



# The Hitchhiker's Guide to Distributed Transactions

@irfansharifm

# Abstract



Recent years have seen a flurry of research, both from academia and industry, enabling distributed transactions while minimizing the performance overhead. There are various proposals afoot; some systems constrain the transaction model, others will trade-off behavior under contention. Some systems sacrifice isolation guarantees, others will incur higher read latencies.

This talk is a survey of the various approaches that academic (Carousel, MDCC, SLOG, TAPIR) and industrial (Spanner, CockroachDB, OceanVista) systems use to achieve atomicity in their transactions. We'll define a shared terminology (ranges, replicas, txn records, etc.) and use it to explore how system unique compose their transaction models with the underlying replication protocol to achieve the theoretical minimum latency for atomic commitment: one round-trip between data centers.



# Agenda

1. Foundations
2. Transactions
3. Implementations



# Agenda

1. **Foundations**
2. Transactions
3. Implementations

- I. Keyspace and Sharding
- II. Replication and Fault Tolerance
- III. APIs



# Agenda

1. Foundations
- 2. Transactions**
3. Implementations

- I. ACID & Isolation Levels
- II. Transaction Basics
- III. Unpipelined Transactions



# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

- I. Spanner/Pipelined Transactions
- II. Parallel Commits
- III. Replicated Commit
- IV. Carousel
- V. MDCC
- VI. SLOG/OceanVista
- VII. TAPIR



# Agenda

1. **Foundations**
2. Transactions
3. Implementations

- I. Keyspace and Sharding
- II. Replication and Fault Tolerance
- III. APIs



## Monolithic, sorted, logical key space

We're ignoring non-ordered databases  
(think consistent hashing), though most of  
the same principles will apply

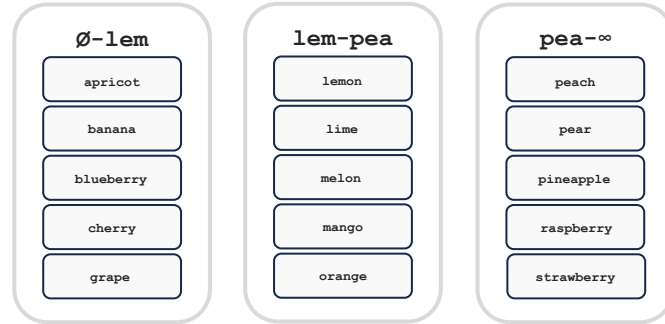






# Ranges

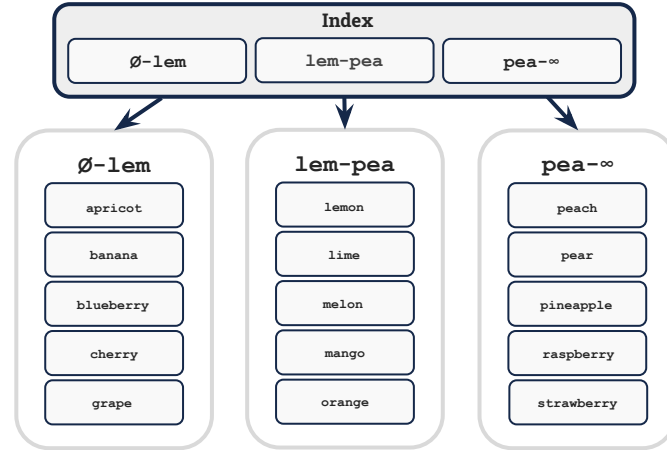
Each *Range* (think “shards”) holds a contiguous span of the key-space





# Indexes

We could use indexing structures to locate shards/ranges

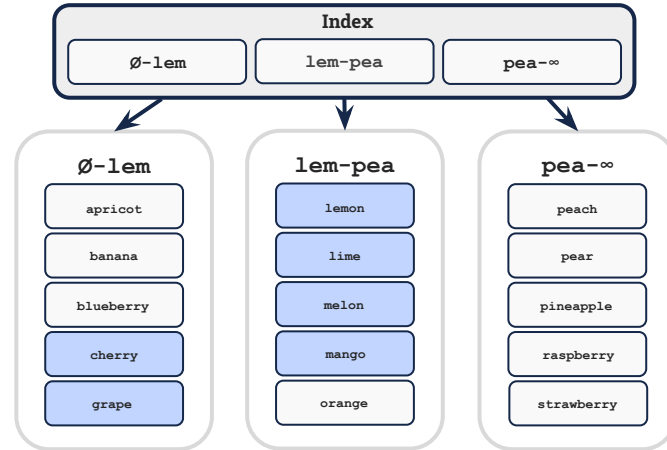




## Range scans

Ordered keys makes range scans efficient

```
fruits >= "cherry" AND <= "mango"
```



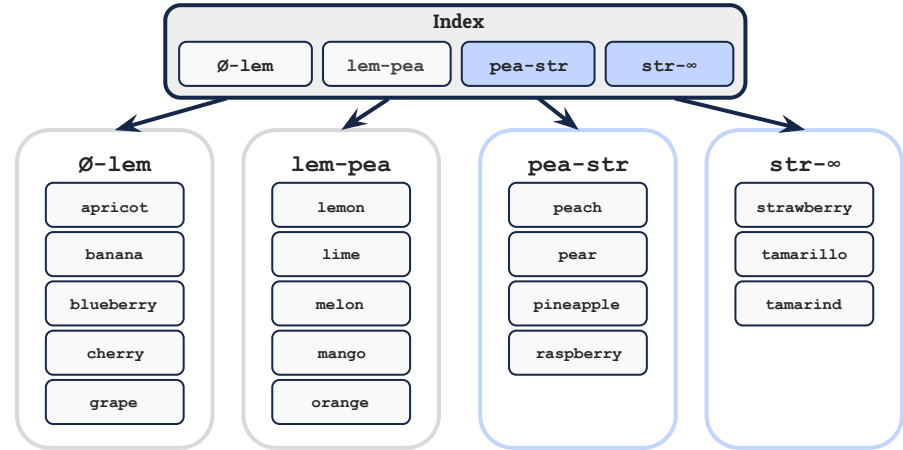


## Range splits/merges

Could split ranges when they get too large,  
merge when they get too small

We want range sizes to be:

- small enough to be moved quickly
- large enough to amortize indexing overhead





## Other systems

Bigtable calls these *tablets*, Hbase calls these *regions*, CRDB calls these *ranges*. Similar structures found in Spanner, YugaByte, SLOG, etc.



# Agenda

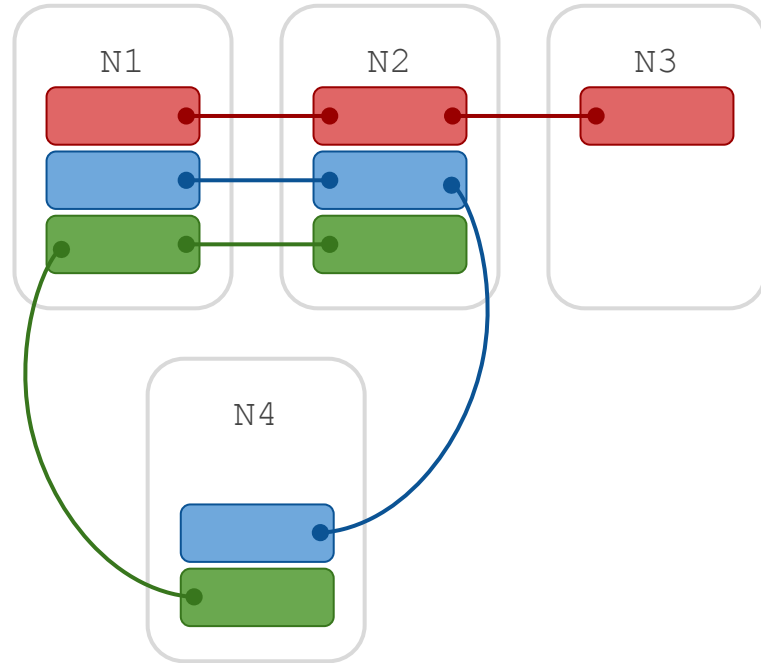
1. **Foundations**
2. Transactions
3. Implementations

- I. Keyspace and Sharding
- II. Replication and Fault Tolerance
- III. APIs



## Consensus replication

Ranges are the unit of replication, each copy is a *Replica*. A single node could hold one or more replicas.

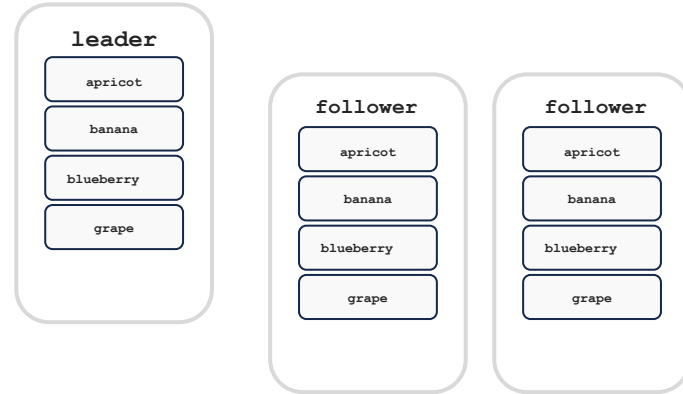




# Consensus replication

There are several variants we could use, leader-based (raft, multi-paxos) and leaderless (epaxos) ones

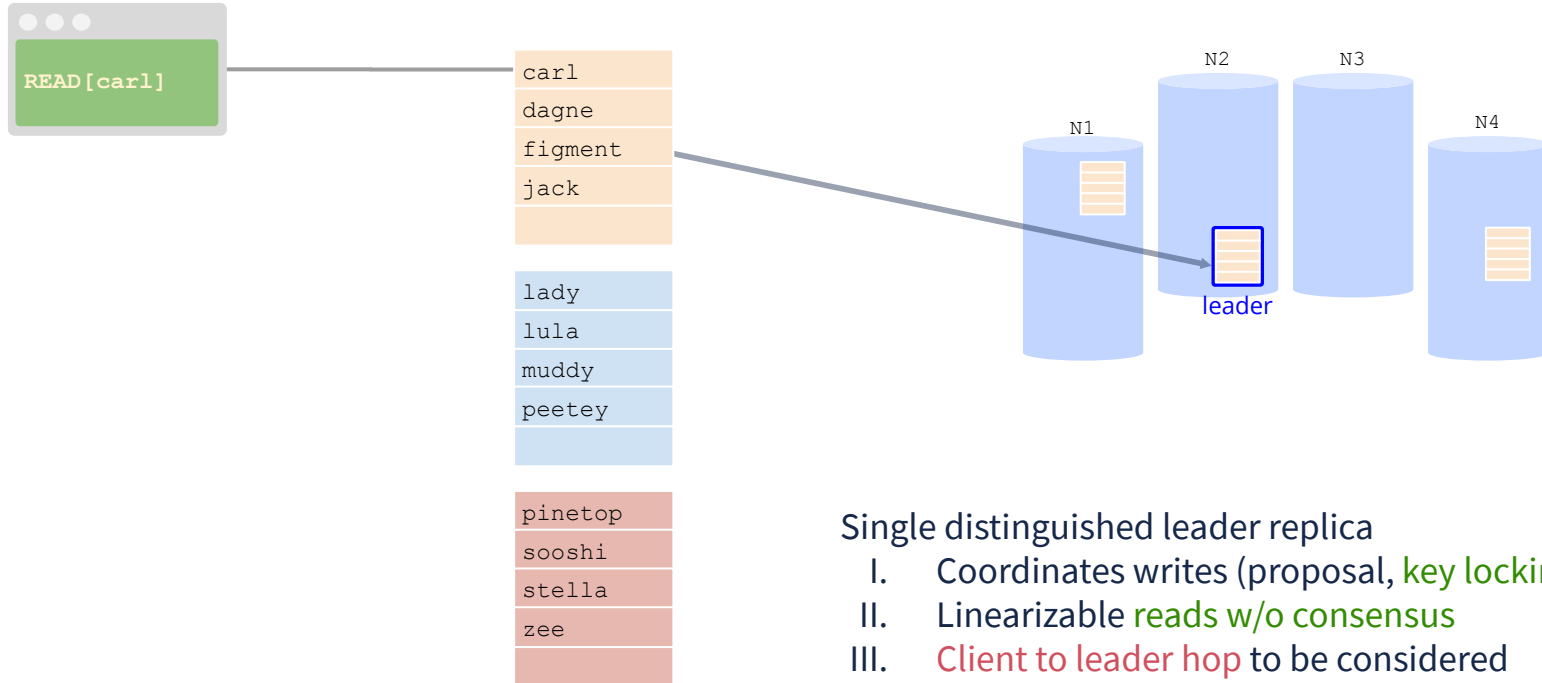
*NB: A replication factor of 2 doesn't make sense for us (but is somewhat akin to primary/secondary replication)*







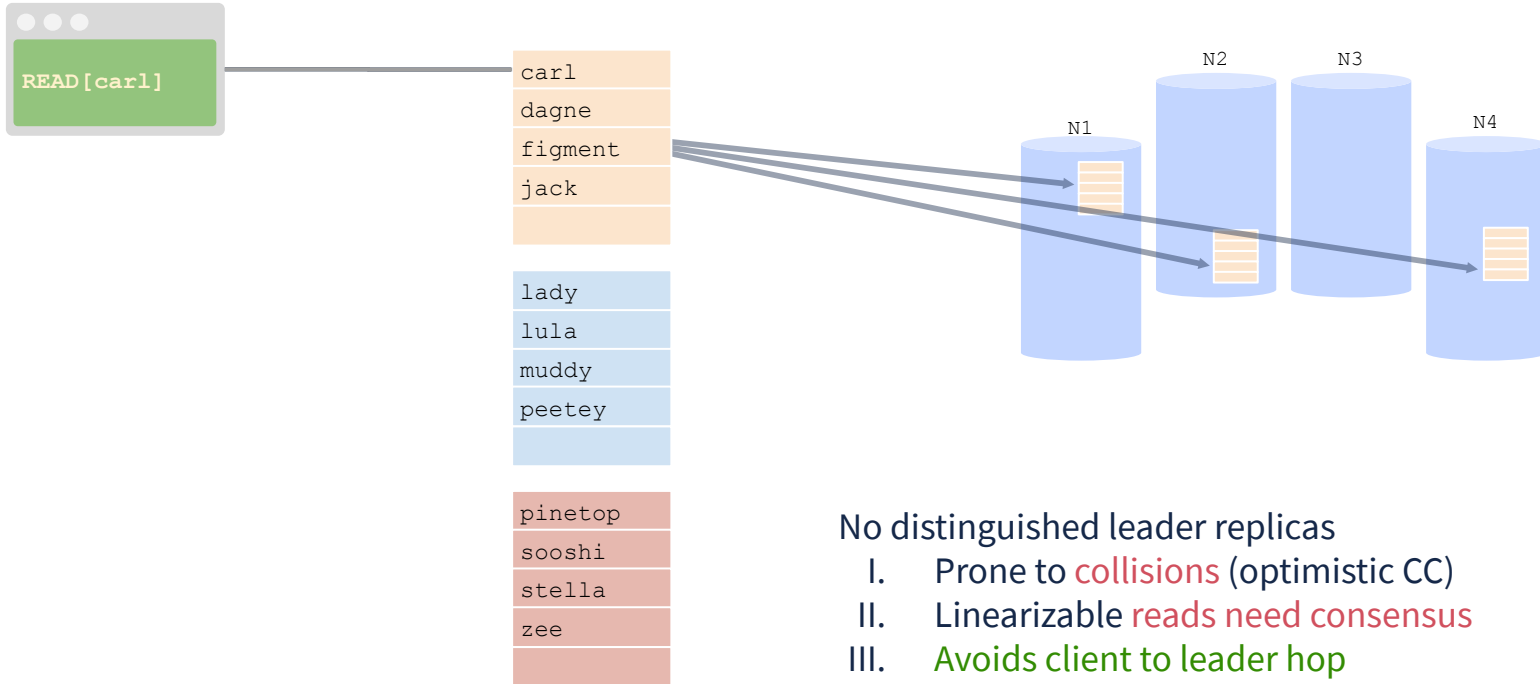
# Leader-based consensus



Single distinguished leader replica

- I. Coordinates writes (proposal, **key locking**)
- II. Linearizable **reads w/o consensus**
- III. **Client to leader hop** to be considered
- IV. **Leader node failures** can cause blips
- V. **Disproportionate** load on leader

# Leaderless consensus



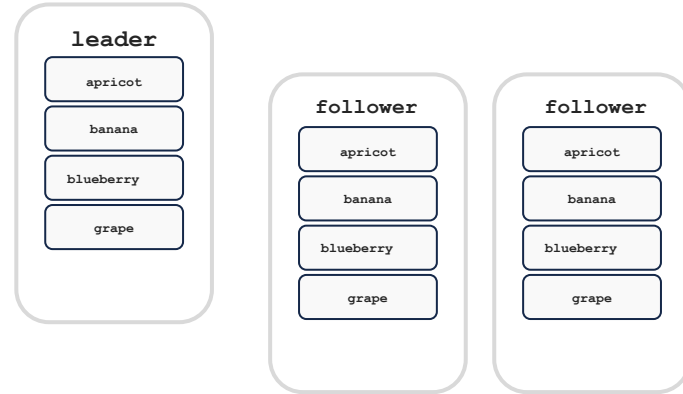
No distinguished leader replicas

- I. Prone to **collisions** (optimistic CC)
- II. Linearizable **reads need consensus**
- III. **Avoids client to leader hop**
- IV. **No blips** caused by node failures
- V. More **evenly distributed load** on nodes



## Consensus replication (Writes)

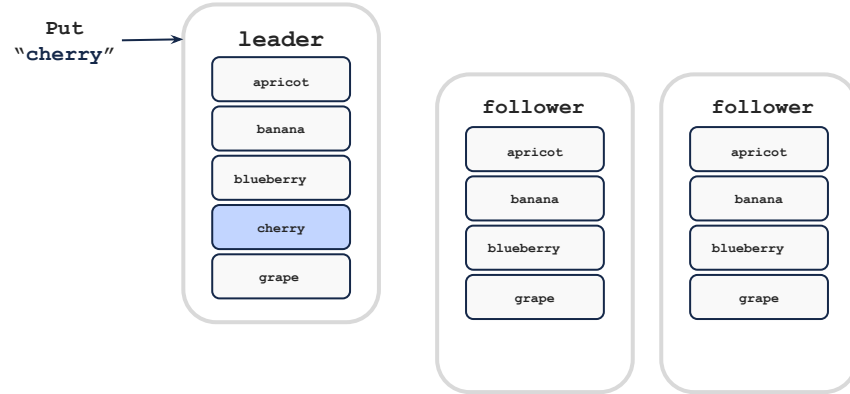
Replicate N-way, committed when acked  
by a quorum of replicas





## Consensus replication (Writes)

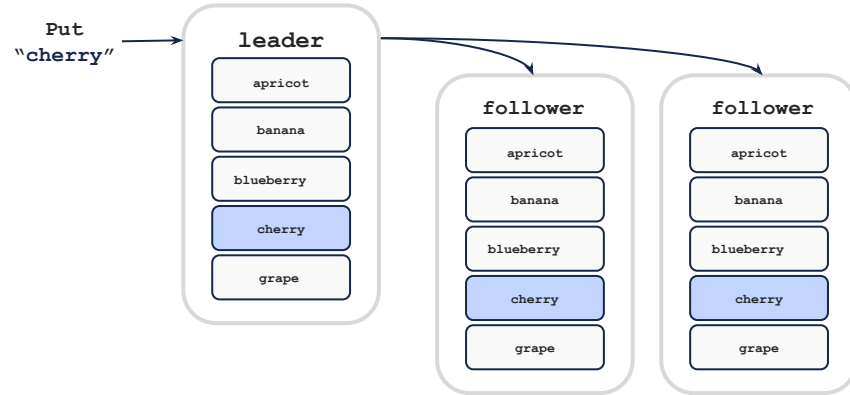
Replicate N-way, committed when acked  
by a quorum of replicas





## Consensus replication (Writes)

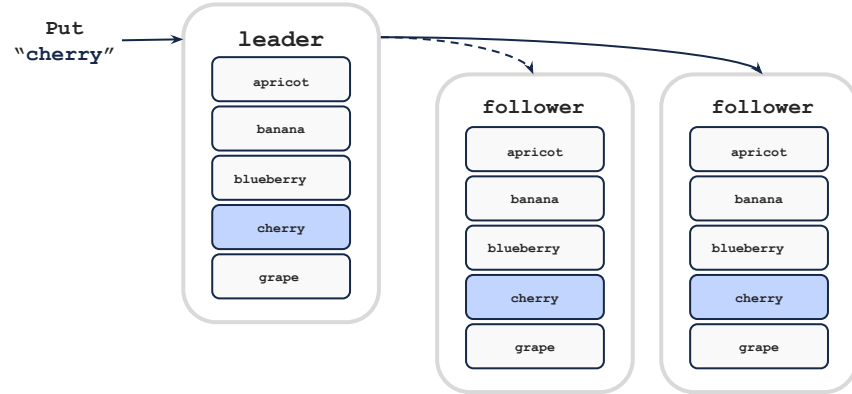
Replicate N-way, committed when acked  
by a quorum of replicas





