# DataStax Enterprise Architecture
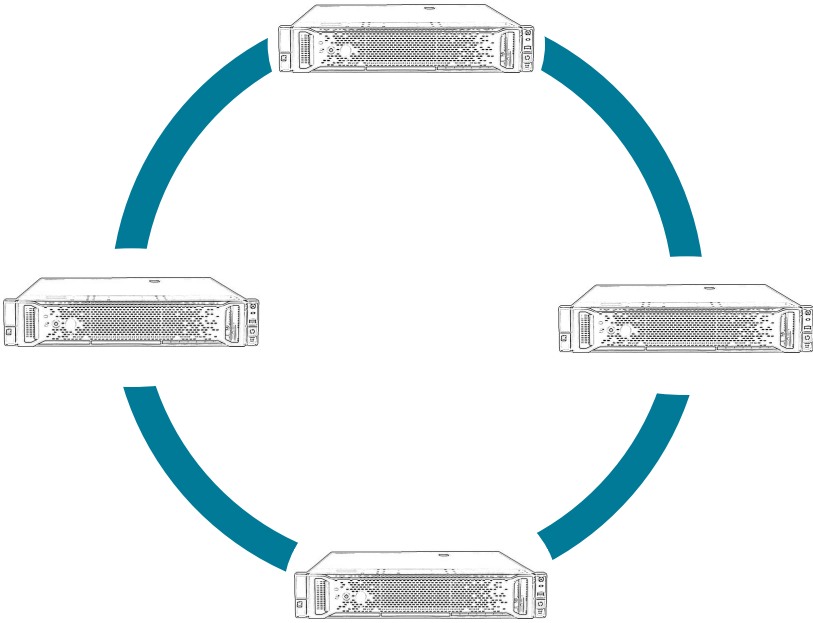
Negib Marhoul, Solution Engineer, DataStax

8. November 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

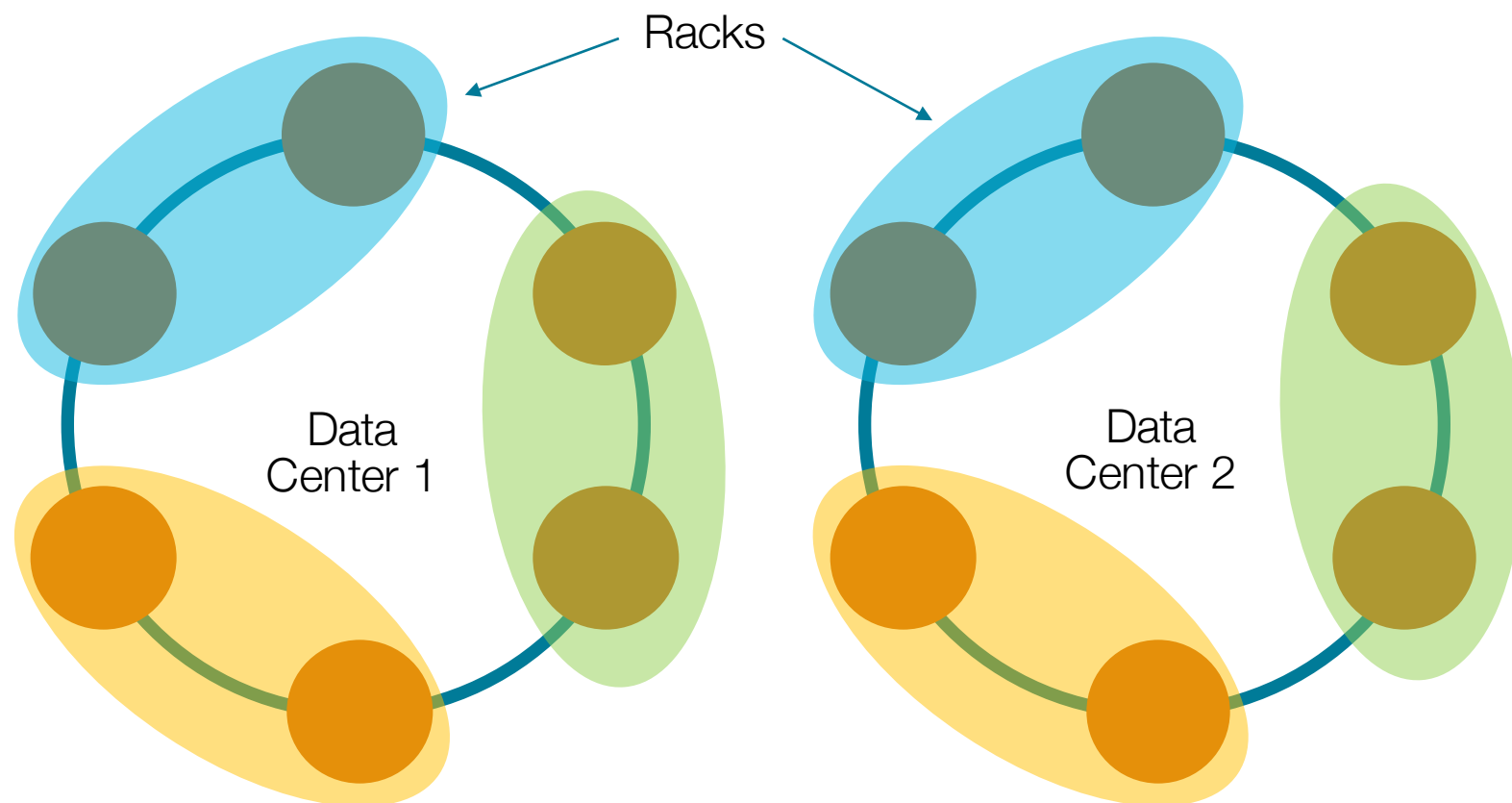DATASTAX

