

CockroachDB: Geo-distributed SQL Database

Make Data Easy

- Distributed
 - Horizontally scalable to grow with your application
- Geo-distributed
 - Handle datacenter failures
 - Place data near usage
 - Push computation near data
- SQL
 - Lingua-franca for rich data storage
 - Schemas, indexes, and transactions make app development easier

Distributed, Replicated, Transactional KV*

- Keys and values are strings
 - Lexicographically ordered by key
- Multi-version concurrency control (MVCC)
 - Values are never updated “in place”, newer versions shadow older versions
 - Tombstones are used to delete values
 - Provides snapshot to each transaction
- Monolithic key-space

* Not exposed for external usage

Monolithic Key Space

DOGS	
carl	
dagne	
figment	
jack	
lady	
lula	
muddy	
peetey	
pinetop	
sooshi	
stella	
zee	

Monolithic logical key space

- Ordered lexicographically by key

Ranges

DOGS	
carl	
dagne	
figment	
jack	
lady	
lula	
muddy	
peetey	
pinetop	
sooshi	
stella	
zee	

Key space divided into contiguous ~64MB ranges

carl	lady	pinetop
dagne	lula	sooshi
figment	muddy	stella
jack	peetey	zee

Ranges are small enough to be moved/split quickly

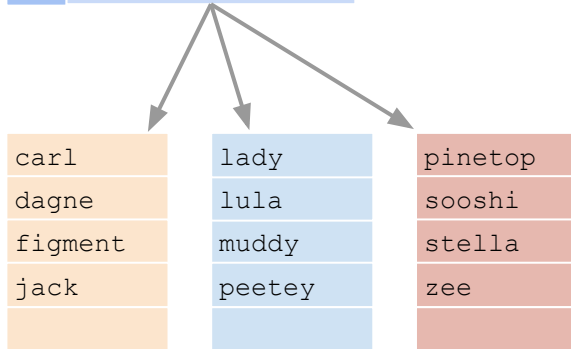
Ranges are large enough to amortize indexing overhead

Range Indexing

DOGS	
carl	
dagne	
figment	
jack	
lady	
lula	
muddy	
peetey	
pinetop	
sooshi	
stella	
zee	

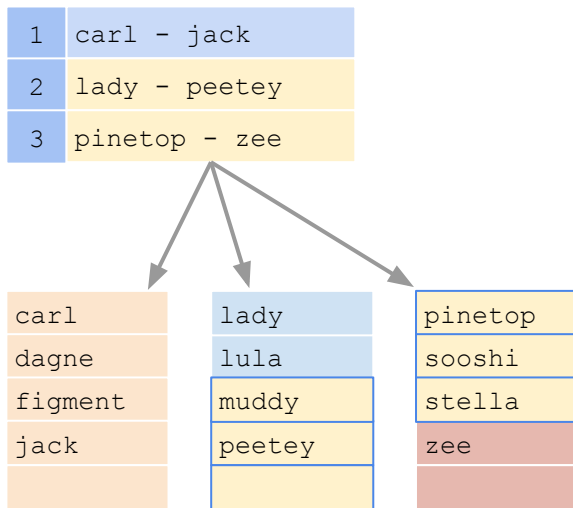
1	carl - jack
2	lady - peetey
3	pinetop - zee

Index structure used to
locate ranges
(very much like a B-tree)



Ordered Range Scans

DOGS	
carl	
dagne	
figment	
jack	
lady	
lula	
muddy	
peetey	
pinetop	
sooshi	
stella	
zee	



Ordered keys enable
efficient range scans

`dogs >= "muddy" AND <= "stella"`

Transactional Updates

INSERT [sunny]

DOGS

carl
dagne
figment
jack
lady
lula
muddy
peetey
pinetop
sooshi
stella
zee

1	carl - jack
2	lady - peetey
3	pinetop - zee

carl
dagne
figment
jack

lady
lula
muddy
peetey

pinetop
sooshi
stella
zee

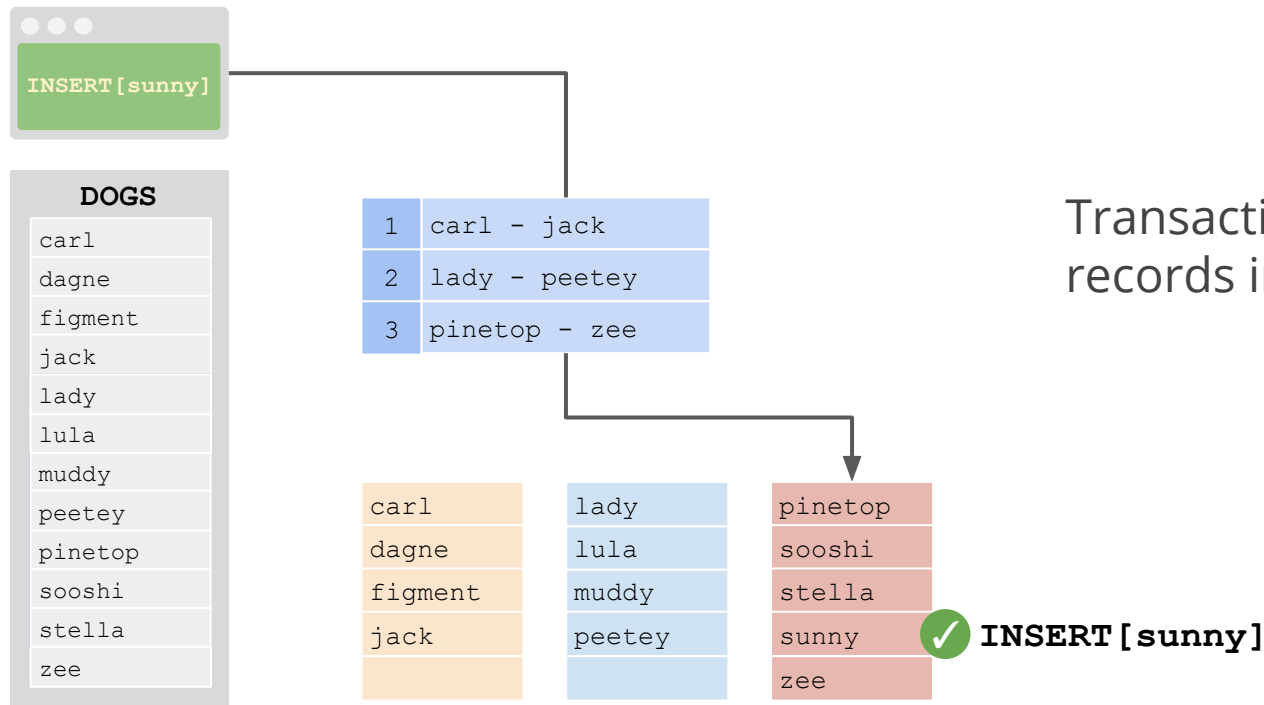
Transactions used to insert records into ranges



INSERT [sunny]

Space available in range? - YES

Transactional Updates



Transactions used to insert records into ranges

Range Splits

INSERT [rudy]

DOGS	
carl	
dagne	
figment	
jack	
lady	
lula	
muddy	
peetey	
pinetop	
sooshi	
stella	
zee	

1	carl - jack
2	lady - peetey
3	pinetop - zee

carl
dagne
figment
jack

lady
lula
muddy
peetey

pinetop
sooshi
stella
sunny
zee

BUT... what happens when a range is full?



INSERT [rudy]

Space available in range? - **NO**

Range Splits

`INSERT [rudy]`

DOGS	
carl	
dagne	
figment	
jack	
lady	
lula	
muddy	
peetey	
pinetop	
sooshi	
stella	
zee	

1	carl - jack
2	lady - peetey
3	pinetop - sooshi
4	stella - zee

carl
dagne
figment
jack

lady
lula
muddy
peetey

pinetop
rudy
sooshi

stella
sunny
zee

Ranges are automatically split, a new range index is created & order maintained

INSERT [rudy]
split range and insert

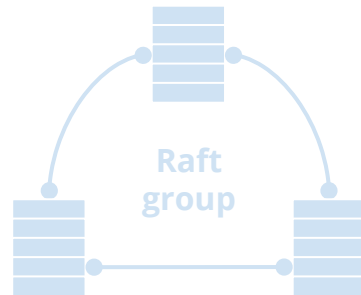
Raft and Replication

Ranges (~64MB) are the unit of replication

Each range is a Raft group
(Raft is a consensus replication protocol)

Default to 3 replicas, though this is configurable

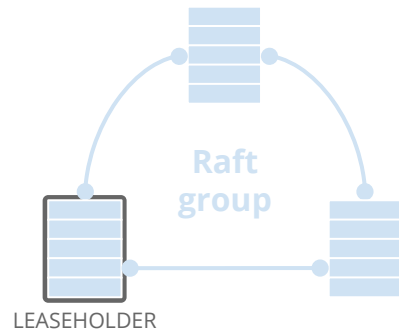
- Important system ranges default to 5 replicas
- Note: 2 replicas doesn't make sense in consensus replication