# Consensus replication (Writes)

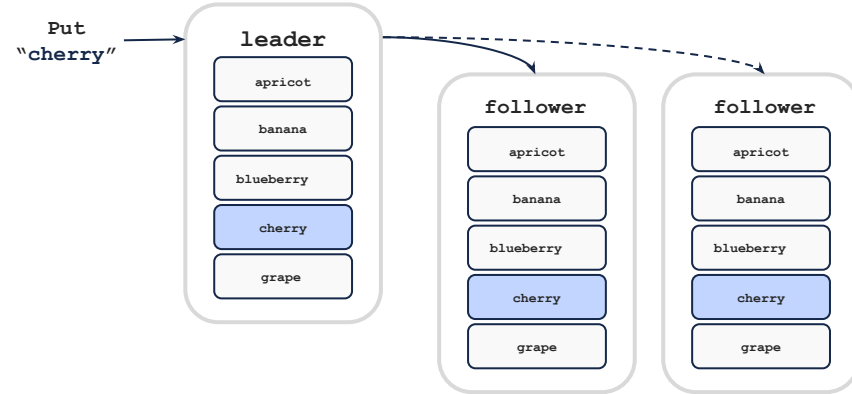
Need ack from just one other node





# Consensus replication (Writes)

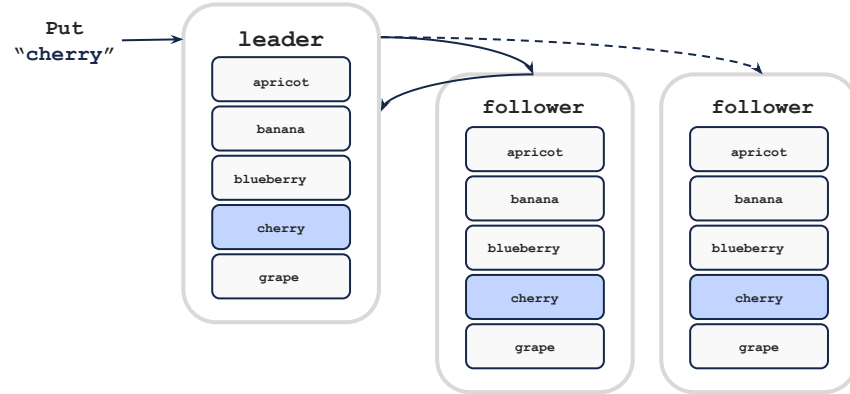
Need ack from just one other node





# Consensus replication (Writes)

Need ack from just one other node

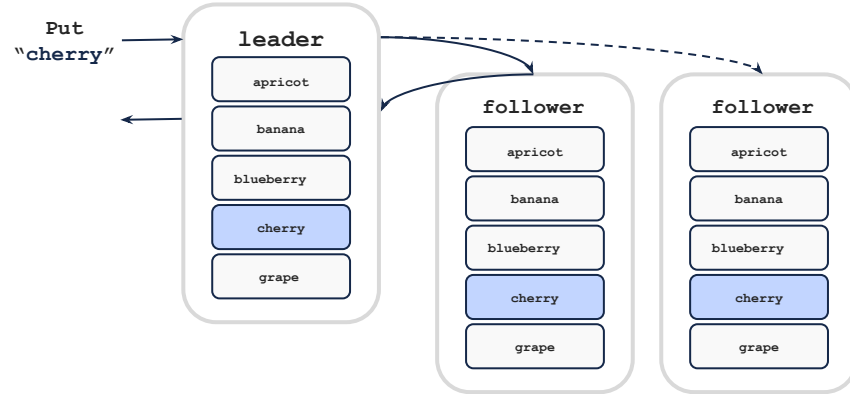






# Consensus replication (Writes)

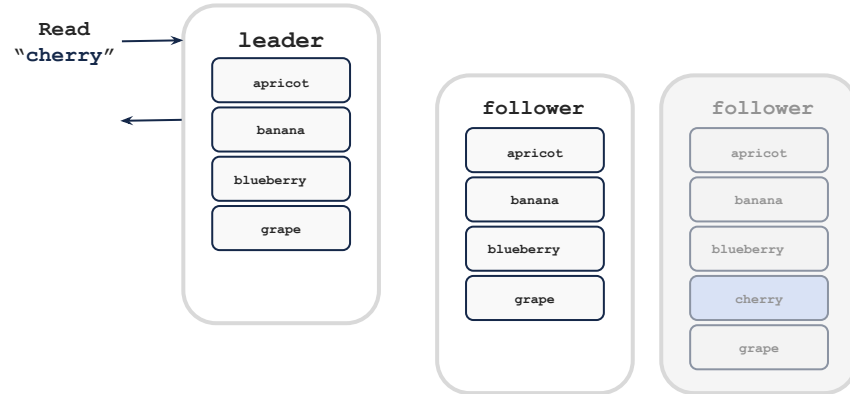
Need ack from just one other node





## Consensus replication (Reads)

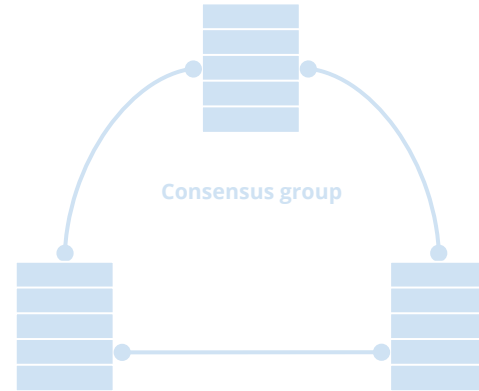
Only data written to a quorum is considered present





# Consensus replication

Consensus provides “durable, atomic replication” of commands





# Agenda

1. **Foundations**
2. Transactions
3. Implementations

- I. Keyspace and Sharding
- II. Replication and Fault Tolerance
- III. APIs

# Data Mapping: SQL



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

```
/<table>/<index>/<key>
```

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

# Data Mapping: SQL



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

```
/<table>/<index>/<key>
```

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



# Agenda

1. Foundations
- 2. Transactions**
3. Implementations

- I. ACID & Isolation Levels
- II. Transaction Basics
- III. Unpipelined Transactions



# Agenda

1. Foundations
- 2. Transactions**
3. Implementations

- I. ACID & Isolation Levels
- II. Transaction Basics
- III. Unpipelined Transactions



# ACID



- I. Atomicity (“all or nothing”; achieved using staged writes and txn records)
- II. Consistency (“db is internally consistent”)
- III. Isolation (“effects of concurrent txns on each other”; determined by locks)
- IV. Durability (“don’t lose committed data”; determined by synced writes)

# ACID



- I. Atomicity (“all or nothing”; achieved using staged writes and txn records)
- II. Consistency (“db is internally consistent”)
- III. Isolation (“effects of concurrent txns on each other”; determined by locks)
- IV. Durability (“don’t lose committed data”; determined by synced writes)

# ACID



- I. Atomicity (“all or nothing”; achieved using staged writes and txn records)
- II. Consistency (“db is internally consistent”)
- III. Isolation (“effects of concurrent txns on each other”; determined by locks)
- IV. Durability (“don’t lose committed data”; determined by synced writes)

# ACID



- I. Atomicity (“all or nothing”; achieved using staged writes and txn records)
- II. Consistency (“db is internally consistent”)
- III. Isolation (“effects of concurrent txns on each other”; determined by locks)
- IV. Durability (“don’t lose committed data”; determined by synced writes)



# Isolation

Isolation Level (“effects of concurrent txns on one another”; determined by locking granularity)

- I. read uncommitted (could read ongoing-but-uncommitted writes)
- II. read committed (could read different values in the same txn)
- III. repeatable read (could range-read different values in the same txn)
- IV. snapshot isolation (could make write decisions based on stale reads)
- V. serializable (none of the above, as if run in serial order)



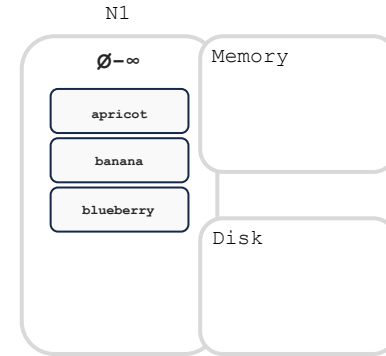
# Agenda

1. Foundations
- 2. Transactions**
3. Implementations

- I. ACID & Isolation Levels
- II. Transaction Basics
- III. Unpipelined Transactions



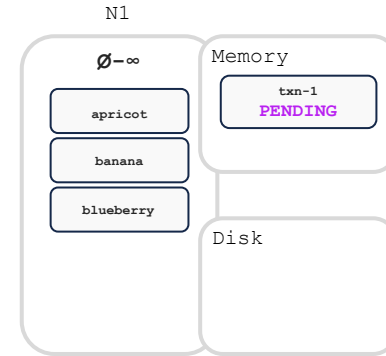
# Single-node transactions





# Single-node transactions

`BEGIN TXN;`

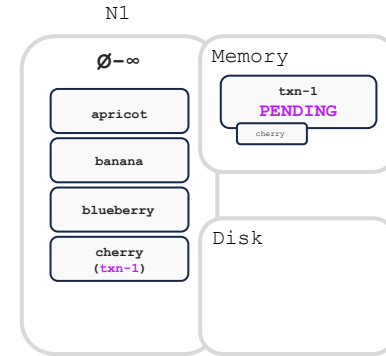






# Single-node transactions

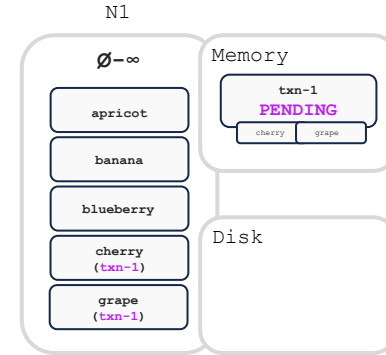
```
BEGIN TXN;  
  PUT "cherry";
```





# Single-node transactions

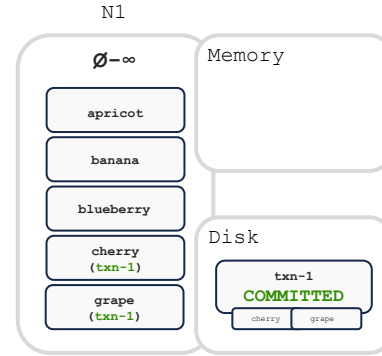
```
BEGIN TXN;  
  PUT "cherry";  
  PUT "grape";
```





# Single-node transactions

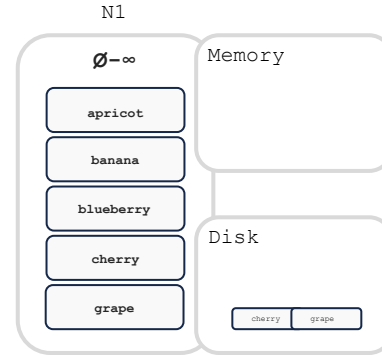
```
BEGIN TXN;  
  PUT "cherry";  
  PUT "grape";  
COMMIT;
```





# Single-node transactions

```
BEGIN TXN;  
    PUT "cherry";  
    PUT "grape";  
COMMIT;  
-- clean-up txn record
```



# Single-node (unreplicated) transactions



Atomicity and durability are achieved by bootstrapping off a lower-level atomic/durable primitive: **disk writes (fsync)**

- PENDING transaction record, with transaction ID
- Staged (**in-memory**) writes, tagged with transaction ID
- Durably persist COMMITTED **transaction record and staged writes**, atomically

# Multi-node (replicated) transactions

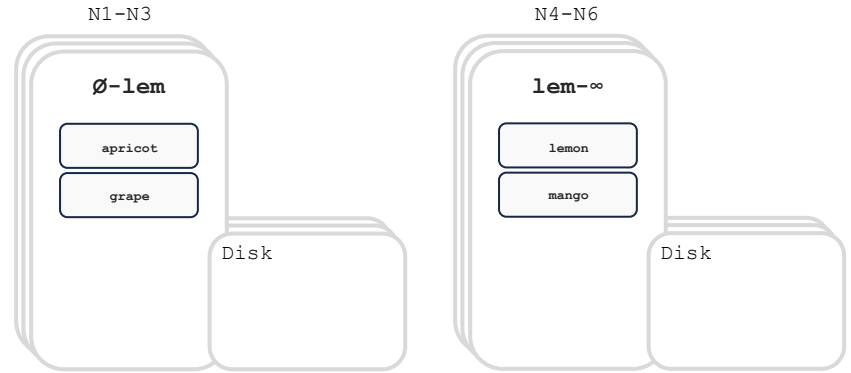


Atomicity and durability are achieved by bootstrapping off a lower-level atomic/durable primitive: `consensus writes (RTT + fsync)`

- `PENDING` transaction record, with transaction ID
- Staged (`consensus`) writes, tagged with transaction ID
- Durably persist `COMMITTED` `transaction record`, atomically



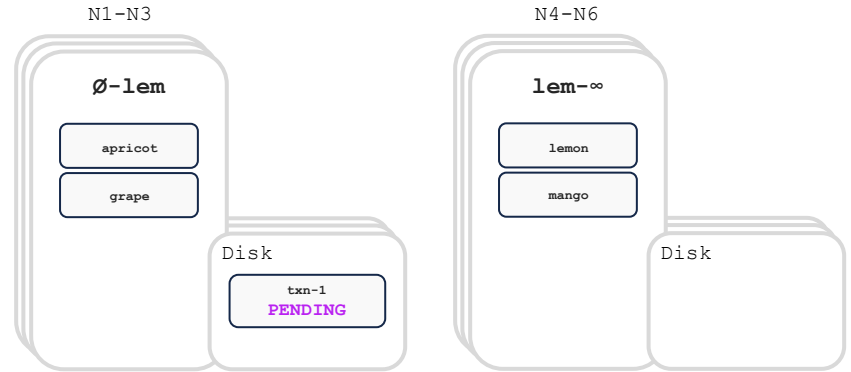
# Multi-node transactions





# Multi-node transactions

`BEGIN TXN;`

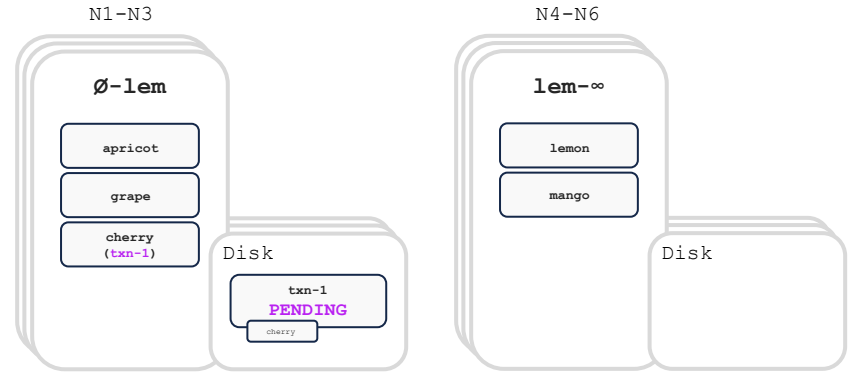






# Multi-node transactions

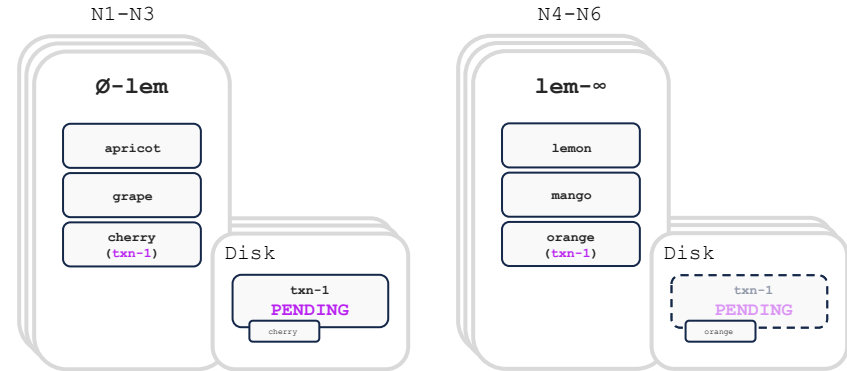
```
BEGIN TXN;  
  PUT "cherry";
```





# Multi-node transactions

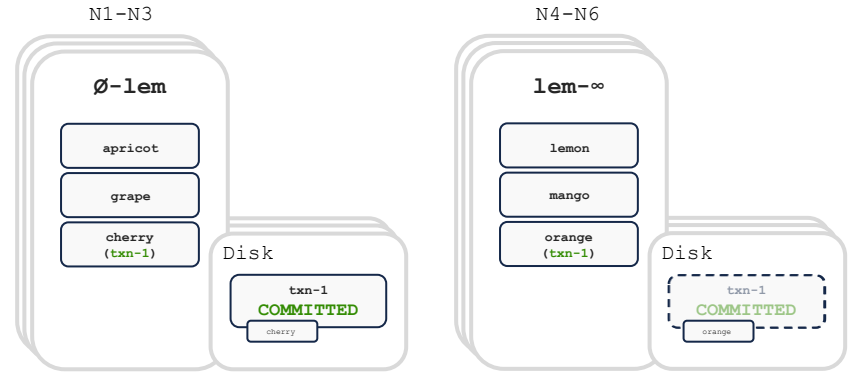
```
BEGIN TXN;  
  PUT "cherry";  
  PUT "orange";
```





# Multi-node transactions

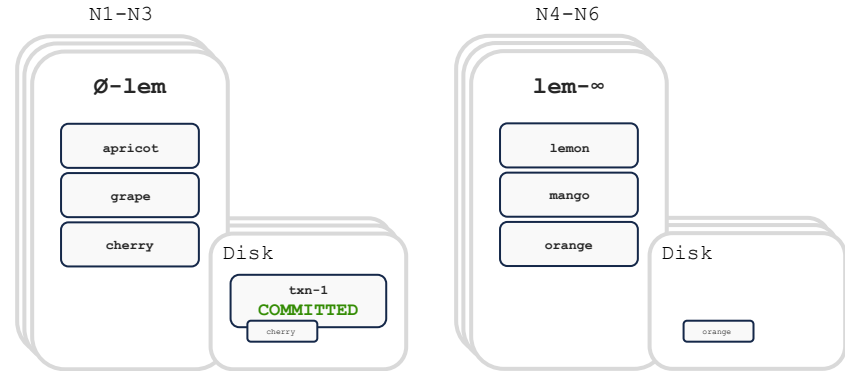
```
BEGIN TXN;  
  PUT "cherry";  
  PUT "orange";  
COMMIT;
```





## Multi-node transactions

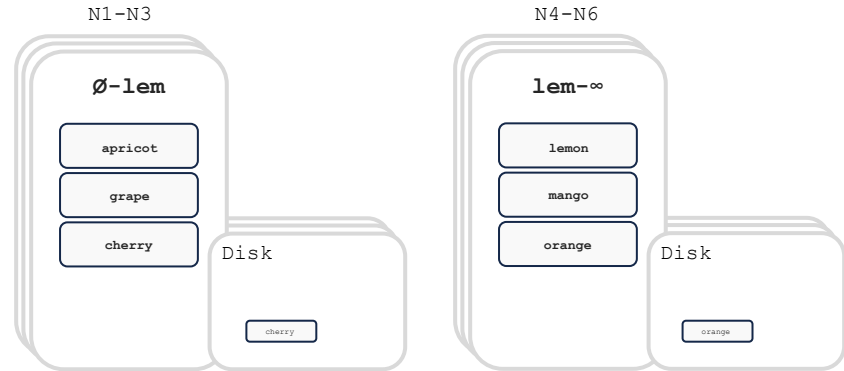
```
BEGIN TXN;  
  PUT "cherry";  
  PUT "orange";  
COMMIT;  
-- clean-up staged writes
```





## Multi-node transactions

```
BEGIN TXN;  
  PUT "cherry";  
  PUT "orange";  
COMMIT;  
-- clean-up staged writes  
-- clean-up txn record
```





# Agenda

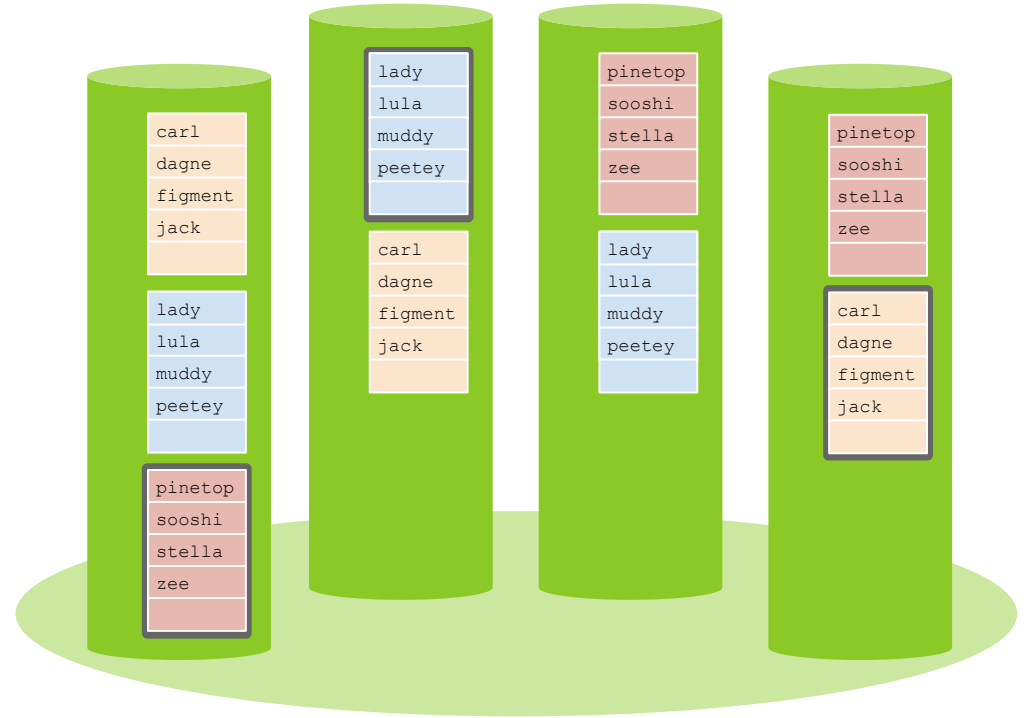
1. Foundations
- 2. Transactions**
3. Implementations

- I. ACID & Isolation Levels
- II. Transaction Basics
- III. Unpipelined Transactions

