



Understanding Cassandra's internal architecture

Apache Cassandra:
Core Concepts, Skills, and Tools

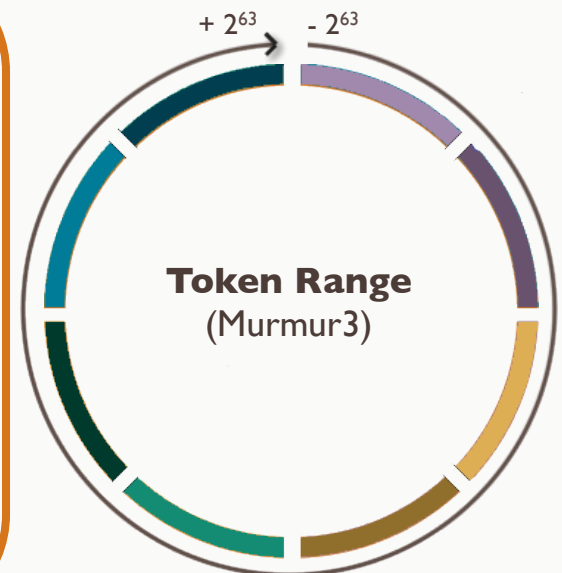
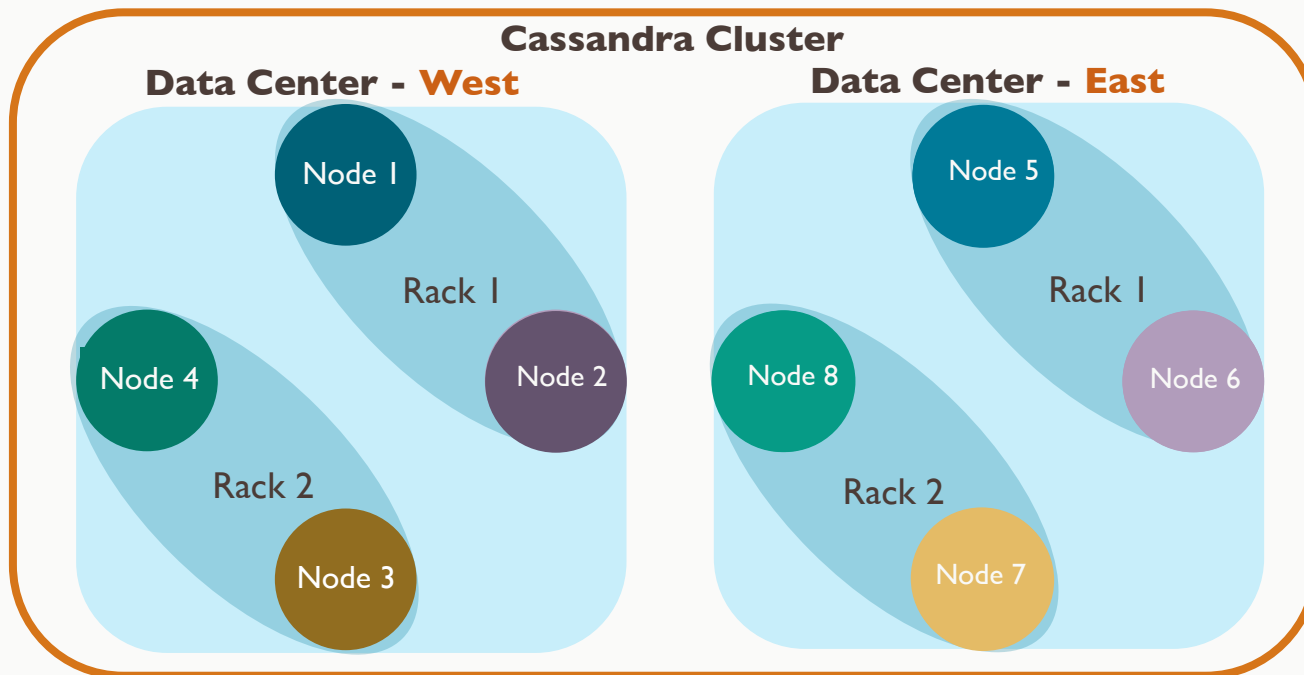
Leo Schuman, Joe Chu

Learning Objectives

- **Understand how requests are coordinated**
- Understand replication
- Understand and tune consistency
- Introduce anti-entropy operations
- Understand how nodes communicate
- Understand the *System* keyspace

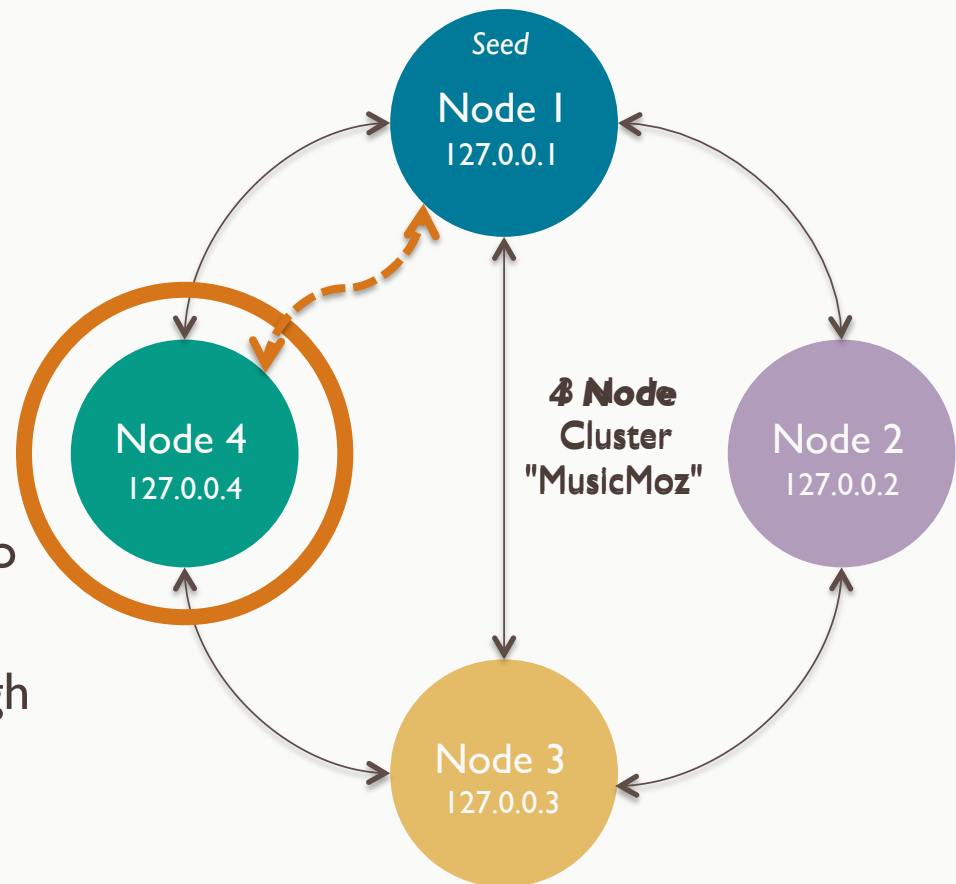
What is a cluster?

- A peer to peer set of nodes
 - **Node** – one Cassandra instance
 - **Rack** – a logical set of nodes
 - **Data Center** – a logical set of racks
 - **Cluster** – the full set of nodes which map to a single complete token ring



