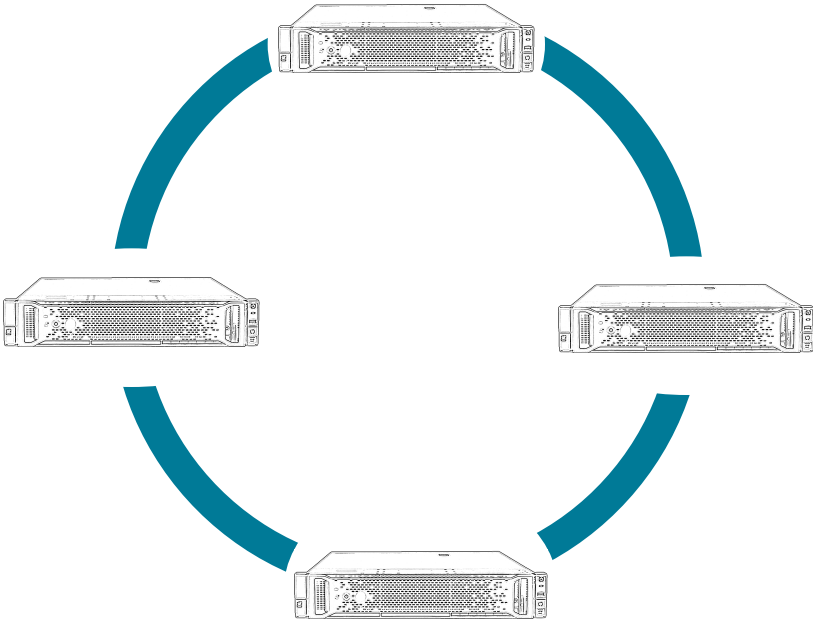# DataStax Enterprise Architecture

Negib Marhoul, Solution Engineer, DataStax

8. June 2017

# Agenda

| | |
|---|---|
| 1 | Topology and Data Structure |
| 2 | Request Handling |
| 3 | Lab1: Cassandra Access and Cassandra Stress |