# Architecture

# Apache Cassandra™ Architecture

- Cluster layer
- Amazon DynamoDB paper
- masterless architecture


- Data-store layer
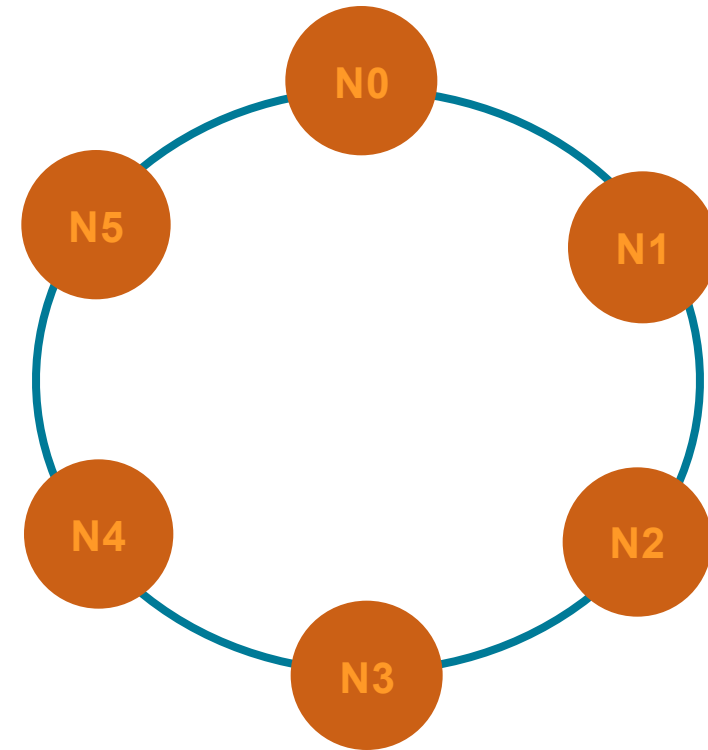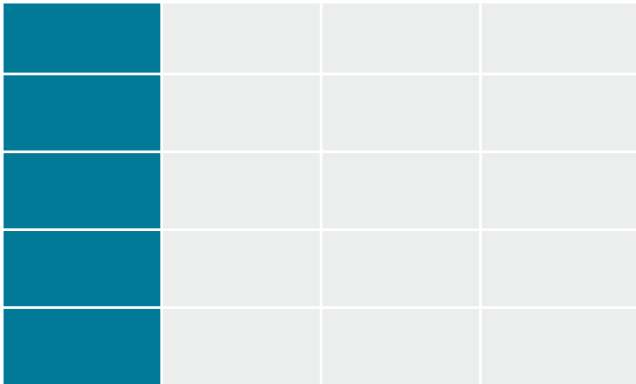- Google Big Table paper
- Columns/columns family