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```



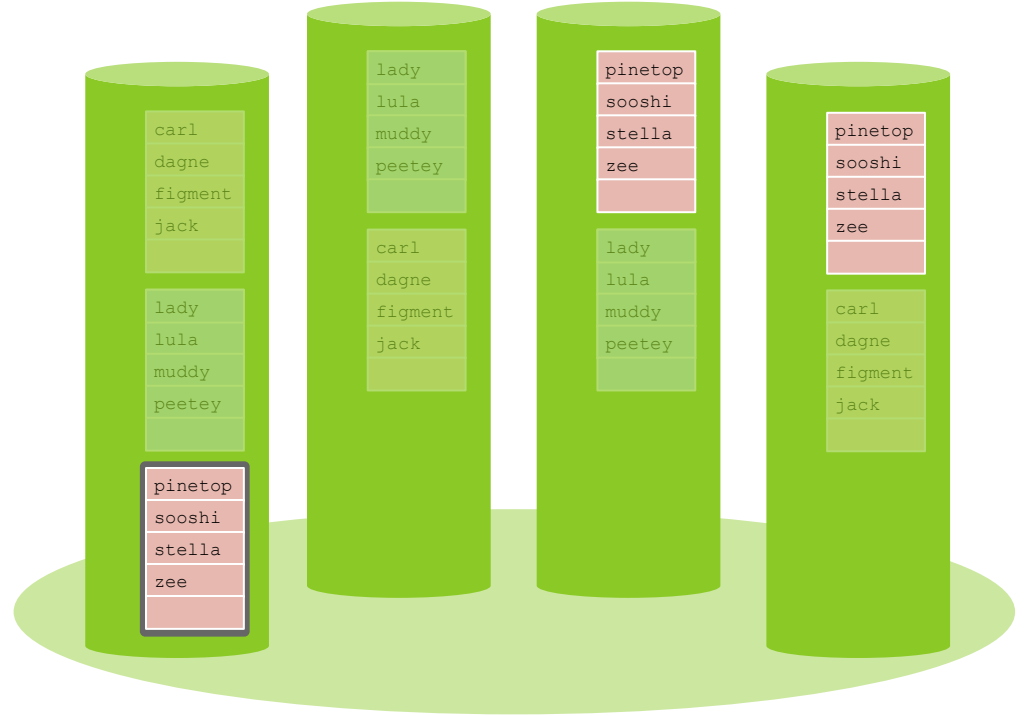
t = 0

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
  
GATEWAY
```



$t = 0$



# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]
```

GATEWAY

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

transactions  
TXN1: PENDING

lady
lula
muddy
peetey
carl
dagne
figment
jack

pinetop
sooshi
stella
zee
lady
lula
muddy
peetey

pinetop
sooshi
stella
zee
carl
dagne
figment
jack

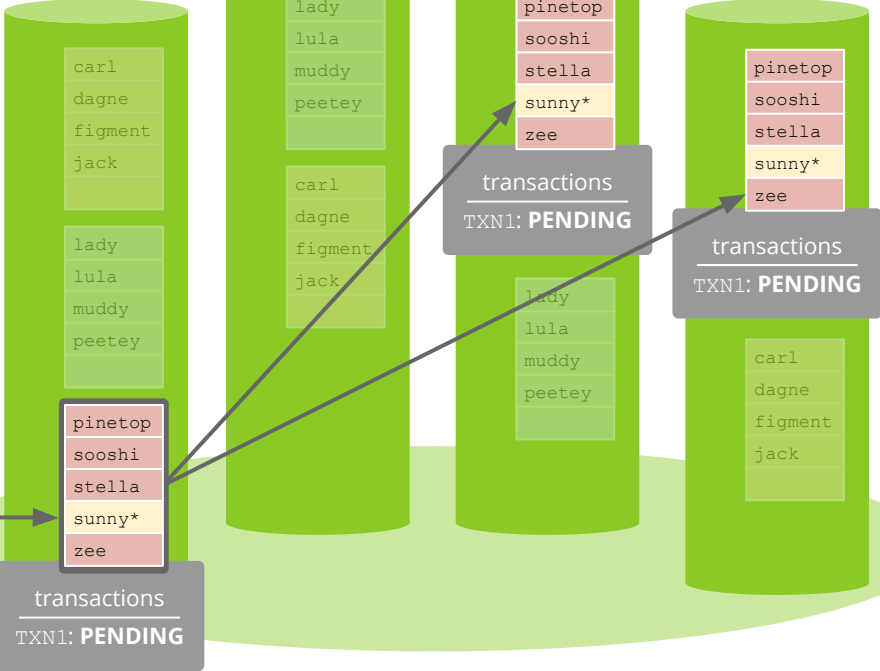
t = 0

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE[sunny]  
  
GATEWAY
```



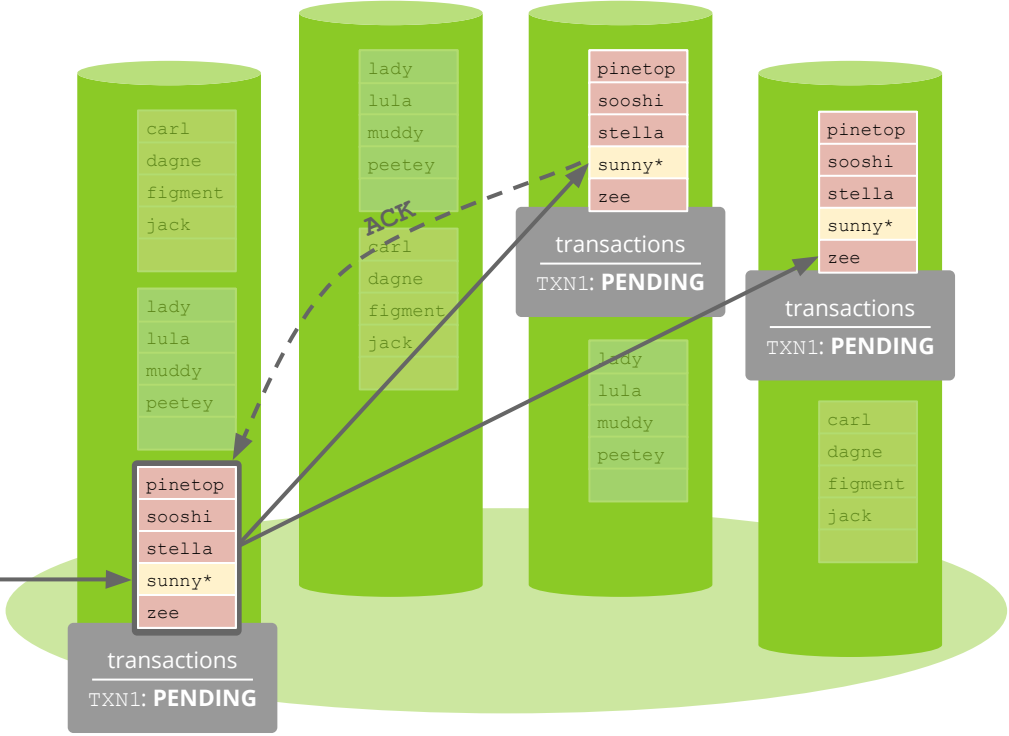
$$t = \frac{1}{2} \text{RTT}$$

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
  
GATEWAY
```



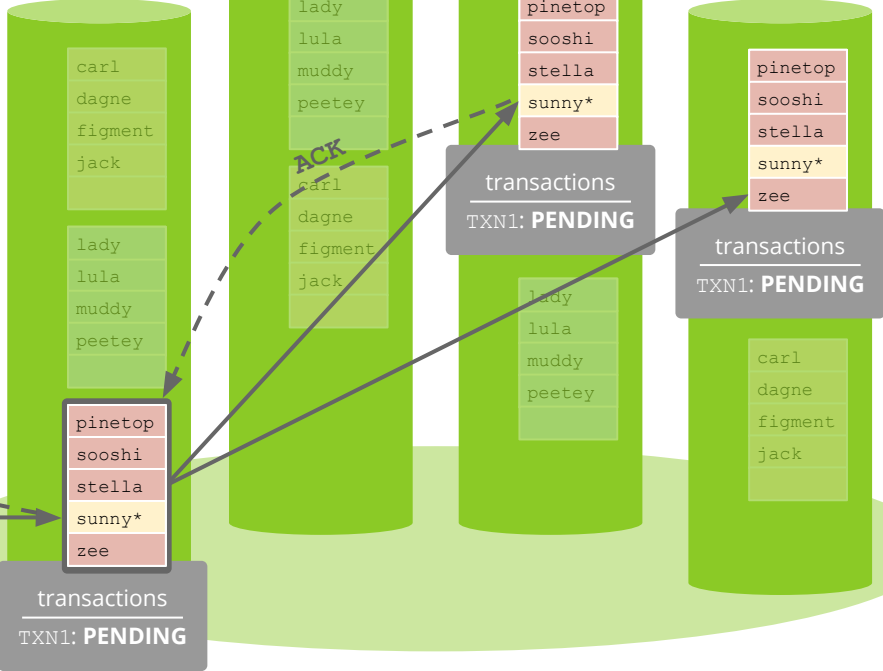
t = 1 RTT

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE[sunny]  
  
GATEWAY
```



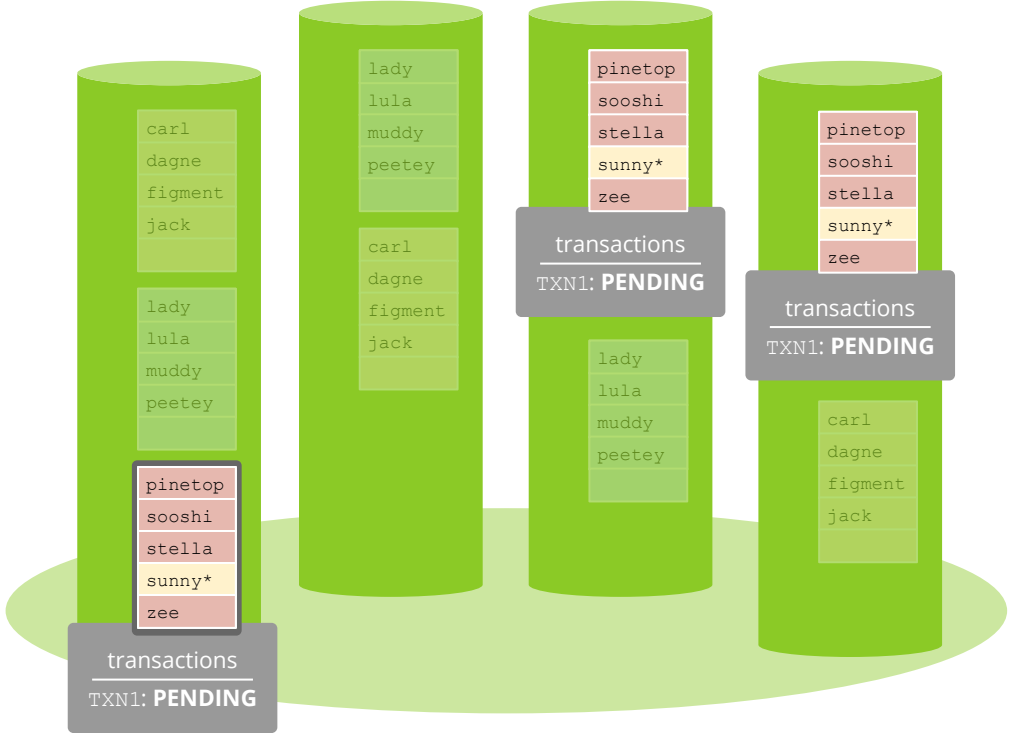
$t = 1 \text{ RTT}$

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
  
GATEWAY
```



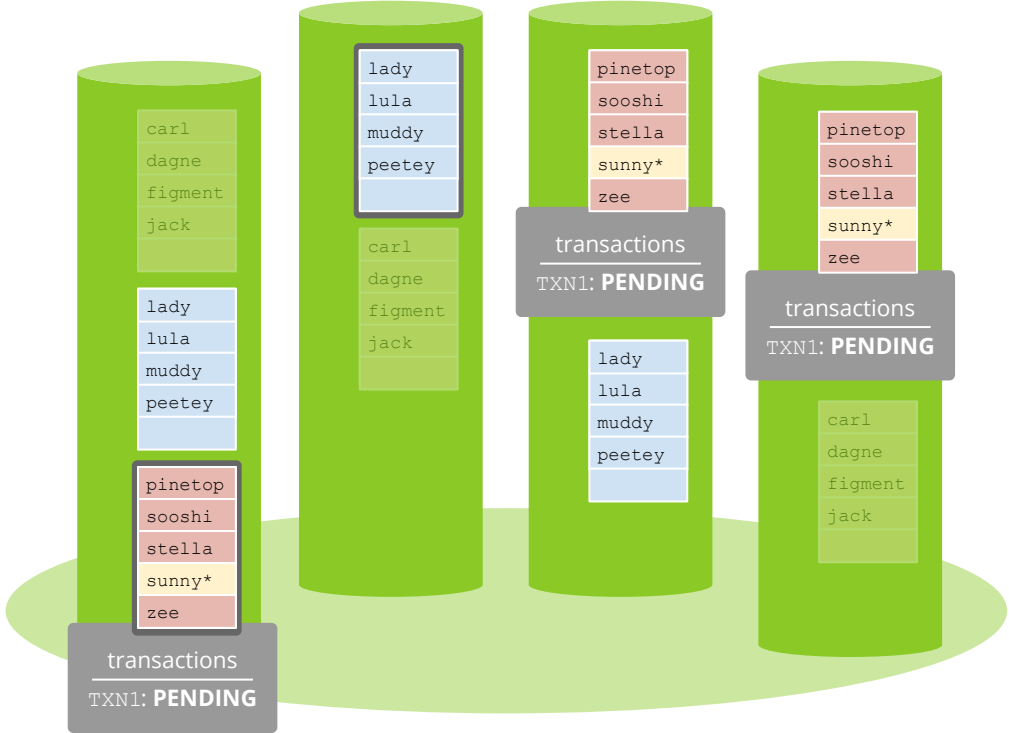
t = 1 RTT

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
GATEWAY
```



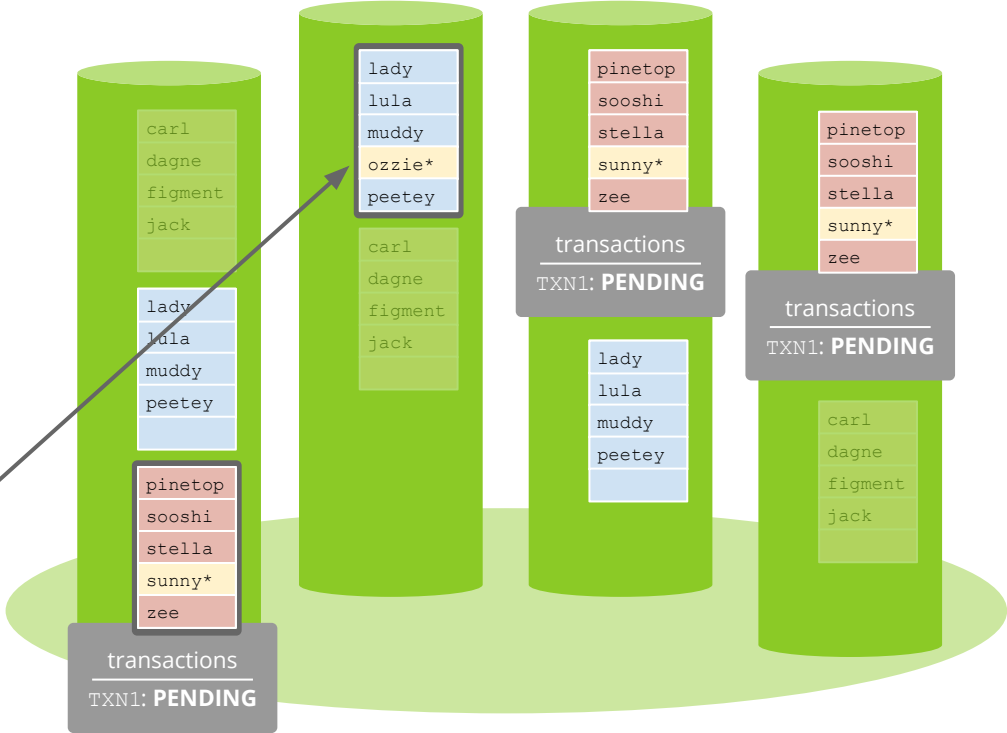
t = 1 RTT

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
GATEWAY
```



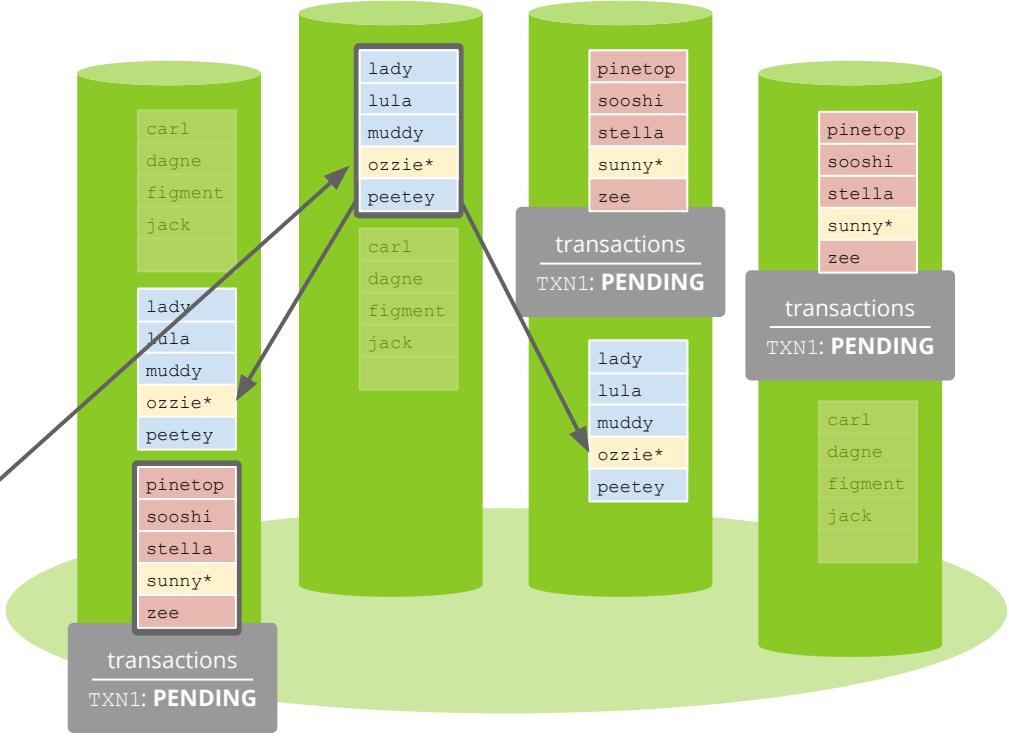
t = 1 RTT

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
GATEWAY
```



$t = 1\frac{1}{2} \text{ RTT}$



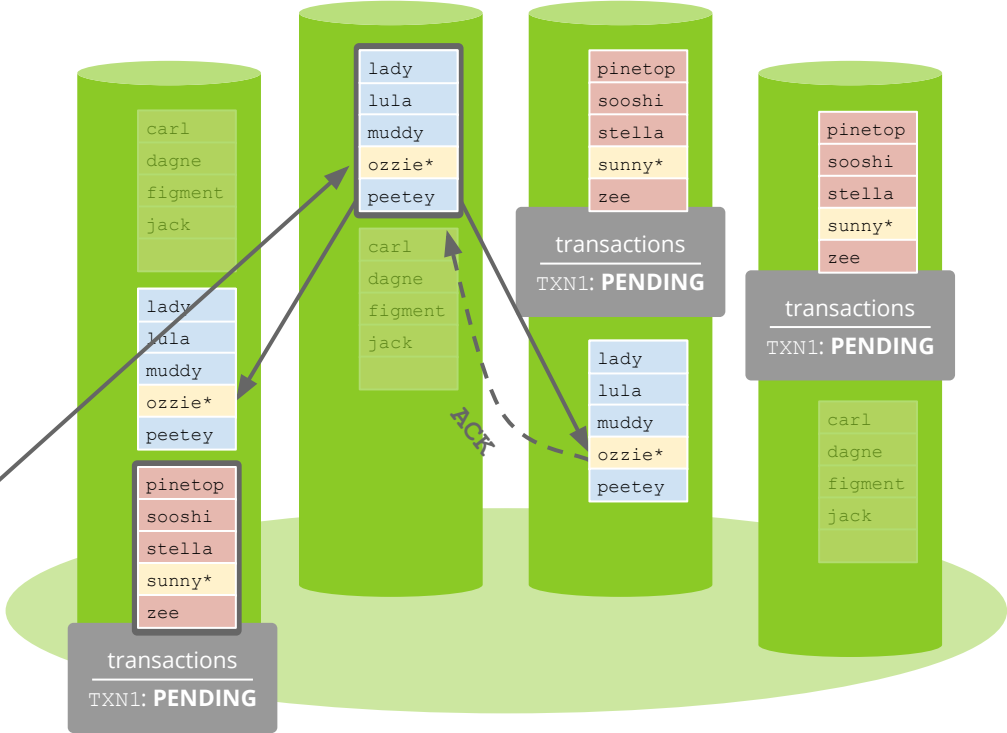
# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]
```