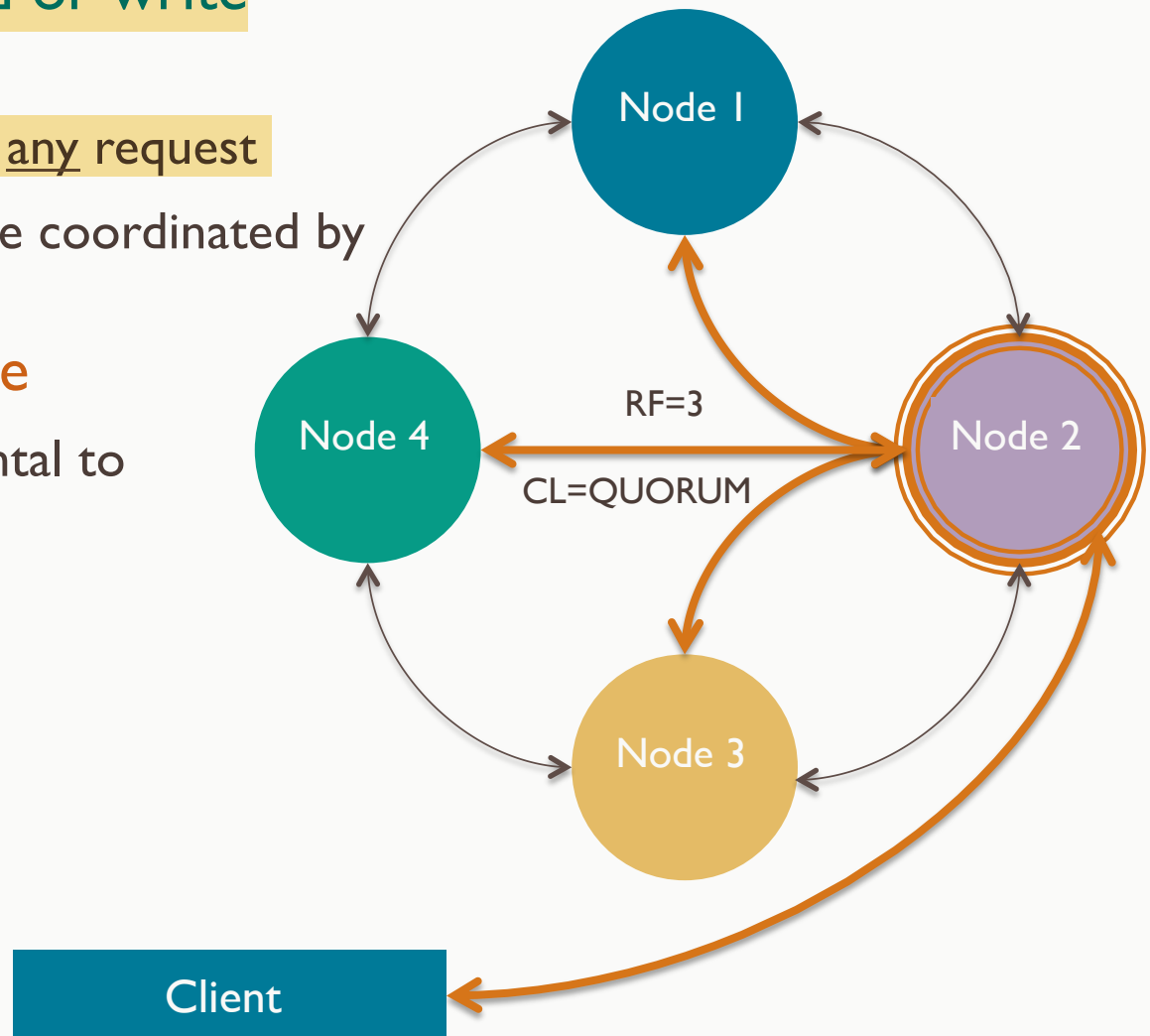
What is a cluster?

- Nodes join a cluster based on the configuration of their own *conf/cassandra.yaml* file
- Key settings include
 - **cluster_name** – shared name to logically distinguish a set of nodes
 - **seeds** – IP addresses of initial nodes for a new node to contact and discover the cluster topology (best practice to use the same two per data center)
 - **listen_address** – IP address through which this particular node communicates



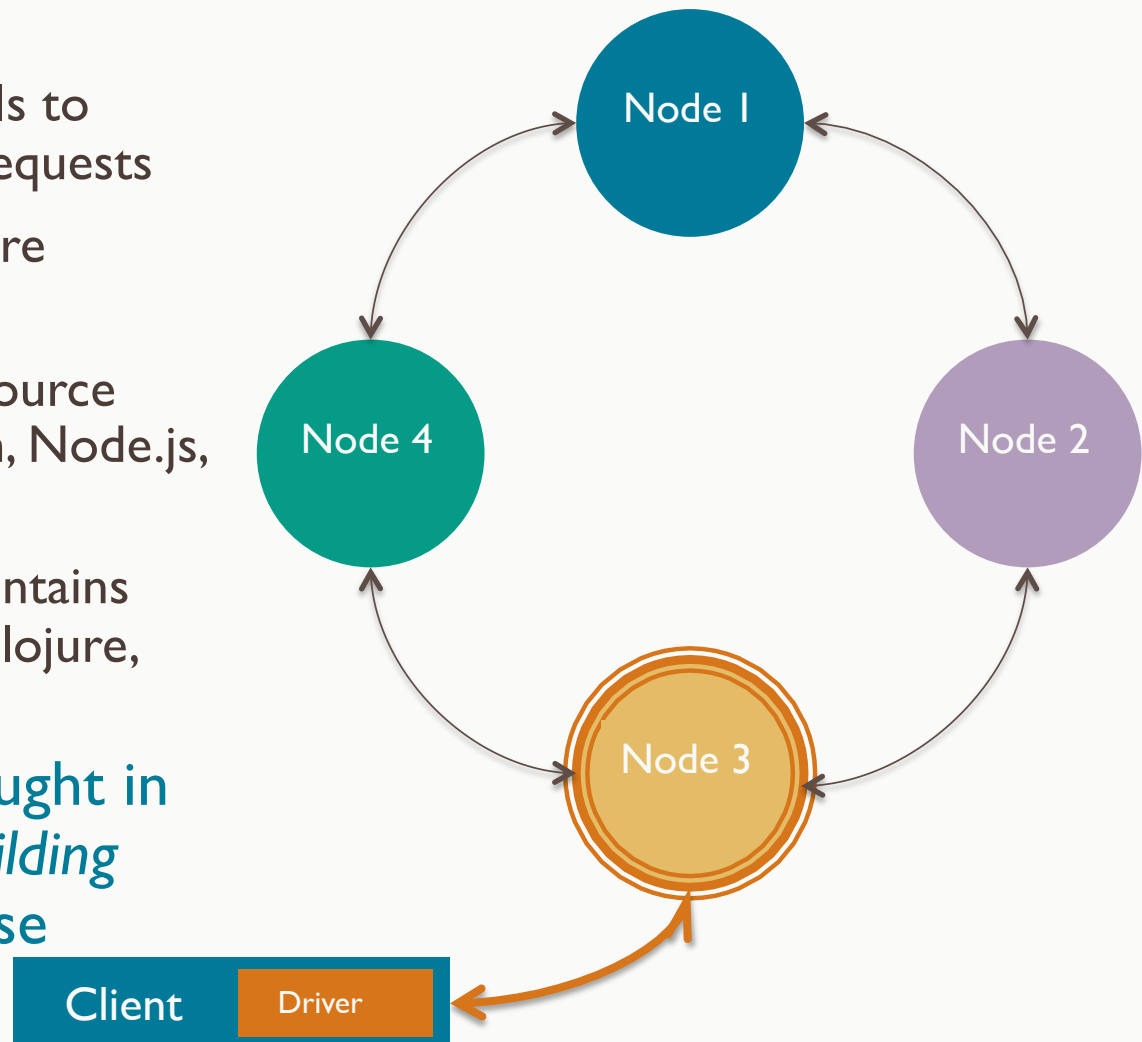
What is a coordinator?

- The *node* chosen by the client to receive a particular read or write request to its *cluster*
 - Any node can coordinate any request
 - Each client request may be coordinated by a different node
- No single point of failure
 - This principle is fundamental to Cassandra's architecture



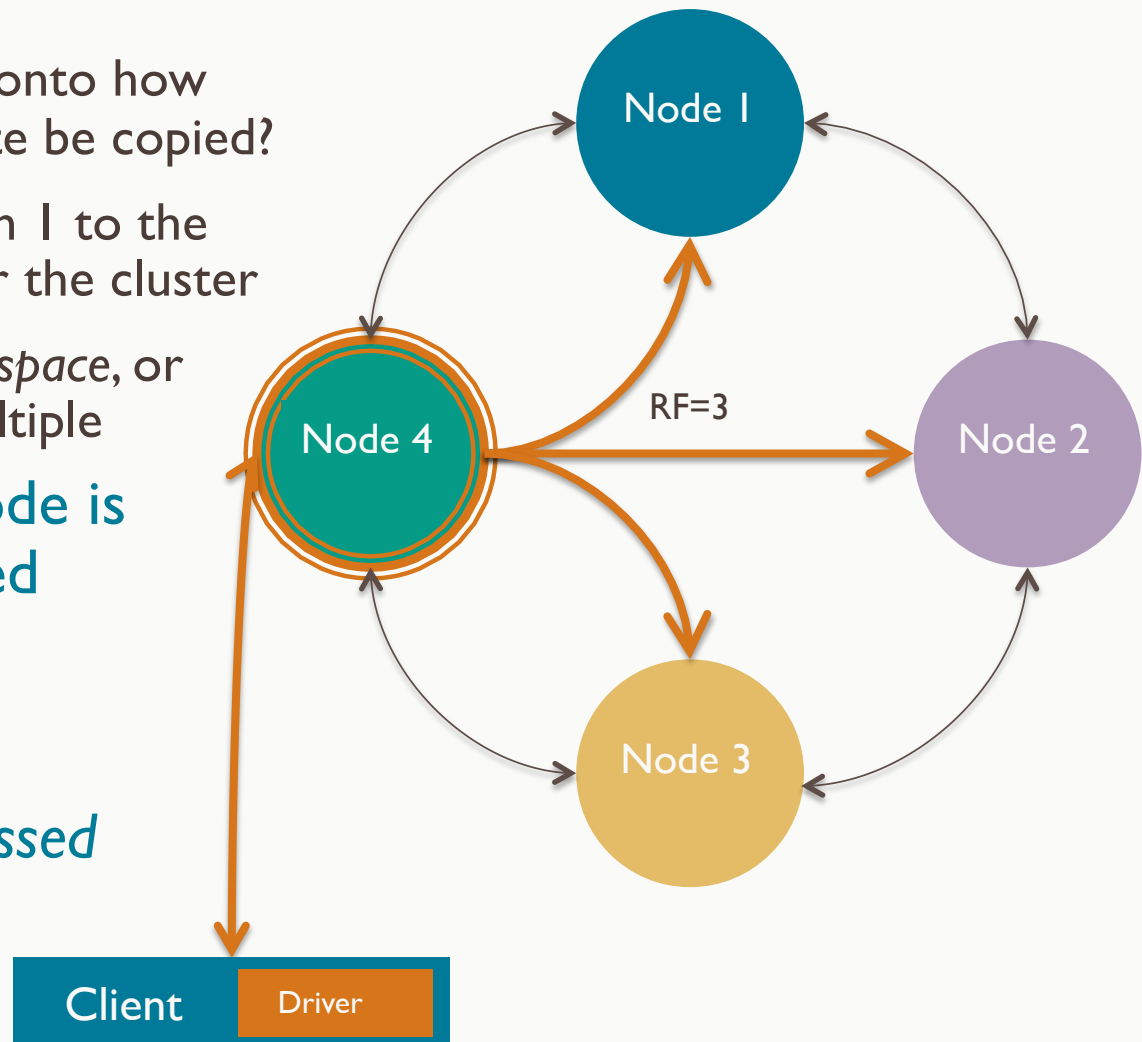
How are client requests coordinated?

