DuyHai DOAN, Technical Advocate

# Shameless self-promotion



Duy Hai DOAN

-       ,       ,
-     -           (Achilles,   )                           :
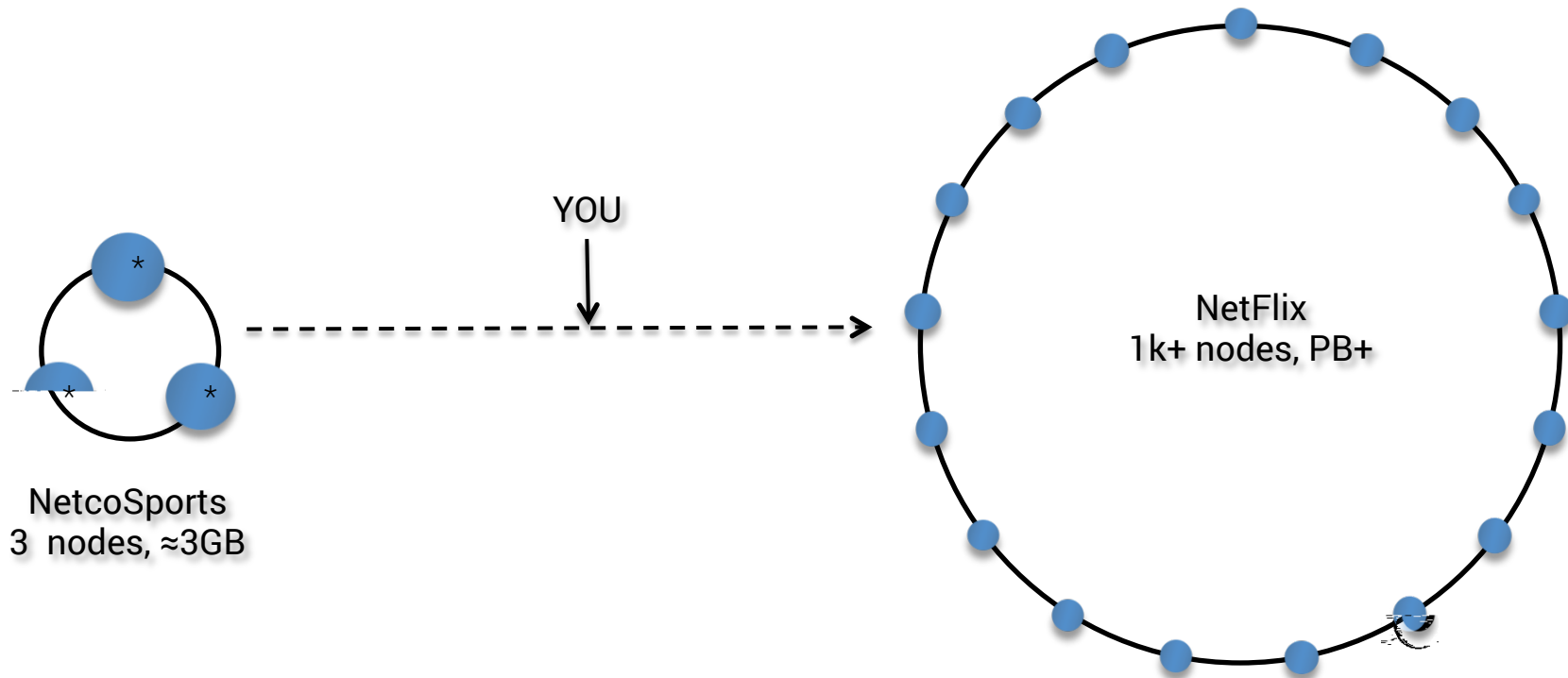- 

      ☞ **duy_hai.doan@datastax.com**

- production

# Datastax

-

# Agenda

## Architecture

- 
- 

## Data model

-       ( ),
-       ( , , , )
-

# Cassandra 5 key facts

YOU

NetcoSports
3 nodes, ≈3GB

NetFlix
1k+ nodes, PB+

# Cassandra 5 key facts

(≈100%    -    )