GATEWAY



t = 2 RTT

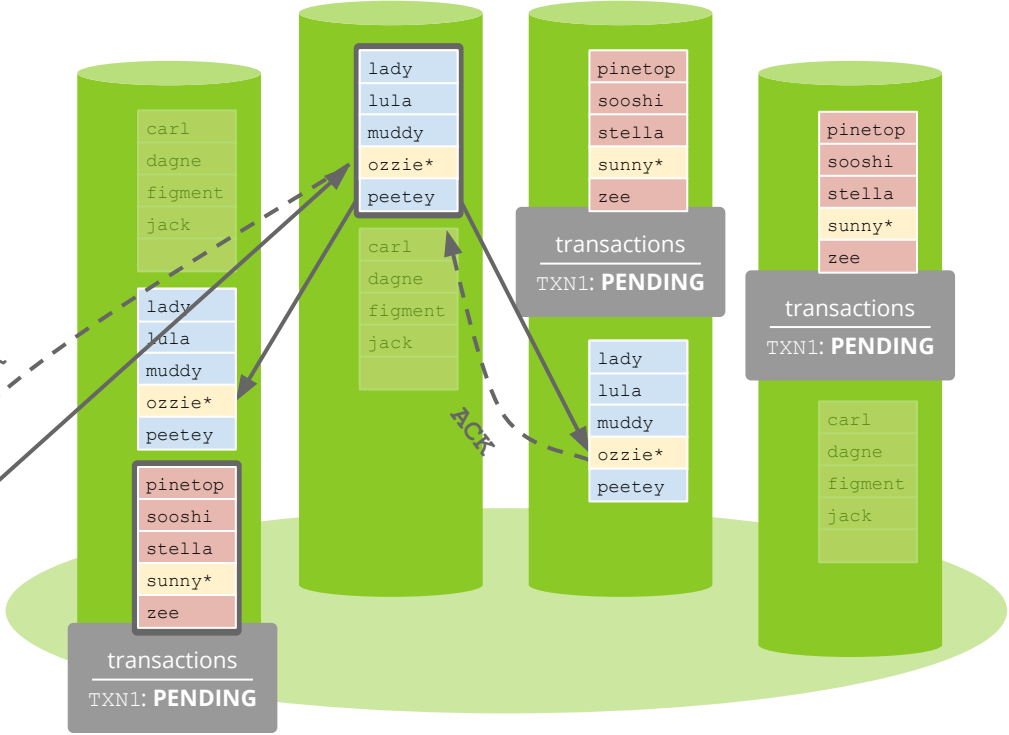
# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]
```

GATEWAY



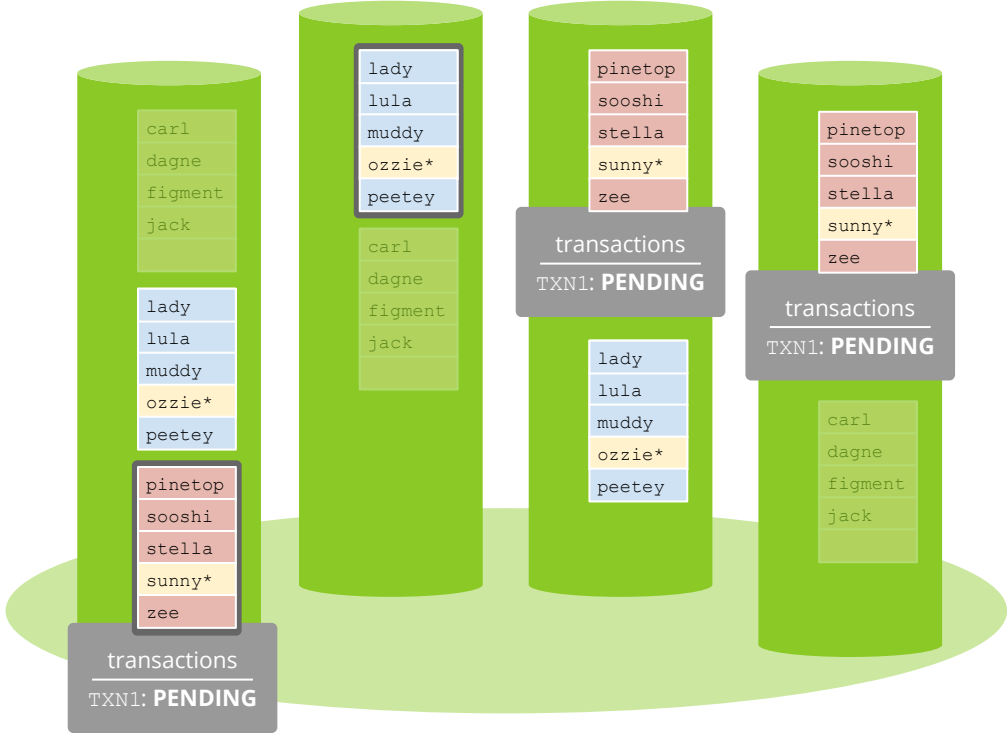
t = 2 RTT

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
GATEWAY
```