- The *Cassandra driver* chooses the node to which each read or write request is sent
 - Client library providing APIs to manage client read/write requests
 - Default policy is TokenAware DCAWARE
 - DataStax maintains open source drivers for Java, C#, Python, Node.js, Ruby, C/C++ (beta)
 - Cassandra Community maintains drivers for PHP, Perl, Go, Clojure, Haskell, R, Scala
- Client development is taught in the *Apache Cassandra: Building Scalable Applications* course



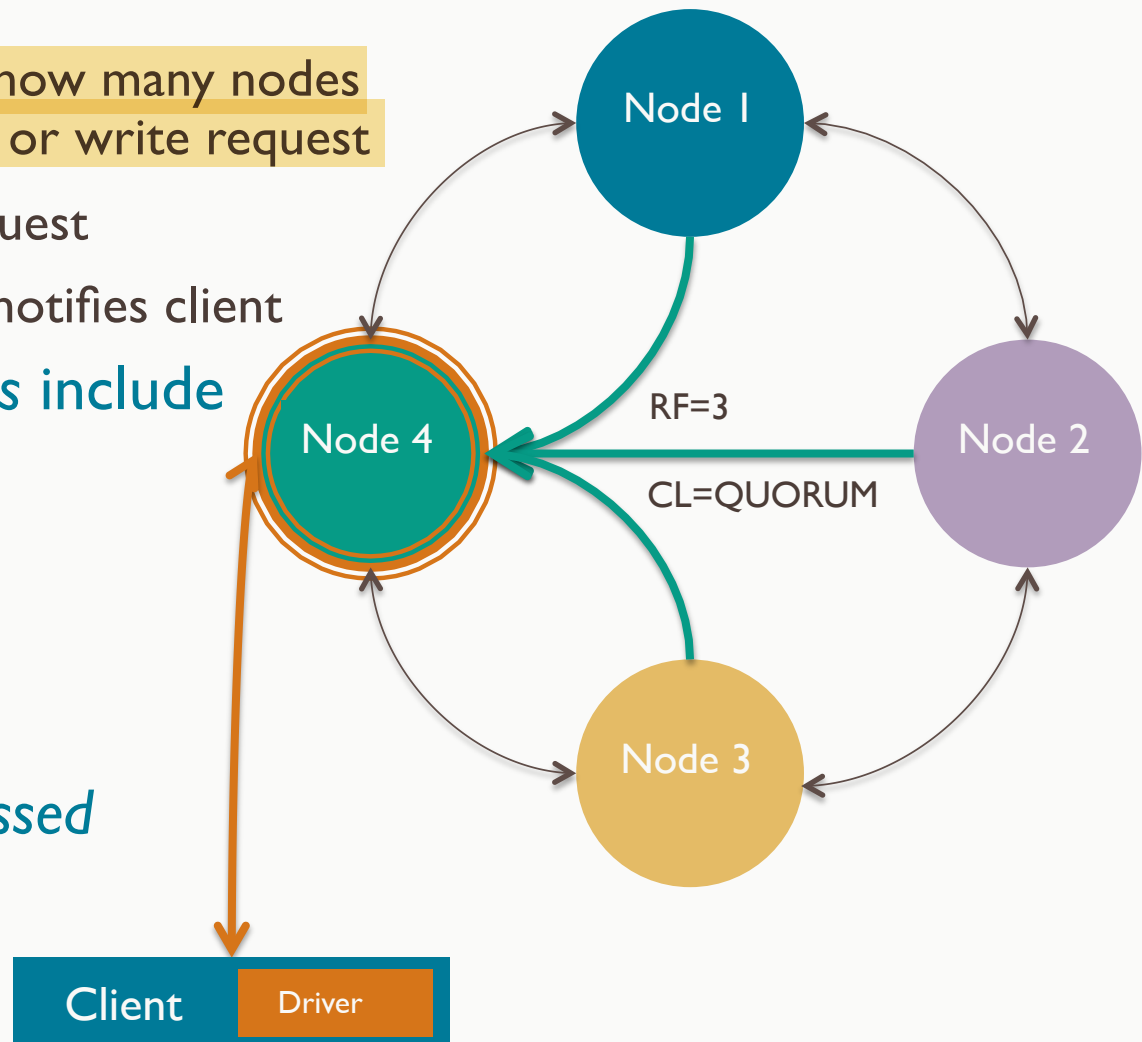
How are client requests coordinated?

- The coordinator manages the Replication Factor (RF)
 - Replication factor (RF) – onto how many nodes should a write be copied?
 - Possible values range from 1 to the total of planned nodes for the cluster
 - RF is set for an entire *keyspace*, or for each *data center*, if multiple
- Every write to every node is individually time-stamped
- Replication factor is discussed further ahead



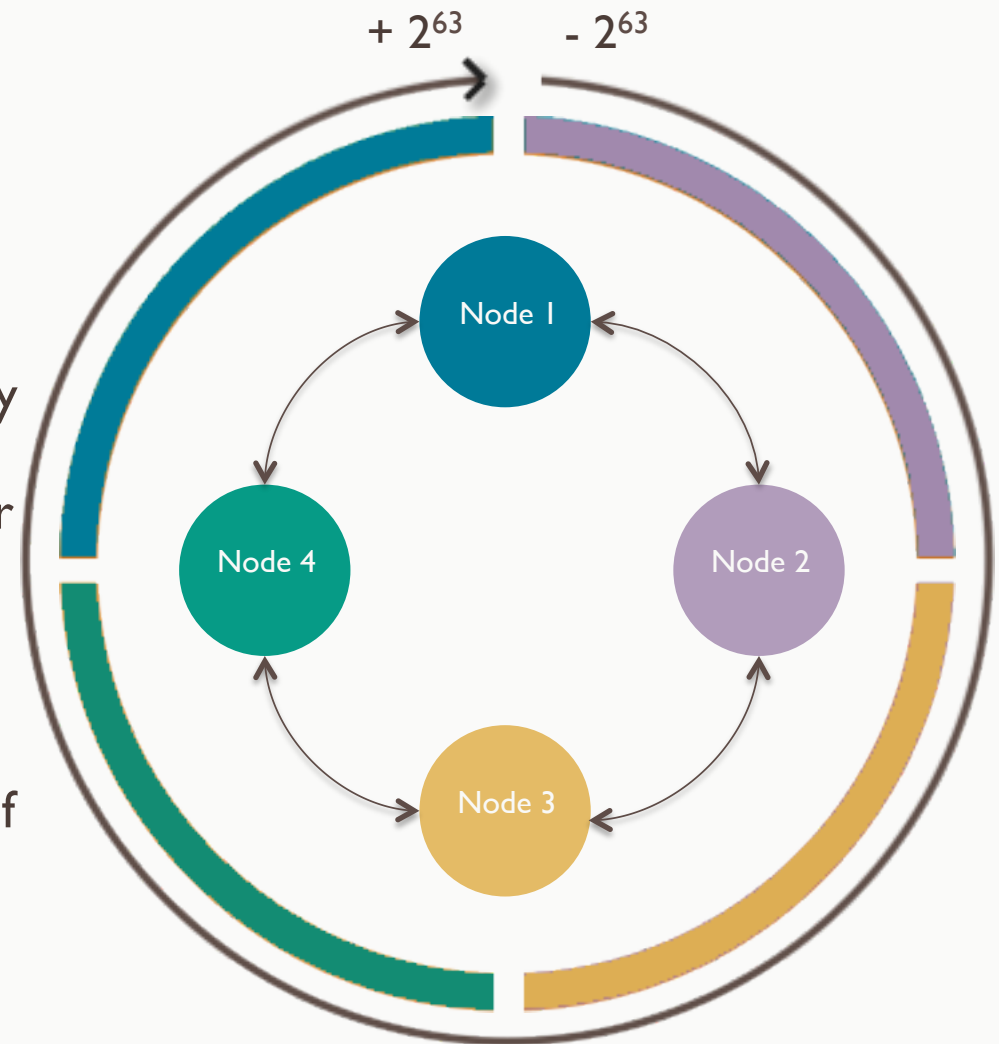
How are clients requests coordinated?