DATASTAX

- All nodes are peers
  - Including seed nodes
  - No master
  - Discovery through gossip

- Built-in replication
  - Simplify your architecture!
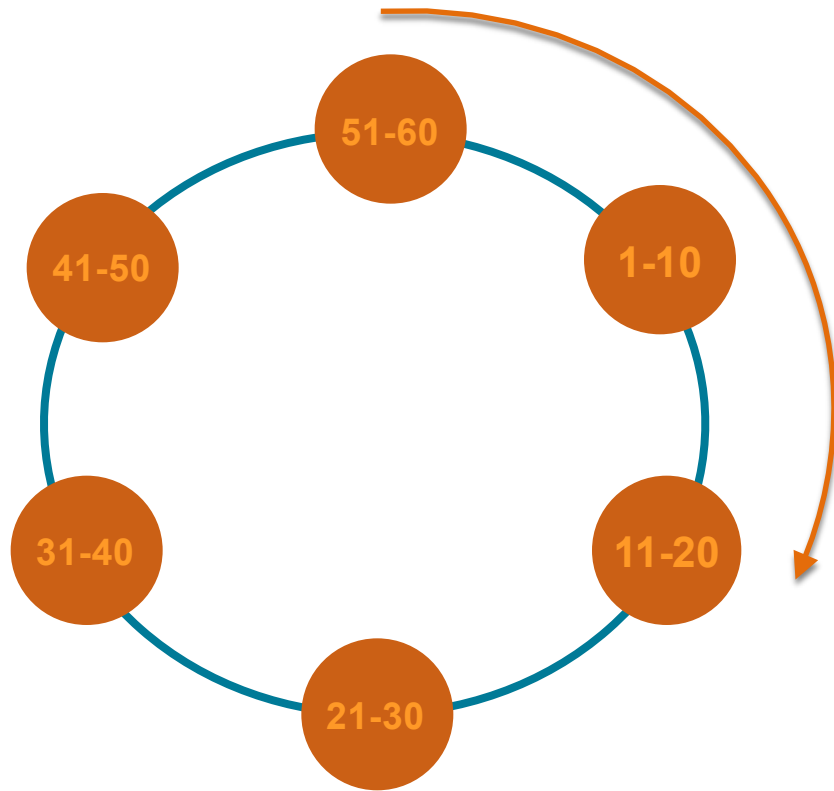
DATASTAX