Raft and Replication

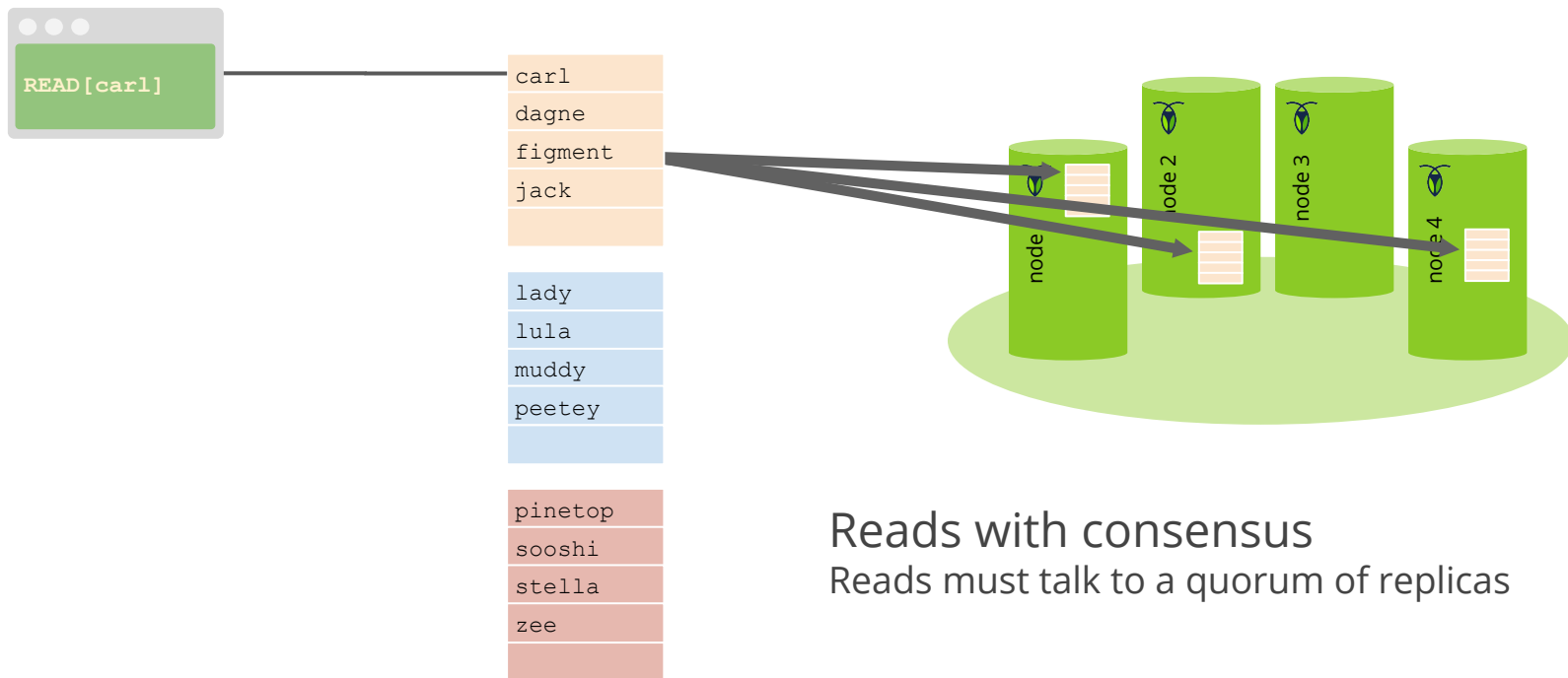
Raft provides “atomic replication” of commands

Commands are proposed by the leaseholder replica and distributed to the follower replicas, but only accepted when a quorum of replicas have acknowledged receipt



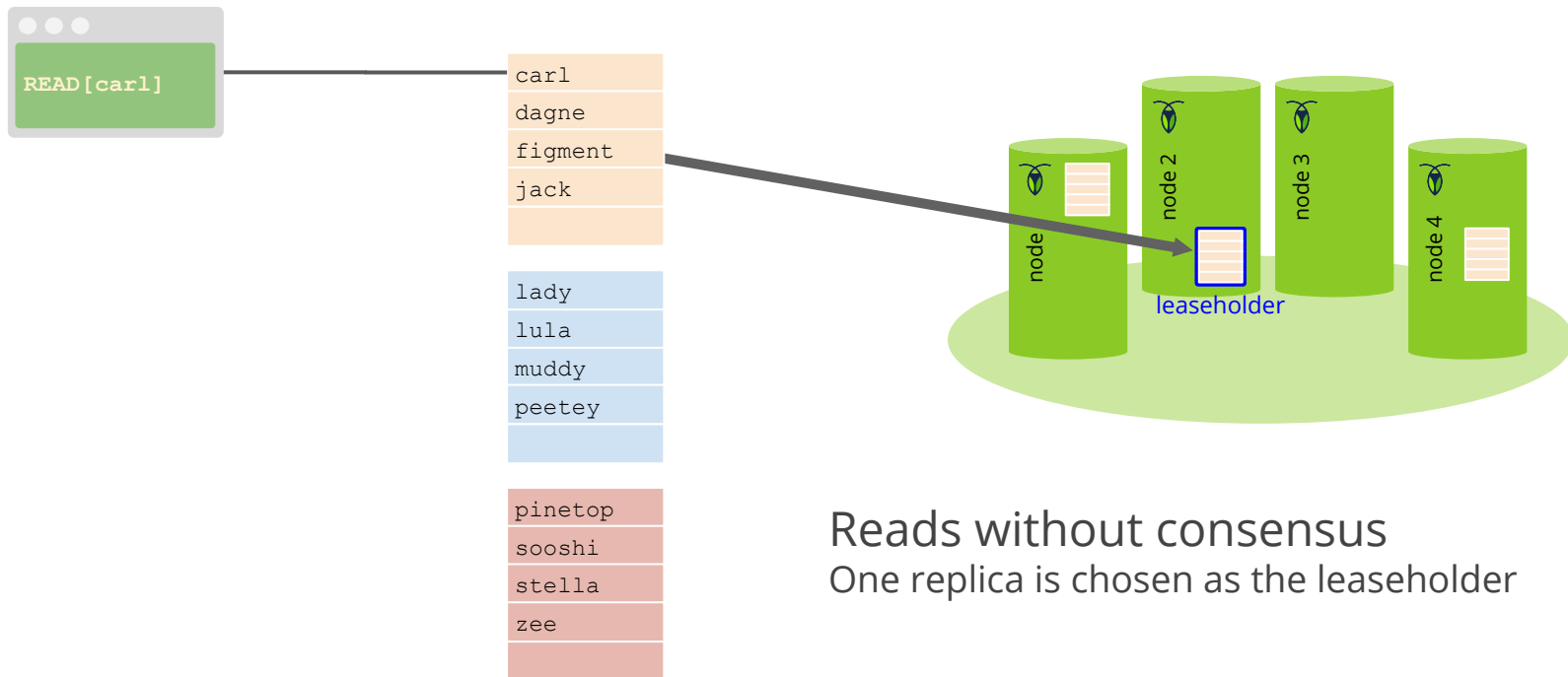
* Leaseholder == Raft leader

Range Leases



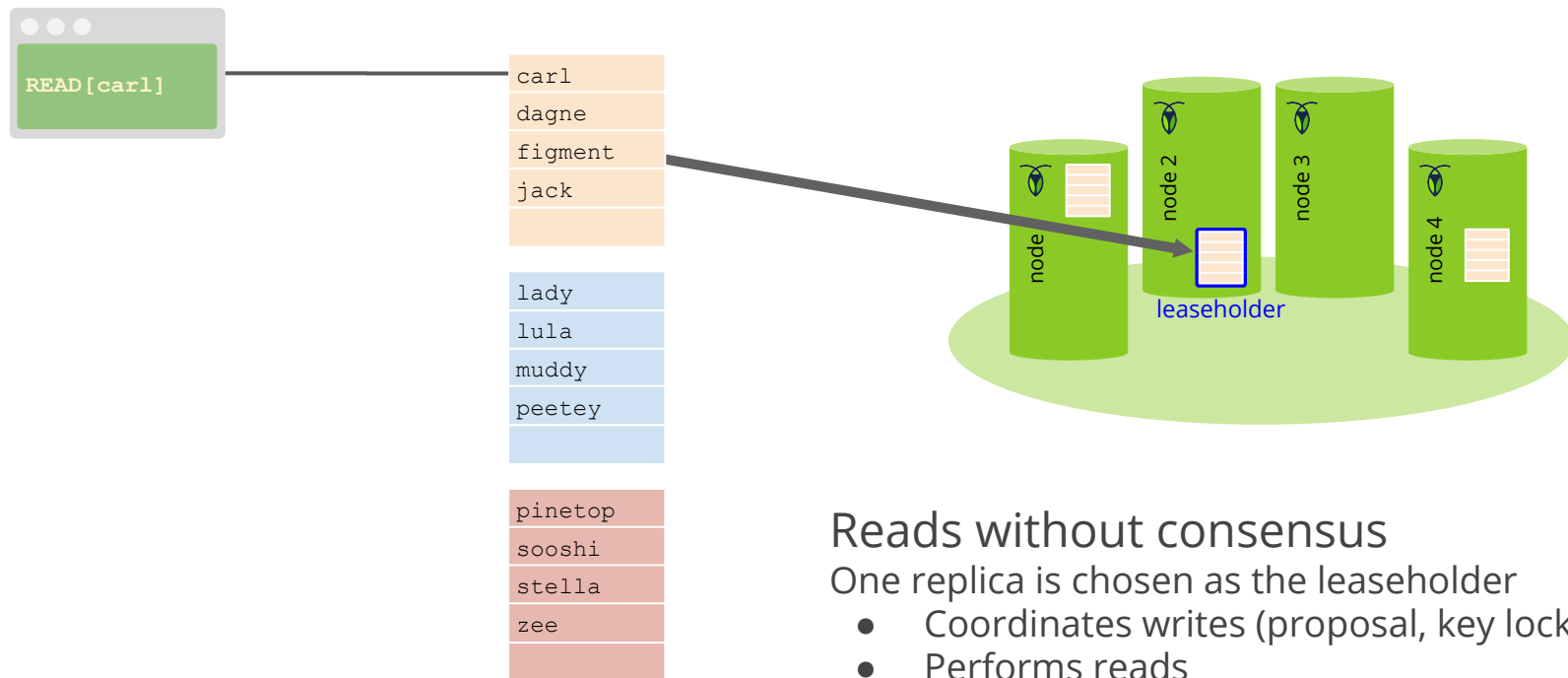
Reads with consensus
Reads must talk to a quorum of replicas