- The coordinator also applies the Consistency Level (CL)
 - Consistency level (CL) – how many nodes must acknowledge a read or write request
 - CL may vary for each request
 - On success, coordinator notifies client
- Possible consistency levels include
 - ANY
 - ONE
 - QUORUM ($RF / 2$) + 1
 - ALL
- Consistency Level is discussed further ahead



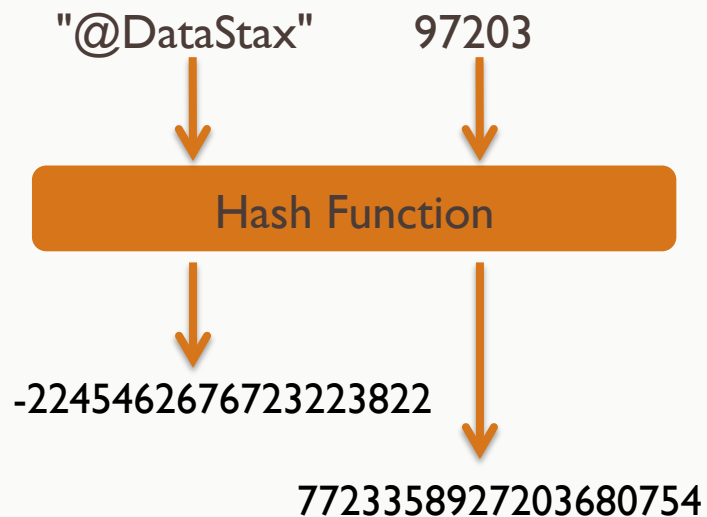
What is consistent hashing?

- Data is stored on *nodes* in *partitions*, each identified by a unique *token*
 - **Partition** – a storage location on a node (analogous to a "table row")
 - **Token** – integer value generated by a hashing algorithm, identifying a partition's location within a cluster
- The 2^{64} value *token range* for a cluster is used as a single ring
 - So, any partition in a cluster is locatable from one *consistent* set of hash values, regardless of its node
 - Specific token range varies by choice of *partitioner*
 - Partitioner options discussed ahead

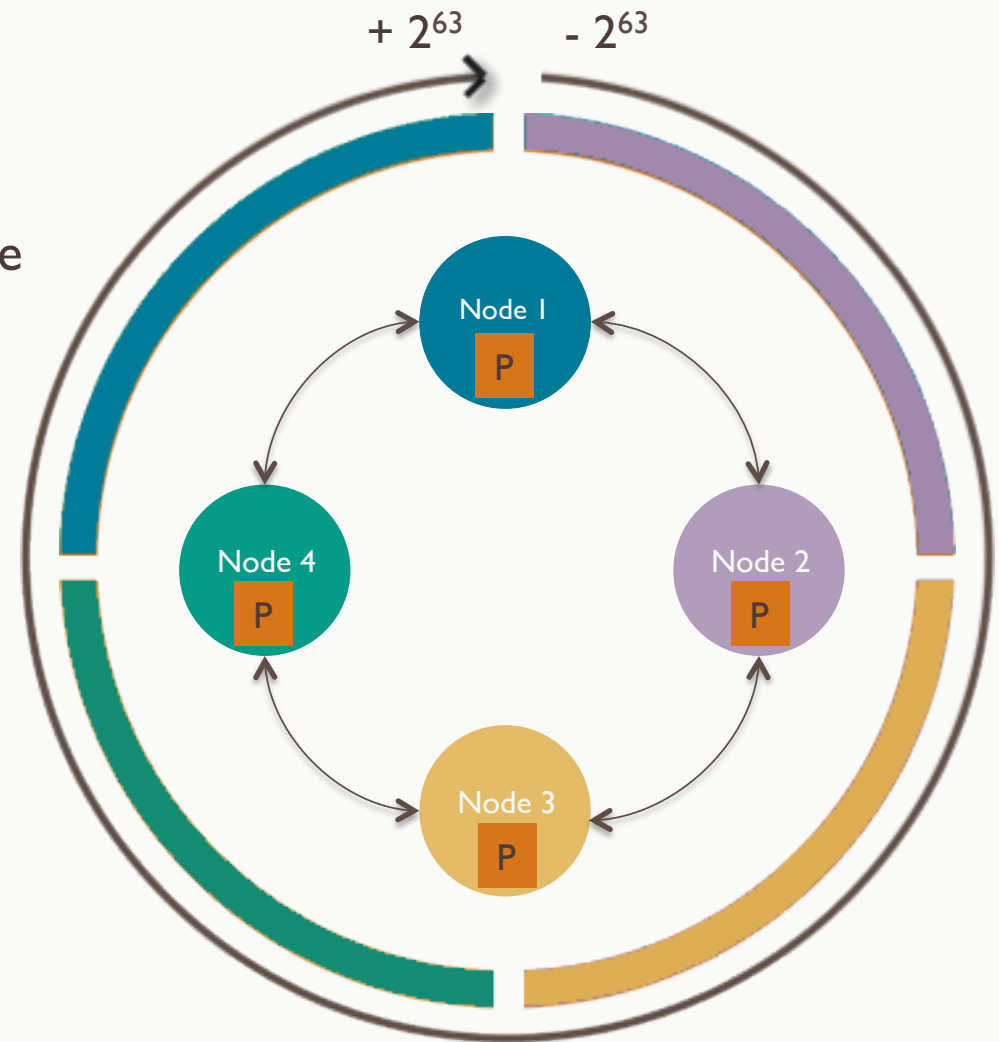


What is the partitioner?

- A system on each node which hashes tokens from designated values in rows being added
 - Hash function – converts a variable length value to a corresponding fixed length value

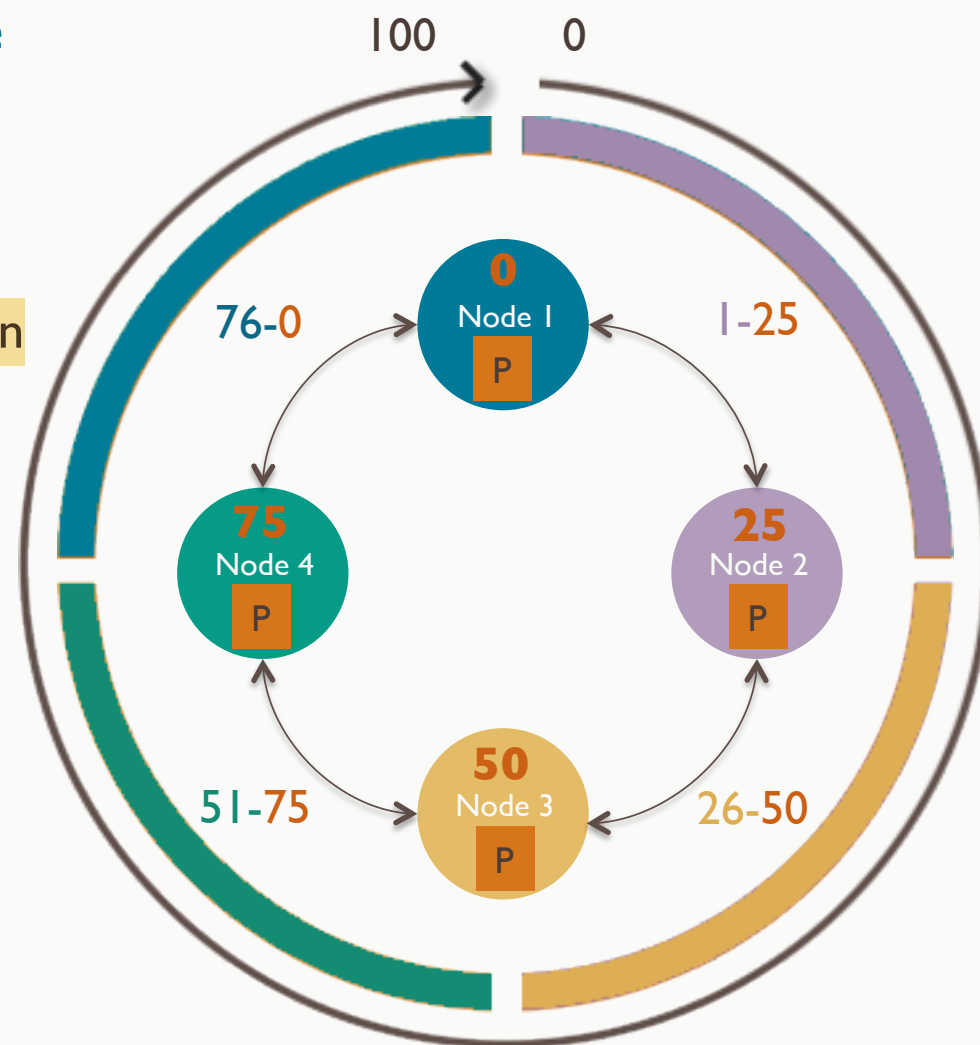


- Various partitioners available



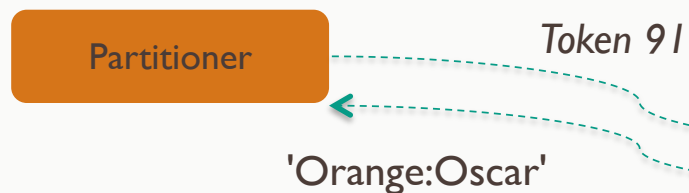
What is the partitioner?