# Tokens

Data is partitioned after its partition key

A unique token is allocated to a partition

Token = random hash of #partition

# Token Ranges



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

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
```
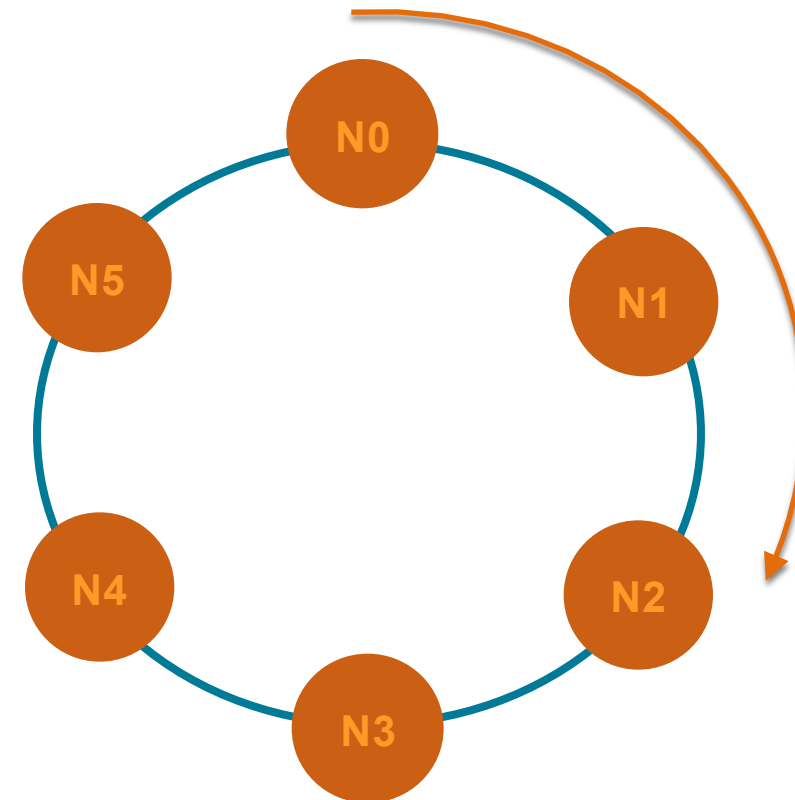
- Primary key
  - Partition key
  - Clustering columns

- Partitioner
  - Generates unique hash from partition key

- Replication strategy
  - Token hash determines starting point
  - Determines replica placement

DATASTAX

# Data Distribution

Token = hash of #partition → #node

| Token1 | custid 1 | | |
|--------|----------|--|--|
| Token2 | curstid 2 | | |
| Token3 | curstid 3 | | |
| Token4 | curstid 4 | | |
| Token4 | curstid 5 | | |



Data is evenly distributed and clock wise replicated

# Data Distribution

# Cassandra Query Language

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


CREATE TABLE retailer.sales_by_customer (
