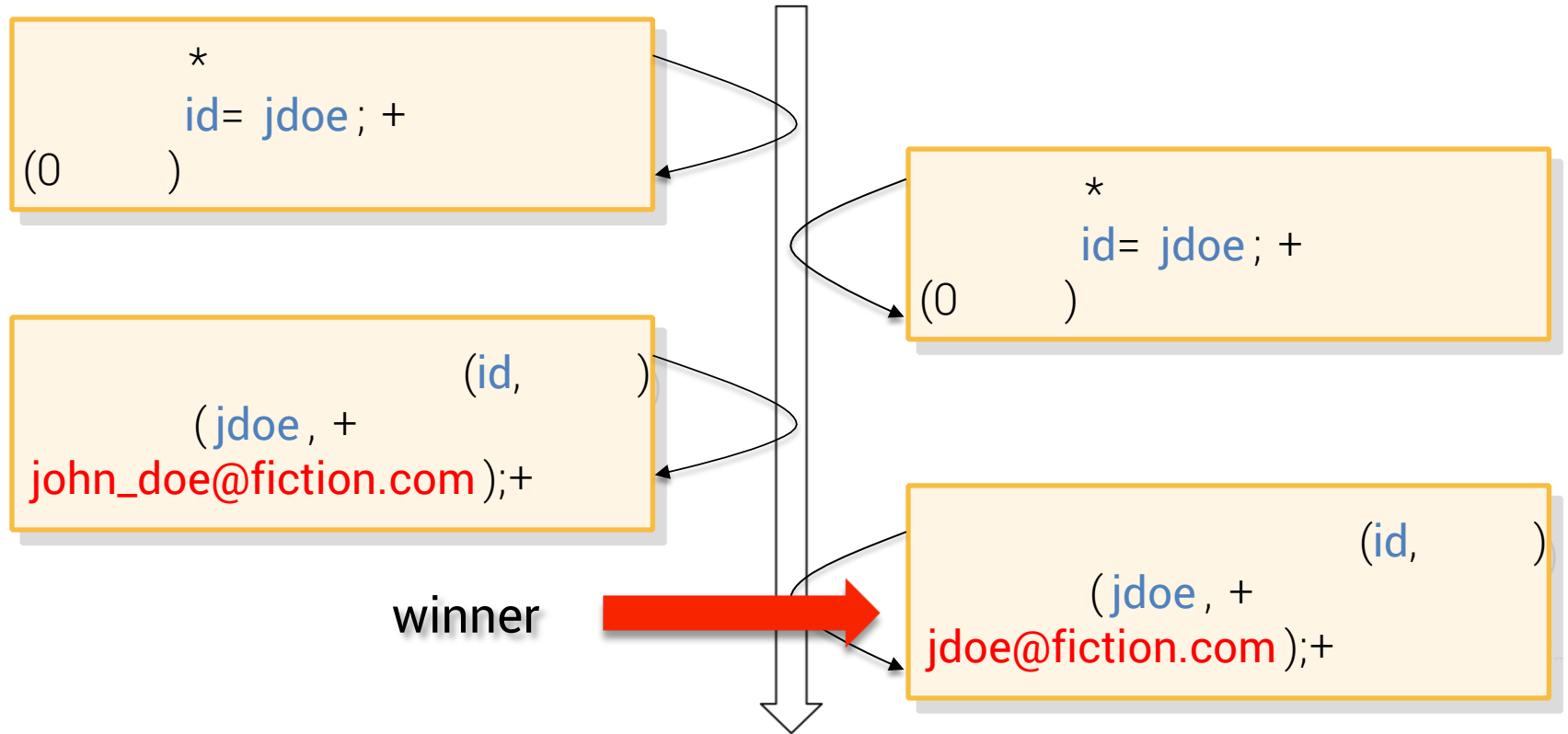
? 

linearizable

? 

---

# Lightweight Transaction (LWT)



# Lightweight Transaction (LWT)

? 

Paxos

?

```
IF NOT EXISTS; ( , ) ( , + . )+  
= . +  
IF email = 'john_doe@fiction.com' = ; +
```