Range Leases



Reads without consensus
One replica is chosen as the leaseholder

Range Leases



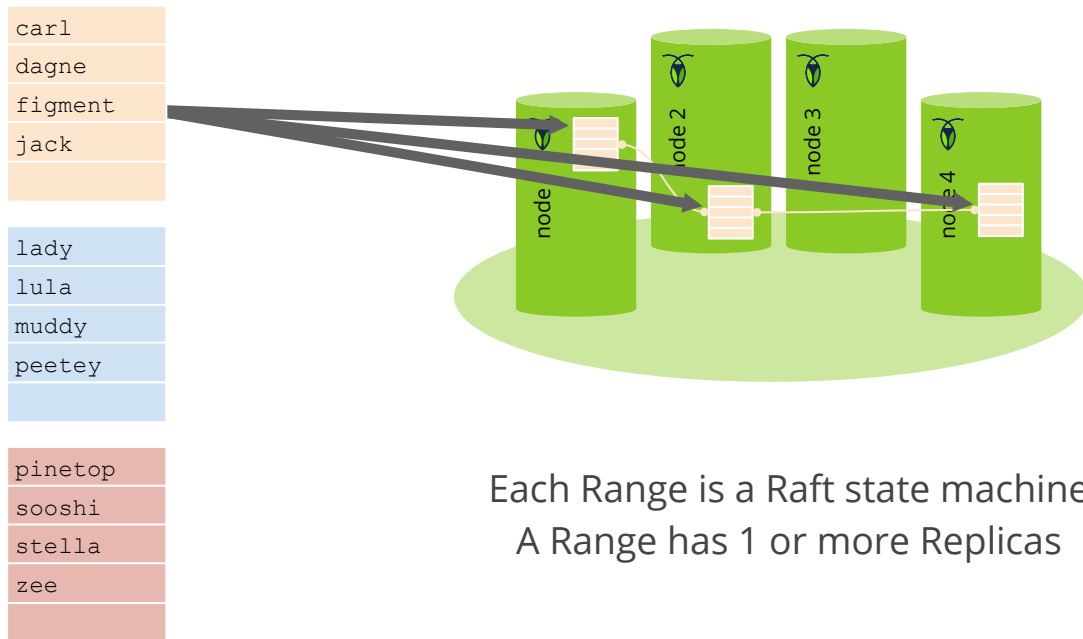
Reads without consensus

One replica is chosen as the leaseholder

- Coordinates writes (proposal, key locking)
- Performs reads

Replica Placement

- Space
- Diversity
- Load
- Latency

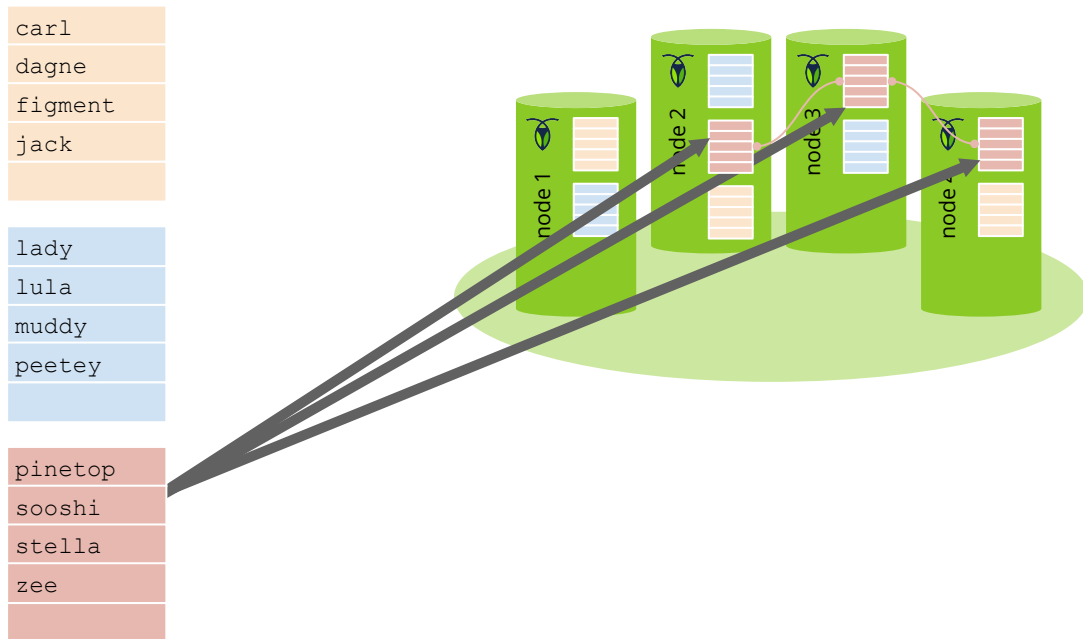


Replica Placement: Diversity

Diversity

optimizes placement of replicas across “failure domains”

- Disk
- Single machine
- Rack
- Datacenter
- Region



Replica Placement: Load

Load

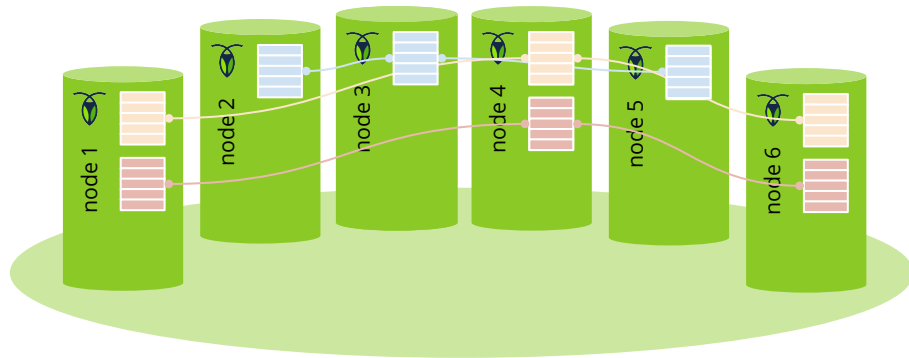
Balances placement using heuristics that considers real-time usage metrics of the data itself

carl
dagne
figment
jack

lady
lula
muddy
peetey

This range is high load as it is accessed more than others

pinetop
sooshi
stella
zee



While we show this for ranges within a single table, this is also applicable across all ranges across ALL tables, which is the more typical situation

Replica Placement: Latency & Geo-partitioning

carl
dagne
figment
jack

EU /carl
EU /lula
EU /sooshi
EU /zee

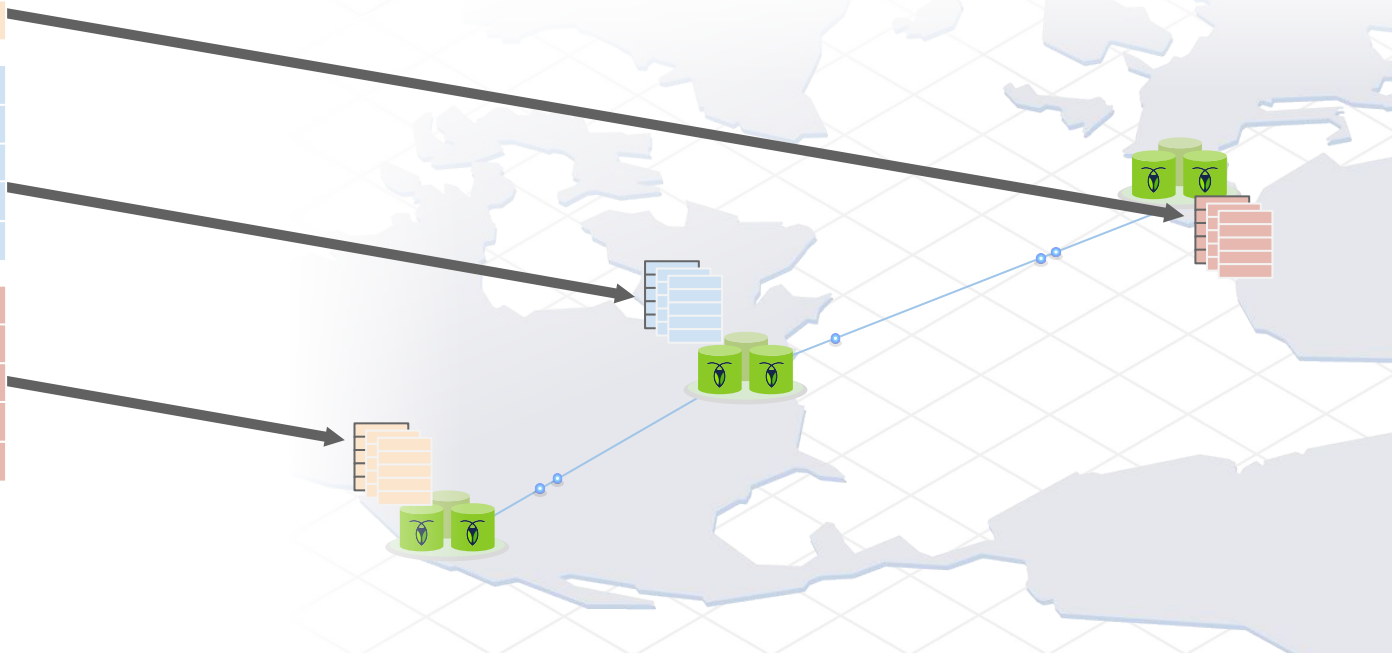
lady
lula
muddy
peetey

USE /dagne
USE /figment
USE /muddy
USE /stella

pinetop
sooshi
stella
zee

USW /jack
USW /lady
USW /peetey
USW /pinetop

We apply a constraint that indicates regional placement so we can ensure low latency access or jurisdictional control of data



Rebalancing Replicas

Scale: Add a node

If we add a node to the cluster, CockroachDB automatically redistributed replicas to even load across the cluster

Uses the replica placement heuristics from previous slides to decide which node to add to and which to remove from



Rebalancing Replicas

Scale: Add a node

If we add a node to the cluster, CockroachDB automatically redistributed replicas to even load across the cluster

Uses the replica placement heuristics from previous slides



Movement is decomposed into adding a replica followed by removing a replica

Rebalancing Replicas

Scale: Add a node

If we add a node to the cluster, CockroachDB automatically redistributed replicas to even load across the cluster

Uses the replica placement heuristics from previous slides



Movement is decomposed into adding a replica followed by removing a replica

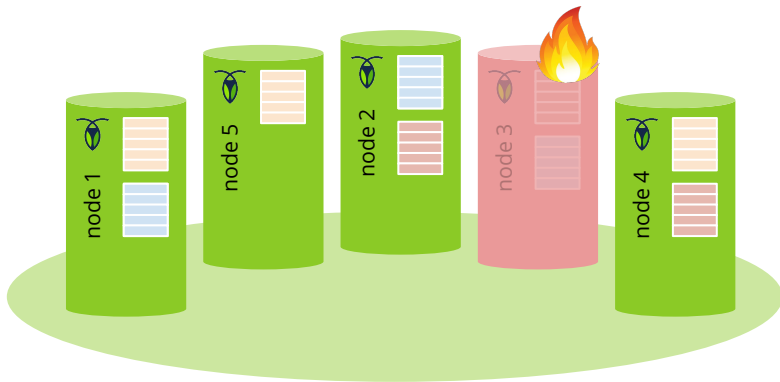
Rebalancing Replicas

Loss of a node

Permanent Failure

If a node goes down, the Raft group realizes a replica is missing and replaces it with a new replica on an active node