DATASTAX

# Design Goals and Objectives

- Continuously Available
- Master Less
- Fully Distributed
- Shared-Nothing Architecture
- Build In Replication
- Linear Scalability
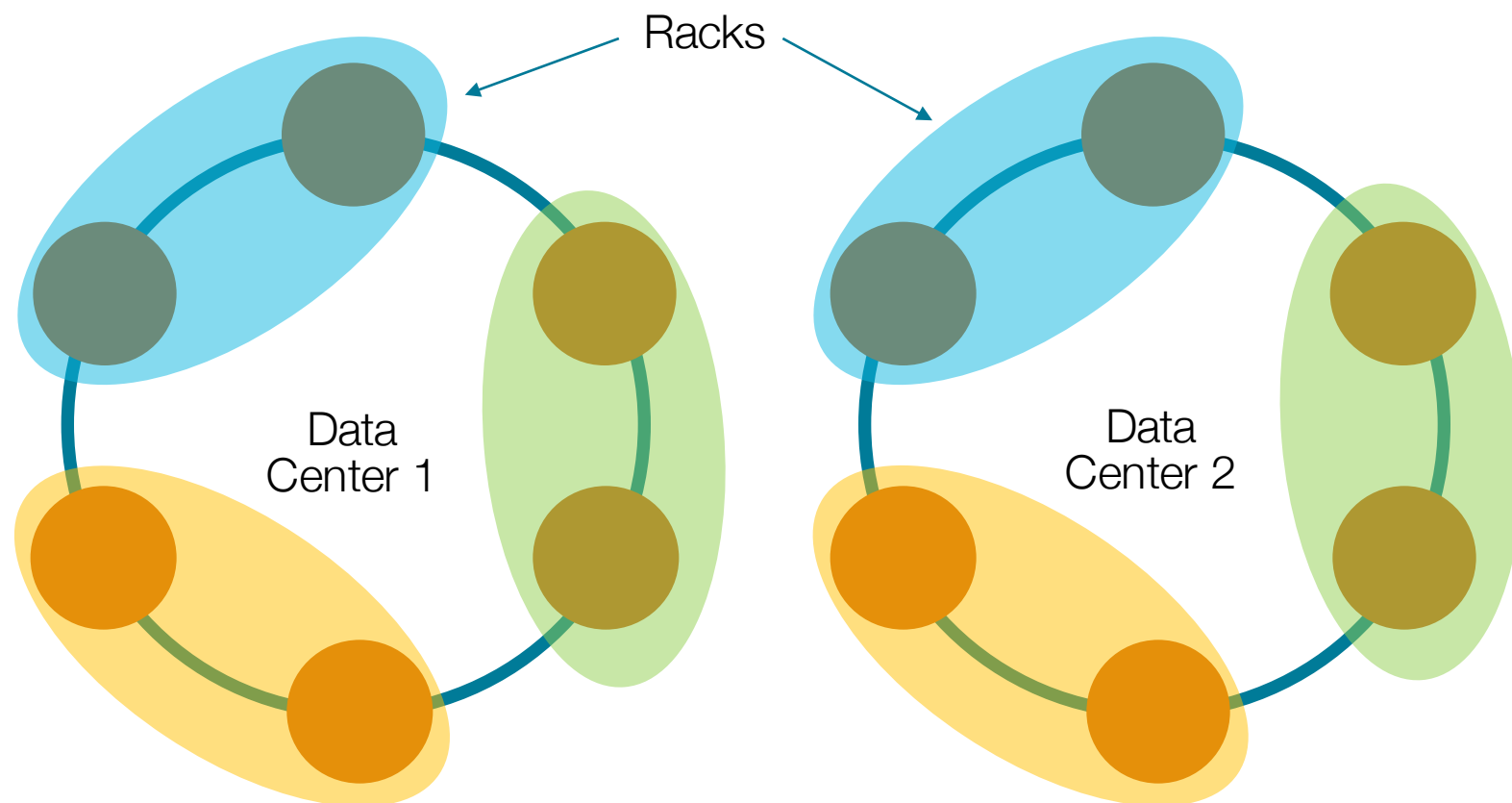- Scale out

# Architecture

# Apache Cassandra™ Architecture

- Cluster layer

- Amazon DynamoDB paper

- masterless architecture


- Data-store layer

- Google Big Table paper
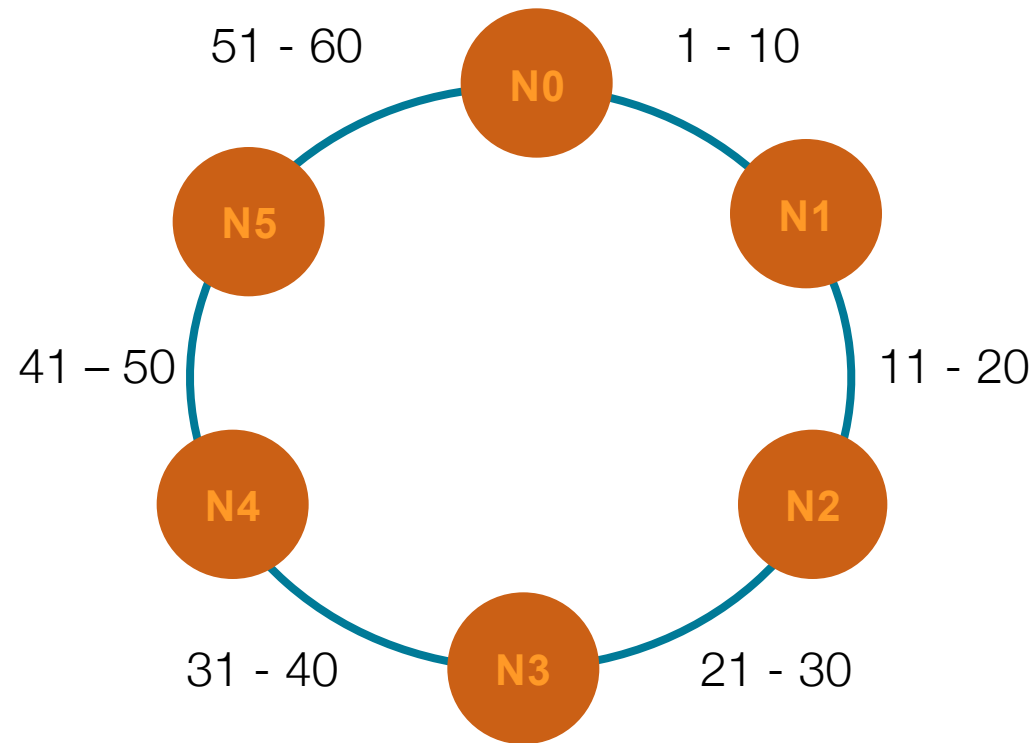
- Columns/columns family