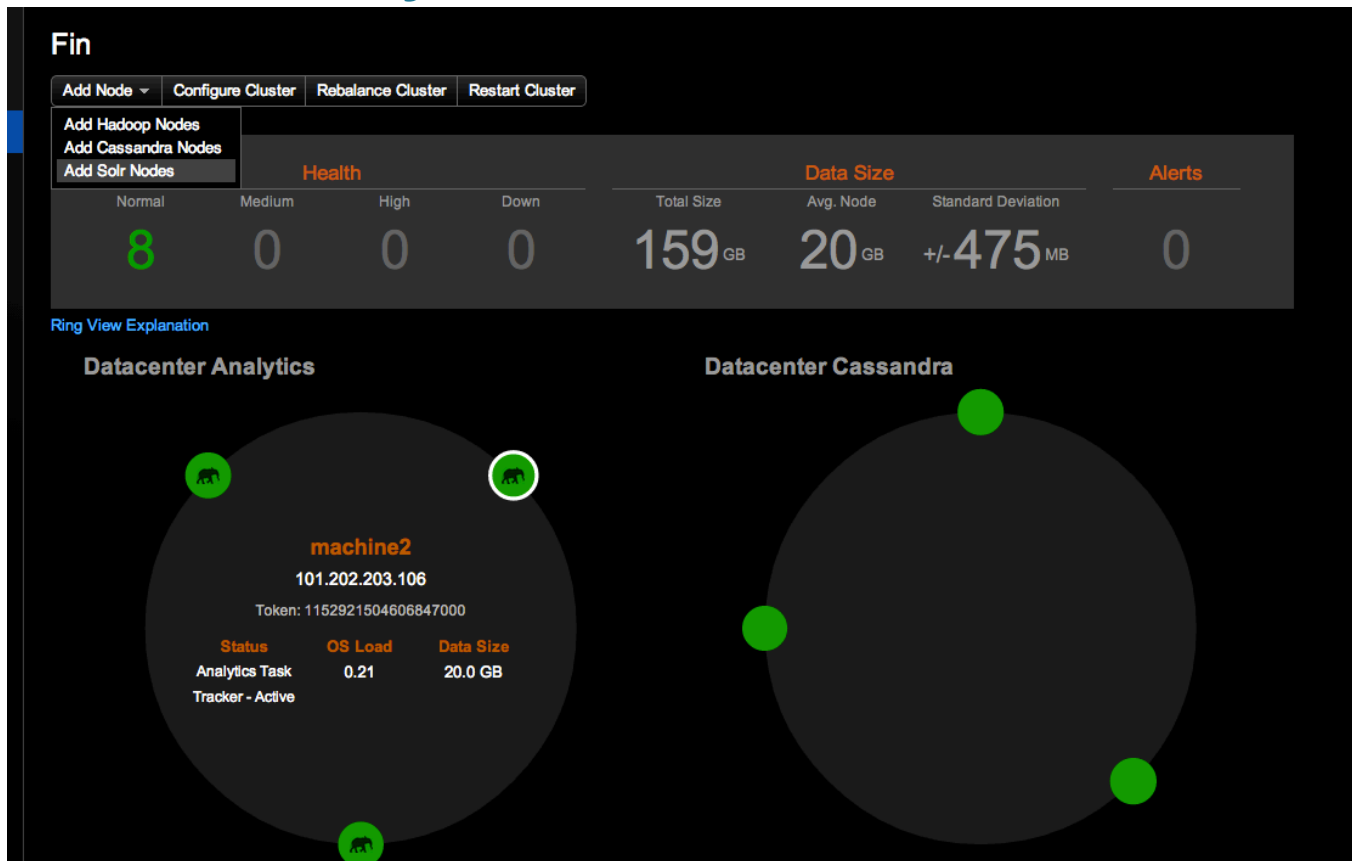- (          )

# Cassandra 5 key facts

-

- - - - ( )
- -
- /

# Cassandra 5 key facts

- 1 = 1 + 1
- 
-

# Cassandra 5 key facts

# Cassandra 5 key facts

- Cassandra + Spark = awesome !
- realtime streaming

# Cassandra architecture

Cluster
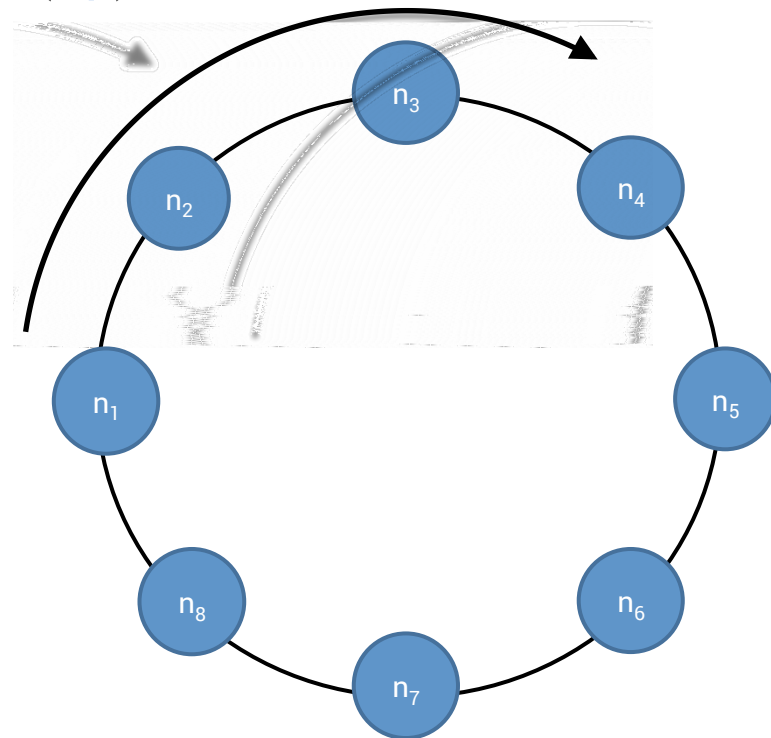
Replication

# Cassandra architecture

DATASTAX

- DynamoDB
- masterless

    -

- Big Table
- /

# Data distribution

: #partition → token = (#p)