    custid int,
    salesdt text,
    comment text,
    discount double,                    Partition Key
    revenue double,
    PRIMARY KEY (custid, salesdt));


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

# Lab 1 : Accessing the cluster

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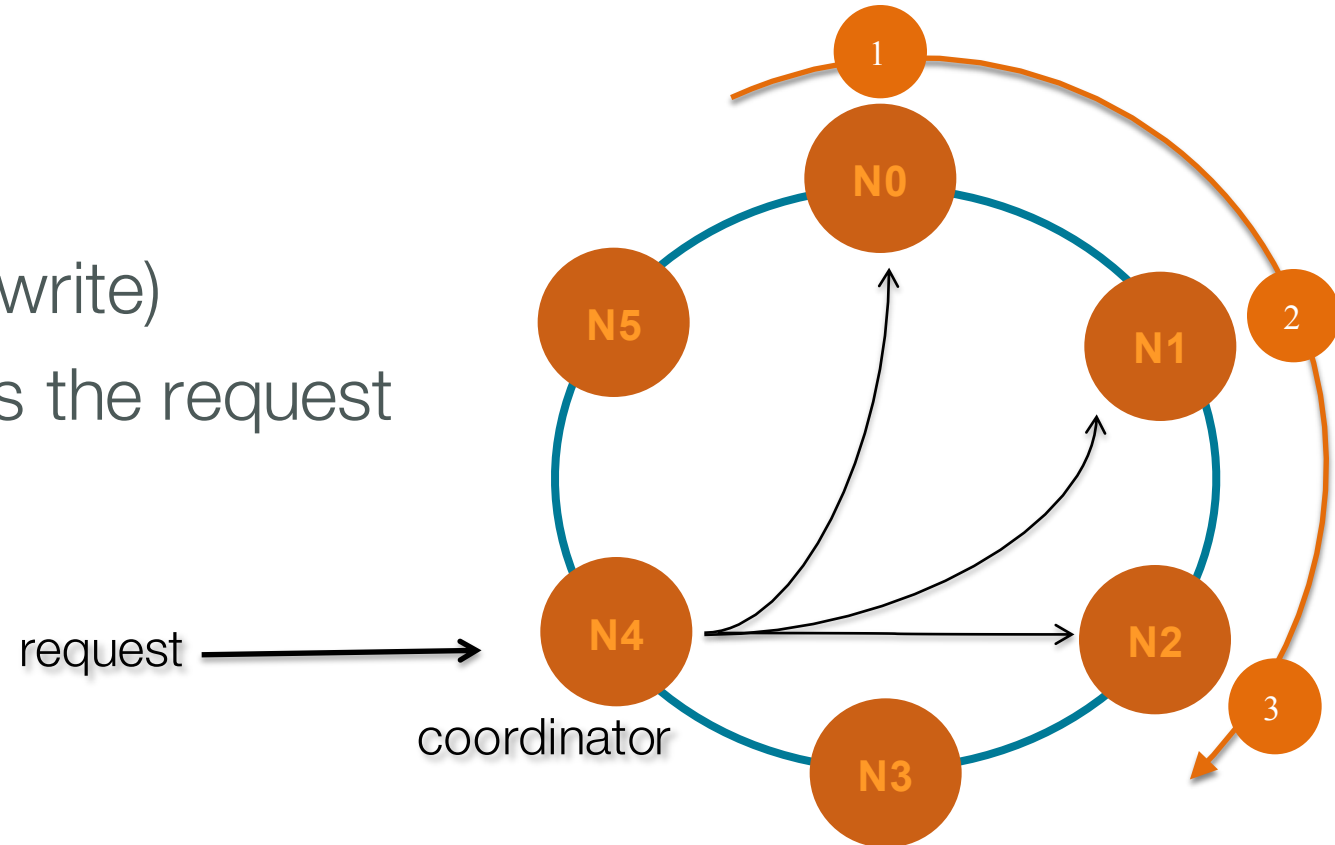
Tunable Consistency
Read and write request handling

# Coordinator node

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



request

coordinator

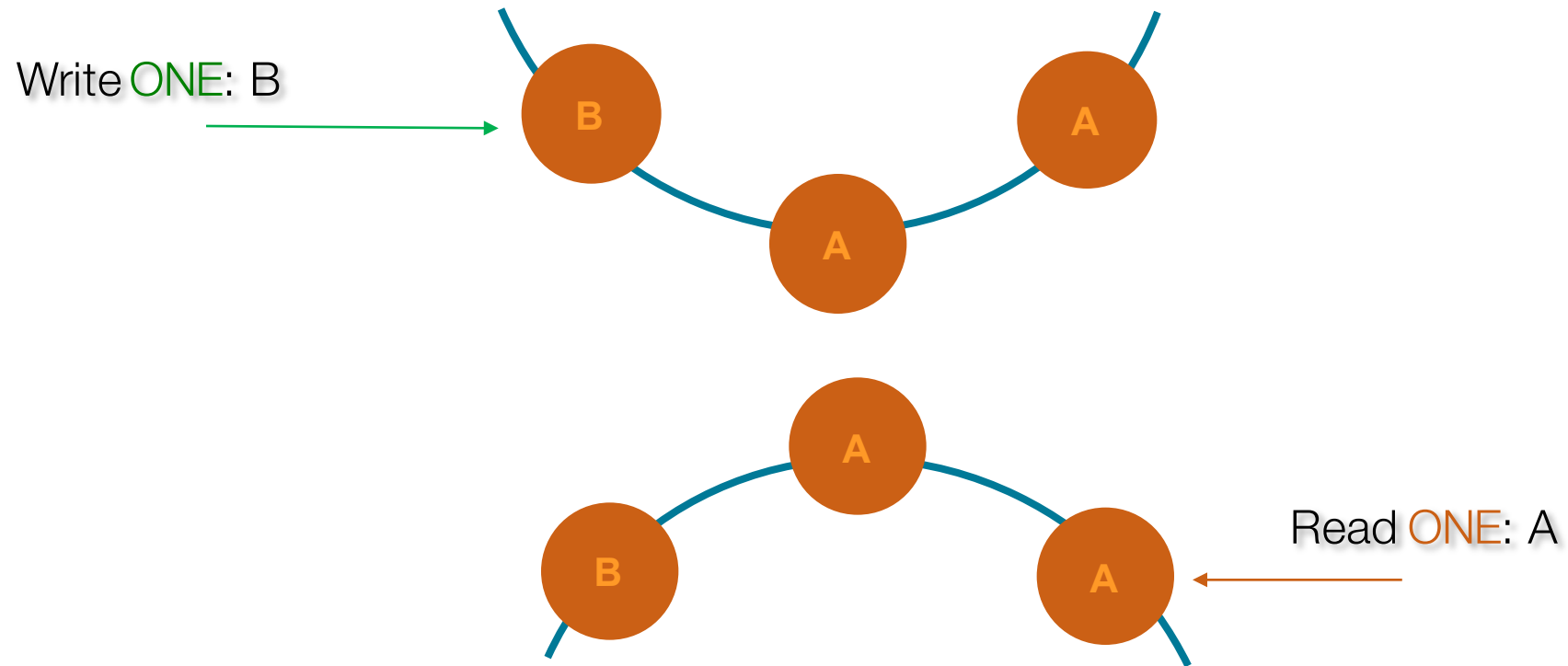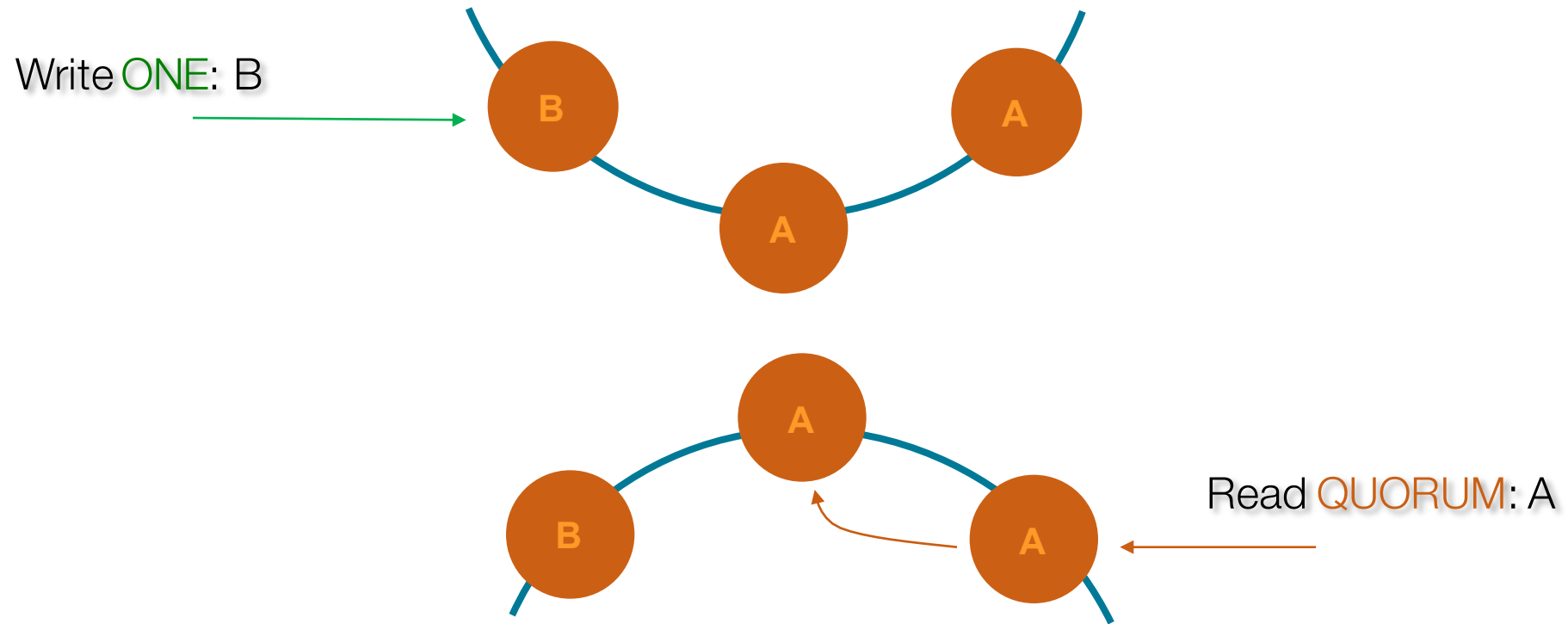Every node can be coordinator → masterless

# Consistency

Tunable at runtime

- ONE
- QUORUM (strict majority w.r.t. RF)
- ALL

Apply both to read & write

- RF = 3, Write ONE, Read ONE



Write ONE: B

Read ONE: A

- RF = 3, Write ONE, Read QUORUM



Write ONE: B

Read QUORUM: A

- RF = 3, Write ONE, Read ALL



Write ONE: B

Read ALL: B