- Imagine a 0 to 100 token range (instead of -2^{63} to $+2^{63}$)
 - Each node is assigned a token, just like each of its partitions
 - Node tokens are the highest value in the segment owned by that node*
- This segment is the *primary token range* of replicas owned by this node
 - Nodes also store *replicas* keyed to tokens outside this range ("secondary range")



How does a partitioner work?

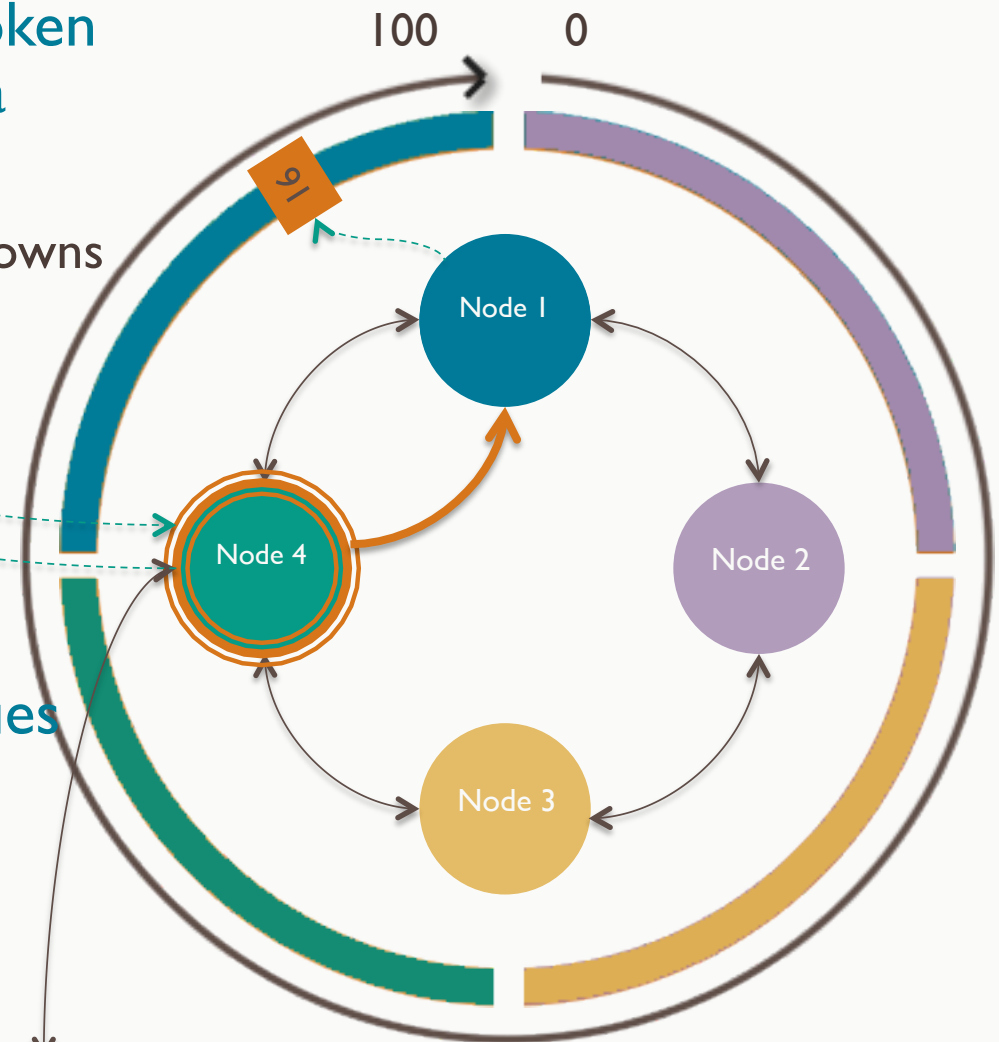
- A node's *partitioner* hashes a token from the *partition key* value of a write request
 - First replica written to node that owns the primary range for this token



- The *primary key* of a table determines its *partition key* values

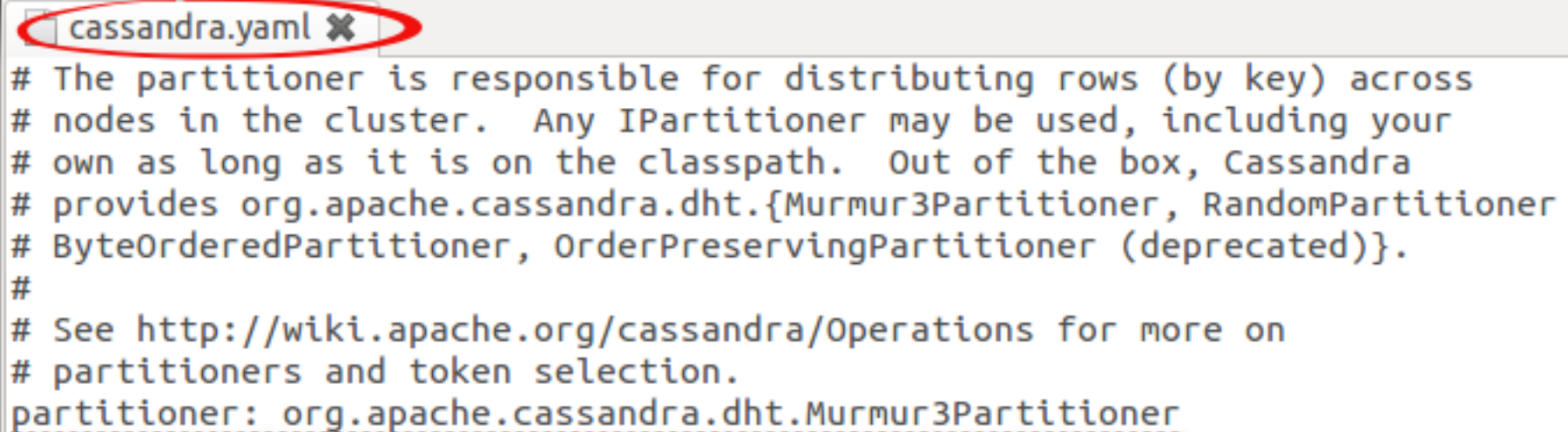
```
CREATE TABLE Users (
  firstname text, lastname text, level text,
  PRIMARY KEY ((lastname, firstname))
);
```

```
INSERT INTO Users (firstname, lastname, level)
VALUES ('Oscar', 'Orange', 42);
```



What partitioners does Cassandra offer?

- Cassandra offers three partitioners
 - **Murmur3Partitioner** (default) – uniform distribution based on Murmur3 hash
 - **RandomPartitioner** – uniform distribution based on MD5 hash
 - **ByteOrderedPartitioner** (legacy only) – lexical distribution based on key bytes
- **Murmur3Partitioner** is the default and best practice
- The partitioner is configured in the *cassandra.yaml* file
 - Must be the same across all nodes in the cluster

A screenshot of a text editor window titled 'cassandra.yaml' with a red circle around the title bar. The window displays the configuration for the partitioner in the cassandra.yaml file. The text is as follows:

```
# The partitioner is responsible for distributing rows (by key) across
# nodes in the cluster. Any IPartitioner may be used, including your
# own as long as it is on the classpath. Out of the box, Cassandra
# provides org.apache.cassandra.dht.{Murmur3Partitioner, RandomPartitioner
# ByteOrderedPartitioner, OrderPreservingPartitioner (deprecated)}.
#
# See http://wiki.apache.org/cassandra/Operations for more on
# partitioners and token selection.
partitioner: org.apache.cassandra.dht.Murmur3Partitioner
```

The line `partitioner: org.apache.cassandra.dht.Murmur3Partitioner` is underlined with a red dotted line.

Learning Objectives

- Understand how requests are coordinated
- **Understand replication**
- Understand and tune consistency
- Introduce anti-entropy operations
- Understand how nodes communicate
- Understand the *System* keyspace

How does the keyspace impact replication?