: - ,

= $(2^{64}/2)$

# Token Ranges
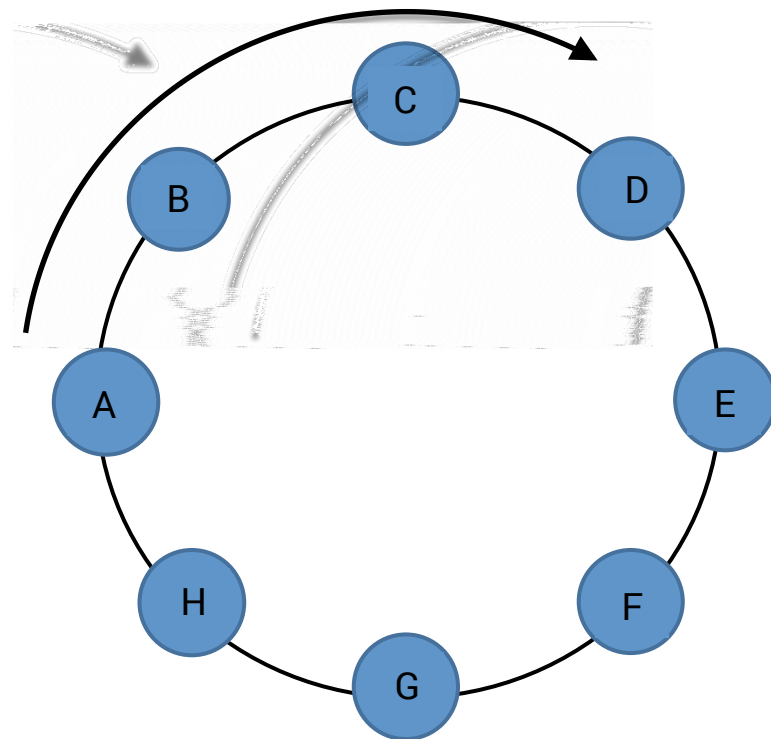
A: 0,  /8

B:   /8, 2  /8

C: 2  /8, 3  /8

D: 3  /8, 4  /8

E: 4  /8, 5  /8
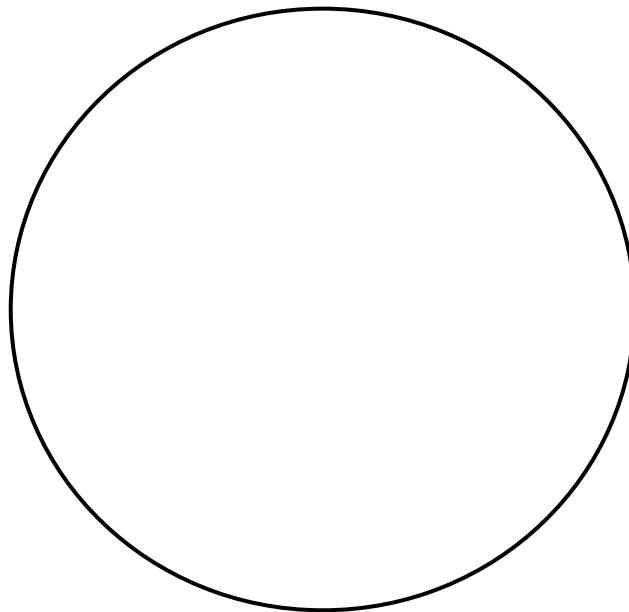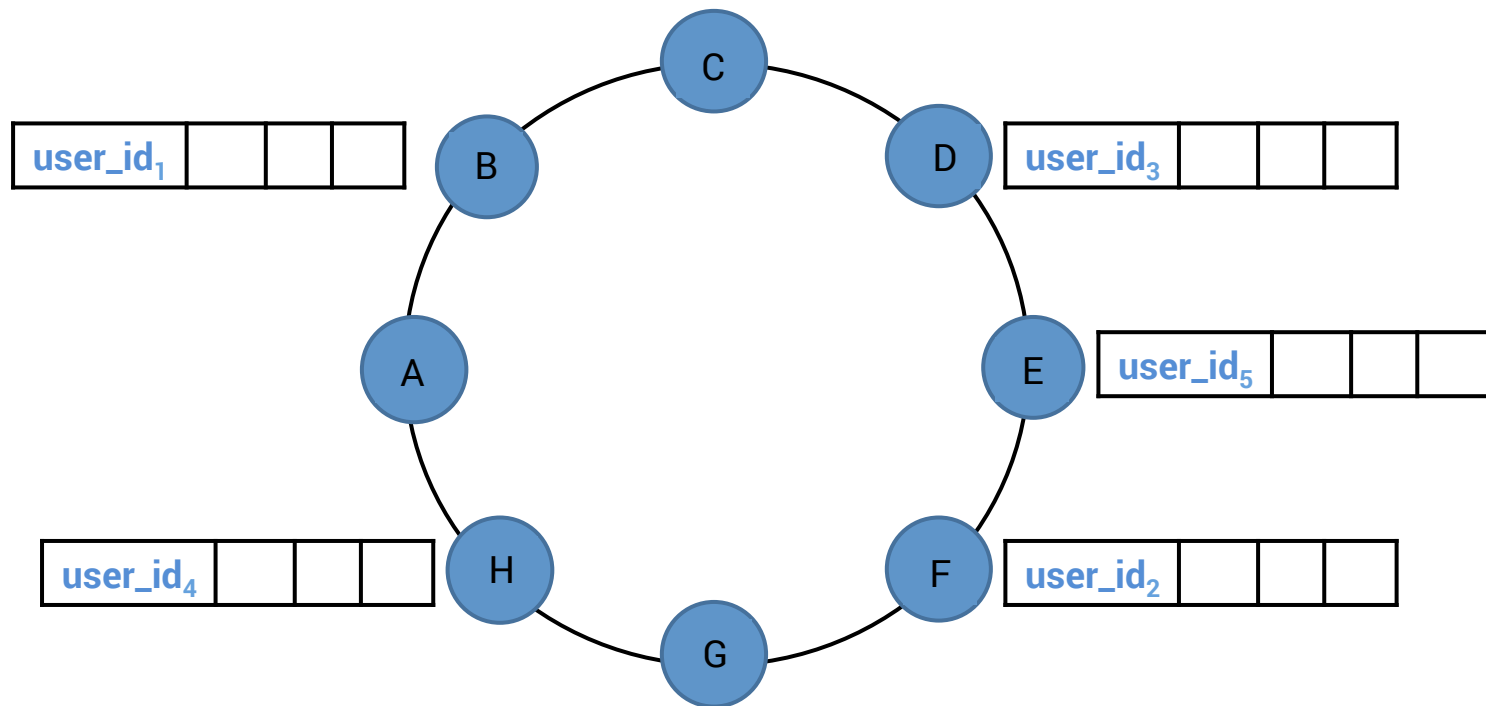
F: 5  /8, 6  /8