# Lightweight Transaction (LWT)

-  must

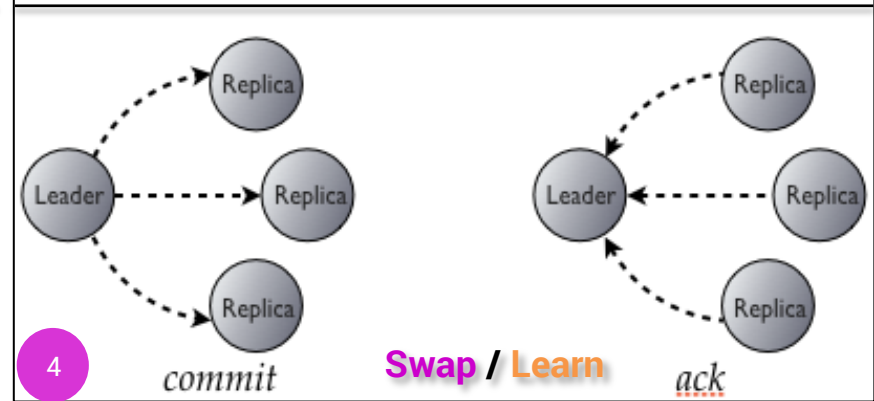
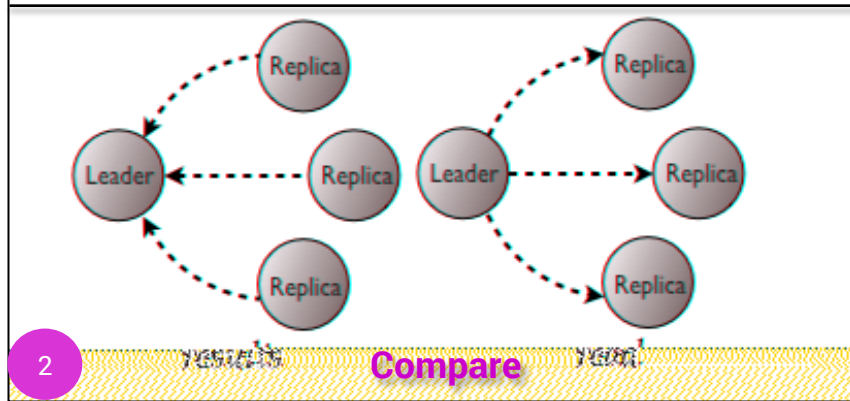
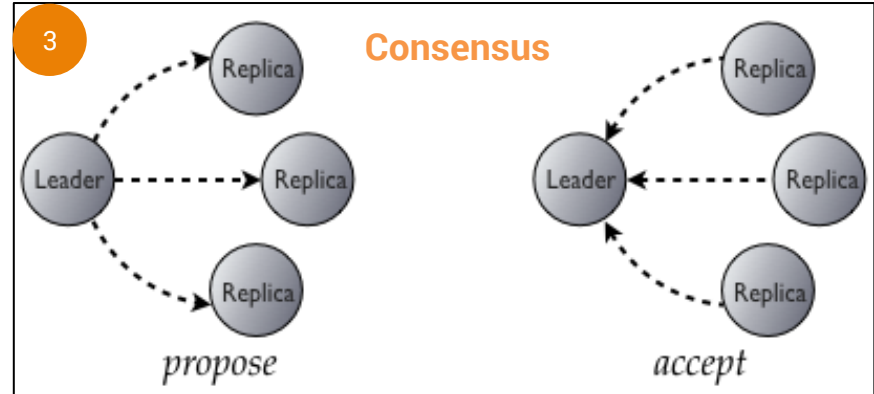
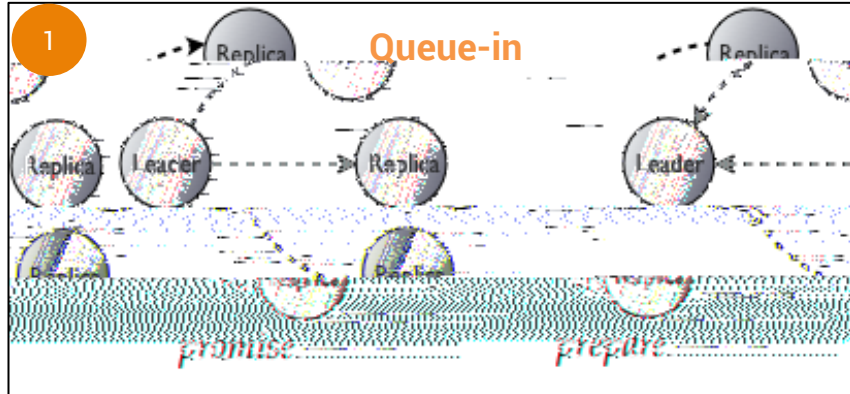
**IF NOT EXISTS  
IF EXISTS**

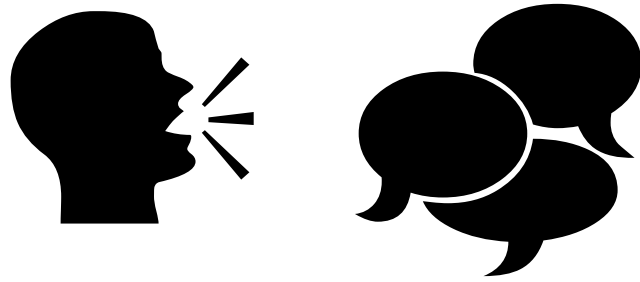


# Lightweight Transaction (LWT)

- (4 - ), **do not abuse**
- 1% - 5%

# Lightweight Transaction (LWT)





Q & R

# Thank You



**@doanduyhai**



**duy\_hai.doan@datastax.com**

**<https://academy.datastax.com/>**