DATASTAX

- **All nodes are peers**

  – Including seed nodes

  – No master

  – Discovery through gossip


- **Built-in replication**

  – Simplify your architecture!

DATASTAX

# Token Ranges



51 - 60    **N0**    1 - 10

**N5**        **N1**

41 – 50            11 - 20

**N4**        **N2**

31 - 40    **N3**    21 - 30

**Token Range : $2^{-63} - 2^{63}$**

Example with **Replication Factor 3**

N3 will own data for tokens 1 – 30

```
Token Range :  1-10, owned by N1,N2,N3
Token Range : 11-20, owned by N2,N3,N4
Token Range : 21-30, owned by N3,N4,N5
```

DATASTAX

- Primary key
  - Partition key
  - Clustering columns
- Partitioner
  - Generates unique hash from partition key
- Replication strategy
  - Token hash determines starting point
  - Determines replica placement
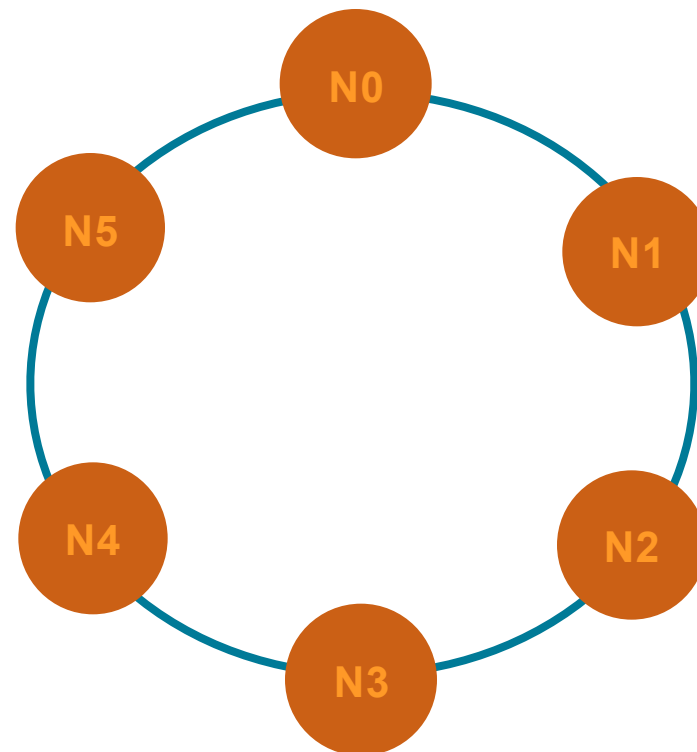
# Cassandra Query Language

```
CREATE KEYSPACE retailer WITH replication =
{'class': 'NetworkTopologyStrategy', 'DC1': '3'}          ←——— Replication Factor
AND durable_writes = true;



CREATE TABLE retailer.sales_by_customer (
    orderid  int,
    salesdt  date,
    revenue  double,
    discount double,        Partition Key
    comment  txt,
    PRIMARY KEY(orderid);
) WITH CLUSTERING ORDER BY (dt DESC)


SELECT * FROM sales_by_customer where custid=1 OR custid=2 AND salesdt >=20160401;
```
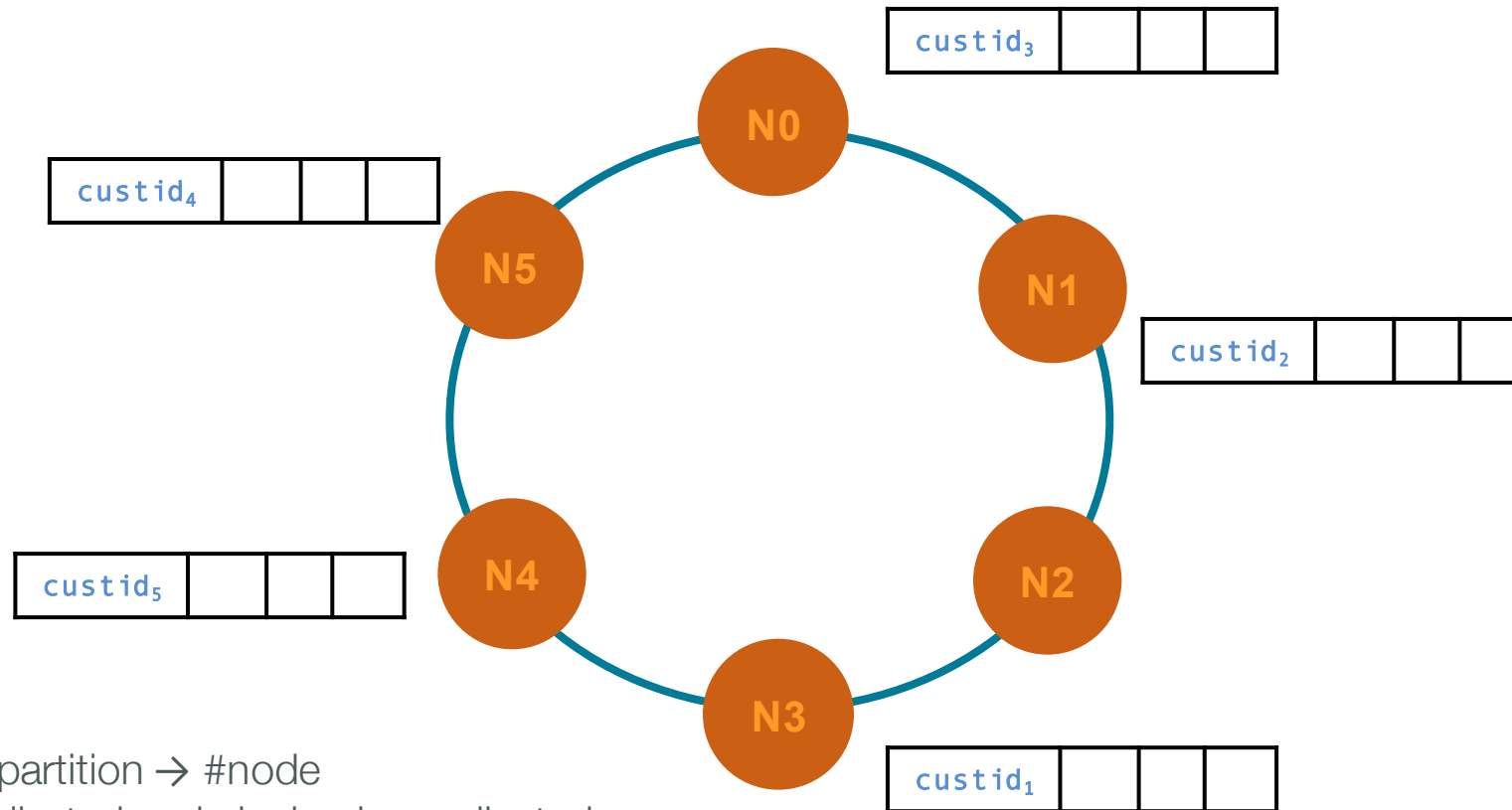
DATASTAX