G: 6  /8, 7  /8

H: 7  /8,

# Distributed Table

# Distributed Table

# Linear scalability

8 nodes

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: ONE: B

B     A     A

:

B     A     A

Read ONE: A ONE: A

# Consistency in action

= 3,   ONE,   QUORUM

Write ONE: B

B    A    A

:

B    A    A

Read QUORUM: A

# Consistency in action

= 3,        ONE,        ALL

Write ONE: B



Read ALL: B

:

# Consistency in action

= 3,          QUORUM,          ONE

QUORUM: B

| B | B | A |
|---|---|---|

:

| B | B | A |
|---|---|---|

ONE: A

# Consistency in action

= 3,     QUORUM,     QUORUM

Write QUORUM: B

| B | B | A |
|---|---|---|
| B | B | A |

Read QUORUM: B

:

# Consistency trade-off

DATASTAX

**Latency**                              **Consistency**

# Consistency level

# ONE

Fast, may not read latest written value

# Consistency level

# QUORUM

Strict majority w.r.t. **R**eplication **F**actor

Good balance

# Consistency level

**DATASTAX**

# ALL

Paranoid

Slow, no high availability

# Consistency summary

$$ONE_{Read} + ONE_{Write}$$

☞ available        /        (N-1)

$$QUORUM_{Read} + QUORUM_{Write}$$

☞ 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 , John DOE , 33);

auto-generated timestamp

| jdoe | (t₁) | (t₁) |
|------|------|------|
|      | 33   | John DOE |

# Last Write Win (LWW)



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

| SSTable$_1$ | | |
|---|---|---|
| jdoe | (t$_1$) | (t$_1$) |
| | 33 | John DOE |

| SSTable$_2$ | |
|---|---|
| jdoe | (t$_2$) |
| | 34 |

# Last Write Win (LWW)

login = jdoe;

tombstone

### SSTable$_1$

| jdoe | ($t_1$) | ($t_1$) |
|------|---------|---------|
|      | 33      | John DOE |

### SSTable$_2$

| jdoe | ($t_2$) |
|------|---------|
|      | 34      |

### SSTable$_3$

| jdoe | ($t_3$) |
|------|---------|
|      | ☒       |

# Last Write Win (LWW)

login = jdoe;

| **SSTable$_1$** | | |
|---|---|---|
| jdoe | $(t_1)$ | $(t_1)$ |
| | 33 | John DOE |

| **SSTable$_2$** | |
|---|---|
| jdoe | $(t_2)$ |
| | 34 |

| **SSTable$_3$** | |
|---|---|
| jdoe | $(t_3)$ |
| | ☒ |

? ? ?

# Last Write Win (LWW)



login = jdoe;

| SSTable$_1$ | | |
|---|---|---|
| jdoe | (t$_1$) | (t$_1$) |
| | 33 | John DOE |

| SSTable$_2$ | |
|---|---|
| jdoe | (t$_2$) |
| | 34 |

| SSTable$_3$ | |
|---|---|
| jdoe | (t$_3$) |
| | ☒ |

# Compaction



SSTable$_1$

SSTable$_2$

SSTable$_3$

jdoe      ($t_3$)

# Historical data

?

- do not
- ☞ time-series

# CRUD operations

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

= 34 login = jdoe;

login = jdoe;

login = jdoe;

# Simple Table

login (

```
                    (
    login        ,
    message_id           ,

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

partition key                clustering column                unicity
                             (sorted)

# Queries

(     )

```
          *                          login = jdoe
     message_id =  2014-09-25 16:00:00 ;
```