- Replication factor is configured when a keyspace is created
 - **SimpleStrategy (learning use only)** – one factor for entire cluster
 - assigned as "replication_factor"

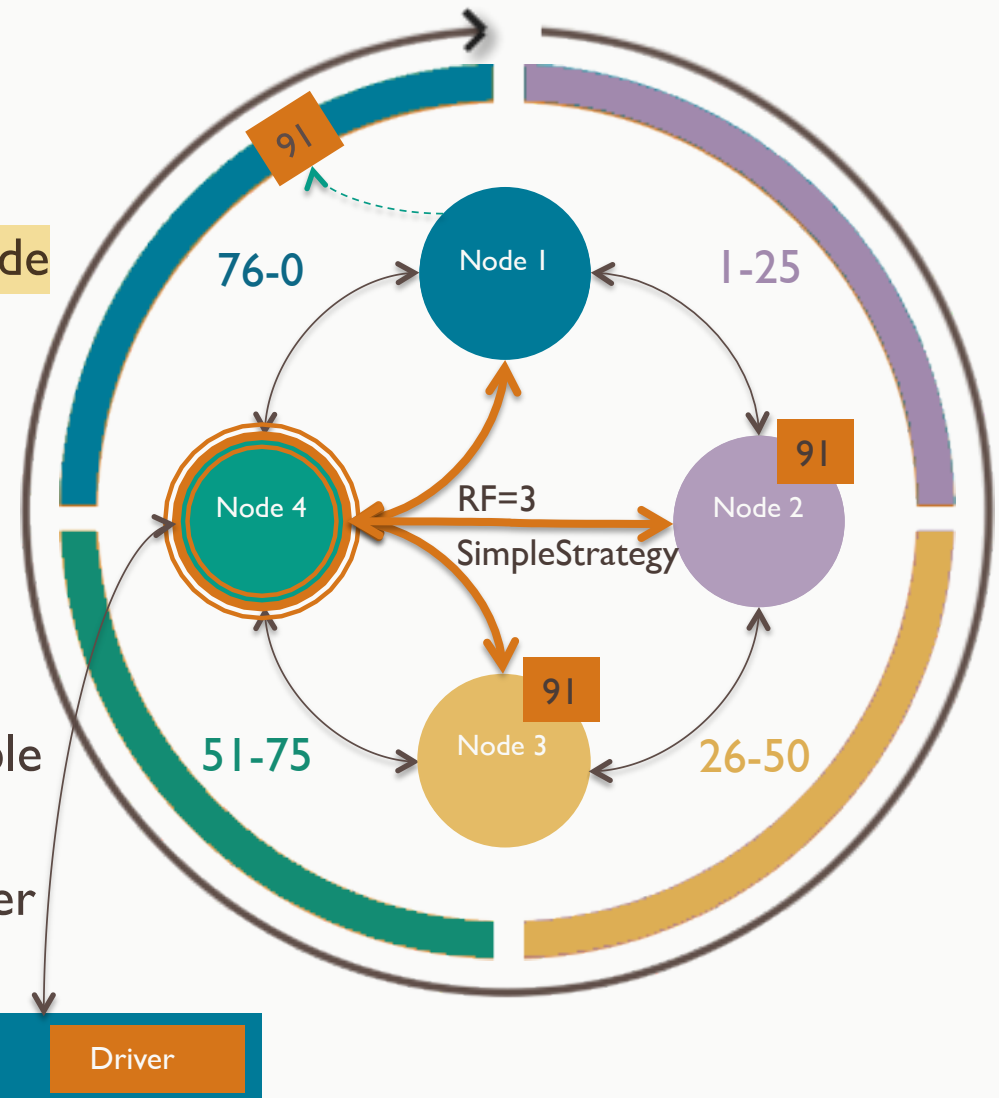
```
CREATE KEYSPACE simple-demo  
WITH REPLICATION =  
{'class':'SimpleStrategy',  
  'replication_factor':2}
```

- **NetworkTopologyStrategy** – separate factor for each data center in cluster
 - assigned by data center id (as also used in *cassandra-rackdc.properties*)

```
CREATE KEYSPACE simple-demo  
WITH REPLICATION =  
{'class':'NetworkTopologyStrategy',  
  'dc-east':2, 'dc-west':3}
```

How does a coordinator forward write requests?

- The target table's *keyspace* determines
 - Replication factor – how many replicas to make of each partition
 - Replication strategy – on which node should each replica be placed
- All partitions are "replicas", there are no "originals"
 - First replica – placed on the node owning its token's primary range
 - Closest node – replicas placed in same rack and data center, if possible
 - (Subsequent) replicas (if $RF > 1$) – placed in "secondary range" of other nodes, per the replication strategy

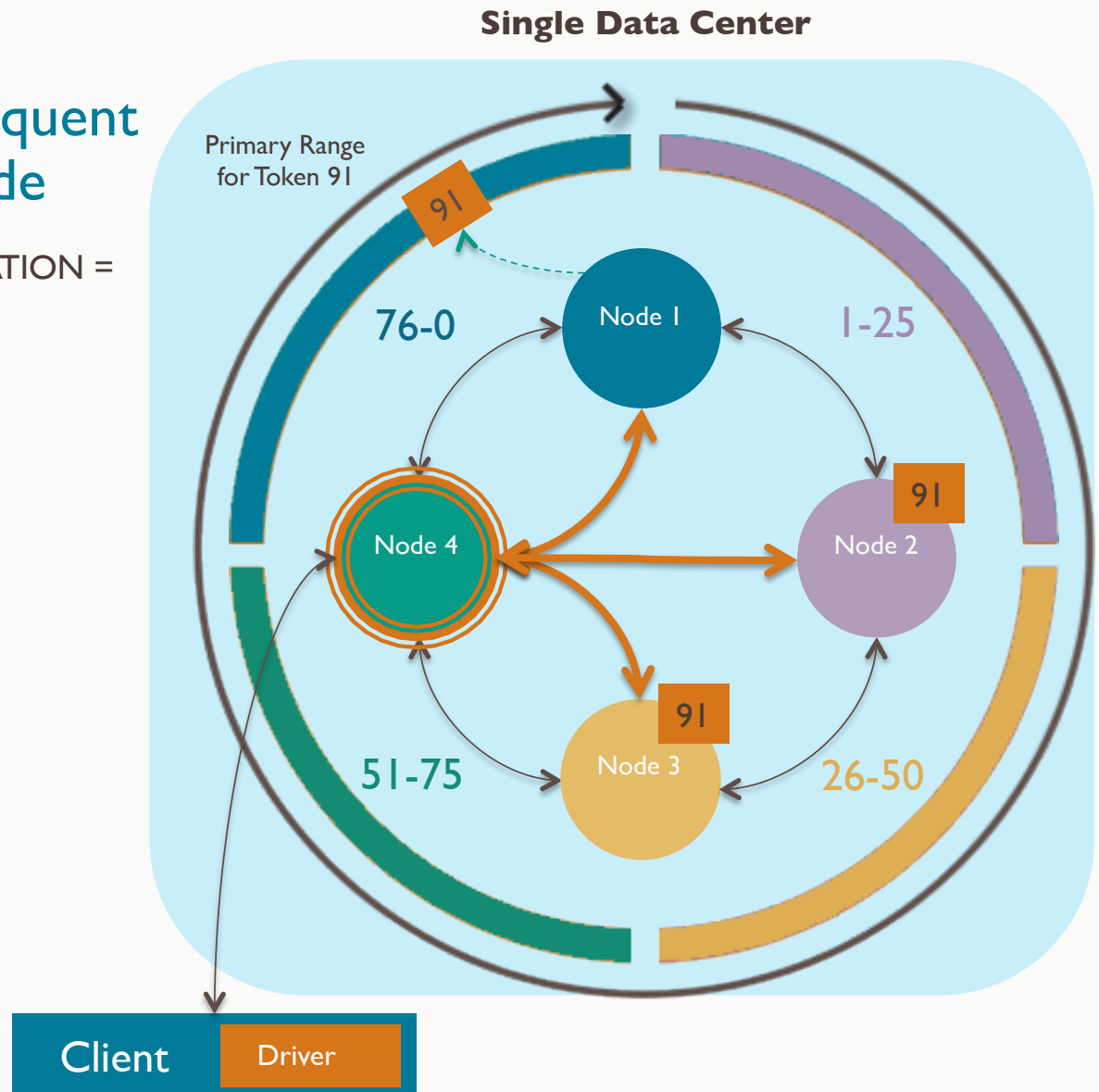


How is data replicated among nodes?

- **SimpleStrategy** – create replicas on nodes subsequent to the *primary range* node

```
CREATE KEYSPACE demo WITH REPLICATION =
{'class': 'SimpleStrategy',
'replication_factor': 3}
```