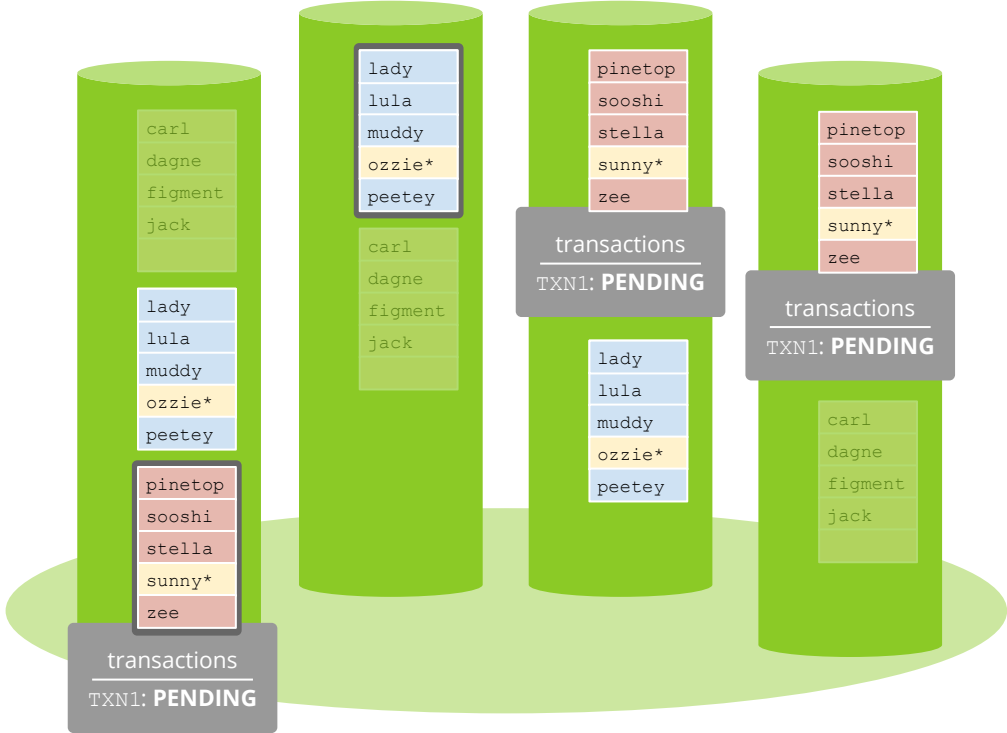
t = 2 RTT

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 2 RTT

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```

transactions  
TXN1: **COMMITTED**

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

lady
lula
muddy
ozzie*
peetey

carl
dagne
figment
jack

transactions  
TXN1: **PENDING**

pinetop
sooshi
stella
sunny*
zee

lady
lula
muddy
ozzie*
peetey

transactions  
TXN1: **PENDING**

pinetop
sooshi
stella
sunny*
zee

carl
dagne
figment
jack

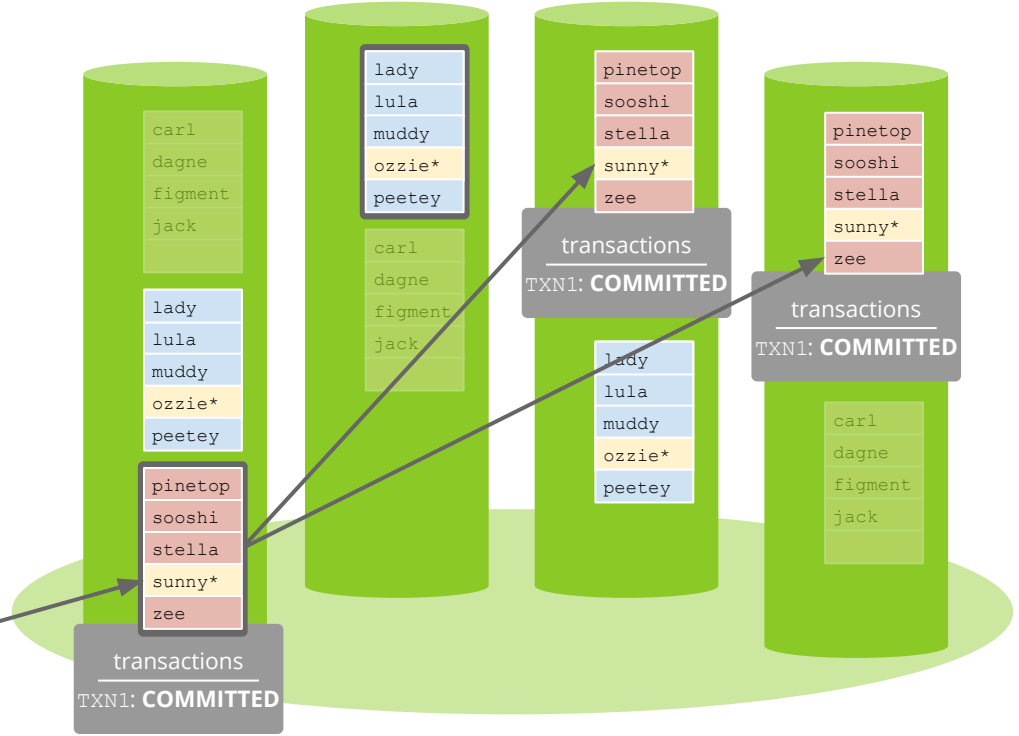
t = 2 RTT

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



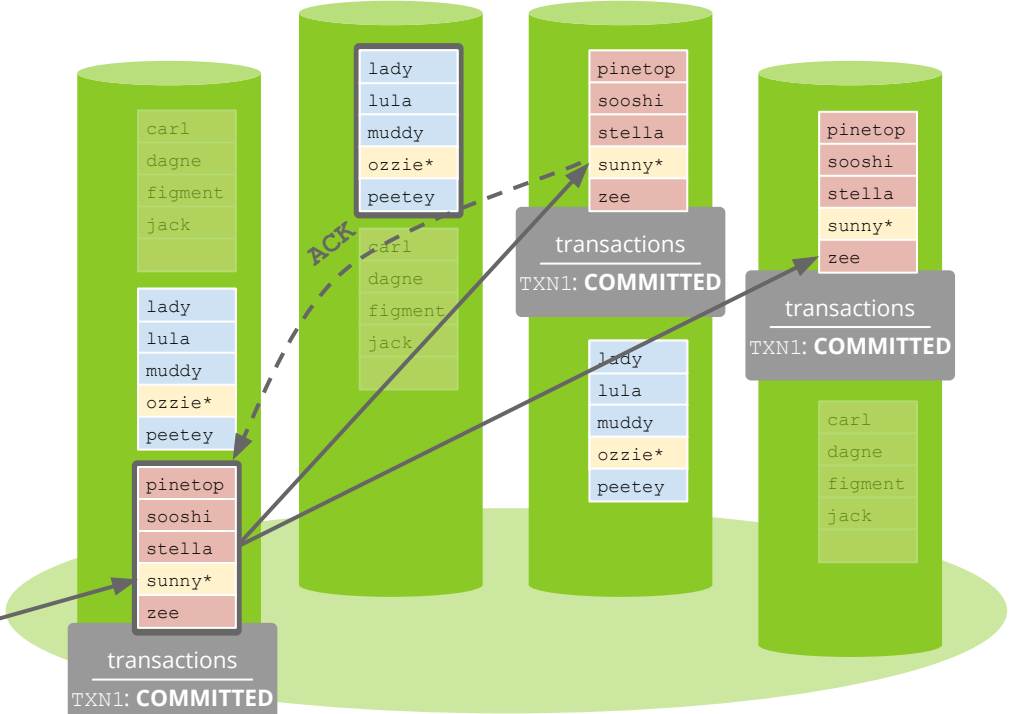
$t = 2\frac{1}{2} \text{ RTT}$

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



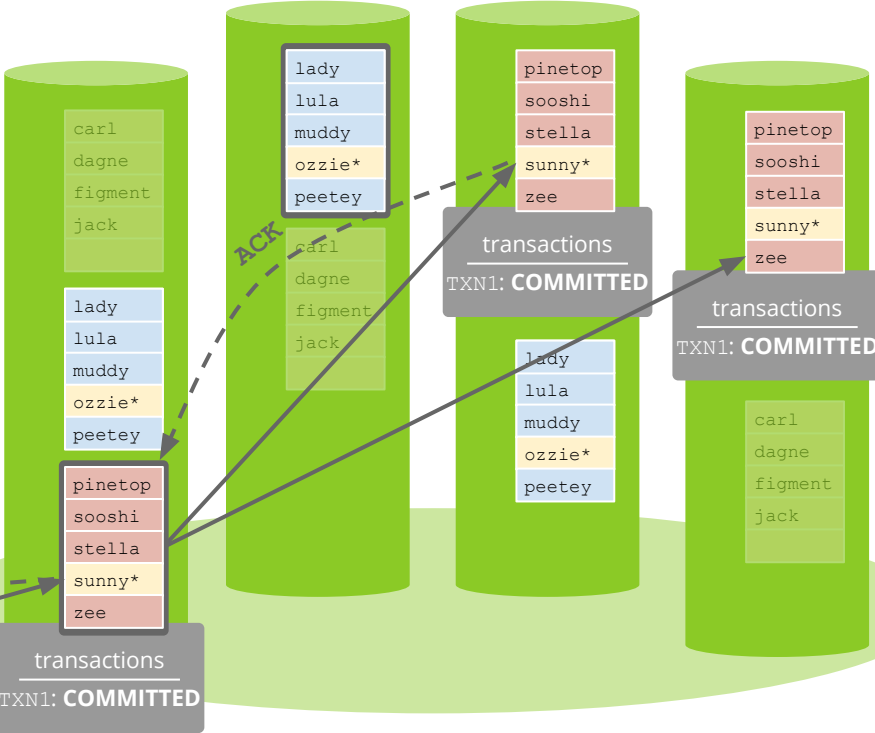
t = 3 RTT

# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 3 RTT

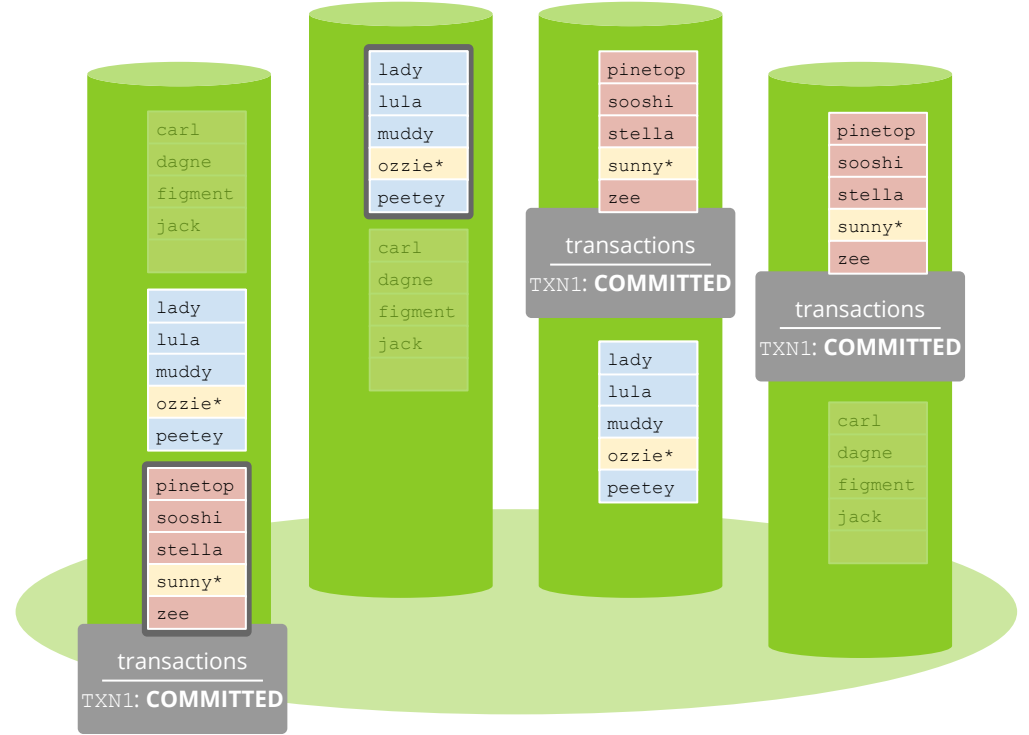


# Unpipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```

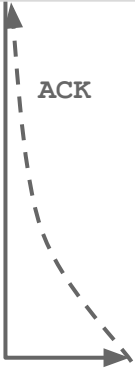


$t = 3 \text{ RTT}$

# Unpipelined Transactions

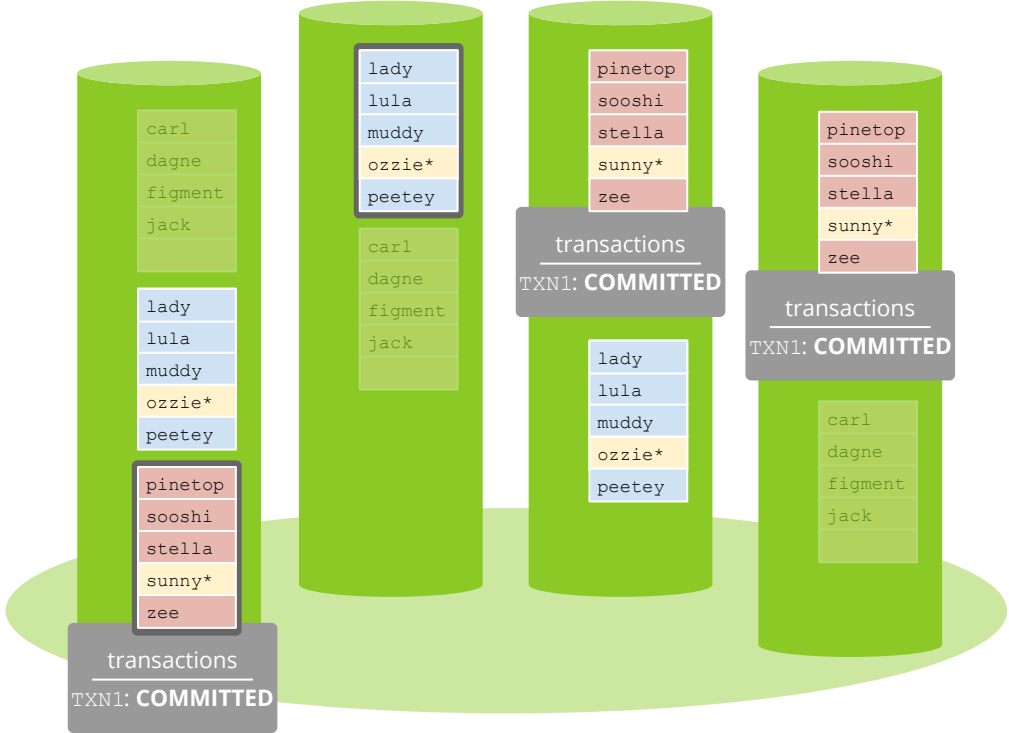


```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```



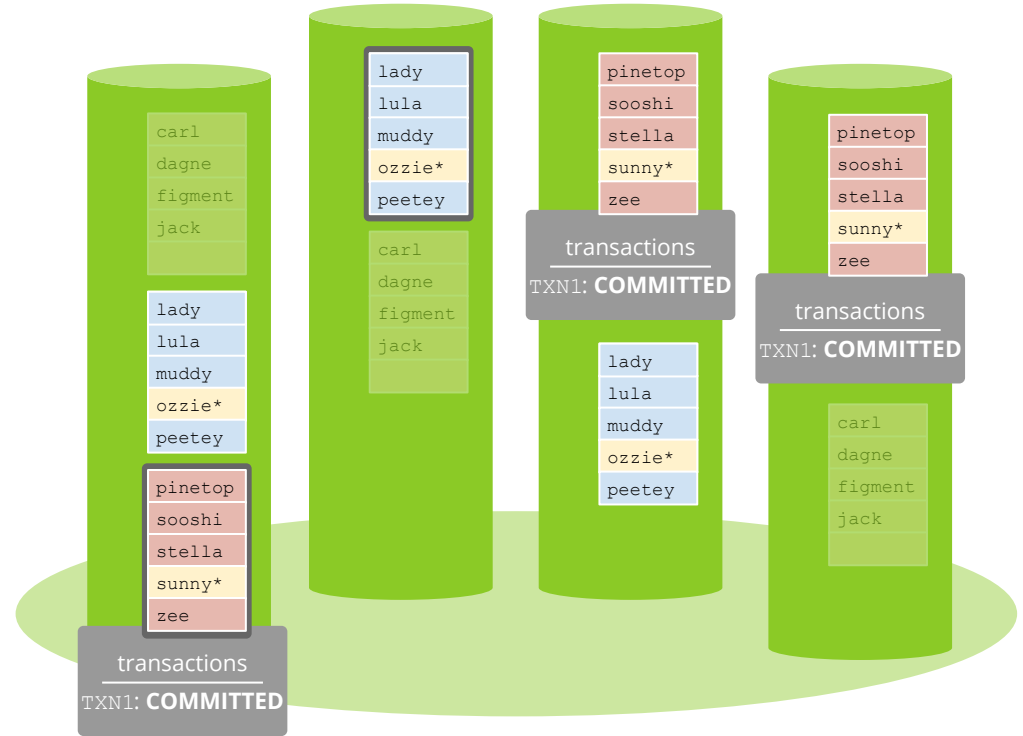
```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT
```

GATEWAY



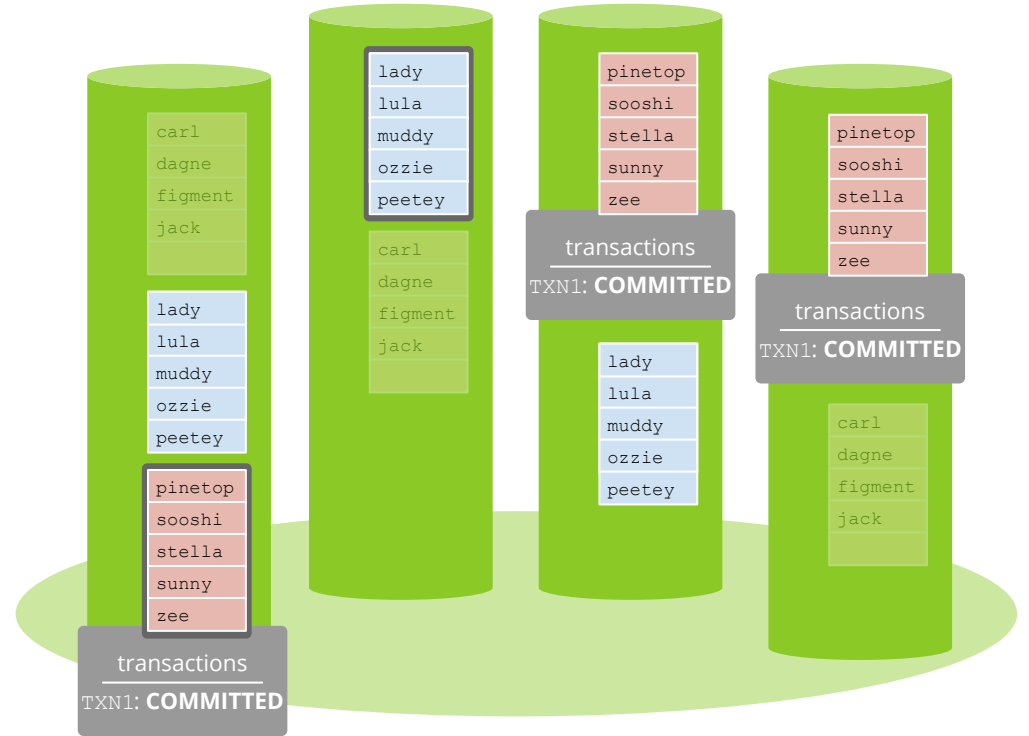
t = 3 RTT

# Unpipelined Transactions

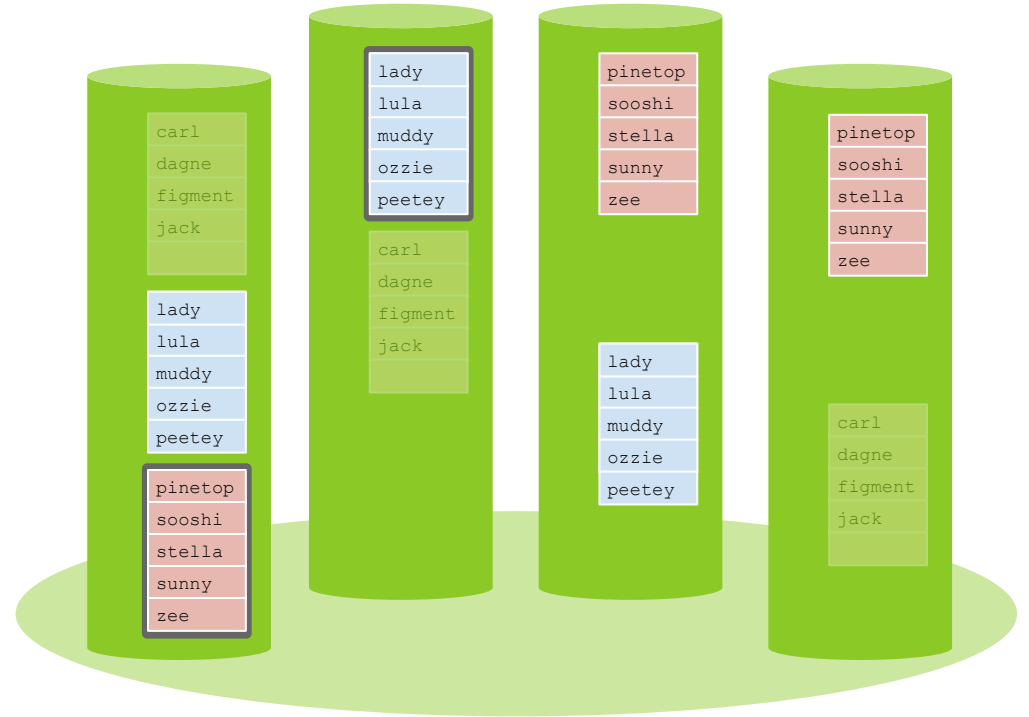


$t = 3 \text{ RTT}$

# Unpipelined Transactions

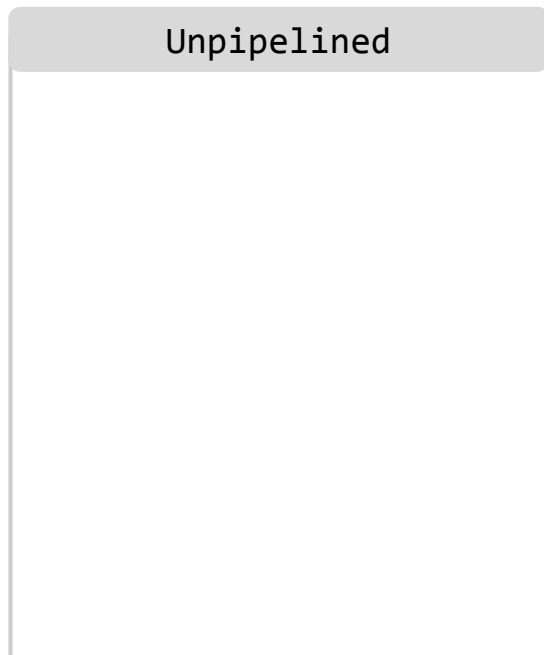
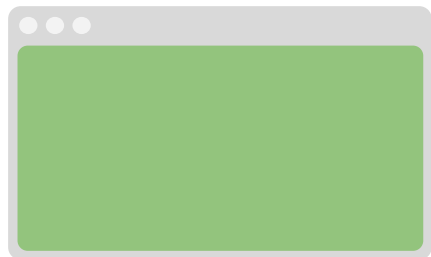


# Unpipelined Transactions

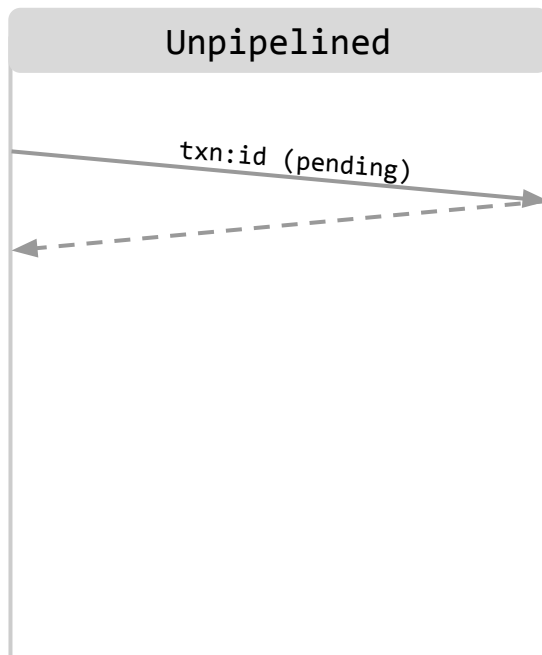
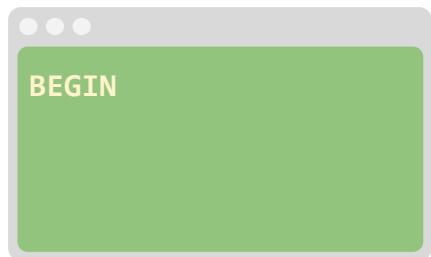


$t = 3 \text{ RTT}$

# Timeline: Unpipelined Transactions



# Timeline: Unpipelined Transactions



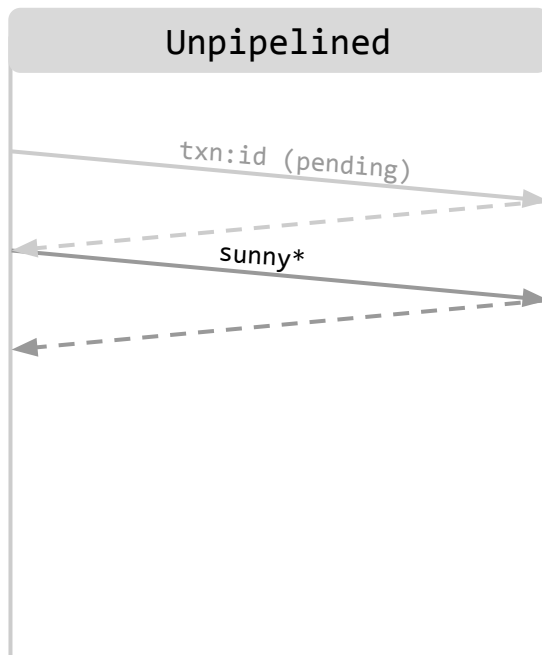
1 WAN RTT



# Timeline: Unpipelined Transactions



```
BEGIN  
WRITE[sunny]
```



1 WAN RTT

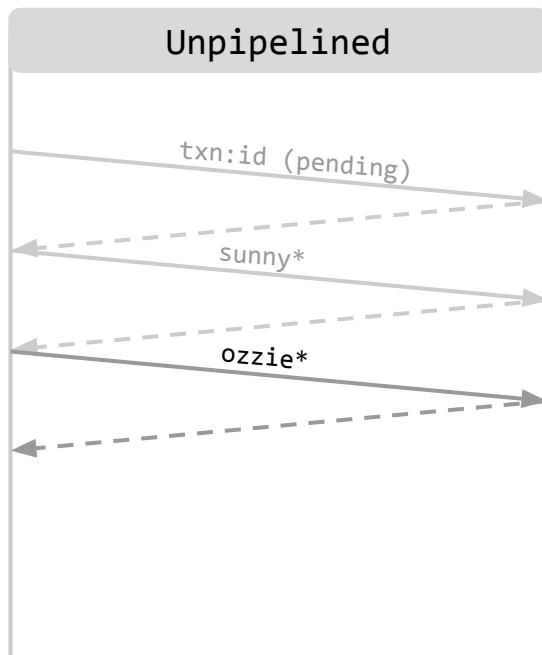




# Timeline: Unpipelined Transactions



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



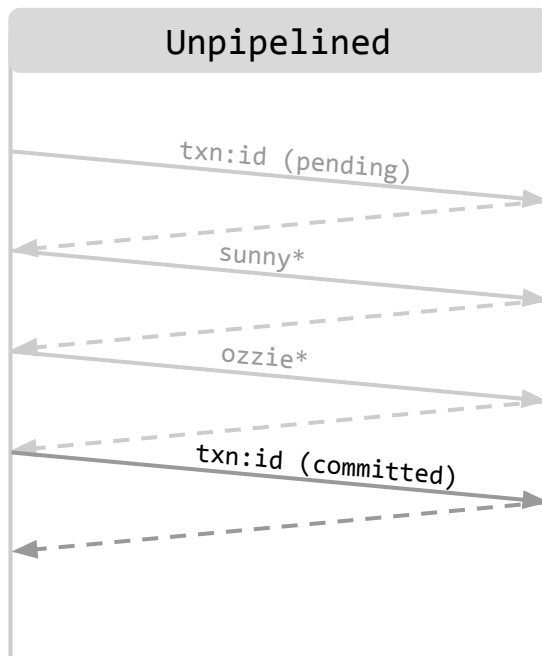
1 WAN RTT



# Timeline: Unpipelined Transactions



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



1 WAN RTT



N+2 WAN RTTS



# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

- I. Spanner/Pipelined Transactions
- II. Parallel Commits
- III. Replicated Commit
- IV. Carousel
- V. MDCC
- VI. SLOG/OceanVista
- VII. TAPIR

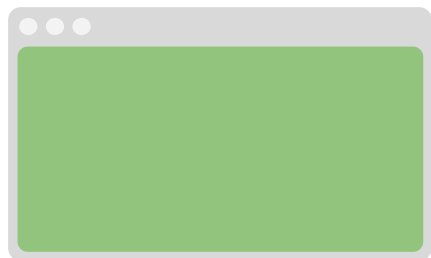


# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

- I. Spanner/Pipelined Transactions
- II. Parallel Commits
- III. Replicated Commit
- IV. Carousel
- V. MDCC
- VI. SLOG/OceanVista
- VII. TAPIR

# Timeline: Pipelined Transactions



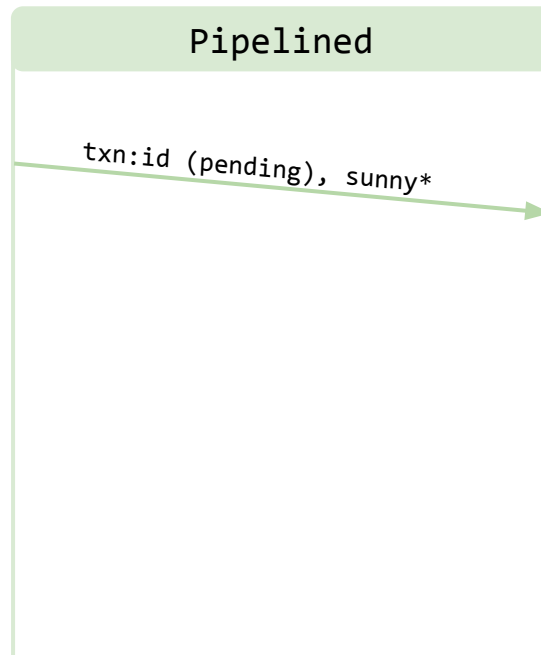
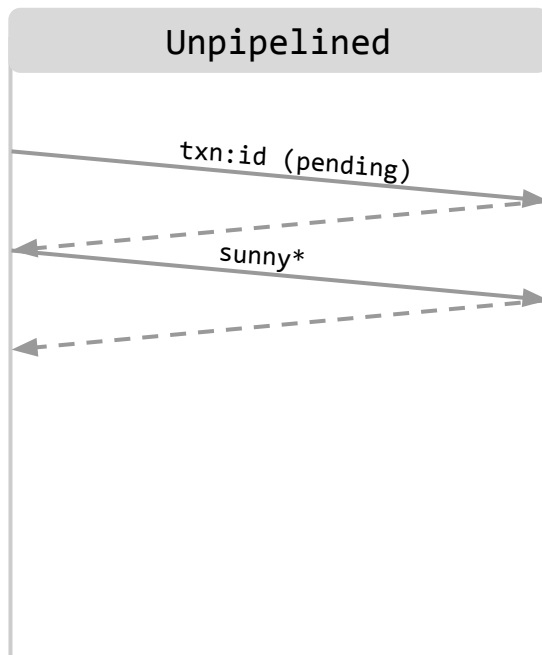
Unpipelined

Pipelined

# Timeline: Pipelined Transactions



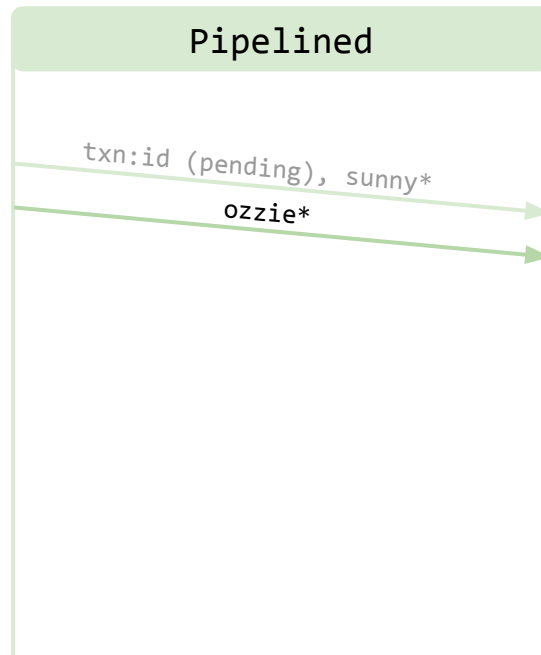
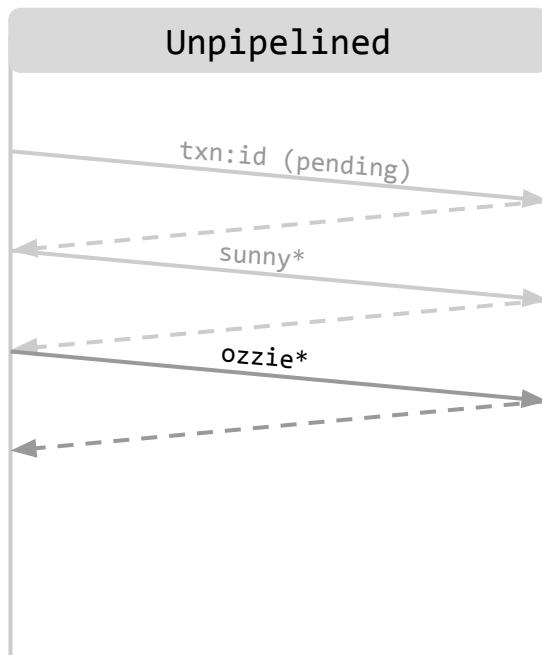
```
BEGIN  
WRITE[sunny]
```



# Timeline: Pipelined Transactions



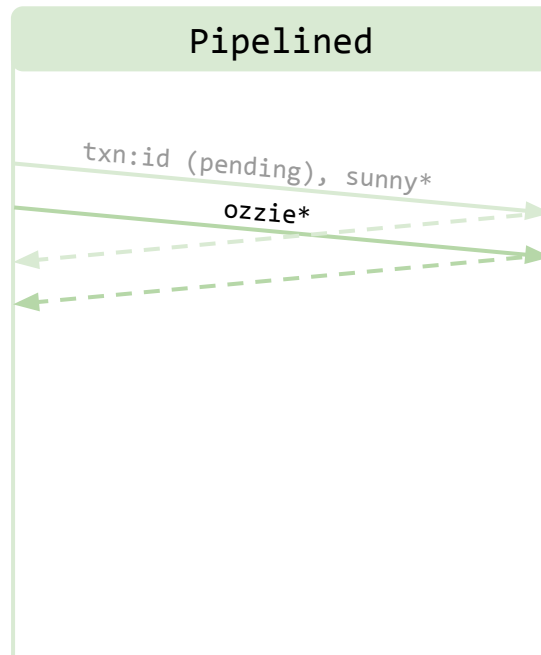
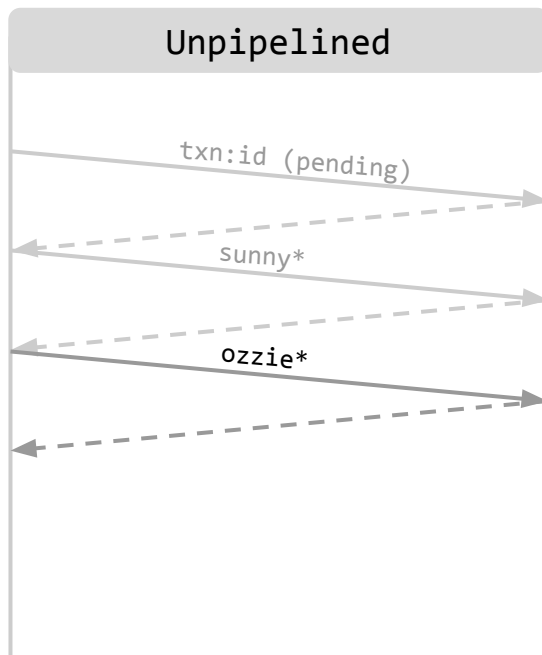
```
BEGIN
WRITE[sunny]
WRITE[ozzie]
```



# Timeline: Pipelined Transactions



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

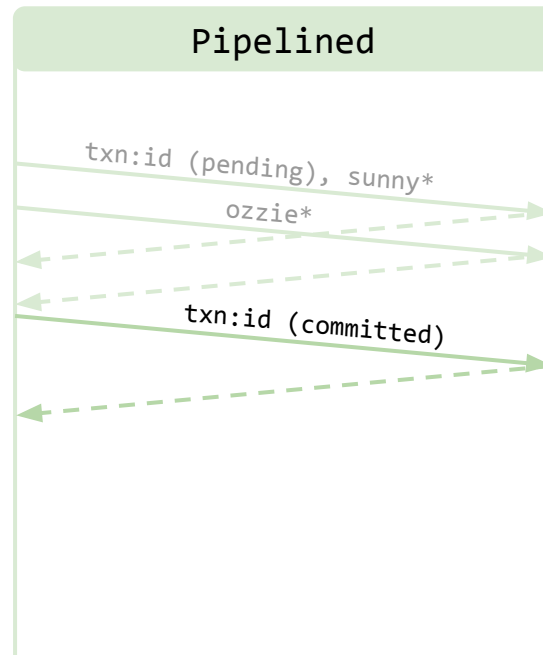
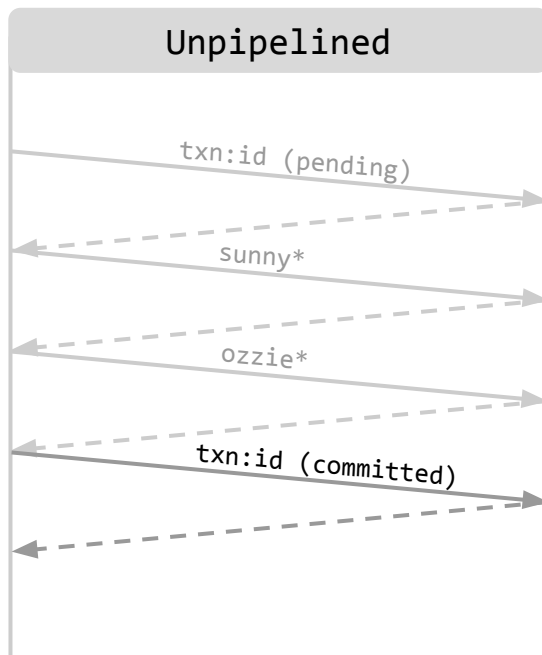




# Timeline: Pipelined Transactions



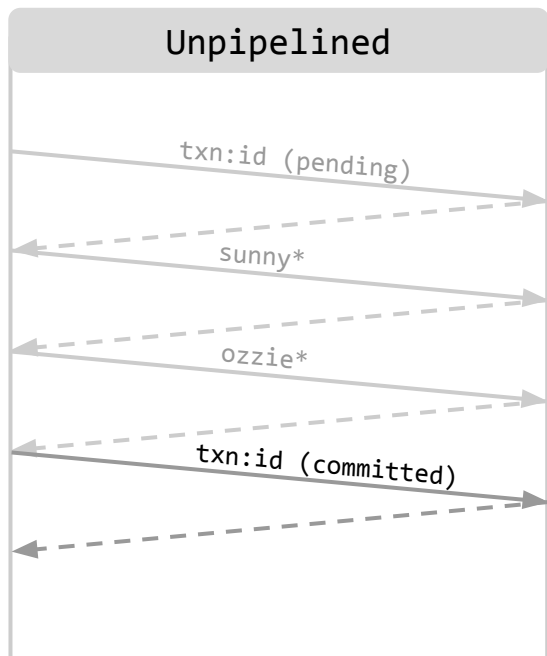
```
BEGIN
WRITE[sunny]
WRITE[ozzie]
COMMIT
```



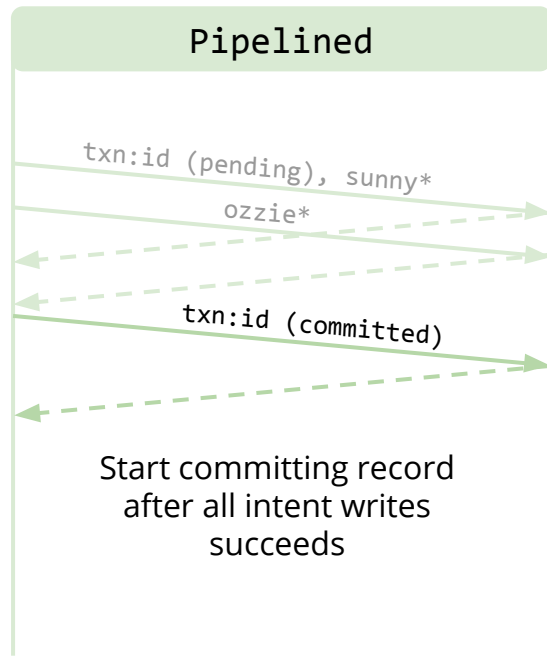
# Timeline: Pipelined Transactions



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



N+2 WAN RTTS

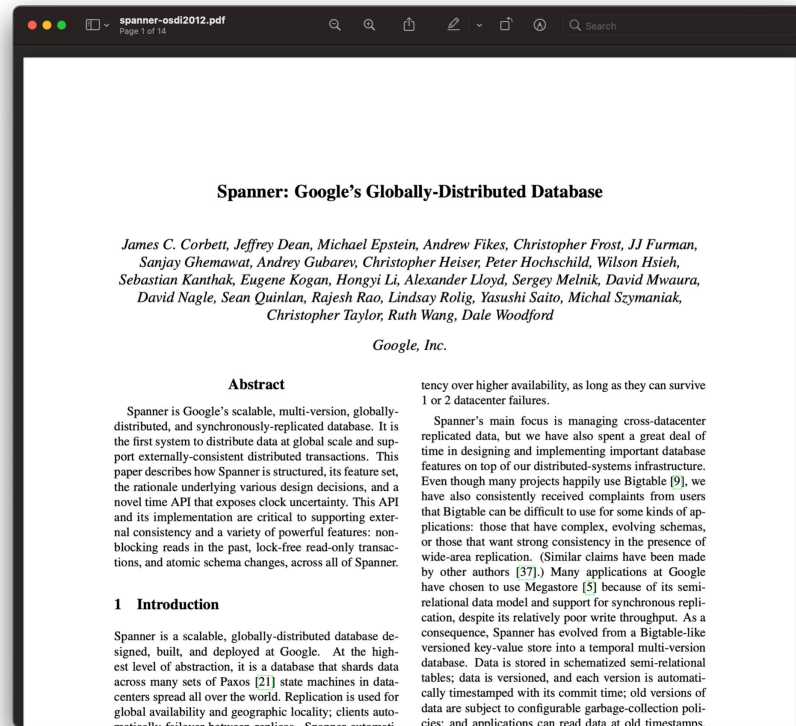


2 WAN RTTS

# Spanner



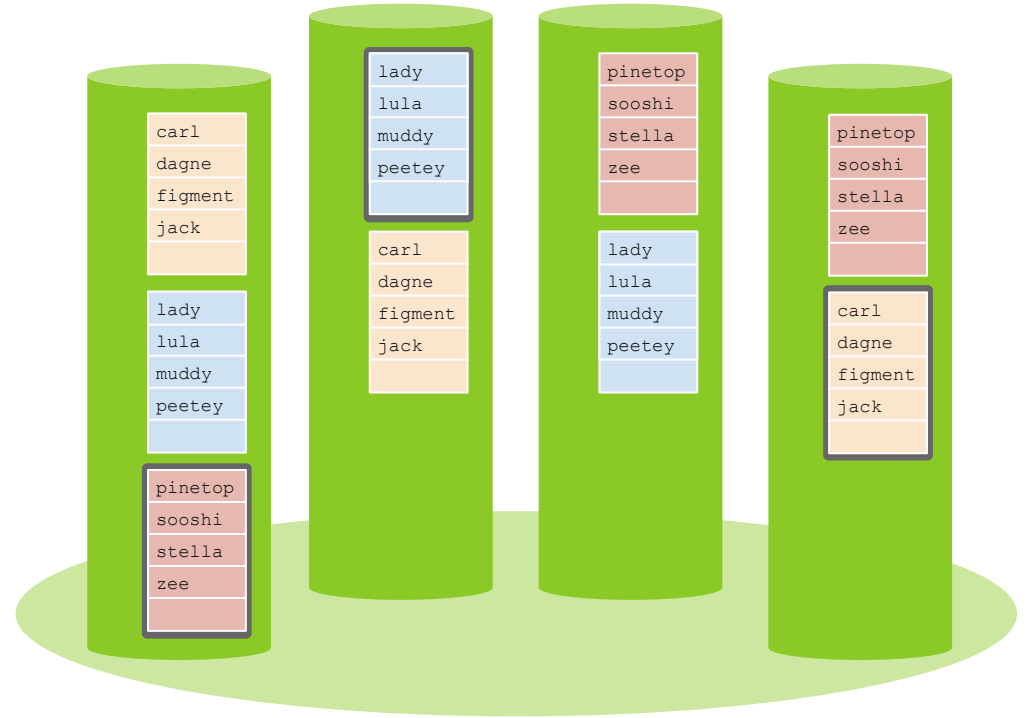
From Google, 2012. Spurred all the following research/derivative systems.



# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```



t = 0

# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
  
GATEWAY
```