```
          *                          login = jdoe
     message_id <=  2014-09-25 16:00:00
     message_id >=  2014-09-20 16:00:00 ;
```

# Queries

(#partition                    )

<div style="border:1px solid orange;background:#fdf6e3;">

    *                    message_id = 2014-09-25 16:00:00 ; ❌

</div>

(#partition                )

<div style="border:1px solid orange;background:#fdf6e3;">

  *

message_id <= 2014-09-25 16:00:00  ❌
message_id >= 2014-09-20 16:00:00 ;

</div>

# Queries

( #partition)

* login >= hsue login <= jdoe; ❌

( #partition)

* login %doe% ; ❌

# On disk layout

SSTable$_1$

| jdoe | message_id$_1$ | message_id$_2$ | ... | message_id$_{104}$ |
|------|----------------|----------------|-----|--------------------|
|      | ...            | ...            | ... | ...                |
| hsue | message_id$_1$ | message_id$_2$ | ... | message_id$_{78}$  |
|      | ...            | ...            | ... | ...                |

SSTable$_2$

| jdoe | message_id$_{105}$ | message_id$_{106}$ | ... | message_id$_{169}$ |
|------|--------------------|--------------------|-----|--------------------|
|      | ...                | ...                | ... | ...                |

DATASTAX

# Clustering order