- *replication factor* of 3 is a recommended minimum



Exercise I: Create a keyspace

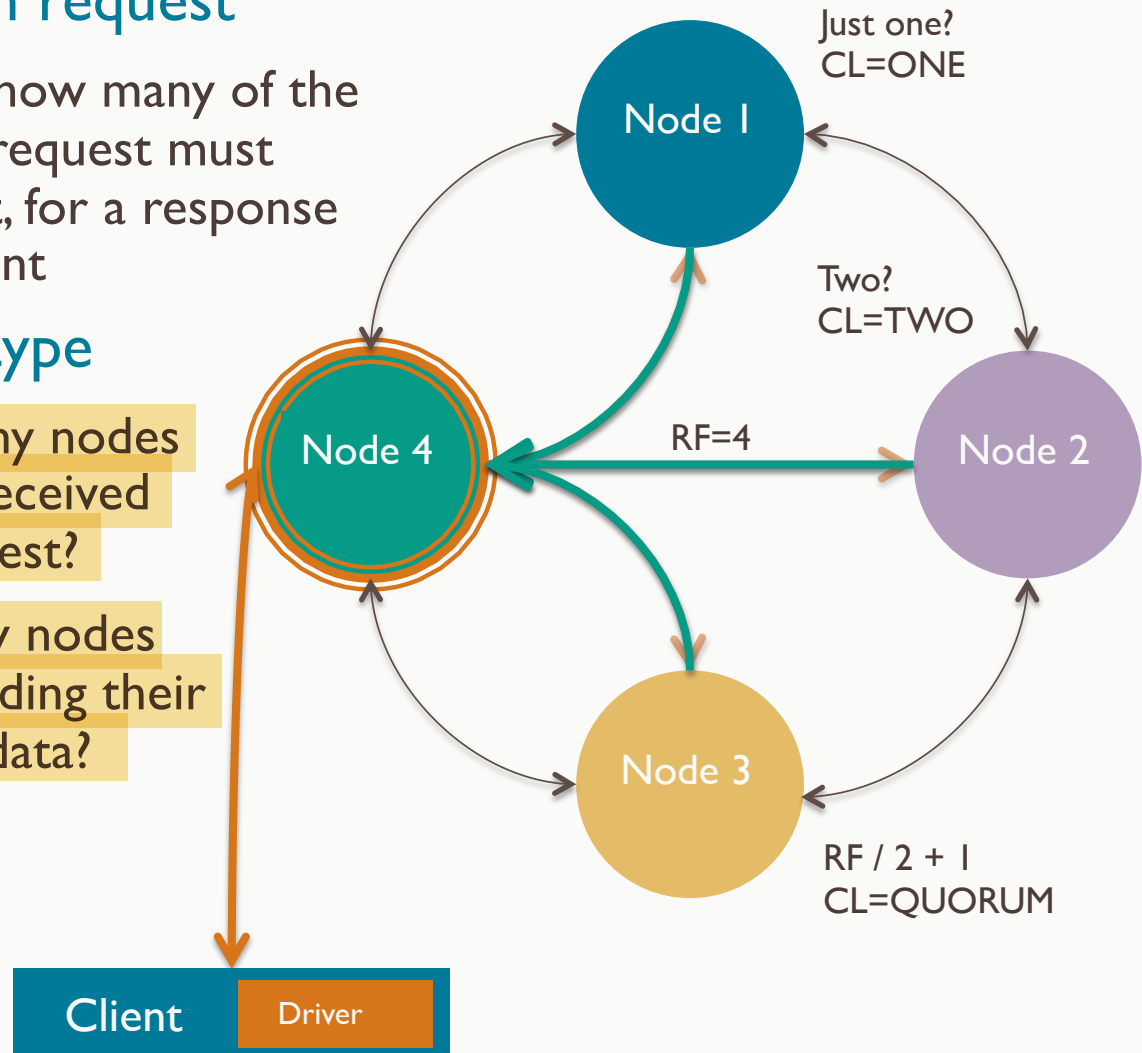


Learning Objectives

- Understand how requests are coordinated
- Understand replication
- **Understand and tune consistency**
- Introduce anti-entropy operations
- Understand how nodes communicate
- Understand the *System* keyspace

What is consistency?

- The *partition key* determines which nodes are sent any given request
 - **Consistency Level** – sets how many of the nodes to be sent a given request must acknowledge that request, for a response to be returned to the client
- The meaning varies by type
 - **Write request** – how many nodes must acknowledge they received and wrote the write request?
 - **Read request** – how many nodes must acknowledge by sending their most recent copy of the data?



What consistency levels are available?

Name	Description	Usage
ANY (writes only)	Write to any node, and store <i>hinted handoff</i> if all nodes are down.	Highest availability and lowest consistency (writes)
ALL	Check all nodes. Fail if any is down.	Highest consistency and lowest availability
ONE (TWO,THREE)	Check closest node to coordinator.	Highest availability and lowest consistency (reads)
QUORUM	Check quorum of available nodes.	Balanced consistency and availability
LOCAL_ONE	Check closest node to coordinator, in the local data center only.	Highest availability, lowest consistency, and no cross-data-center traffic
LOCAL_QUORUM	Check quorum of available nodes, in the local data center only.	Balanced consistency and availability, with no cross-data-center traffic
EACH_QUORUM	Only valid for writes. Check quorum of available nodes, in <u>each</u> data center of the cluster.	Balanced consistency and availability, with cross-data-center consistency
SERIAL	Conditional write to quorum of nodes. Read current state with no change.	Used to support linearizable consistency for lightweight transactions
LOCAL_SERIAL	Conditional write to quorum of nodes in local data center.	Used to support linearizable consistency for lightweight transactions

How do you set consistency per request?

- The default *consistency level* for all requests is **ONE**
 - In *cqlsh*, the *CONSISTENCY* command modifies this value for all subsequent requests during the same *cqlsh* session
 - In *client drivers*, a *ConsistencyLevel* constant is passed as part of each request

```
Cassandra x dstaining@DST: /home/dsc-c
dstaining@DST:/home/dsc-cassandra-2.0.5/bin$ ./cqlsh
Connected to Test Cluster at localhost:9160.
[cqlsh 4.1.1 | Cassandra 2.0.5 | CQL spec 3.1.1 | Thrift protocol 19.39.0]
Use HELP for help.
cqlsh> USE demo;
cqlsh:demo> CONSISTENCY;
Current consistency level is ONE.
cqlsh:demo> CONSISTENCY QUORUM;
Consistency level set to QUORUM.
cqlsh:demo> CONSISTENCY ANY;
Consistency level set to ANY.
cqlsh:demo> CONSISTENCY ALL;
Consistency level set to ALL.
cqlsh:demo> █
```

What is immediate vs. eventual consistency?

- For any given read, how likely is it the data may be stale?
- **Immediate Consistency** – reads always return the most recent data
 - Consistency Level ALL guarantees immediate consistency, because all replica nodes are checked and compared before a result is returned
 - *Highest latency* because all replicas are checked and compared
- **Eventual Consistency** – reads may return stale data
 - Consistency Level ONE carries the highest risk of stale data, because only one replica node is checked before a result is returned
 - *Lowest latency* because the response from one replica is immediately returned

Available Replicas



Consistency Level

How do you choose a consistency level?

- In any given scenario, is the value of immediate consistency worth the latency cost?
 - Netflix uses CL ONE and measures its "eventual" consistency in milliseconds
 - Consistency Level ONE is your friend ...

Consistency Level ONE	Consistency Level QUORUM	Consistency Level ALL
Lowest latency	Higher latency (than ONE)	Highest latency
Highest throughput	Lower throughput	Lowest throughput
Highest availability	Higher availability (than ALL)	Lowest availability
Stale read possible (if read CL + write CL < RF)	No stale reads (if read <u>and</u> write at quorum)	No stale reads (if either read <u>or</u> write at ALL)

- If "stale" is measured in milliseconds, how much are those milliseconds worth?

Exercise 2: Working with consistency level

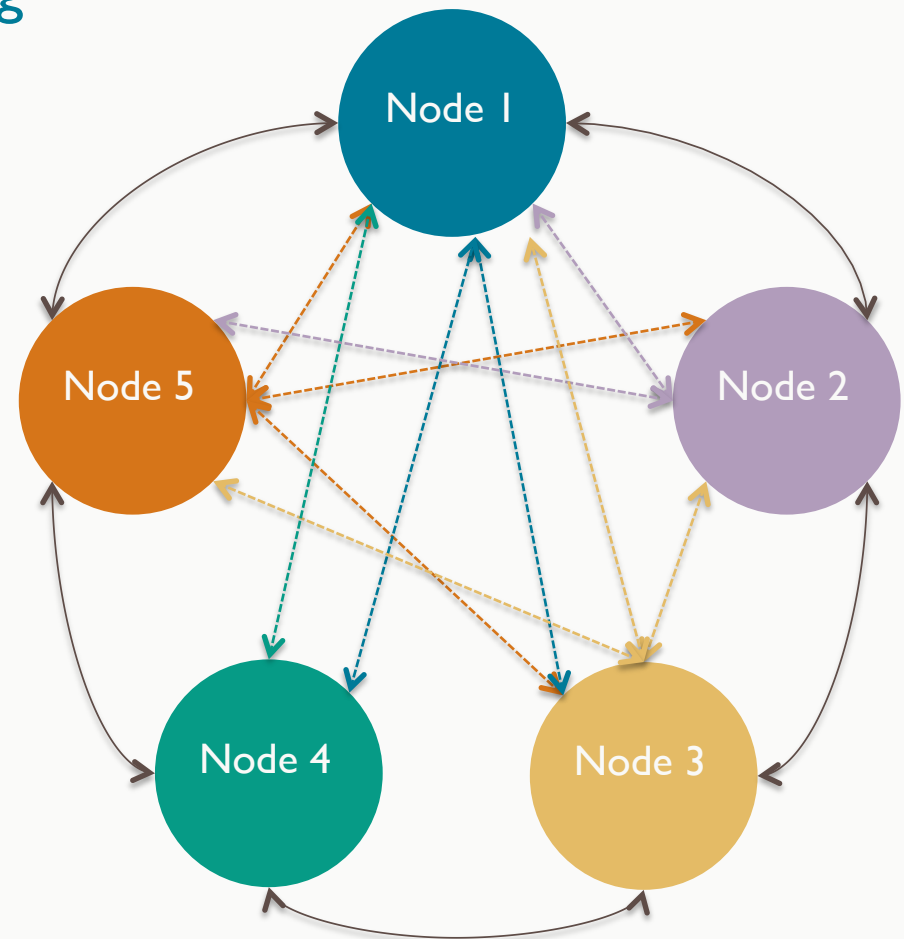


Learning Objectives

- Understand how requests are coordinated
- Understand replication
- Understand and tune consistency
- Introduce anti-entropy operations
- **Understand how nodes communicate**
- Understand the *System* keyspace