# Data Distribution

| | | | |
|---|---|---|---|
| $custid_1$ | | | |
| $custid_2$ | | | |
| $custid_3$ | | | |
| $custid_4$ | | | |
| $custid_5$ | | | |

N0

N5

N1

N4

N2

N3

Token = hash of #partition → #node
Data is evenly distributed and clock wise replicated

DATASTAX

# Data Distribution



custid$_3$

custid$_4$

custid$_2$

N0

N5

N1

N4

N2

N3

custid$_5$

custid$_1$

Token = hash of #partition → #node
Data is evenly distributed and clock wise replicated

DATASTAX

# Lab 1 : Accessing the cluster

The power behind the moment. | DATASTAX
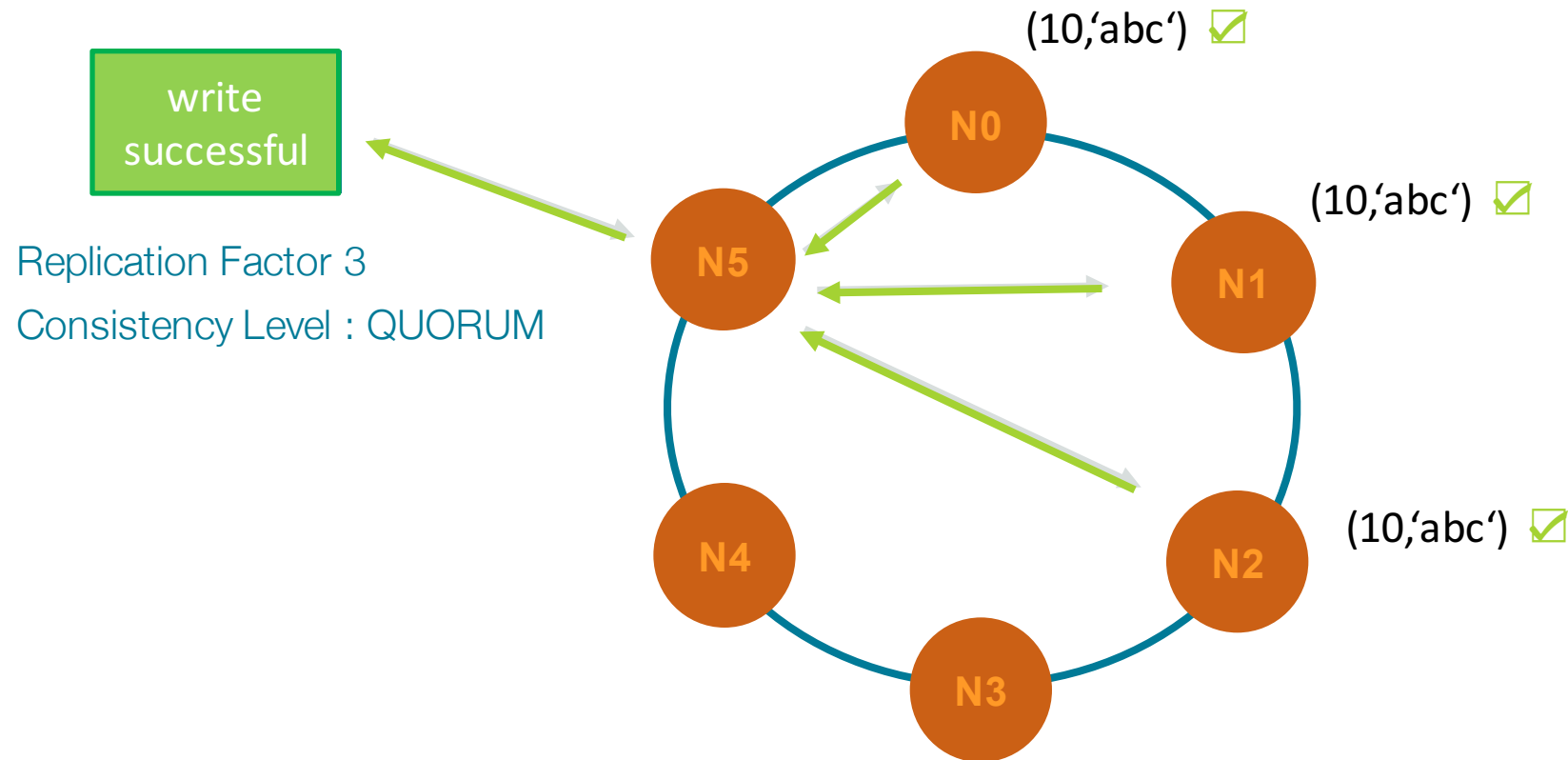
Read and write request handling