Uses the replica placement heuristics from previous slides



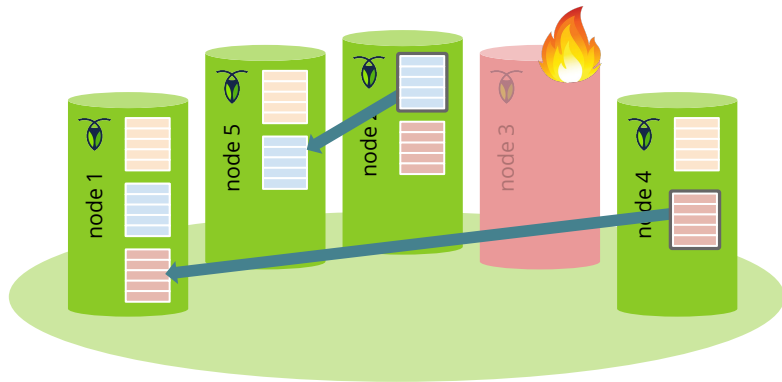
Rebalancing Replicas

Loss of a node

Permanent Failure

If a node goes down, the Raft group realizes a replica is missing and replaces it with a new replica on an active node

Uses the replica placement heuristics from previous slides



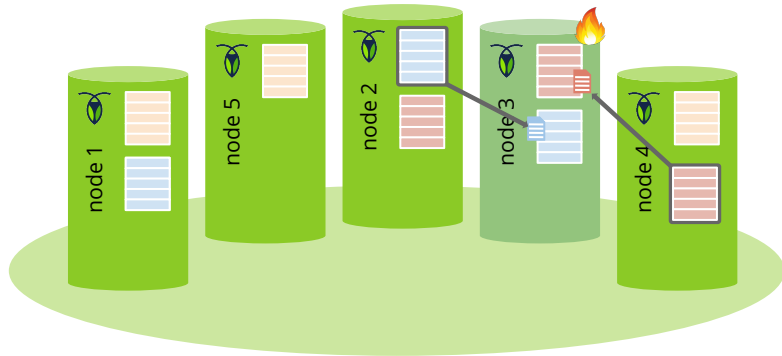
The failed replica is removed from the Raft group and a new replica created. The leaseholder sends a snapshot of the Range's state to bring the new replica up to date.

Rebalancing Replicas

Loss of a node

Temporary Failure

If a node goes down for a moment, the leaseholder can “catch up” any replica that is behind



The leaseholder can send commands to be replayed OR it can send a snapshot of the current Range data. We apply heuristics to decide which is most efficient for a given failure.

Transactions

Atomicity, **C**onsistency, **I**solation, **D**urability

Serializable Isolation

- As if the transactions are run in a serial order
- Gold standard isolation level
- Make Data Easy - weaker isolation levels are too great a burden

Transactions can span arbitrary ranges

Conversational

- The full set of operations is not required up front

Transactions

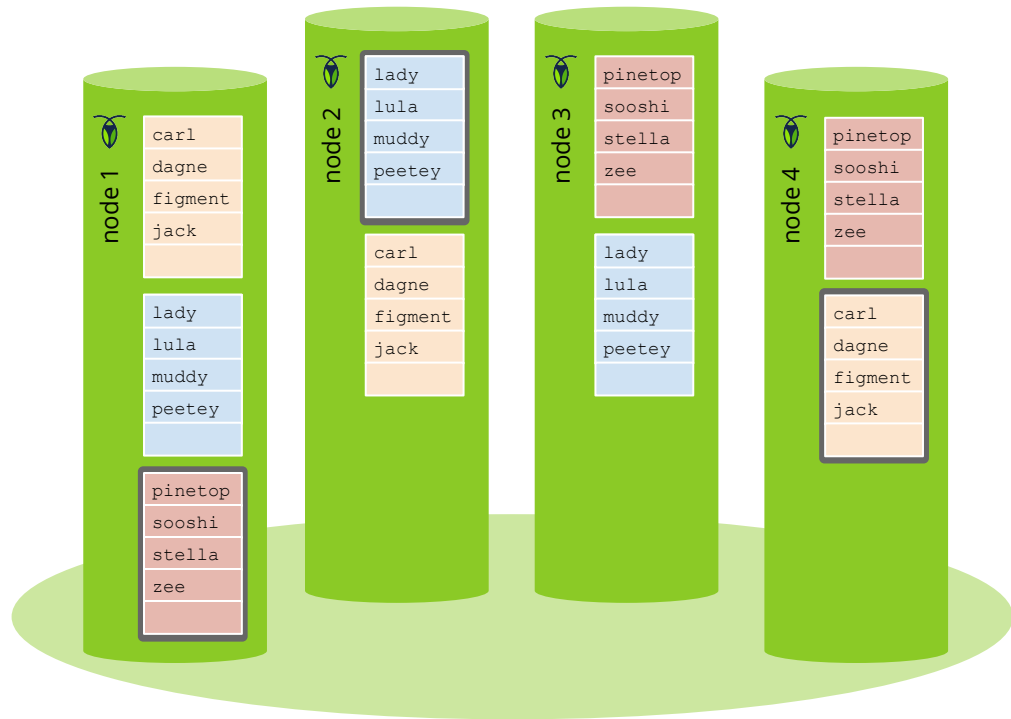
Raft provides atomic writes to individual ranges

Bootstrap transaction atomicity using Raft atomic writes

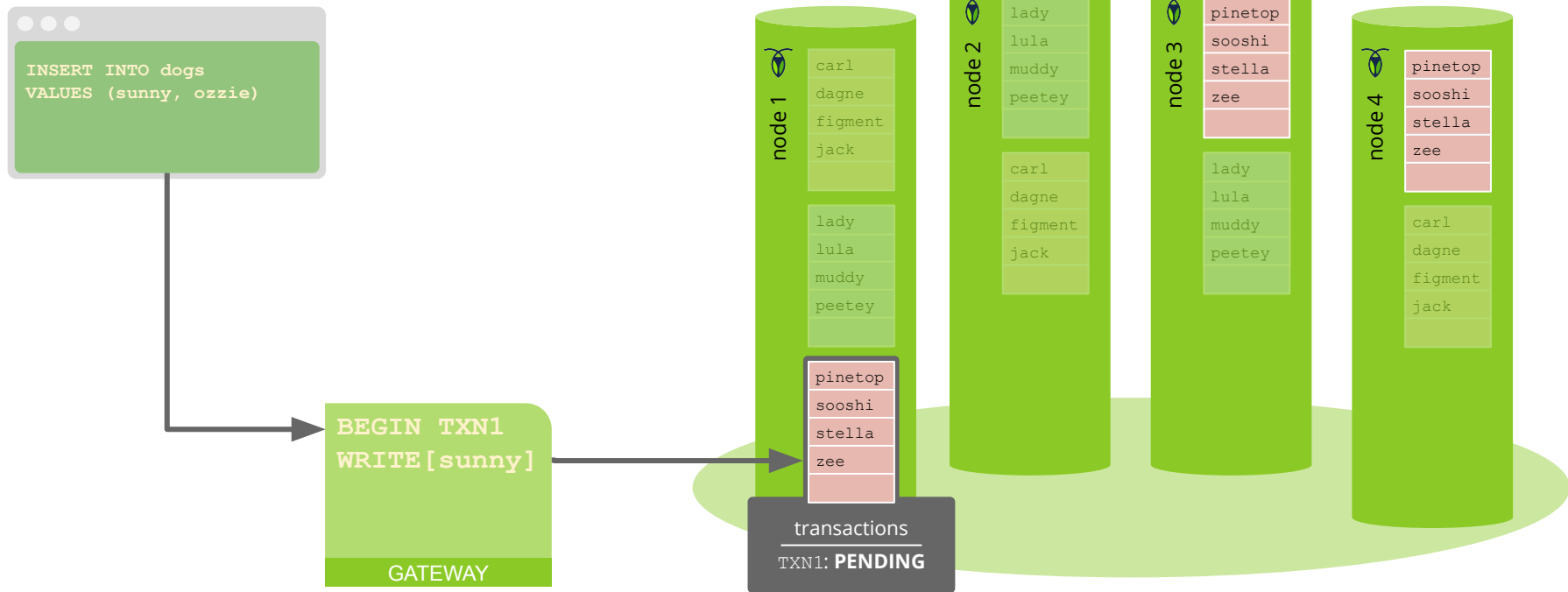
Transaction record atomically flipped from PENDING to COMMIT

Distributed Transactions

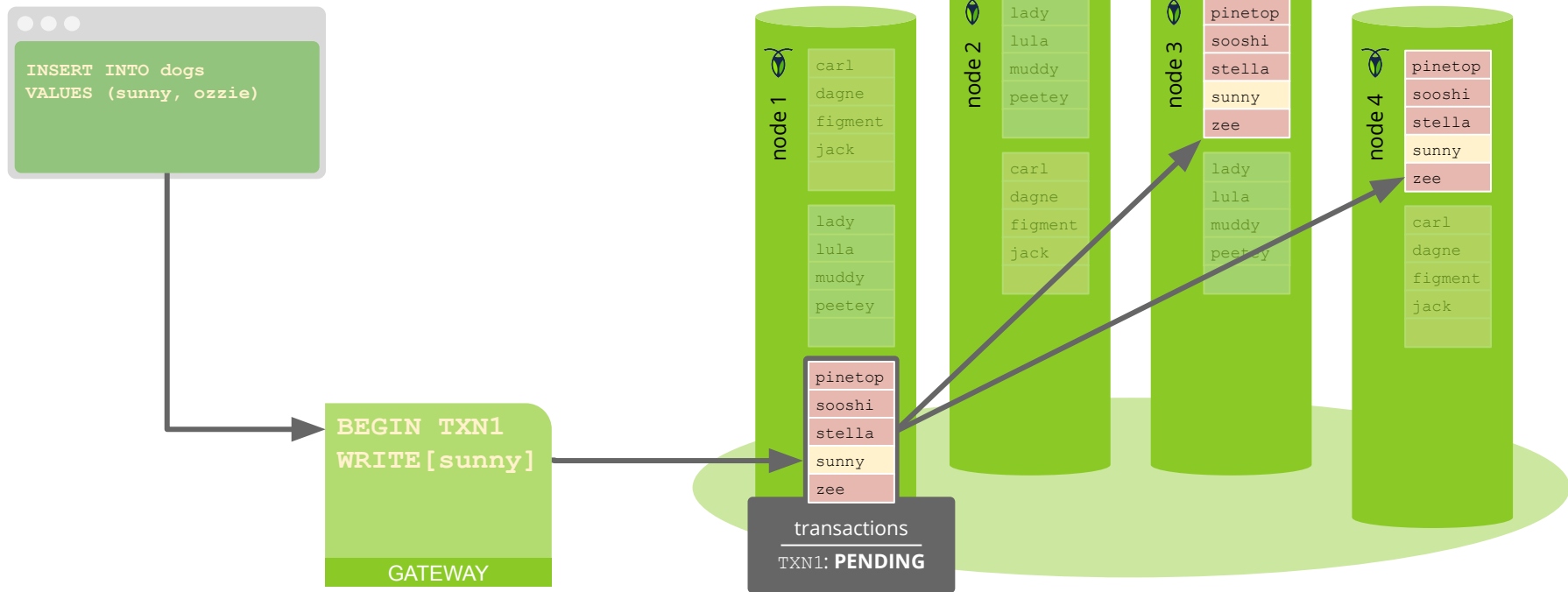
```
INSERT INTO dogs  
VALUES (sunny, ozzie)
```