```
                          (
        login      ,
        message_id              ,
                        ,
                    ,
                    ((login), message_id))
        CLUSTERING ORDER BY (message_id        ) ;
```
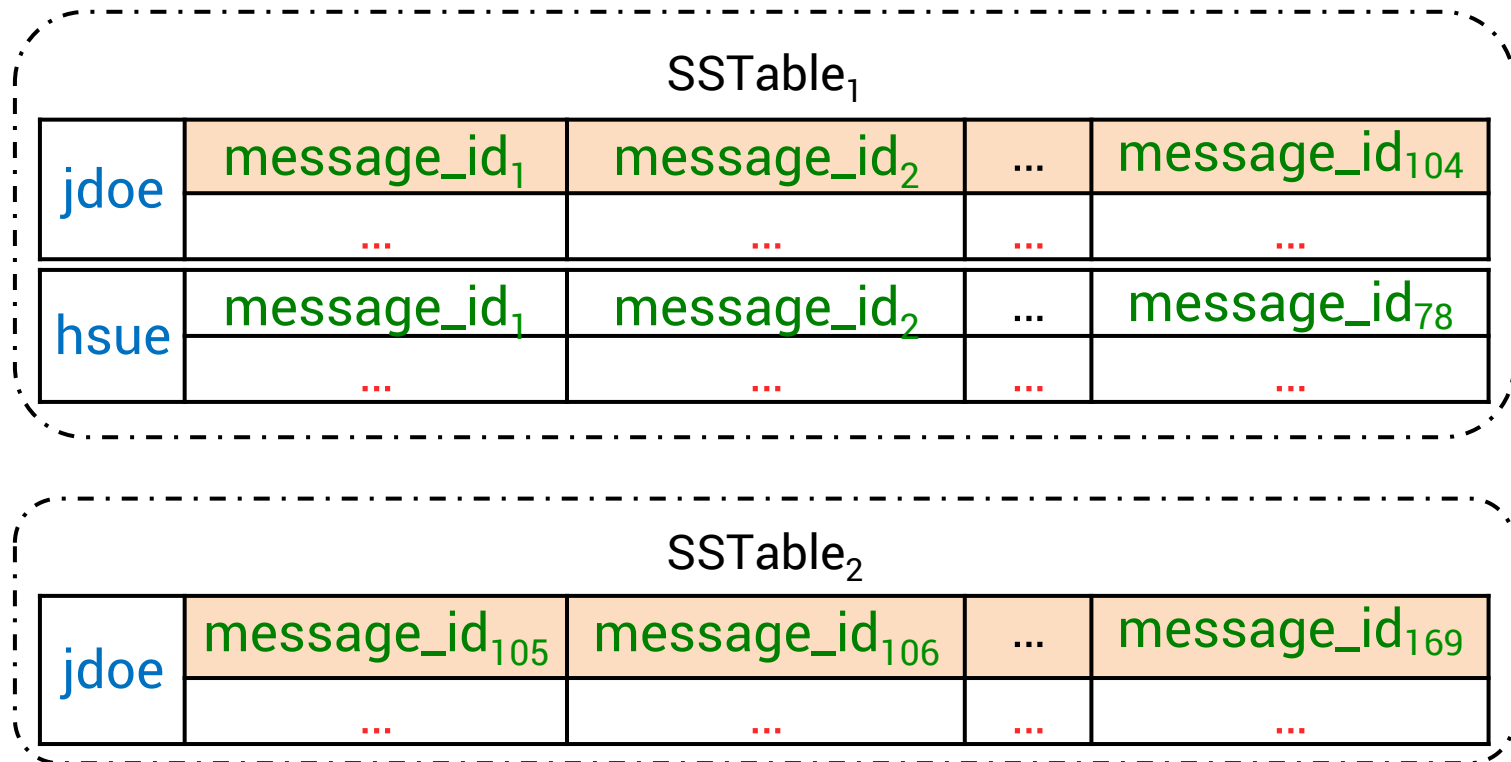
# Reverse on disk layout

SSTable$_1$

| jdoe | message_id$_{169}$ | message_id$_{168}$ | ... | message_id$_{105}$ |
|------|---------------------|---------------------|-----|---------------------|
|      | ... | ... | ... | ... |

SSTable$_2$

| jdoe | message_id$_{104}$ | message_id$_{103}$ | ... | message_id$_1$ |
|------|---------------------|---------------------|-----|-----------------|
|      | ... | ... | ... | ... |

# WHERE clause restrictions

( / / / ) #partition

exact match (=) #partition, (<, ≤, >, ≥)
- ☞ full cluster scan

clustering columns, (<, ≤, >, ≥) exact match

WHERE PRIMARY KEY

# Dynamic search

?

- ,

# Dynamic search

?

- ,

☞ Apache Solr (         )              (Datastax Enterprise)

*              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

:                                   (

:                                 (