# Coordinator node

Incoming requests (read/write)
Coordinator node handles the request

request → N4
coordinator
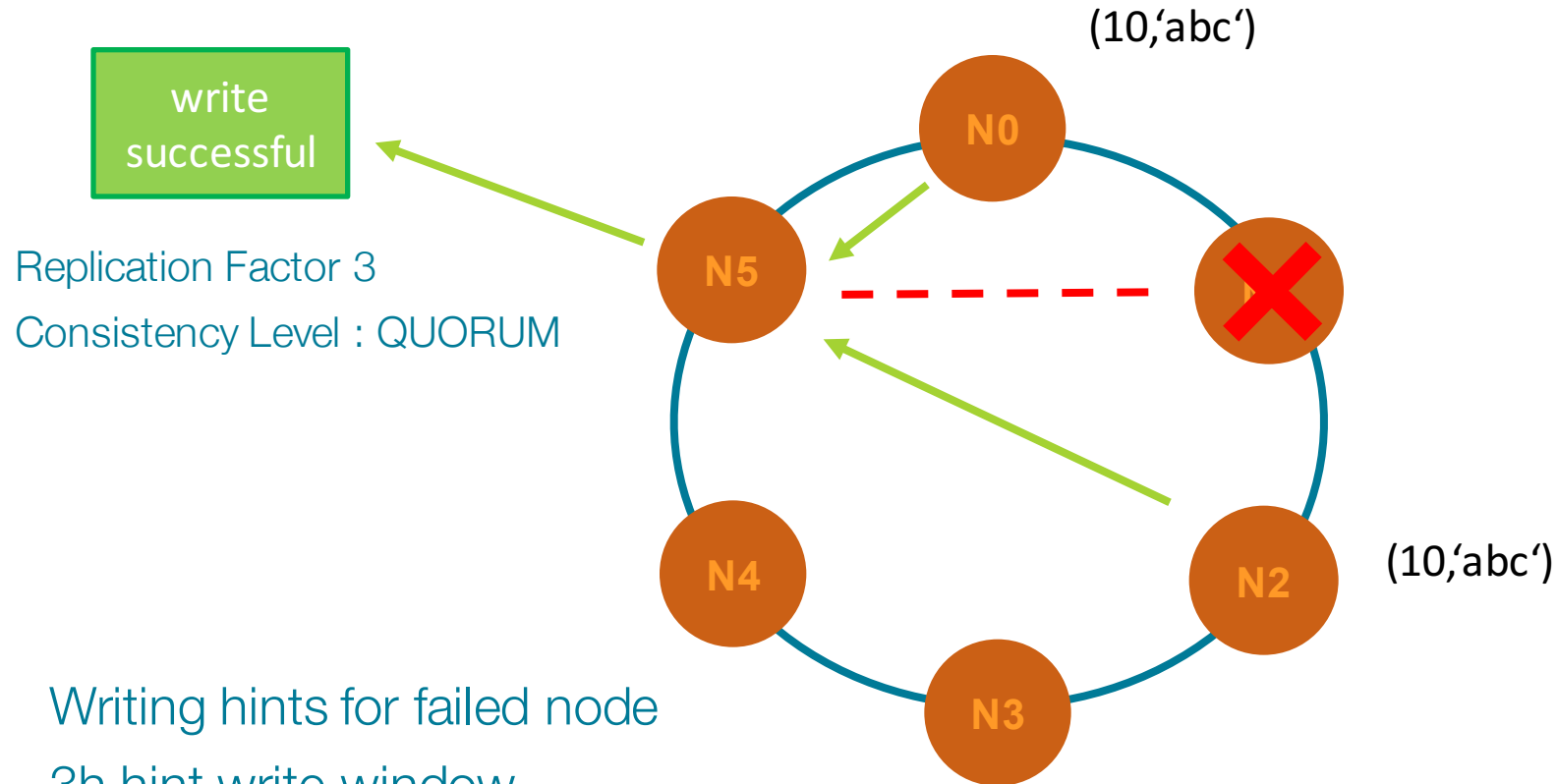
N0
N1
N2
N3
N4
N5
1
2
3