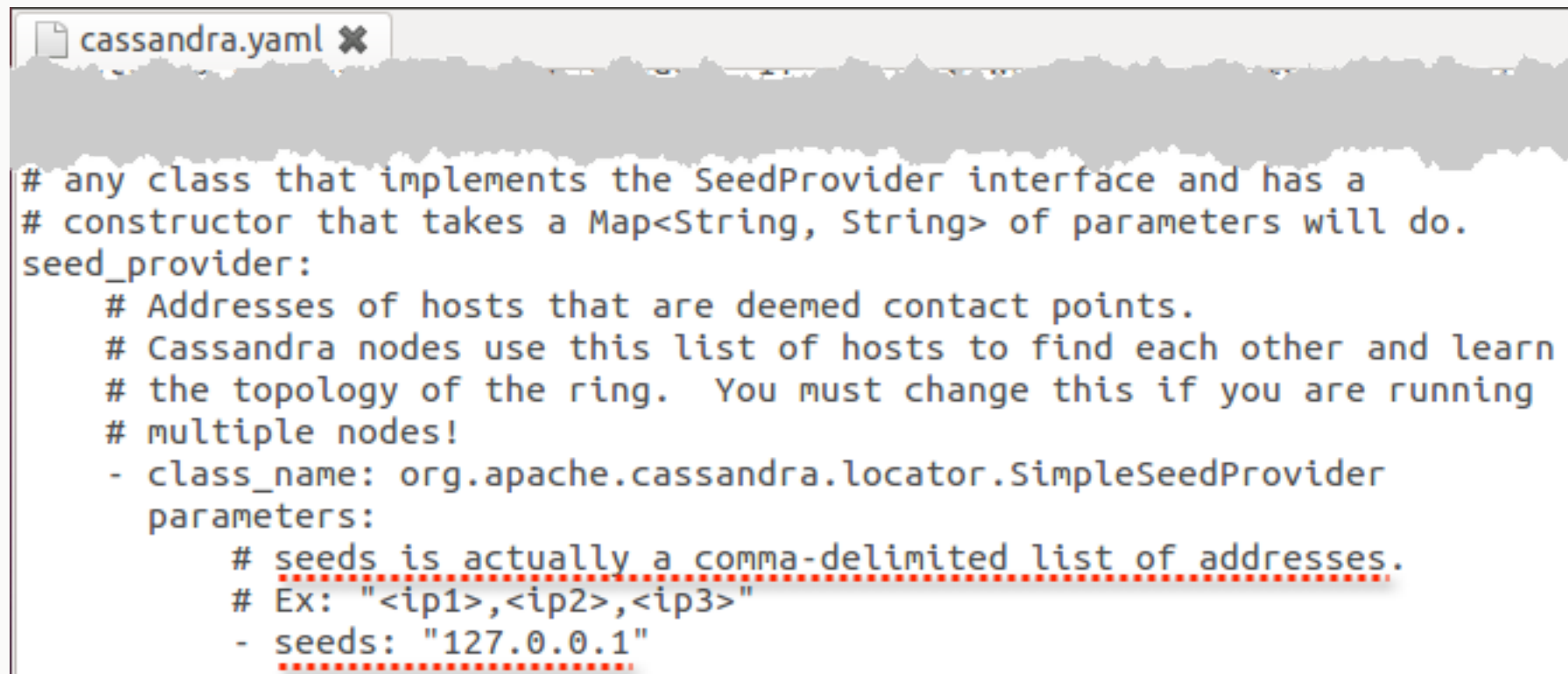
What is the Gossip protocol?

- Once per second, each node contacts 1 to 3 others, requesting and sharing updates about
 - Known node states ("heartbeats")
 - Known node locations
 - Requests and acknowledgments are timestamped, so information is continually updated and discarded



What is the Gossip protocol?

- As a node joins a cluster, it gossips with the *seed nodes* set in its *cassandra.yaml* to learn its cluster's topology
 - Assign the same seed nodes to each node in each data center
 - If more than one data center, include a seed node from each



```
cassandra.yaml ✕  
  
# any class that implements the SeedProvider interface and has a  
# constructor that takes a Map<String, String> of parameters will do.  
seed_provider:  
  # Addresses of hosts that are deemed contact points.  
  # Cassandra nodes use this list of hosts to find each other and learn  
  # the topology of the ring. You must change this if you are running  
  # multiple nodes!  
  - class_name: org.apache.cassandra.locator.SimpleSeedProvider  
    parameters:  
      # seeds is actually a comma-delimited list of addresses.  
      # Ex: "<ip1>,<ip2>,<ip3>"  
      - seeds: "127.0.0.1"
```


Learning Objectives

- Understand how requests are coordinated
- Understand replication
- Understand and tune consistency
- Introduce anti-entropy operations
- Understand how nodes communicate
- **Understand the *System* keyspace**

What is the System keyspace?

- Cassandra stores its state in *System* keyspace tables
 - Examining *System* keyspace tables is useful towards understanding Cassandra
 - But, directly editing any *System* keyspace table is an anti-pattern, so don't

```
dstraining@DST: /home/cassandra
dstraining@DST: /home/dsc-cassandra-2.0.5

dstraining@DST:/home/cassandra$ bin/cqlsh
Connected to Test Cluster at localhost:9160.
[cqlsh 4.1.1 | Cassandra 2.0.5 | CQL spec 3.1.1 | Thrift protocol 19.39.0]
Use HELP for help.
cqlsh> USE SYSTEM;
cqlsh:system> SELECT * FROM system.schema_keyspaces;
```

keyspace_name	durable_writes	strategy_class	strategy_options
test	True	org.apache.cassandra.locator.SimpleStrategy	{"replication_factor":"1"}
system	True	org.apache.cassandra.locator.LocalStrategy	{}
system_traces	True	org.apache.cassandra.locator.SimpleStrategy	{"replication_factor":"2"}
demo	True	org.apache.cassandra.locator.SimpleStrategy	{"replication_factor":"1"}

```
(4 rows)
cqlsh:system>
```

- Use the CQL *DESCRIBE KEYSPACE* command to list all System table schema

What tables are in the System keyspace, and why?

- System keyspace tables include

Table	Purpose
schema_keyspaces	Keyspaces available within this cluster, along with their assigned replication strategy and replication factor
schema_columns	Details of cluster compound primary key columns
schema_columnfamilies	Details of cluster tables and their configuration
peers	Local tracking of cluster-wide gossip for this node
local	Details of the local node's own state
hints	Stores information for hinted handoffs

- Other System keyspace tables

IndexInfo	schema_usertypes	batchlog	compaction_history
compactions_in_progress	paxos	peer_events	range_xfers
schema_triggers	sstable_activity		

Demo 5: Query and examine System tables



Summary

- A Cassandra *cluster* is comprised of peer-to-peer *nodes* logically organized into *racks* within *data centers*
- Any node may *coordinate* any *request* issued by a Cassandra *client*
- Data is organized into *partitions* ("rows") identified by *tokens* in an integer range
- The total *token range* is treated internally as a ring whose segments are owned by nodes
- Nodes are identified by the highest token in their segment of the total range
- A node's *partitioner* hashes a token from the *partition key* of a value being written
- The *first replica* ("copy") of a *partition* is written to the *node* owning the *primary range* containing its *token*
- The *Murmur3Partitioner* is the default and best practice
- *Virtual nodes* are multiple smaller token ring segments managed by a single machine
- Virtual nodes improve performance as nodes are *bootstrapped*, added, and removed
- A *keyspace* defines a cluster's *replication strategy* and *replication factor*
- *Replication factor* (RF) determines how many replicas ("copies") are made of each partition

Summary

- *Replication strategy* determines how replicas are distributed across the cluster
- A *per-request consistency level (CL)* determines how many nodes must acknowledge
- A *hinted handoff* temporarily stores requests issued to unavailable nodes
- *Consistency levels* include *ANY*, *ONE*, *QUORUM*, *LOCAL_QUORUM*, *EACH_QUORUM*, and *ALL*
- *Immediately consistent* data is guaranteed to be current
- Nodes with stale data are updated during each read request through *read repair*
- The *nodetool repair* command makes stale data consistent for a node or set of nodes
- **IF *nodes_written* + *nodes_read* > *replication_factor* THEN results are immediately consistent**
- *Eventually consistent* data may be a few milliseconds stale
- Node clocks must be in sync as the most recently timestamped data returns to the client
- Nodes continually exchange state and location information via the *Gossip* protocol
- Each node includes a *Snitch* which tracks and reports on the current cluster topology
- Cassandra tracks its internal state and structure in *System* keyspace tables

Review Questions

- Describe the relationship of nodes, racks, clusters, and data centers
- What is the function of the partitioner?
- Can a node hold a partition with a token outside its primary range?
- In a 3 node cluster with RF=2, how much total data volume does each node own?
- What is the function of the *nodetool repair* operation?
- What is a remote coordinator?
- How could RF and CL be tuned to ensure immediate consistency?