# CQL    not SQL

**DATASTAX**

:                                        (

# no join

(do you want to scale ?)

:                                    (

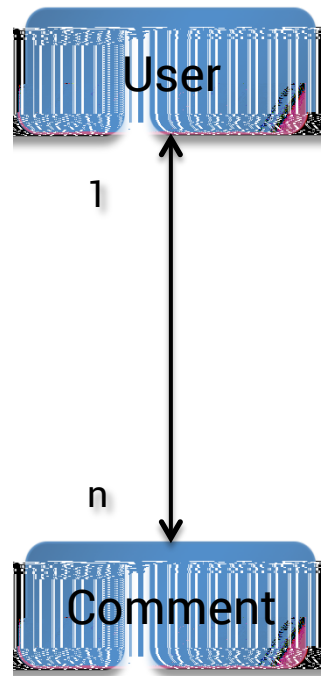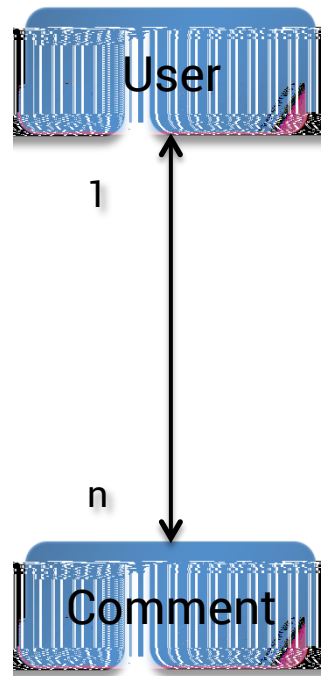# no integrity constraint

(do you want to read-before-write ?)

# From SQL to CQL

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

User

1

n

Comment

# From SQL to CQL

‒

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

User

1

n

Comment

# Data modeling best practices

- 

- 1          $\approx$ 1

# Data modeling best practices

- 
- 1                     ≈ 1



-   ,     necessary & immutable data
-     /   trade-off

# Data modeling best practices

**Article title**

Person **JSON**
- firstname/lastname
- date of birth
- gender
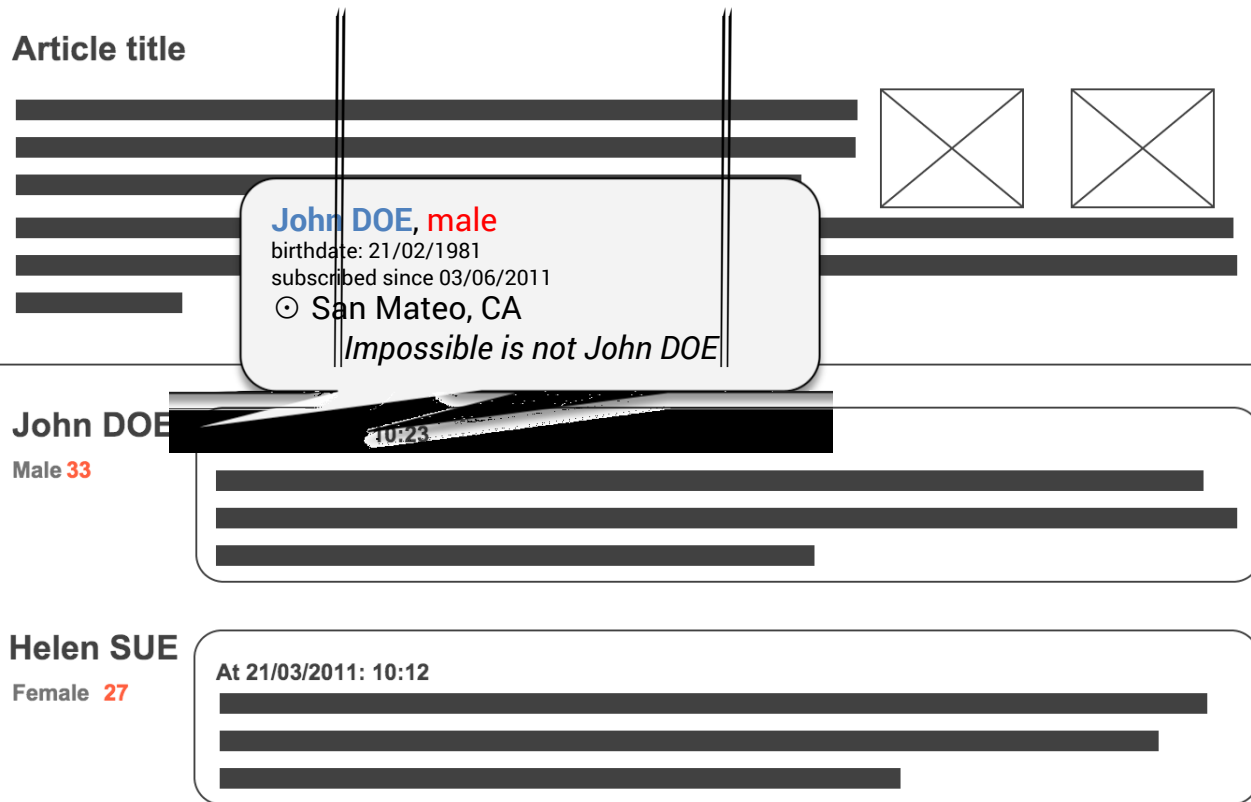- ~~mood~~
- ~~location~~

**John DOE**
Male **33**

At 21/03/2011: 10:23

**Helen SUE**
Female **27**

At 21/03/2011: 10:12

# Data modeling best practices

**Article title**

John DOE, male
birthdate: 21/02/1981
subscribed since 03/06/2011
⊙ San Mateo, CA
*Impossible is not John DOE*