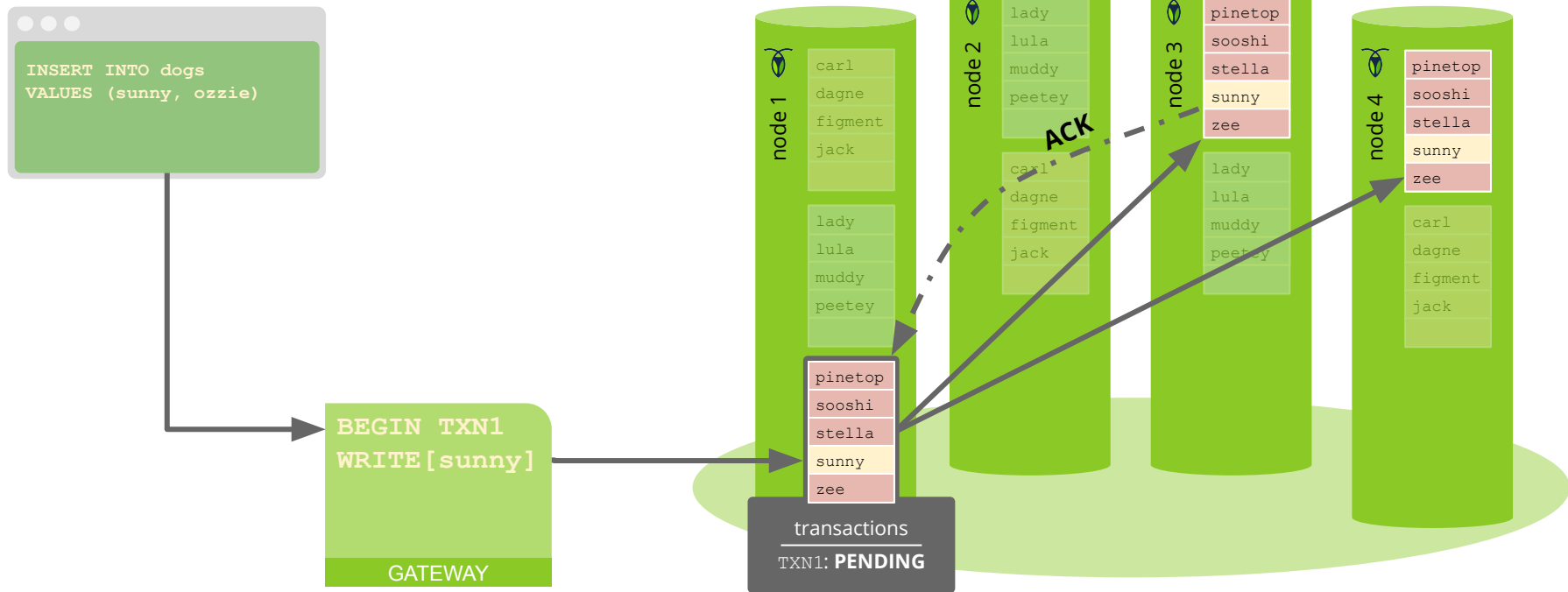
Distributed Transactions



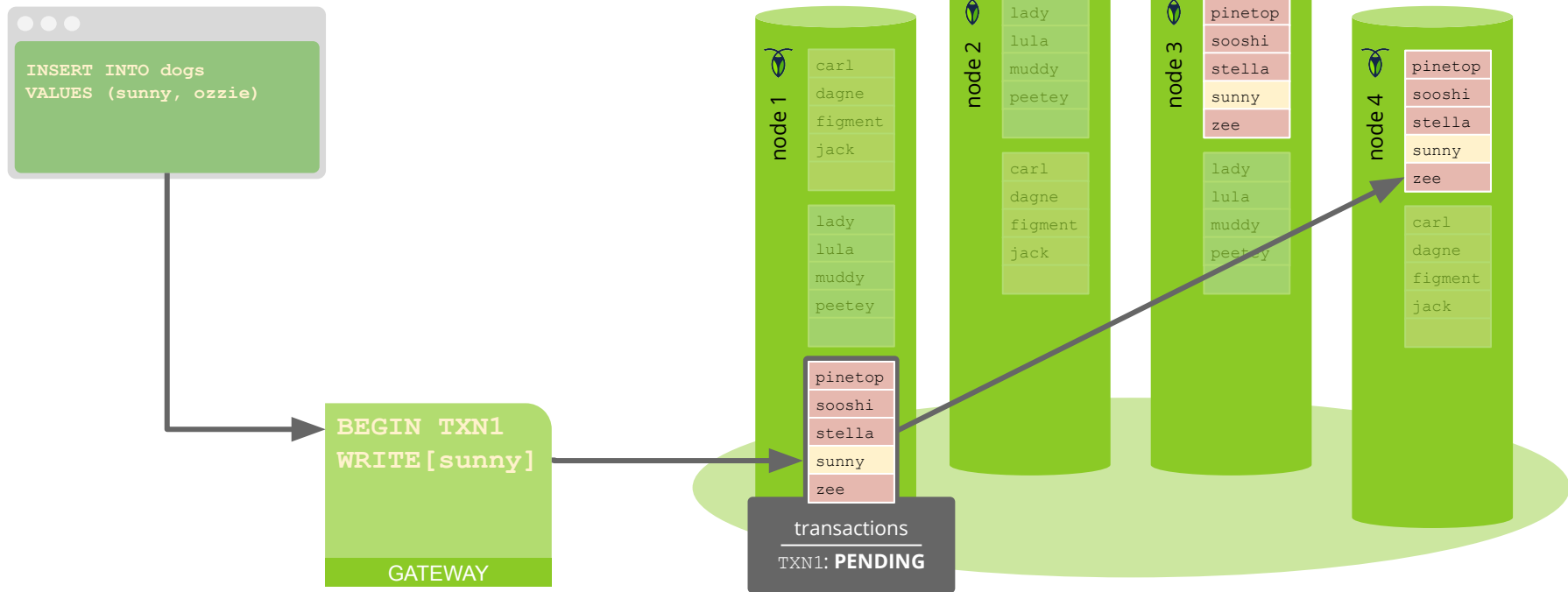
Distributed Transactions



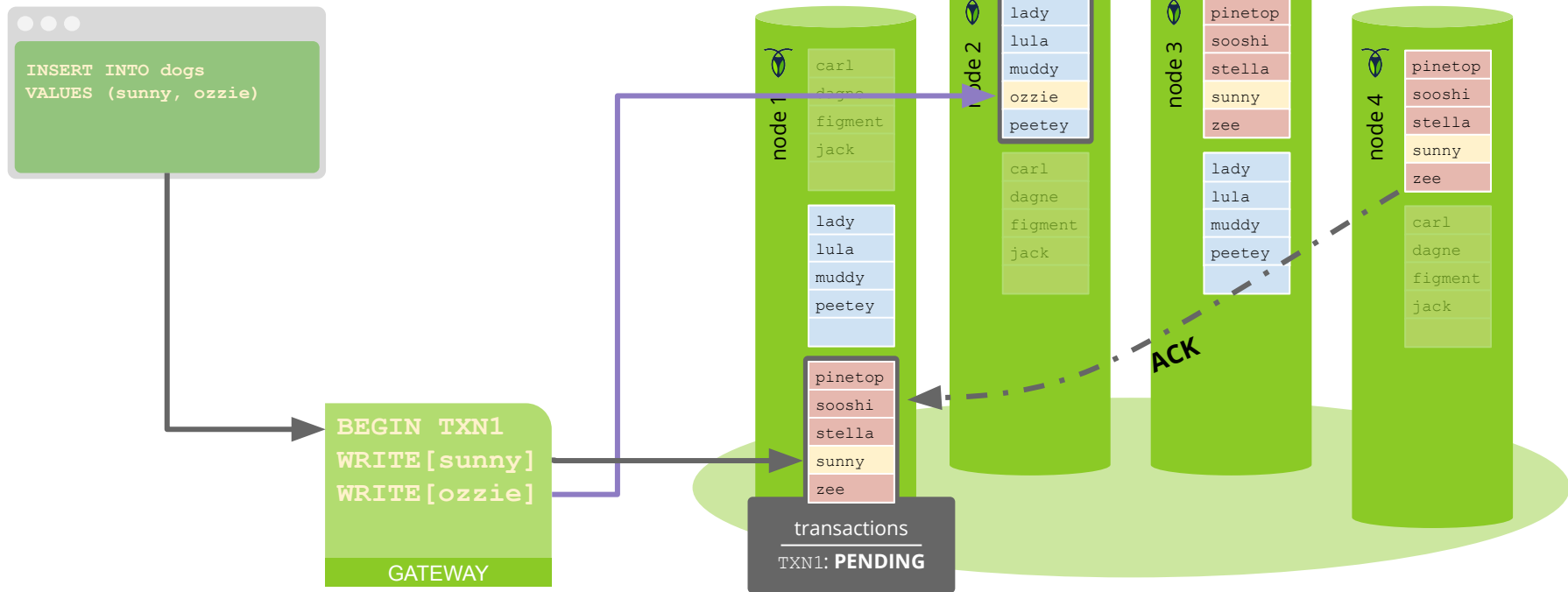
Distributed Transactions



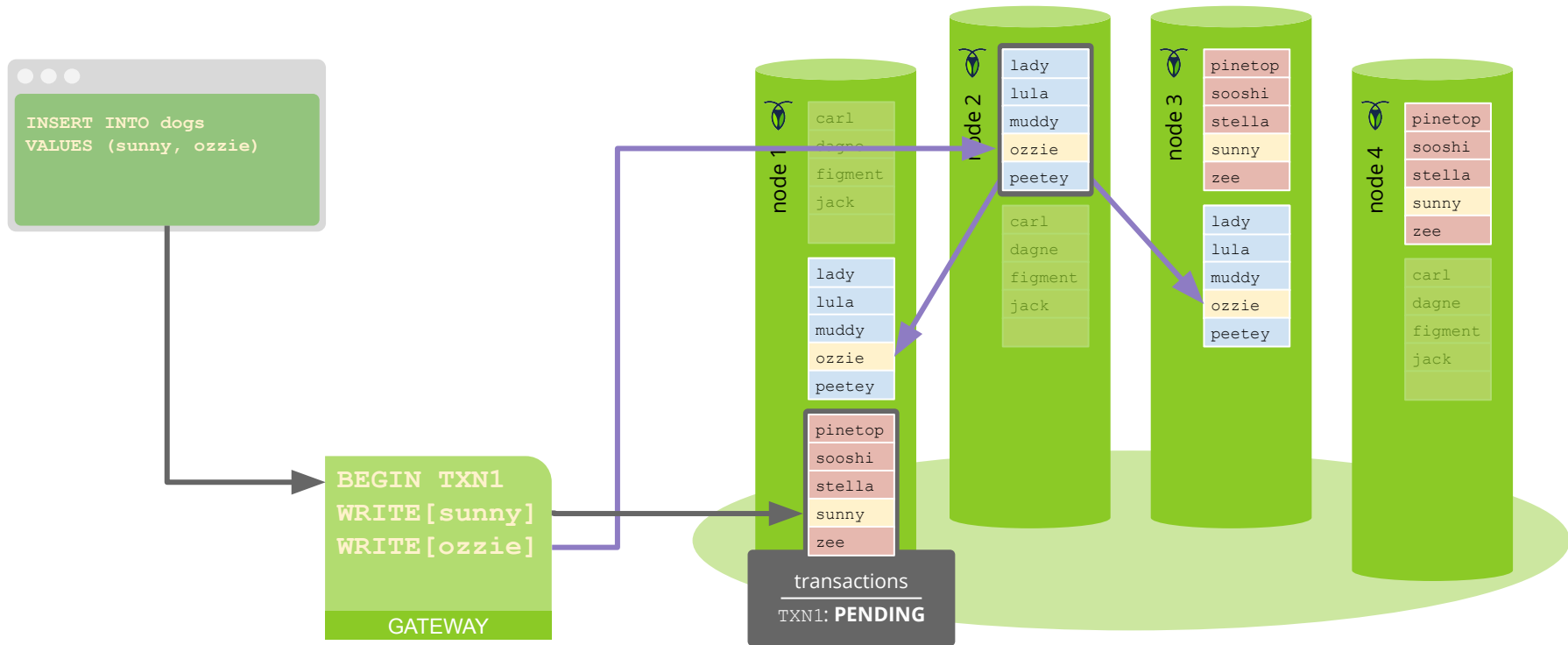
Distributed Transactions