Every node can be coordinator → masterless

# Write request handling

write successful

Replication Factor 3

Consistency Level : QUORUM

(10,'abc') ☑

(10,'abc') ☑

(10,'abc') ☑

N0

N1

N2

N3

N4

N5

# Write request handling

write successful

Replication Factor 3

Consistency Level : QUORUM

(10,'abc')

N0

N5

N4

N2

N3

(10,'abc')
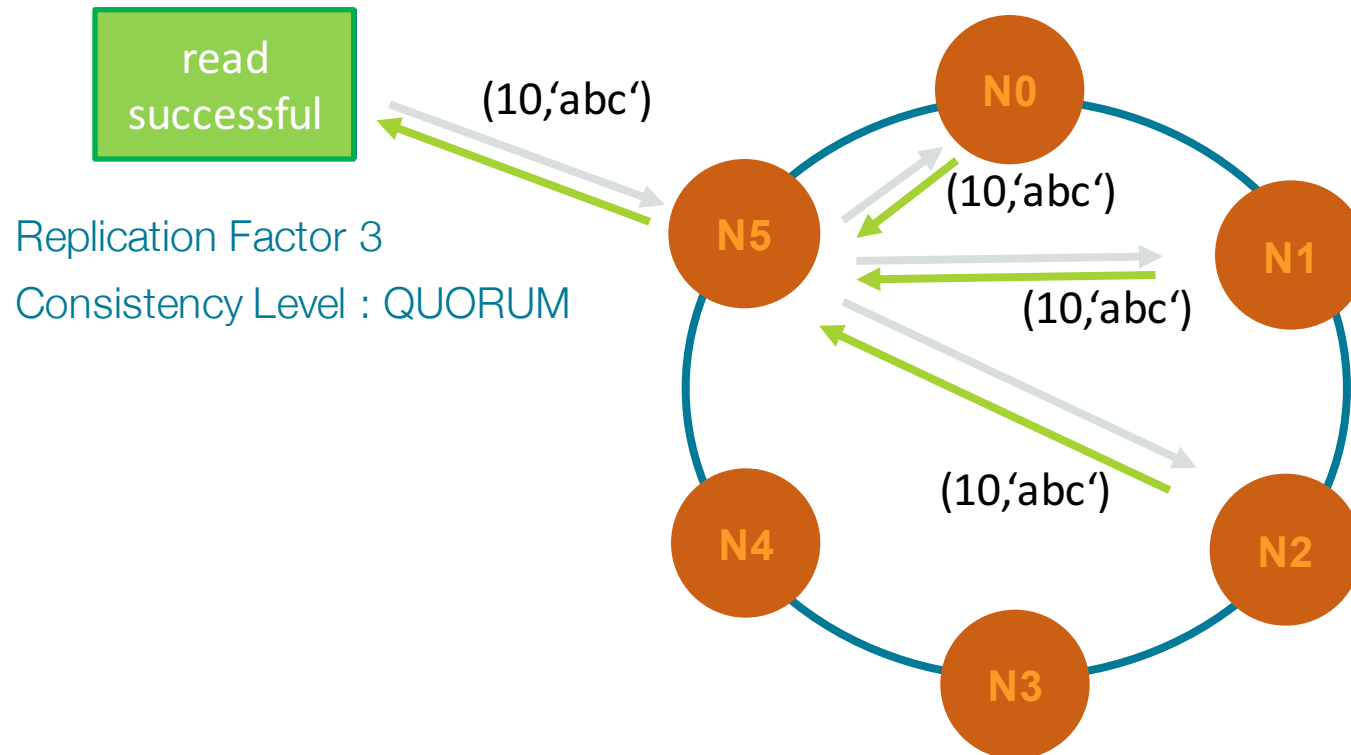
- Writing hints for failed node
- 3h hint write window
- CL: LOCAL_ALL will fail