Full detail read from
**User** table on click

**John DOE**
Male **33**

**Helen SUE**
Female **27**

At 21/03/2011: 10:12

# Lightweight Transaction (LWT)

? ☞                    linearizable

? ☞

_____

# Lightweight Transaction (LWT)



winner
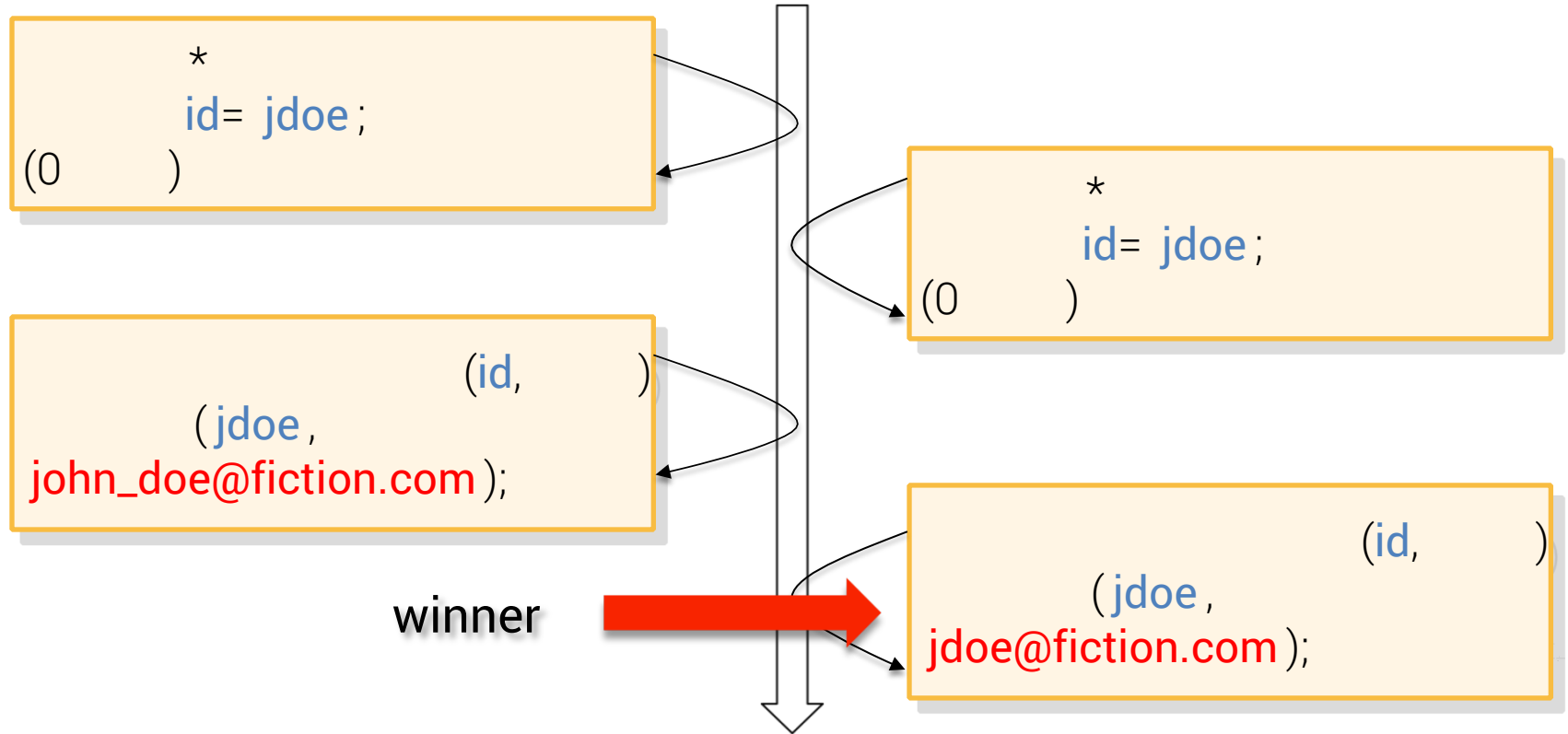
# Lightweight Transaction (LWT)

?  ☞                          Paxos

    ?

IF NOT EXISTS;            (   ,      )          (      ,                              .      )

                                    =                          .
IF email = 'john_doe@fiction.com'              =        ;

# Lightweight Transaction (LWT)

- ☞ **must**

IF NOT EXISTS : (
☞            IF EXISTS : (

# Lightweight Transaction (LWT)

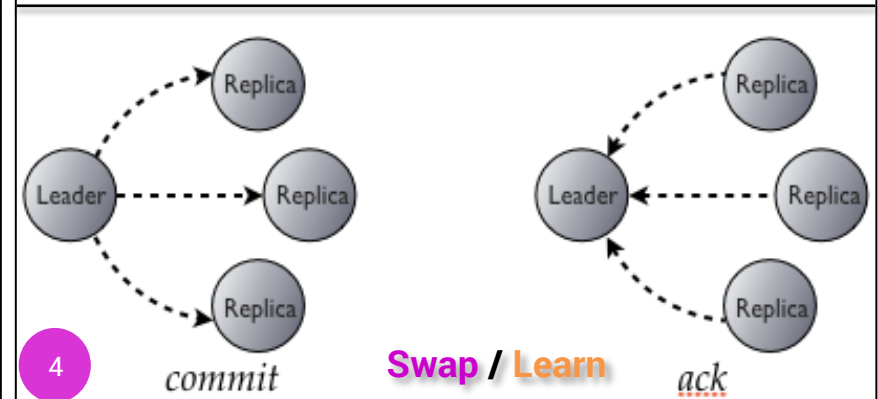- must ☞ _____

> =     <span style="color:red">IF condition_column = yyy</span>
>
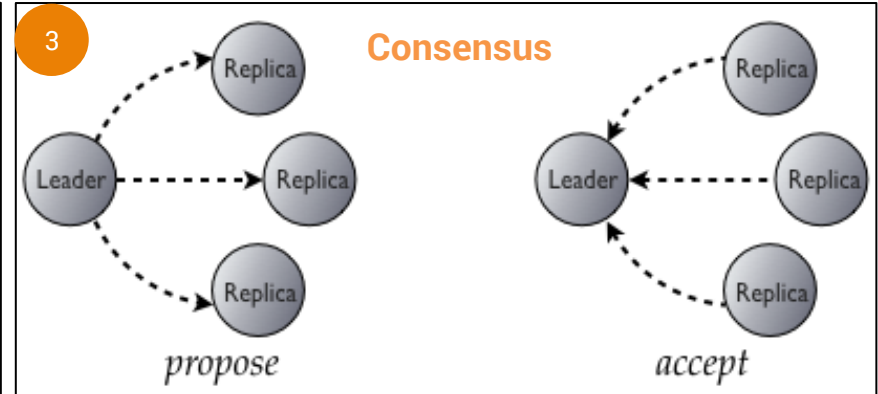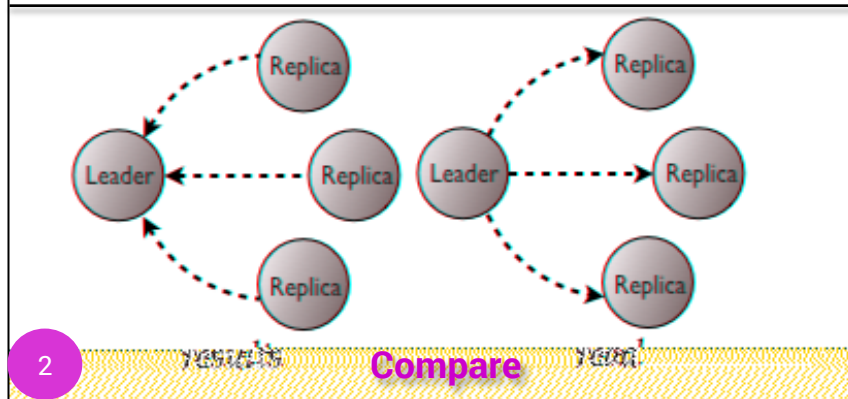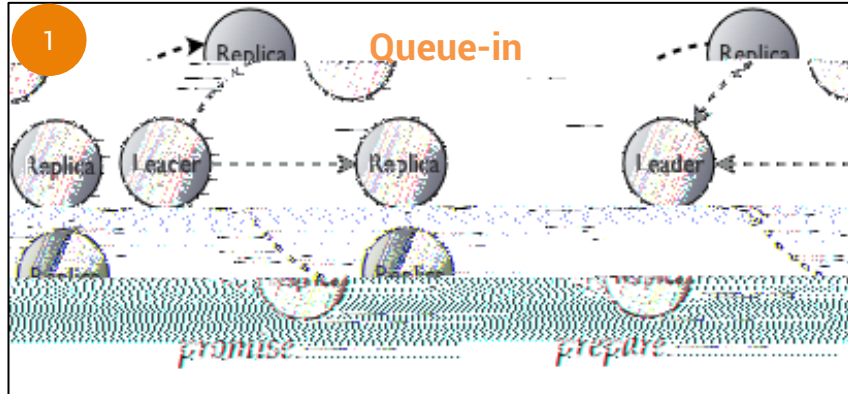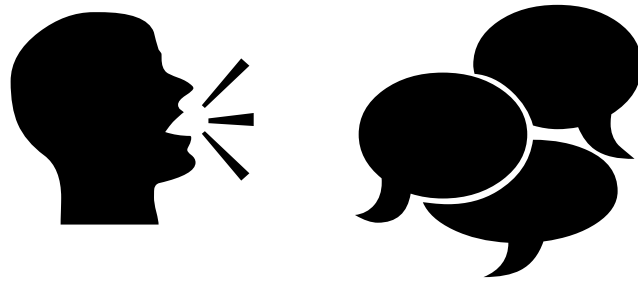> ☞     condition_column =    <span style="color:red">IF</span>       :

# 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/