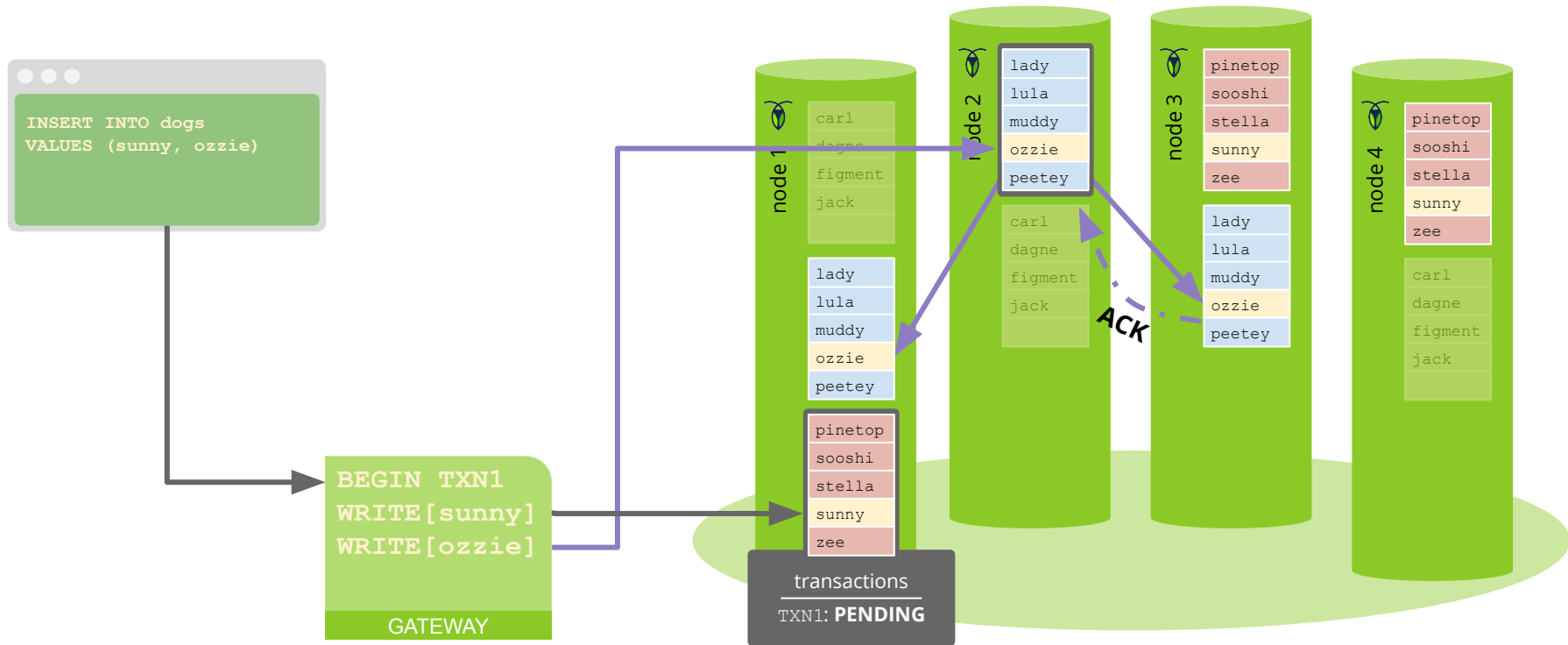
Distributed Transactions



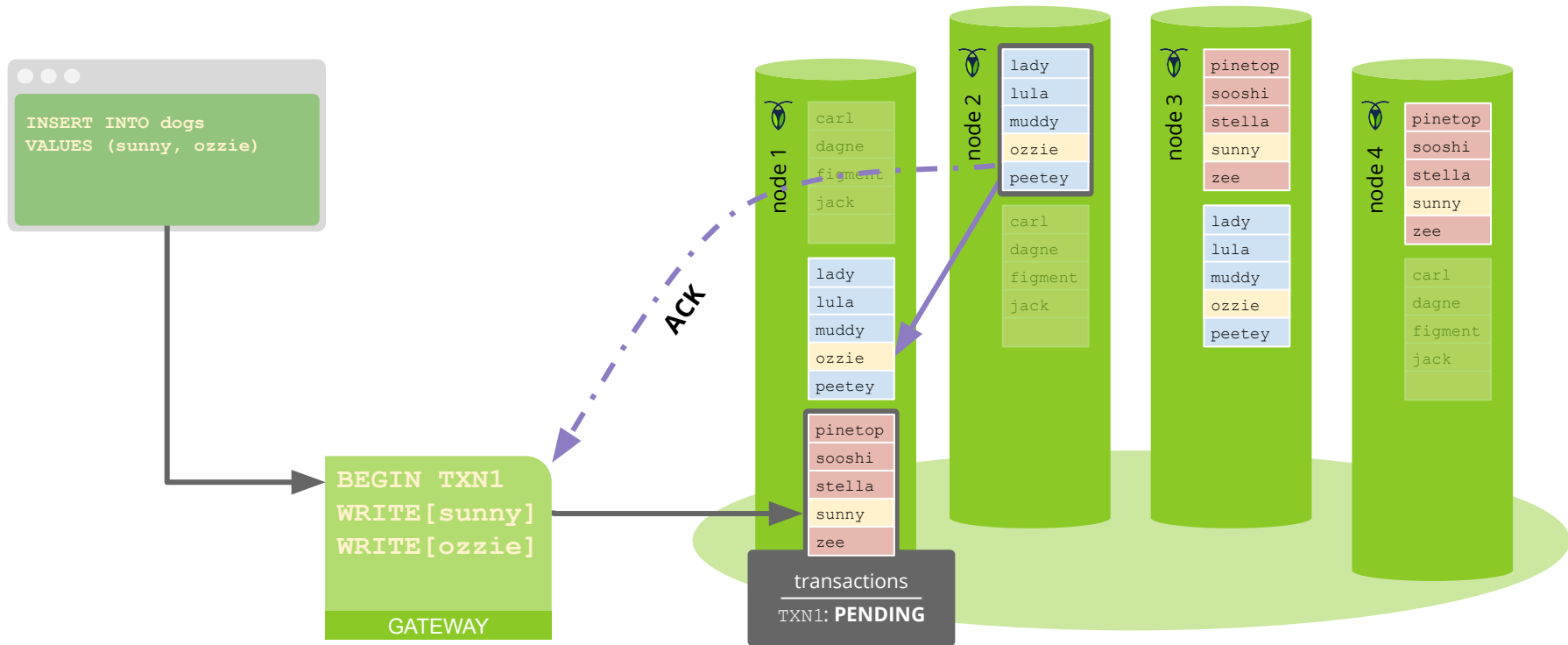
Distributed Transactions



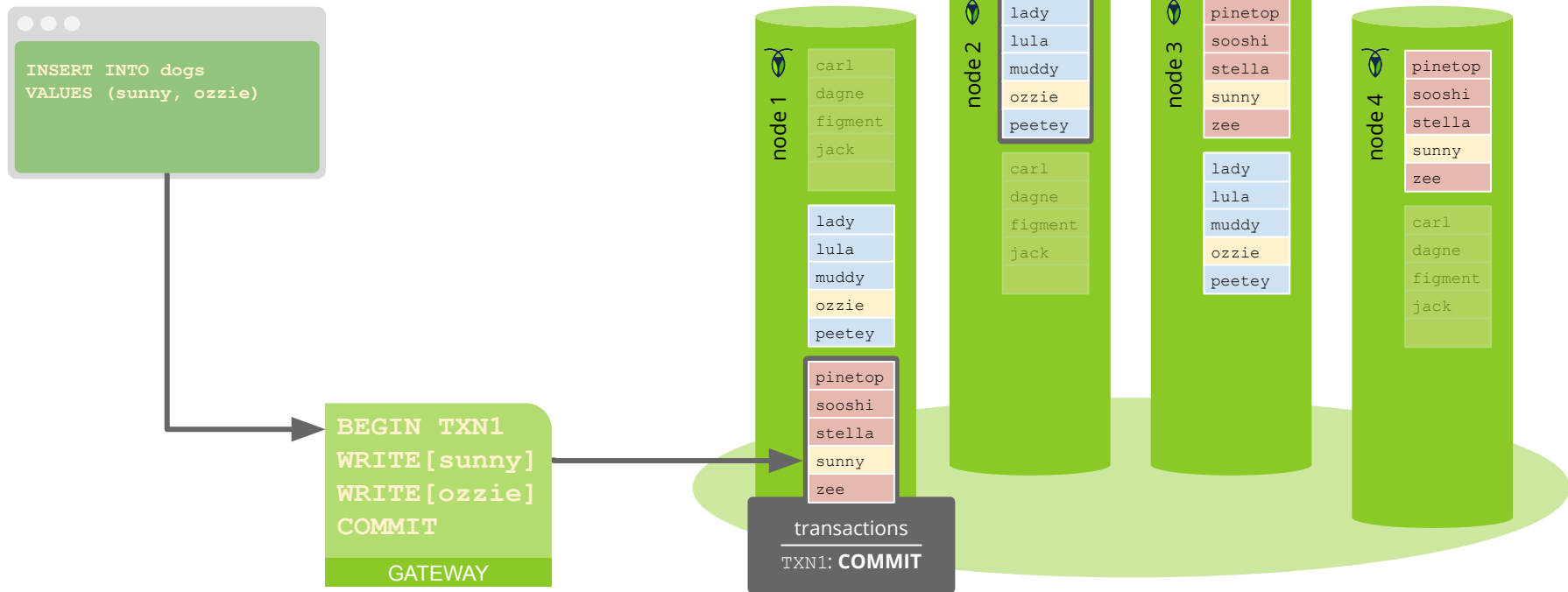
Distributed Transactions



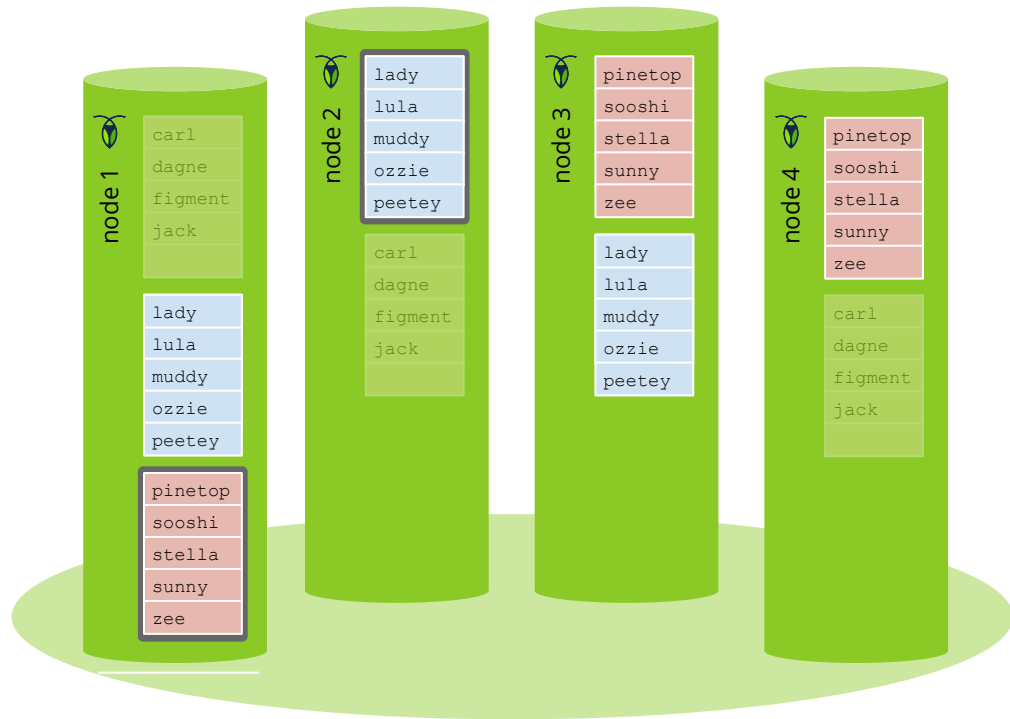