- RF = 3, Write QUORUM, Read ONE



Write QUORUM: B

Read ONE: A

- RF = 3, Write QUORUM, Read QUORUM

Write QUORUM: B

B

A

B

B

B

A

Read QUORUM: B

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

# Consistency trade-off

**Latency**                                       **Consistency**

DATASTAX

# Consistency summary

ONE~Read~ + ONE~Write~

available for read/write even (N-1) replicas down

QUORUM~Read~ + QUORUM~Write~

available for read/write even 1+ replica down

DATASTAX

# 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')

(10,'abc')

N0

N5

N4

N3

N2

- Writing hints for failed node
- 3h hint write window
- CL: LOCAL_ALL will fail

DATASTAX

# Write request handling

(10,'abc')

write
failure

Replication Factor 3

Consistency Level : QUORUM

N0

N5

N4

N3

- quorum not met, failure

- CL: LOCAL_ONE will succeed

DATASTAX

# Read request handling

read successful

(10,'abc')

**N0**

(10,'abc')

**N5**

**N1**

Replication Factor 3

Consistency Level : QUORUM

(10,'abc')

(10,'abc')

**N4**

**N2**

**N3**

DATASTAX

# Read request handling

read
successful

(10,'abc')

N0

(10,'abc')

N5

N1

(10,'abc')

Replication Factor 3

Consistency Level : QUORUM

N4

N3

- Reading LOCAL_QUORUM succeeds
  CL: LOCAL_ALL will fail

DATASTAX

# Read request handling

read
failure

(10,'abc')

Replication Factor 3

Consistency Level : QUORUM

(10,'abc')

N0

N5

N4

N3

- quorum not met, failure

- CL: LOCAL_ONE will succeed. See more...

# Read request handling – Read Repair

Repair
(10,'def')

```
read
successful
```

(10,'def')

**N0**

(10,'abc',t1)

**N5**

**N1**

(10,'def',t2)

(10,'def',t2)

**N4**

**N2**

**N3**

Replication Factor 3

Consistency Level : QUORUM

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

DATASTAX

- 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

DATASTAX

# 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