# Write request handling



write failure

Replication Factor 3

Consistency Level : QUORUM

(10,'abc')

N0

N5

N4

N3

- quorum not met, failure
- CL: LOCAL_ONE will succeed

DATASTAX

# Read request handling



read
successful

(10,'abc')

Replication Factor 3

Consistency Level : QUORUM

(10,'abc')

(10,'abc')

(10,'abc')

N0

N1

N2

N3

N4

N5

DATASTAX

# Read request handling

read
successful

(10,'abc')

N0

(10,'abc')

Replication Factor 3

Consistency Level : QUORUM

N5

N1

(10,'abc')

N4

N3

- Reading LOCAL_QUORUM succeeds
  CL: LOCAL_ALL will fail

DATASTAX

# Read request handling

read
failure

(10,'abc')

**N0**

Replication Factor 3

Consistency Level : QUORUM

**N5**

(10,'abc')

❌

**N4**

❌

**N3**

- quorum not met, failure
- CL: LOCAL_ONE will succeed. See more...

DATASTAX

# Read request handling – Read Repair



Replication Factor 3

Consistency Level : QUORUM

- Last Write Win
- R+W > RF = immediate consistency
- Background vs. foreground Read Repair. See more...

- Background vs. Forground Read Repair
  - Compare digests
    - If any mismatch
  - re-request to same nodes (full data set)
    - compare full data sets, send update
    - block until out-of-date replicas respond
  - Return merged data set to the client

- Consistency Level
  - one, quorum, all
  - local vs. cluster wide

# Driver Code

```
Cluster cluster = Cluster.builder()
    .addContactPoint("127.0.0.1")
    .withLoadBalancingPolicy(new TokenAwarePolicy(DCAwareRoundRobinPolicy.builder()
        .withLocalDc("myLocalDC")
        .build())
    .build();


PreparedStatement prepared = session.prepare
    ( "insert into sales_by_customer(custid, salesdt) values (?, ?)");


BoundStatement bound = prepared.bind("1", "20170102");


session.execute(bound);   // Throws UnavailableException  If consistency doesn't met, downgrade is
                          possible with corresponding RetryPolicy. Read More…
```

DATASTAX

# Take away

- Data distribution (hash, tokens)

- Data replication (RF)

- All nodes are peer nodes , master less

- Background Read Repairs

- RetryPolicy in driver

DATASTAX

# Lab 2 : Hands-on DSE CQL

The power behind the moment. | DATASTAX

Vielen Dank!

# Eventuell Bootstrap, ReBalance, Num Tokens VNodes

- All nodes are peers

  – Including seed nodes

  – No master

  – Discovery through gossip

- Built-in replication

  – Simplify your architecture!

DATASTAX