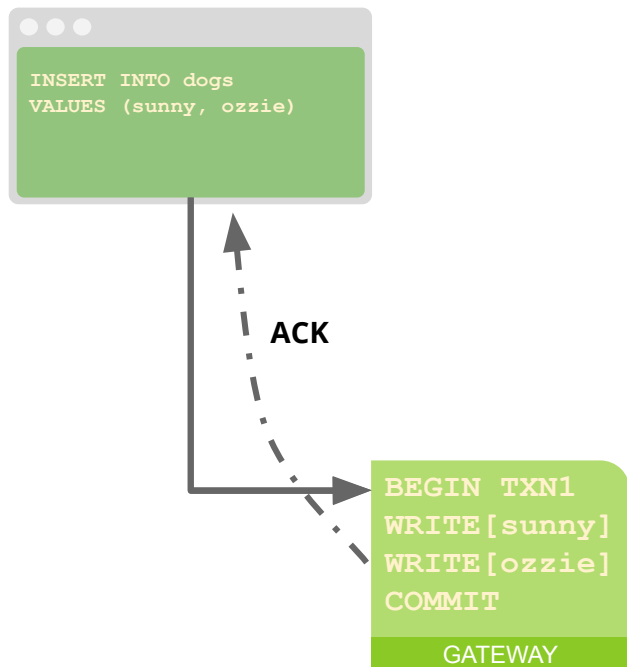
Distributed Transactions



Distributed Transactions



Distributed Transactions



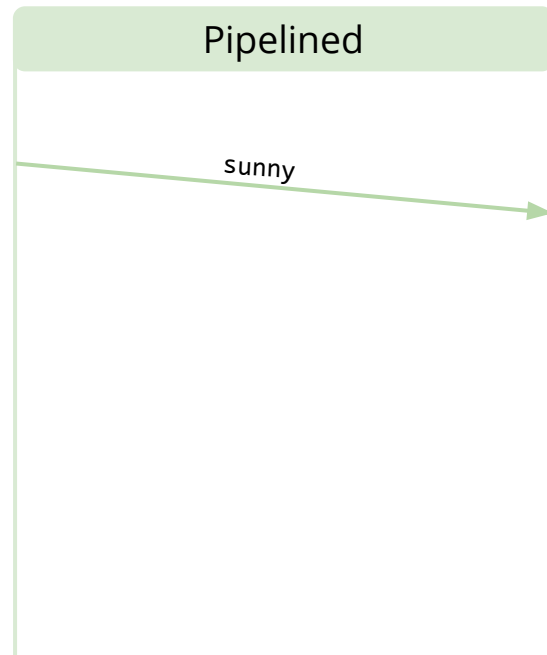
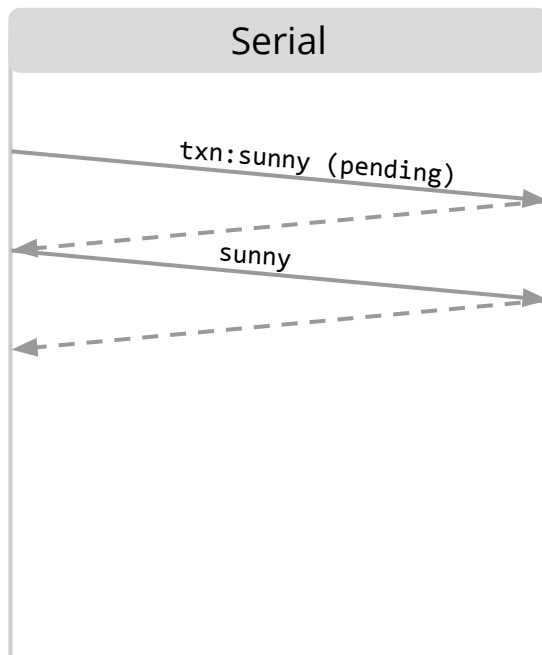
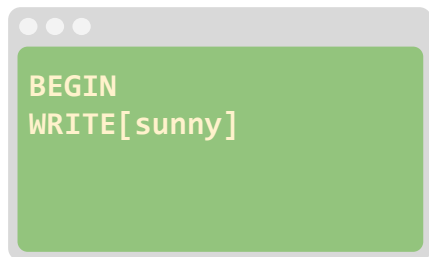
Transactions: Pipelining