t = 0

# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]
```

GATEWAY

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

transactions  
TXN1: PENDING

lady
lula
muddy
peetey
carl
dagne
figment
jack

pinetop
sooshi
stella
zee
lady
lula
muddy
peetey

pinetop
sooshi
stella
zee
carl
dagne
figment
jack

t = 0

# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE[sunny]  
WRITE[ozzie]
```

GATEWAY

carl
dagne
figment
jack
lady
lula
muddy
peetey

pinetop
sooshi
stella
sunny*
zee

transactions  
TXN1: PENDING

lady
lula
muddy
ozzie*
peetey

carl
dagne
figment
jack

pinetop
sooshi
stella
zee

lady
lula
muddy
peetey

pinetop
sooshi
stella
zee

carl
dagne
figment
jack

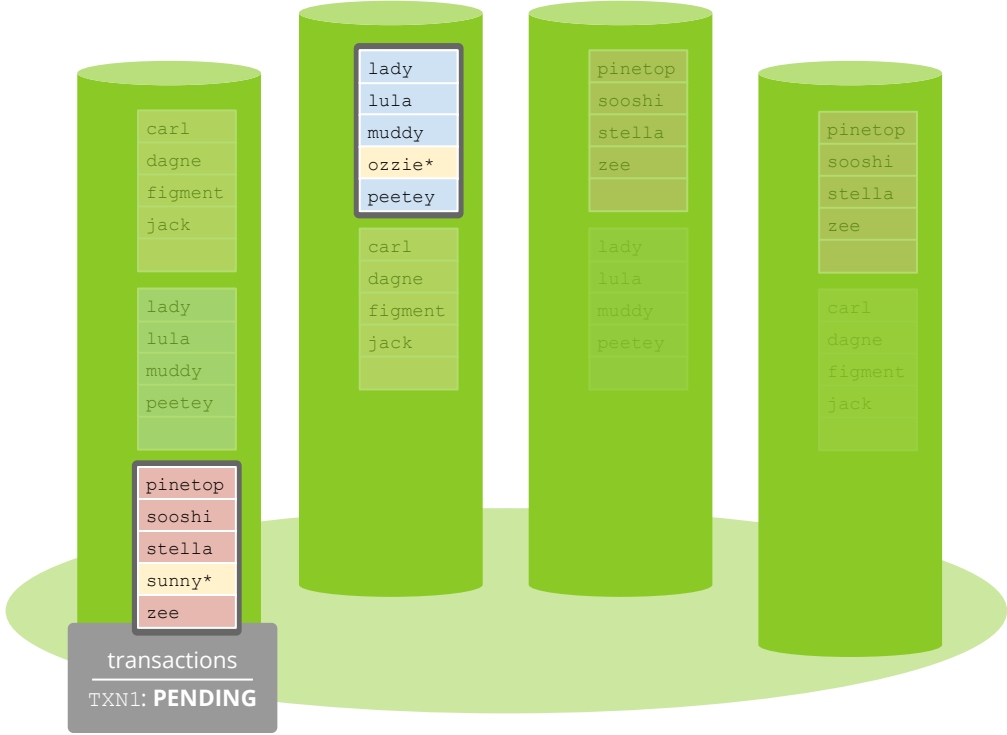
t = 0

# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 0

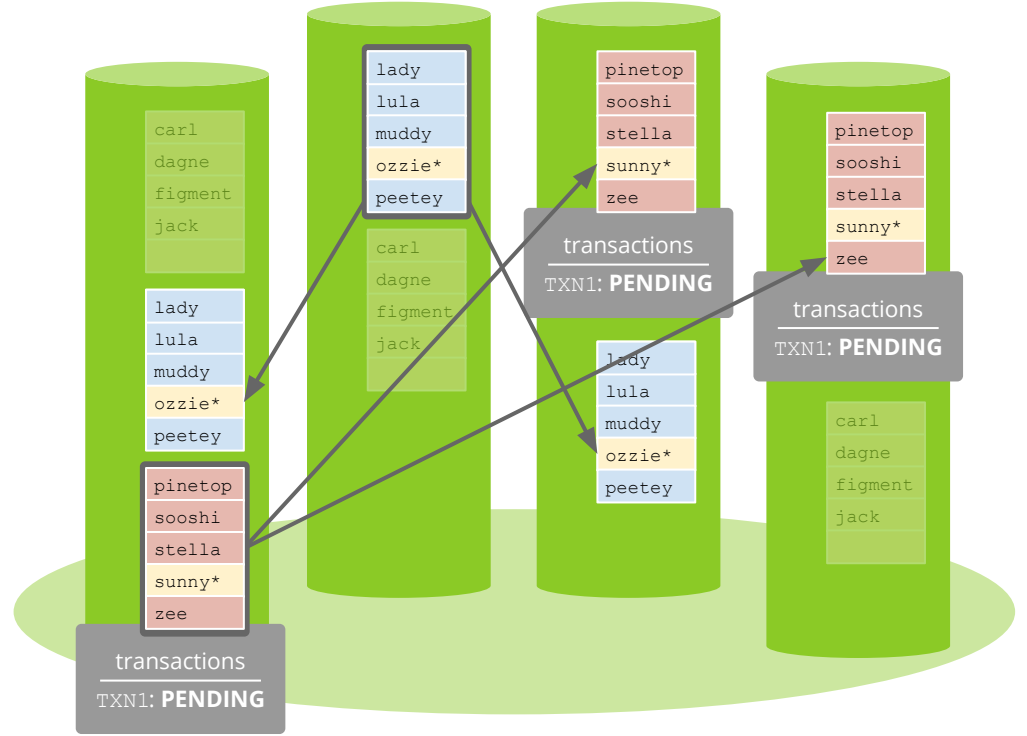


# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



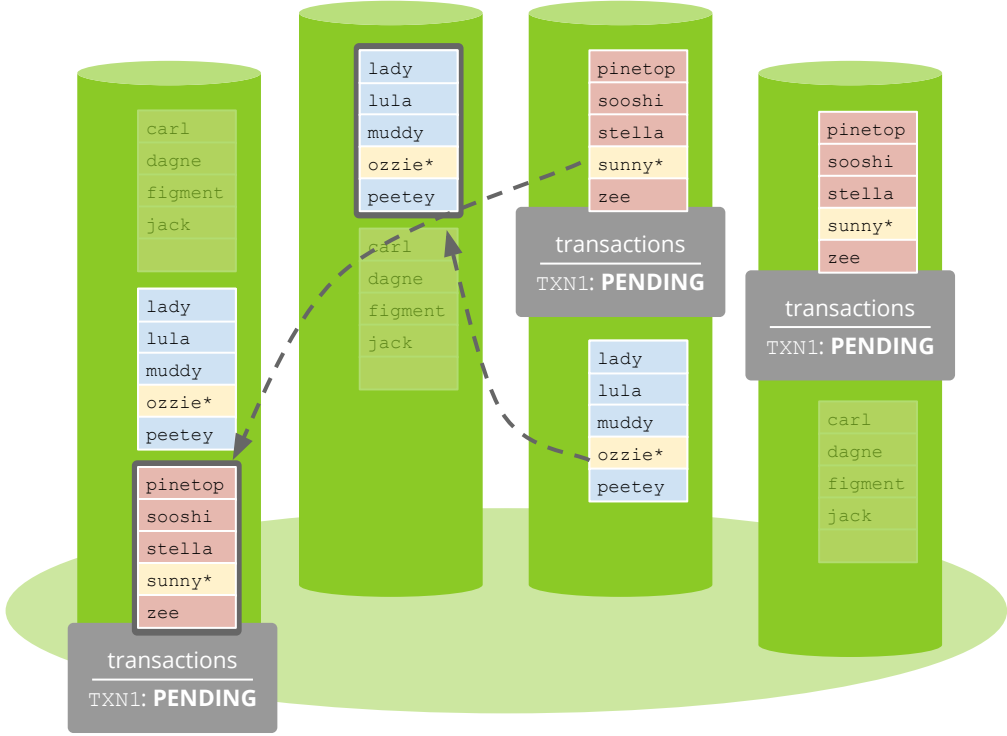
$$t = \frac{1}{2} \text{ RTT}$$

# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 1 RTT

# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```

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

transactions  
TXN1: PENDING

lady
lula
muddy
ozzie*
peetey
carl
dagne
figment
jack

pinetop
sooshi
stella
sunny*
zee

transactions  
TXN1: PENDING

lady
lula
muddy
ozzie*
peetey

pinetop
sooshi
stella
sunny*
zee

transactions  
TXN1: PENDING

carl
dagne
figment
jack

t = 1 RTT

# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```

transactions  
TXN1: **COMMITTED**

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

lady
lula
muddy
ozzie*
peetey

carl
dagne
figment
jack

transactions  
TXN1: **PENDING**

pinetop
sooshi
stella
sunny*
zee

lady
lula
muddy
ozzie*
peetey

transactions  
TXN1: **PENDING**

pinetop
sooshi
stella
sunny*
zee

carl
dagne
figment
jack

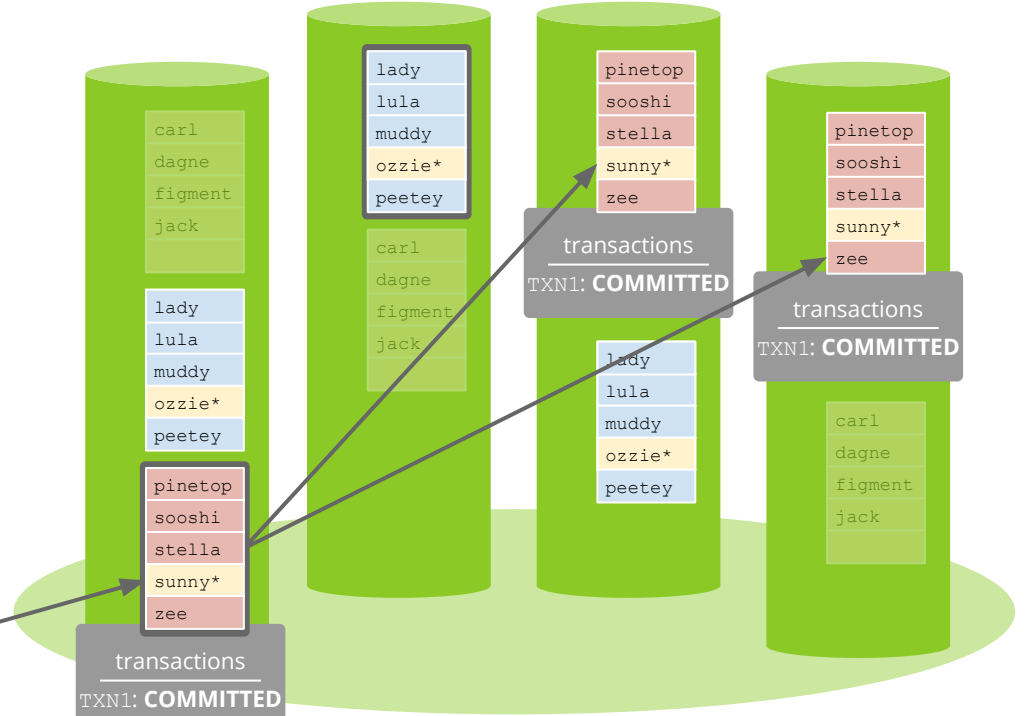
t = 1 RTT

# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



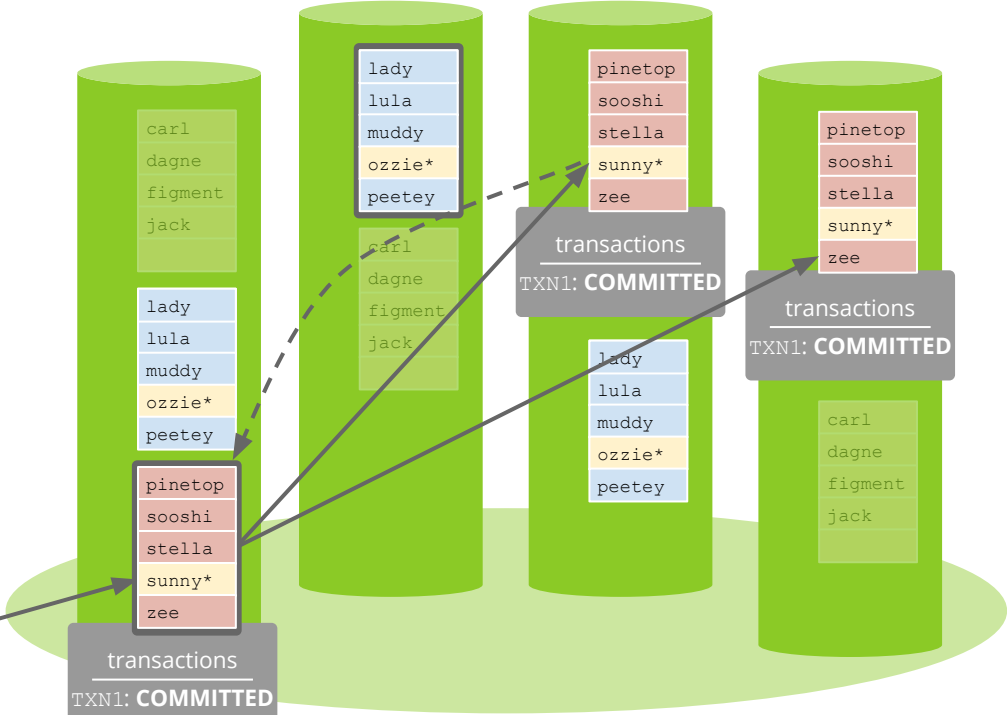
$t = 1\frac{1}{2} \text{ RTT}$

# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 2 RTT

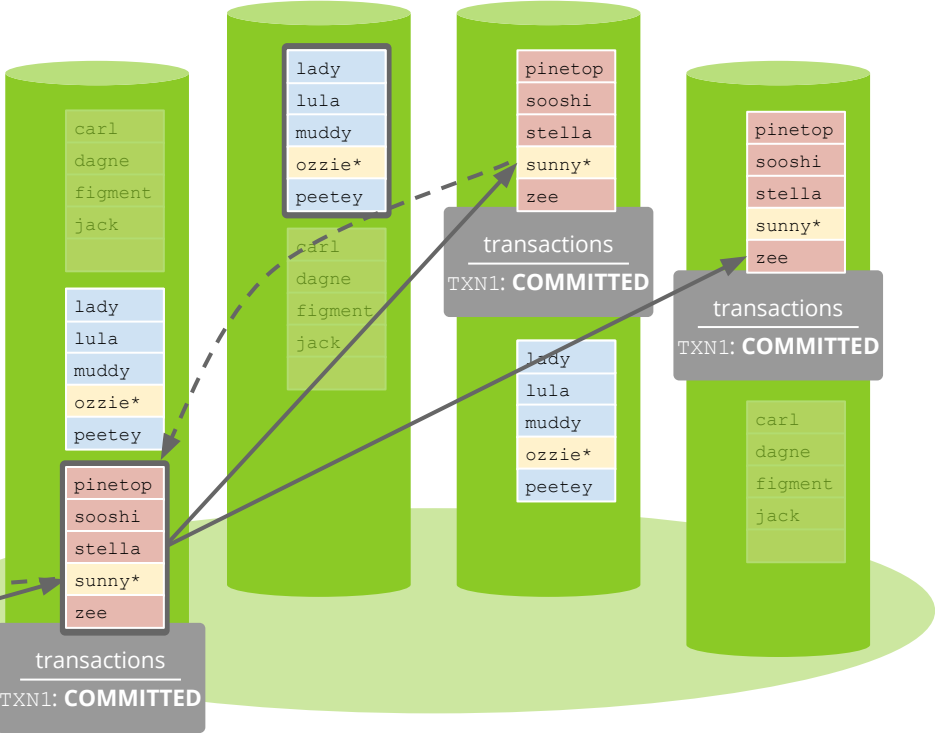
# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```