Serial

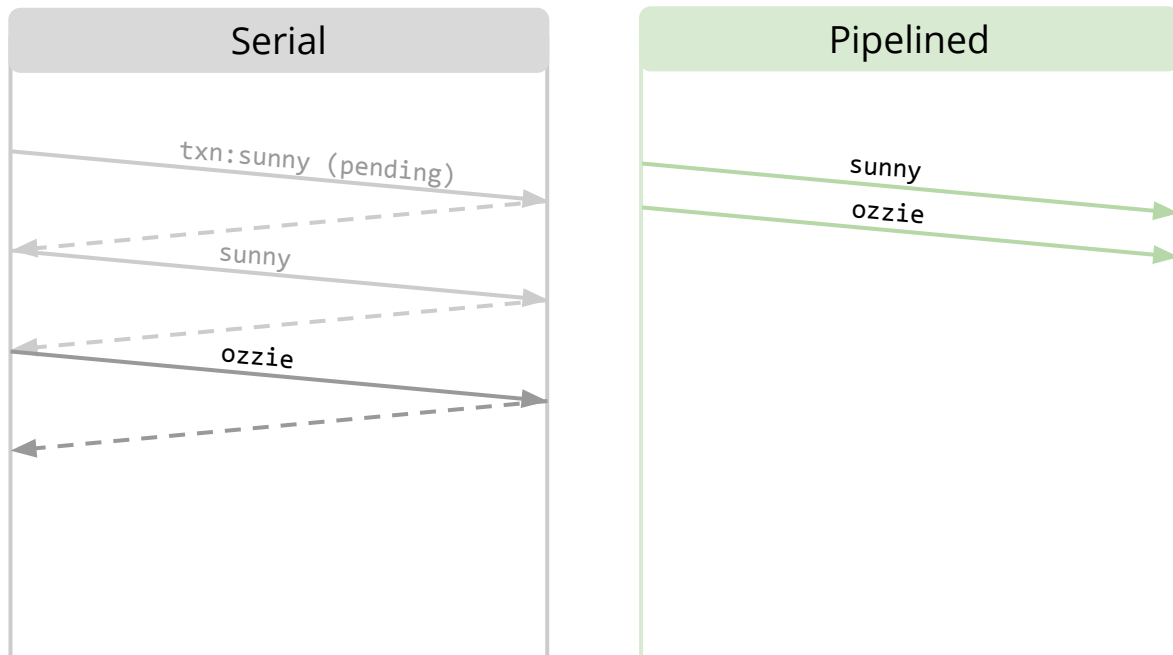
Pipelined

Transactions: Pipelining



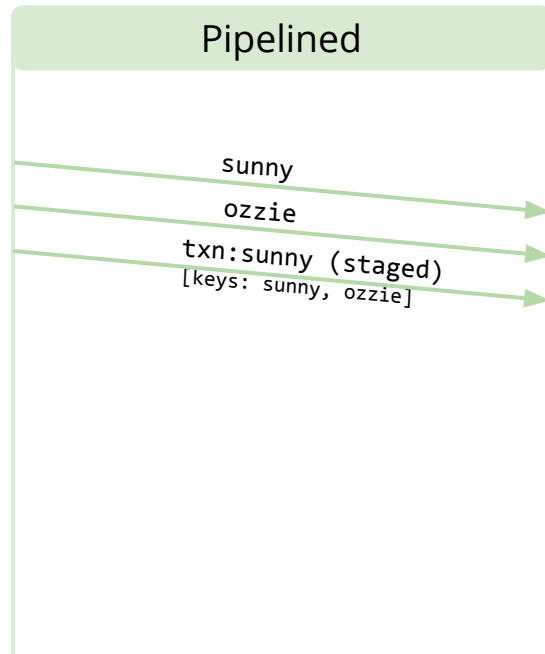
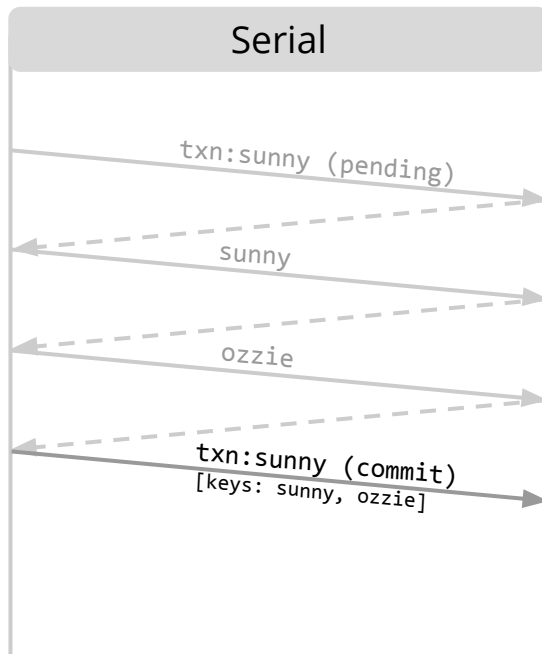
Transactions: Pipelining

```
BEGIN  
WRITE[sunny]  
WRITE[ozzie]
```

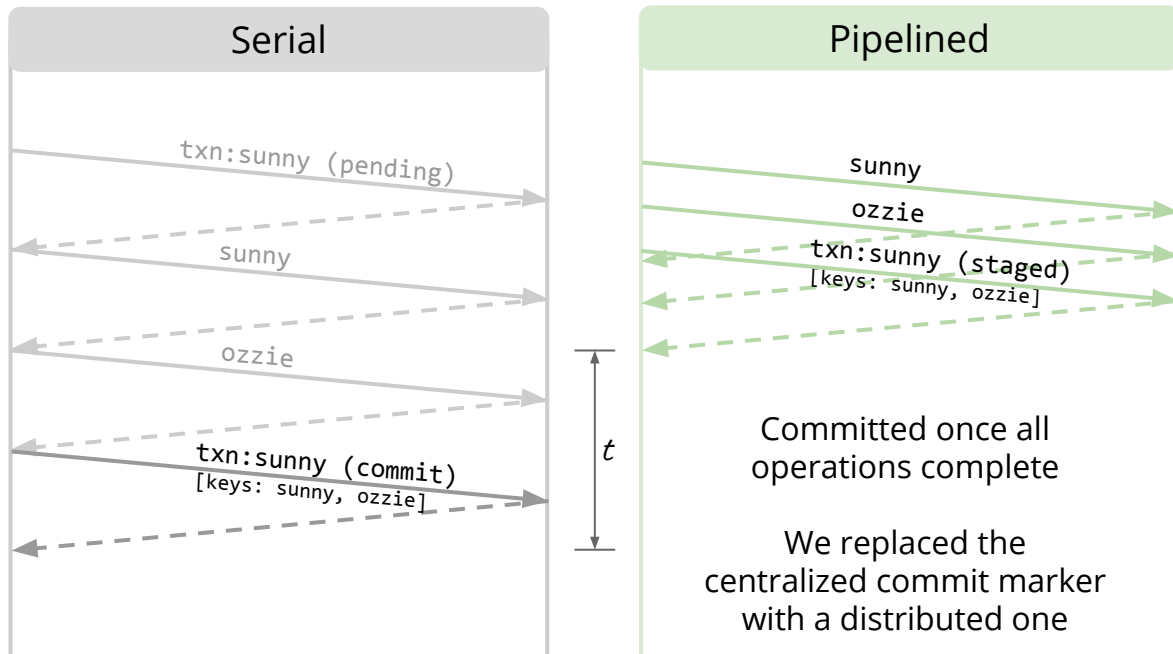
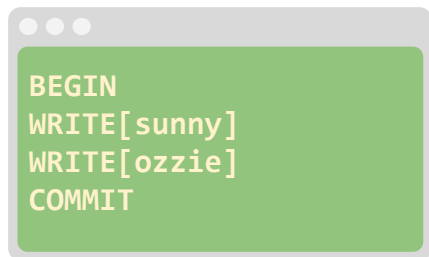


Transactions: Pipelining

```
BEGIN  
WRITE[sunny]  
WRITE[ozzie]  
COMMIT
```



Transactions: Pipelining



* "Proved" with TLA+

SQL

Structured Query Language

Declarative, not imperative

- These are the results I want vs perform these operations in this sequence

Relational data model

- Typed: INT, FLOAT, STRING, ...
- Schemas: tables, rows, columns, foreign keys

SQL: Tabular Data in a KV World

SQL data has columns and types?!?

How do we store typed and columnar data in a distributed, replicated, transactional key-value store?

- The SQL data model needs to be mapped to KV data
- Reminder: keys and values are lexicographically sorted

SQL Data Mapping: Inventory Table

```
CREATE TABLE inventory (  
    id INT PRIMARY KEY,  
    name STRING,  
    price FLOAT  
)
```

ID	Name	Price
1	Bat	1.11
2	Ball	2.22
3	Glove	3.33

Key	Value
/1	"Bat",1.11
/2	"Ball",2.22
/3	"Glove",3.33

SQL Data Mapping: Inventory Table

```
CREATE TABLE inventory (  
    id INT PRIMARY KEY,  
    name STRING,  
    price FLOAT  
)
```

ID	Name	Price
1	Bat	1.11
2	Ball	2.22
3	Glove	3.33

Key	Value
/<Table>/<Index>/1	"Bat",1.11
/<Table>/<Index>/2	"Ball",2.22
/<Table>/<Index>/3	"Glove",3.33

SQL Data Mapping: Inventory Table

```
CREATE TABLE inventory (  
    id INT PRIMARY KEY,  
    name STRING,  
    price FLOAT  
)
```

ID	Name	Price
1	Bat	1.11
2	Ball	2.22
3	Glove	3.33

Key	Value
/inventory/primary/1	"Bat",1.11
/inventory/primary/2	"Ball",2.22
/inventory/primary/3	"Glove",3.33

SQL Data Mapping: Inventory Table

```
CREATE TABLE inventory (  
    id INT PRIMARY KEY,  
    name STRING,  
    price FLOAT,  
    INDEX name_idx (name)  
)
```

ID	Name	Price
1	Bat	1.11
2	Ball	2.22
3	Glove	3.33

Key	Value
/inventory/name_idx/"Bat"/1	Ø
/inventory/name_idx/"Ball"/2	Ø
/inventory/name_idx/"Glove"/3	Ø

SQL Data Mapping: Inventory Table

```
CREATE TABLE inventory (  
    id INT PRIMARY KEY,  
    name STRING,  
    price FLOAT,  
    INDEX name_idx (name)  
)
```

ID	Name	Price
1	Bat	1.11
2	Ball	2.22
3	Glove	3.33
4	Bat	4.44

Key	Value
/inventory/name_idx/"Bat"/1	Ø
/inventory/name_idx/"Ball"/2	Ø
/inventory/name_idx/"Glove"/3	Ø

SQL Data Mapping: Inventory Table

```
CREATE TABLE inventory (  
    id INT PRIMARY KEY,  
    name STRING,  
    price FLOAT,  
    INDEX name_idx (name)  
)
```

ID	Name	Price
1	Bat	1.11
2	Ball	2.22
3	Glove	3.33
4	Bat	4.44

Key	Value
/inventory/name_idx/"Bat"/1	∅
/inventory/name_idx/"Ball"/2	∅
/inventory/name_idx/"Glove"/3	∅
/inventory/name_idx/"Bat"/4	∅