t = 2 RTT

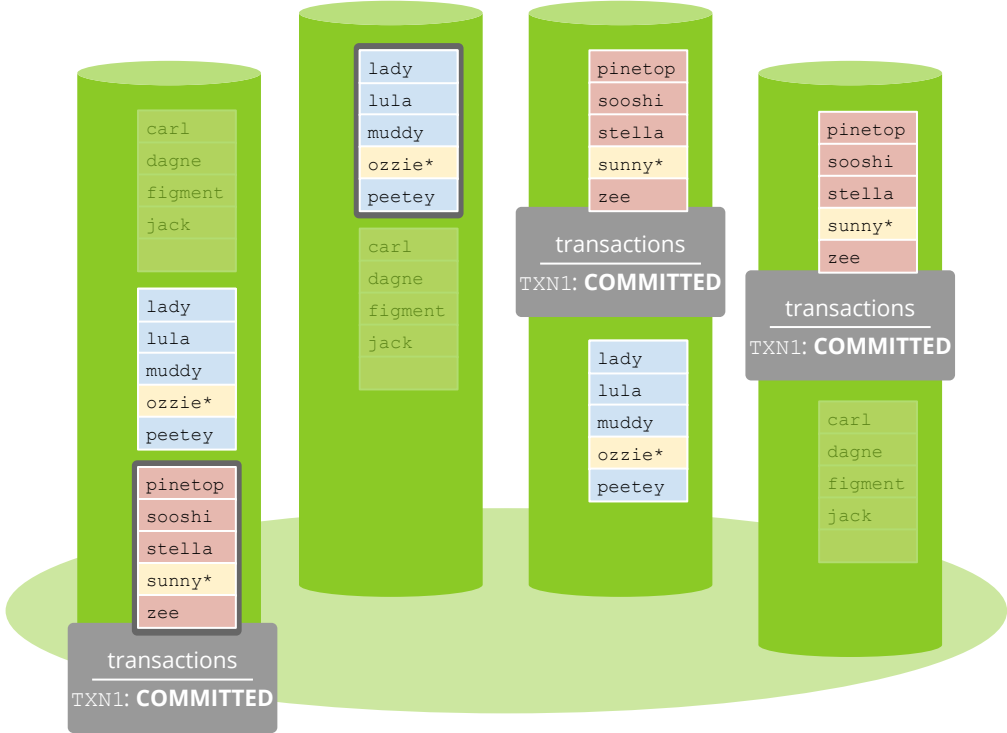


# Pipelined Transactions



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 2 RTT

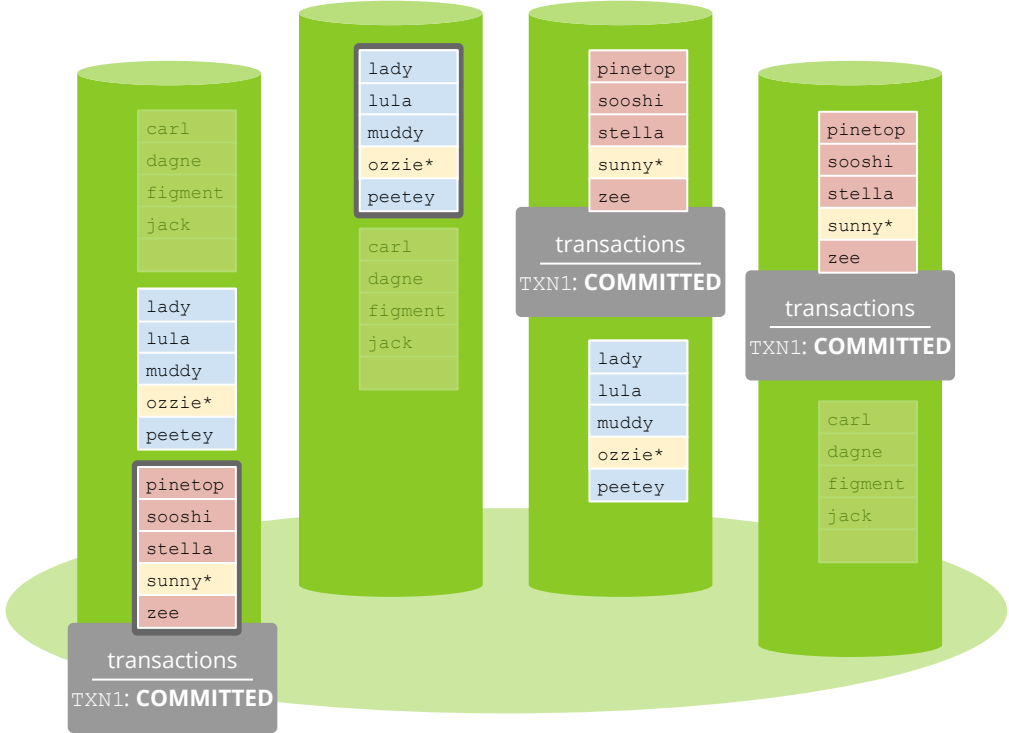


# Pipelined Transactions



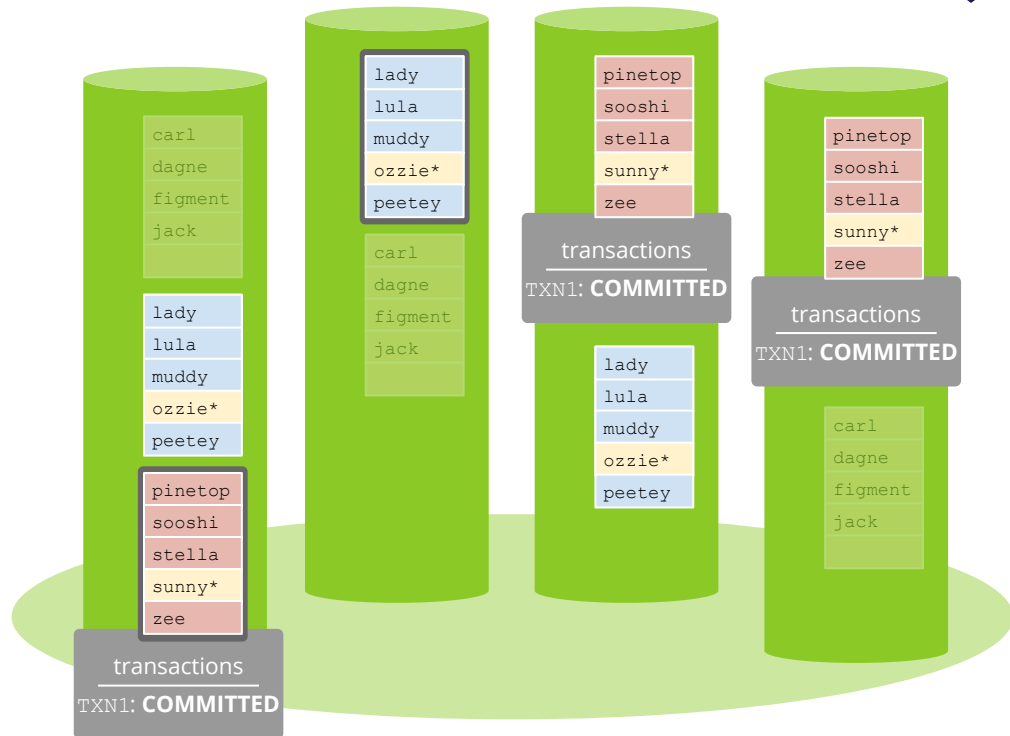
```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



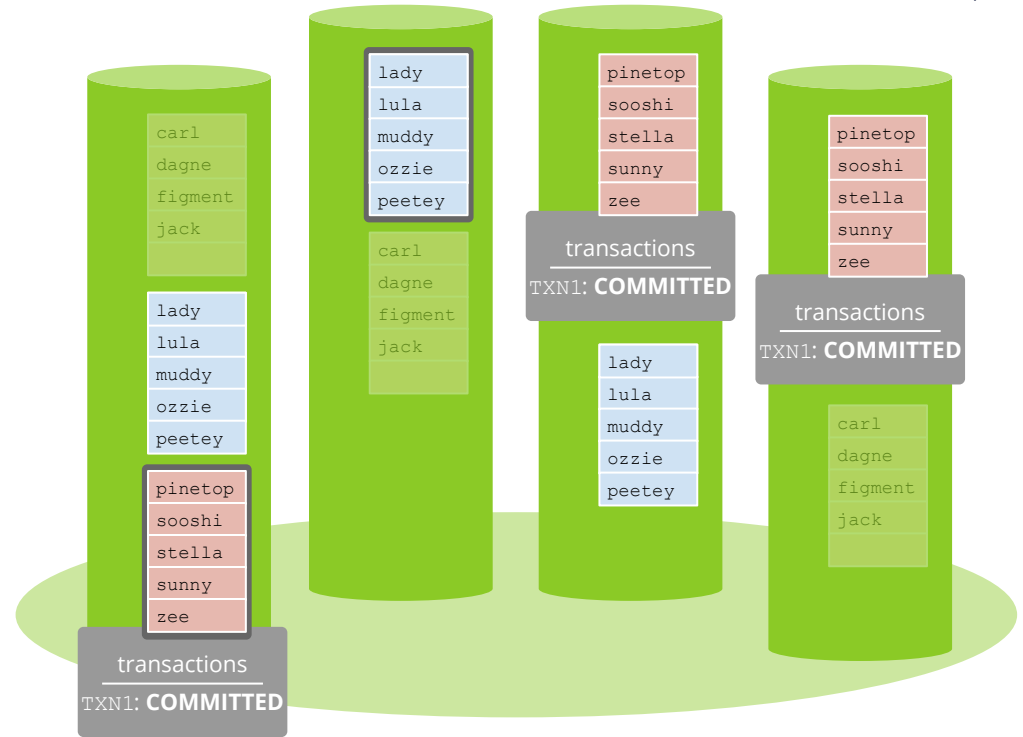
t = 2 RTT

# Pipelined Transactions



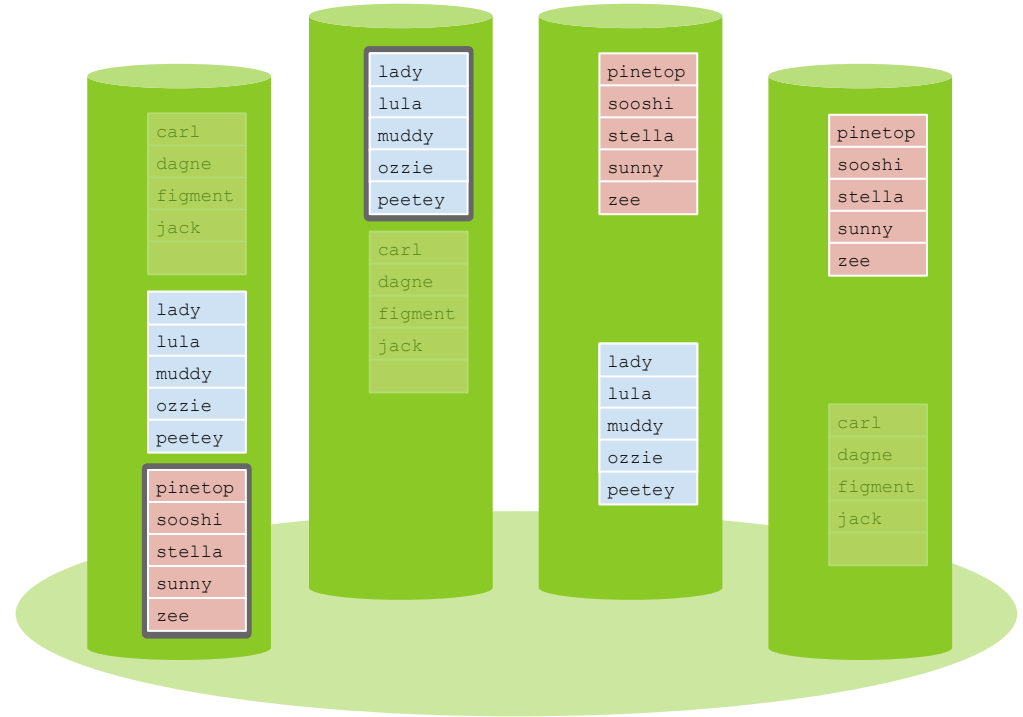
t = 2 RTT

# Pipelined Transactions



$t = 2 \text{ RTT}$

# Pipelined Transactions



$t = 2 \text{ RTT}$

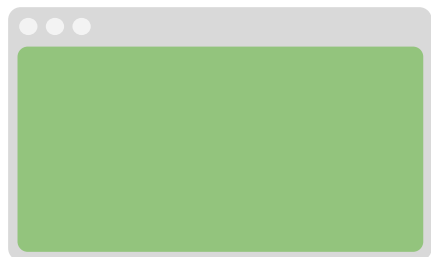


# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

- I. Spanner/Pipelined Transactions
- II. Parallel Commits
- III. Replicated Commit
- IV. Carousel
- V. MDCC
- VI. SLOG/OceanVista
- VII. TAPIR

# Timeline: Parallel Commits



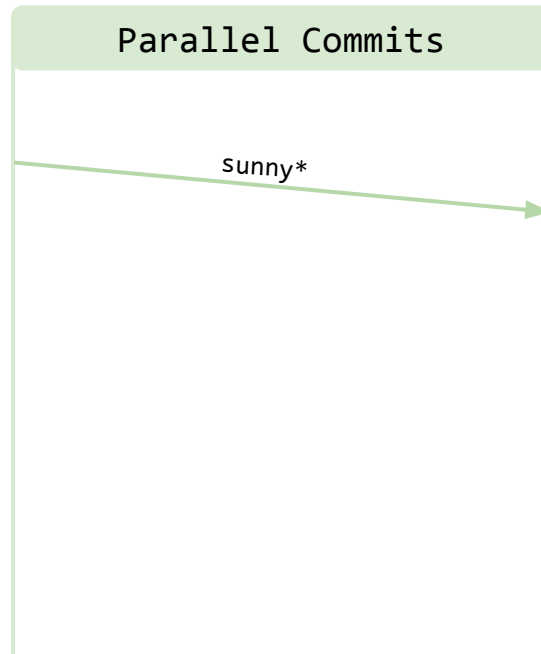
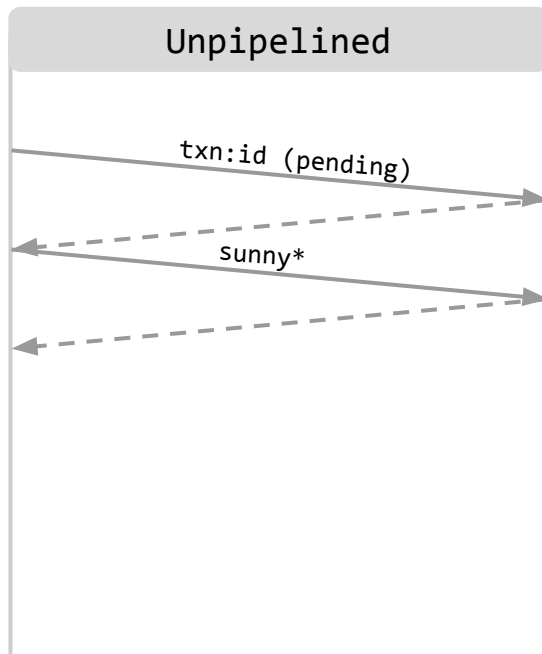
Unpipelined

Parallel Commits

# Timeline: Parallel Commits



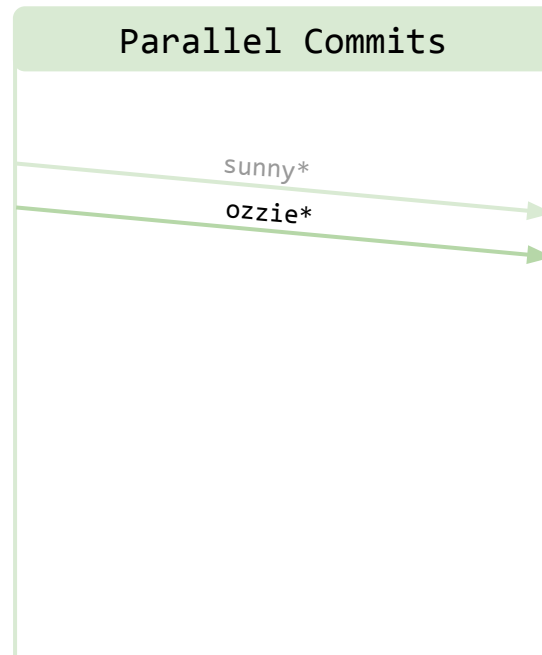
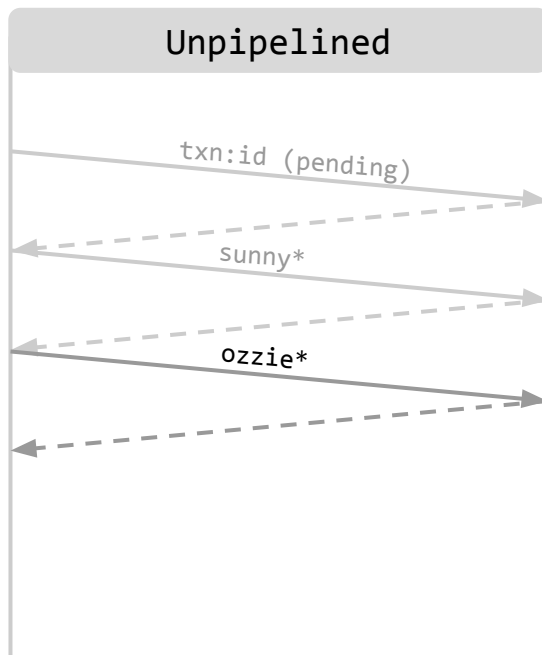
```
BEGIN  
WRITE[sunny]
```



# Timeline: Parallel Commits



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

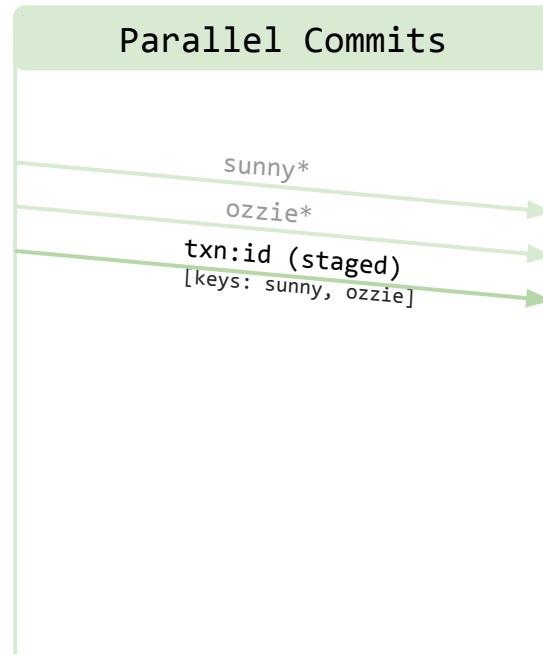
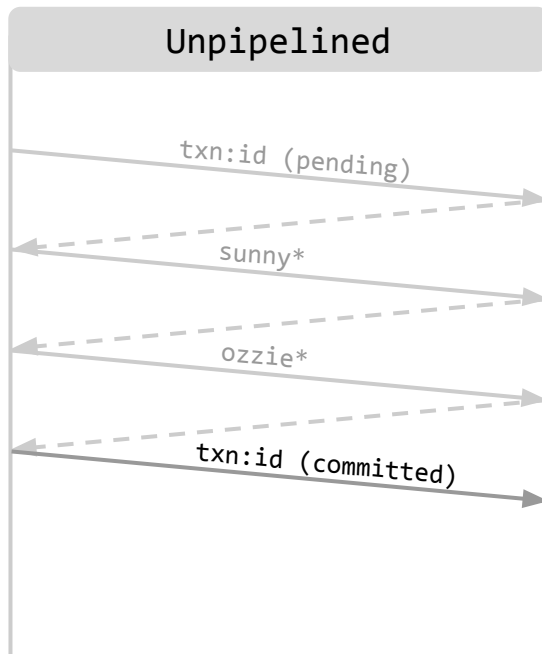




# Timeline: Parallel Commits



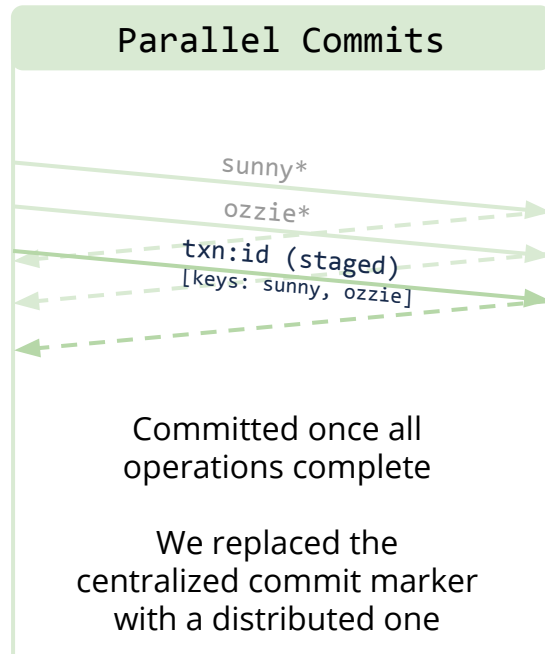
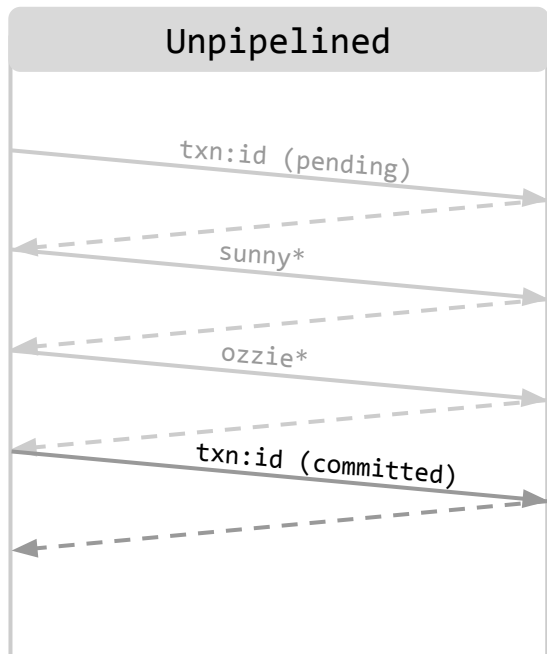
```
BEGIN
WRITE[sunny]
WRITE[ozzie]
COMMIT
```



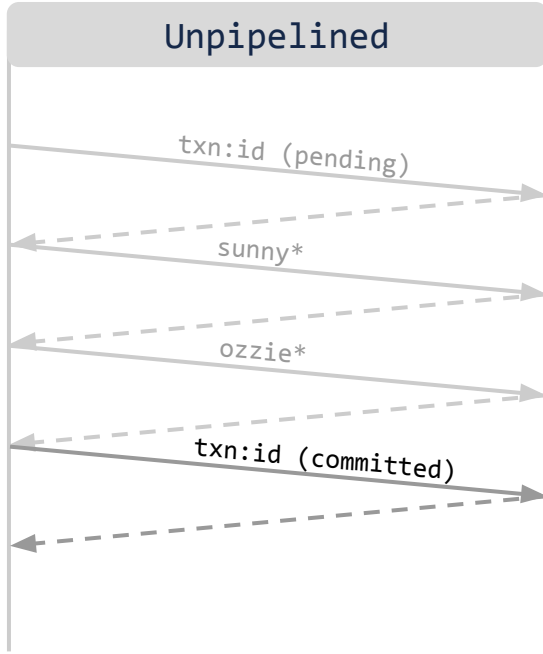
# Timeline: Parallel Commits



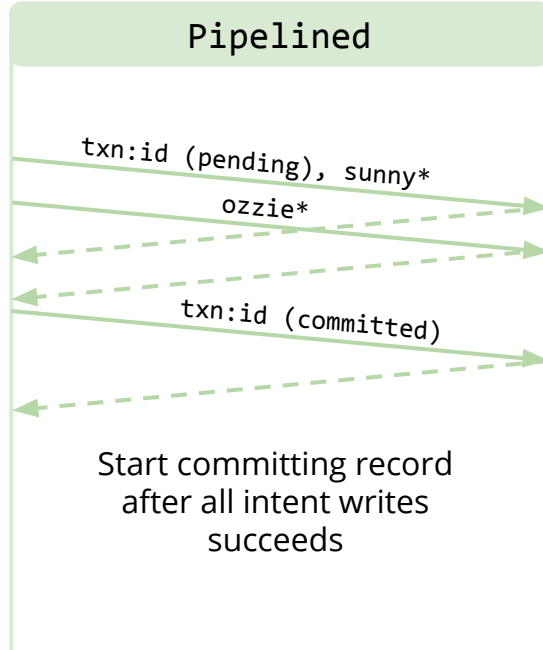
```
BEGIN
WRITE[sunny]
WRITE[ozzie]
COMMIT
```



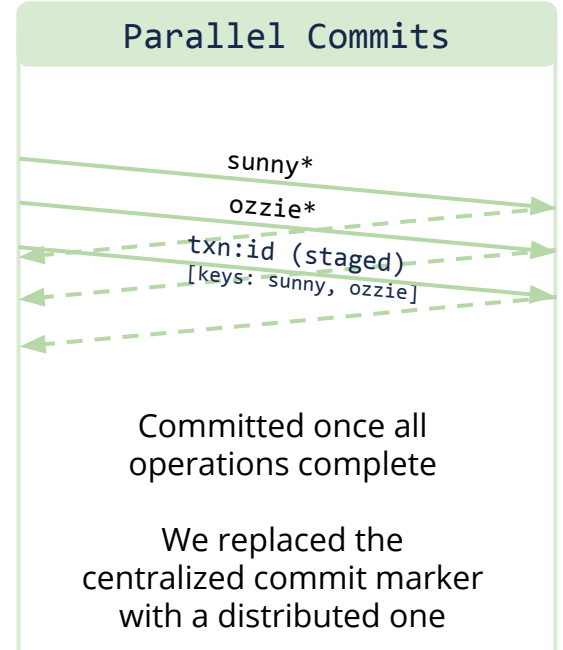
- Proved using TLA+



N+2 WAN RTTS



2 WAN RTTS

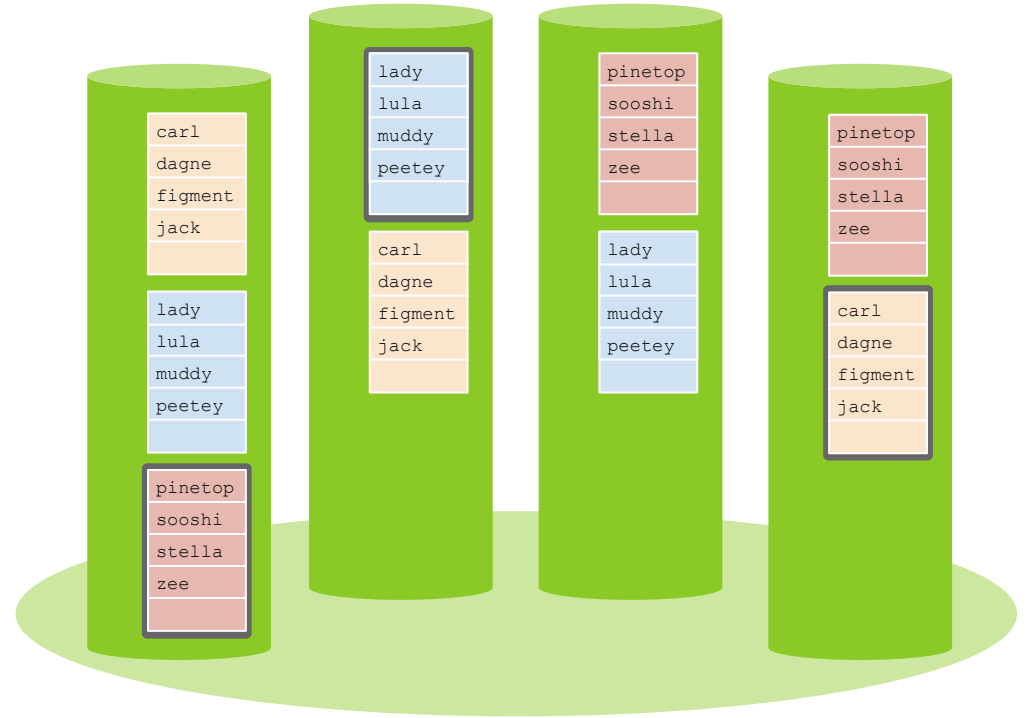


1 WAN RTT

# Parallel Commits



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```



$t = 0$

# Parallel Commits



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
  
GATEWAY
```



t = 0

# Parallel Commits



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
  
GATEWAY
```

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

transactions  
TXN1: **STAGED**

lady
lula
muddy
peetey
carl
dagne
figment
jack

pinetop
sooshi
stella
zee
lady
lula
muddy
peetey

pinetop
sooshi
stella
zee
carl
dagne
figment
jack

t = 0

# Parallel Commits



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]
```

GATEWAY

carl
dagne
figment
jack
lady
lula
muddy
peetey

pinetop
sooshi
stella
sunny*
zee

transactions  
TXN1: **STAGED**

lady
lula
muddy
ozzie*
peetey

carl
dagne
figment
jack

pinetop
sooshi
stella
zee

lady
lula
muddy
peetey

pinetop
sooshi
stella
zee

carl
dagne
figment
jack

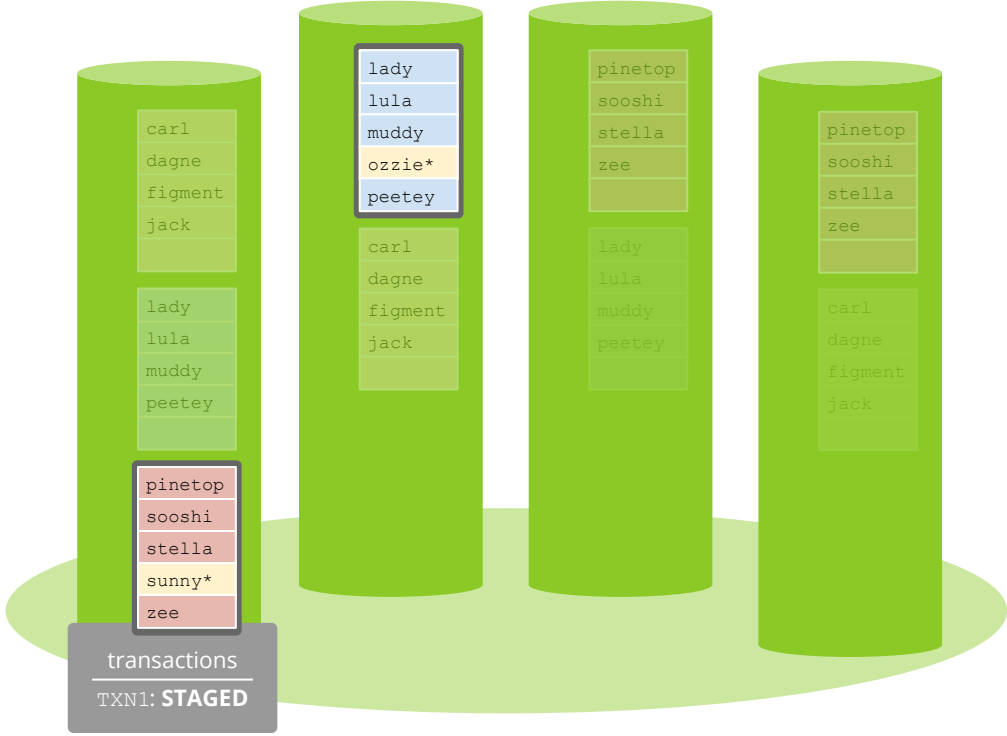
t = 0

# Parallel Commits



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 0

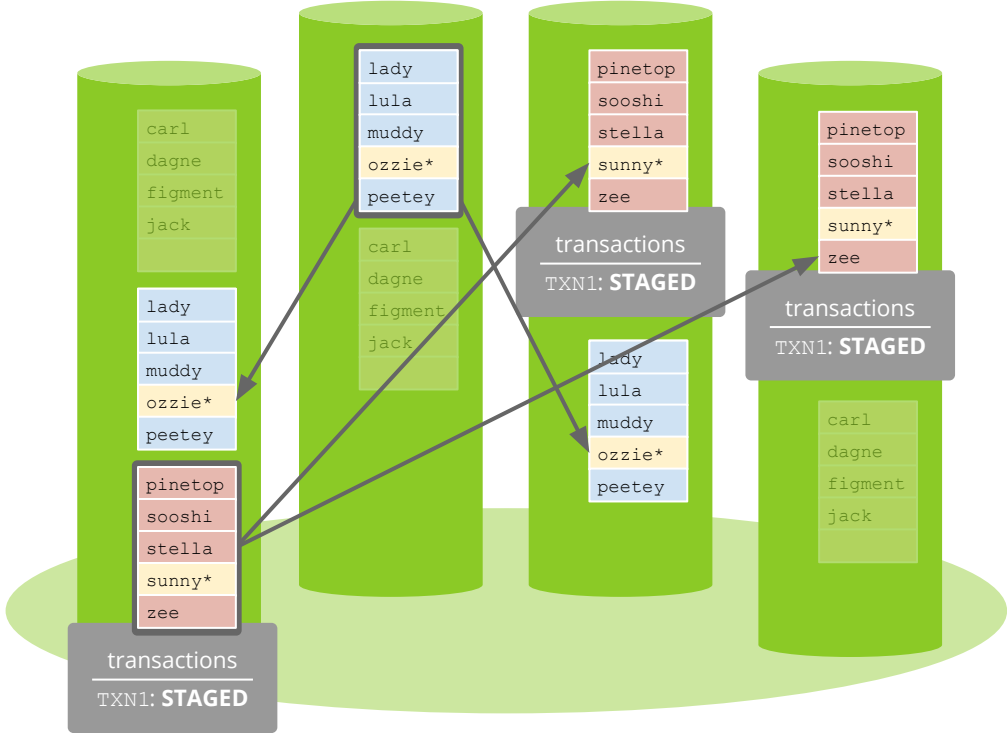


# Parallel Commits



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



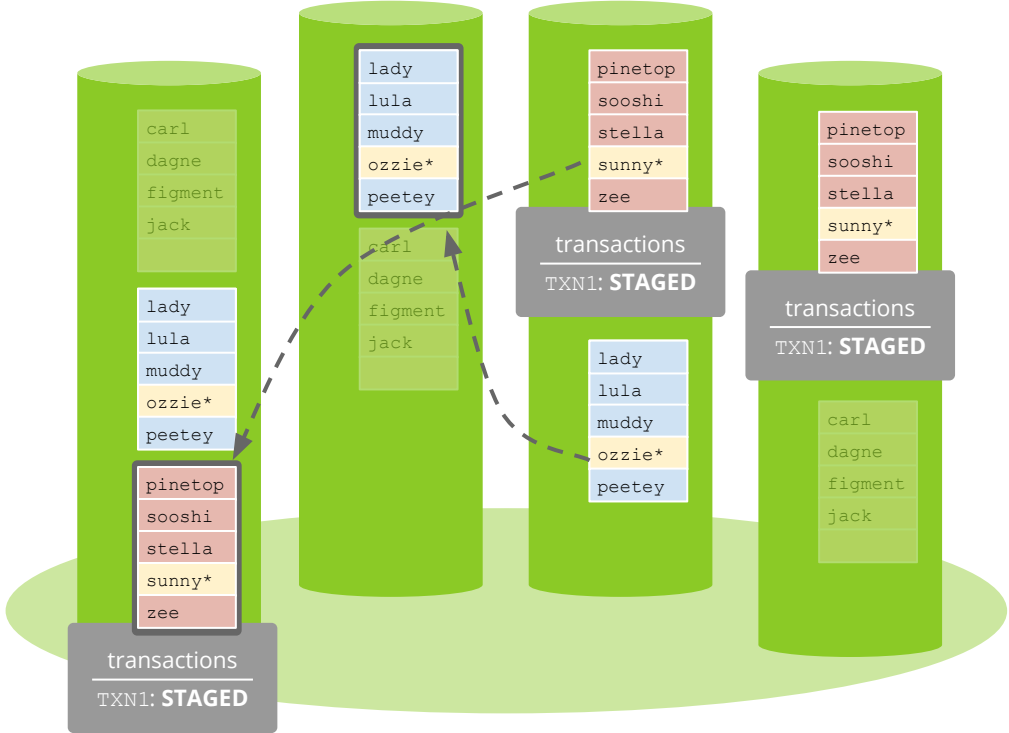
$$t = \frac{1}{2} \text{ RTT}$$

# Parallel Commits



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 1 RTT

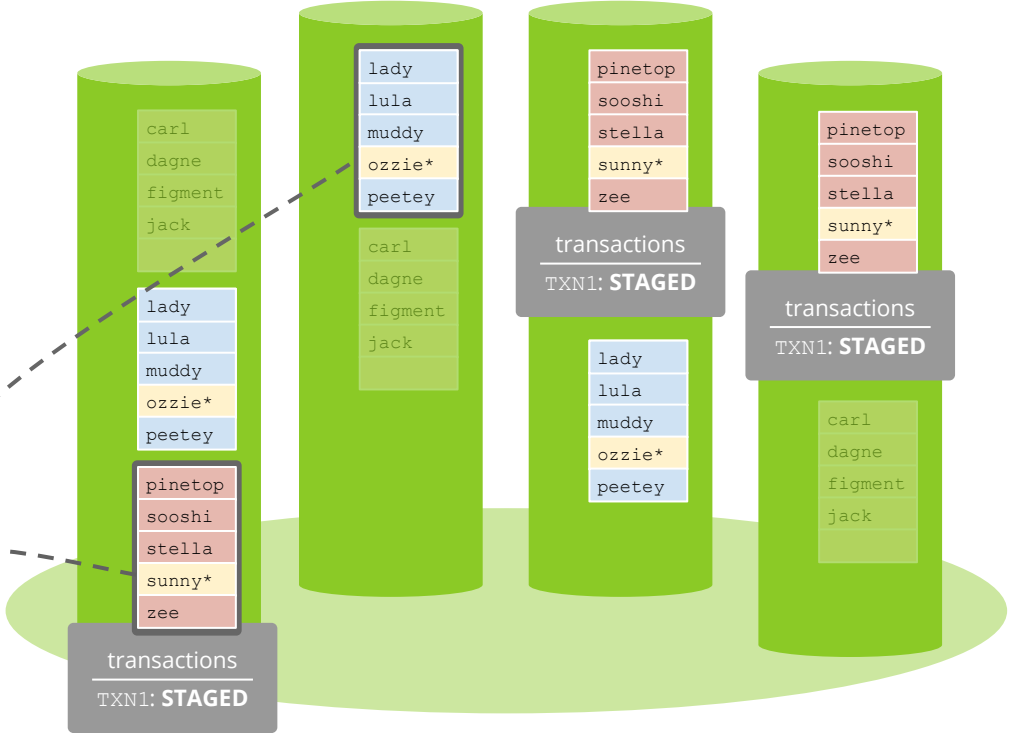
# Parallel Commits



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```

t = 1 RTT

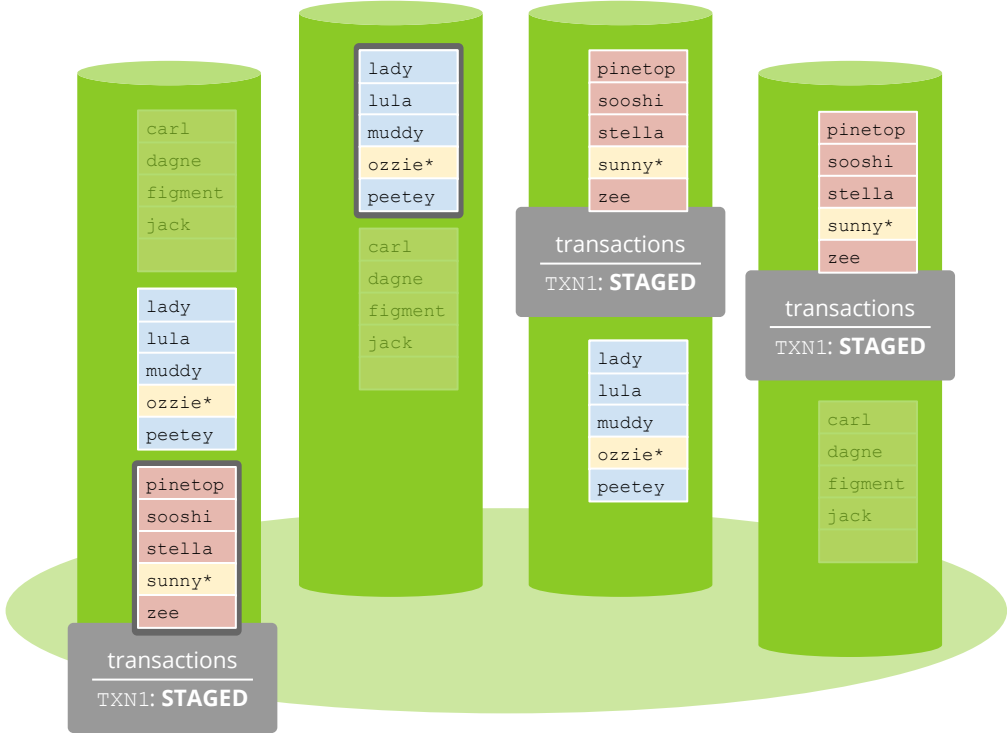


# Parallel Commits



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



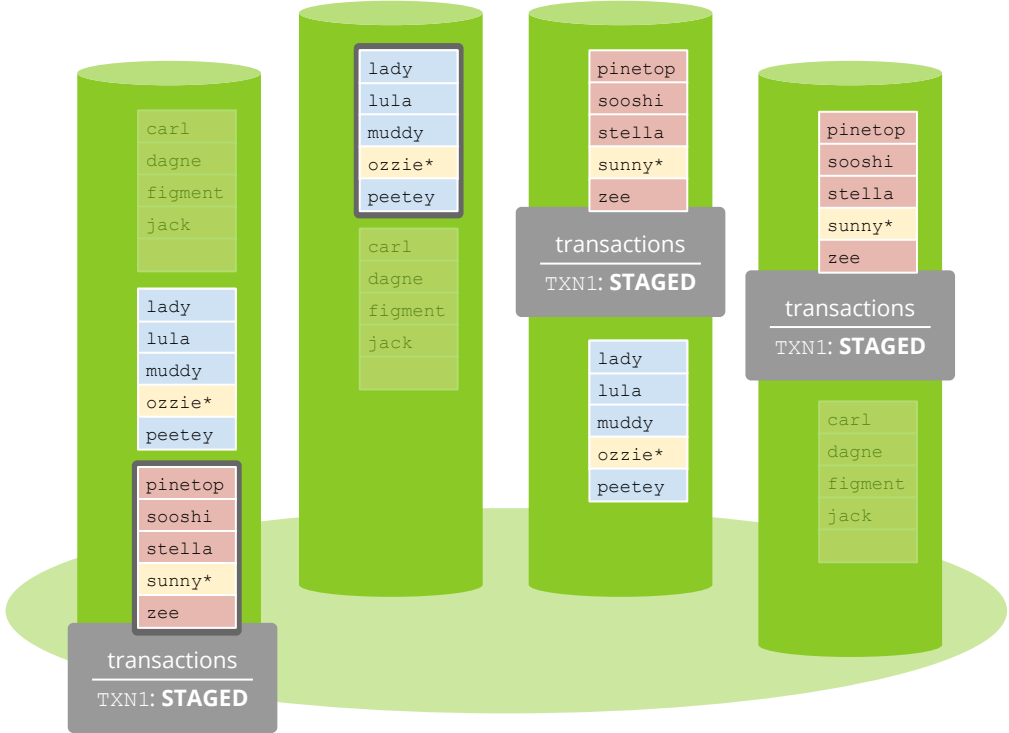
t = 1 RTT

# Parallel Commits



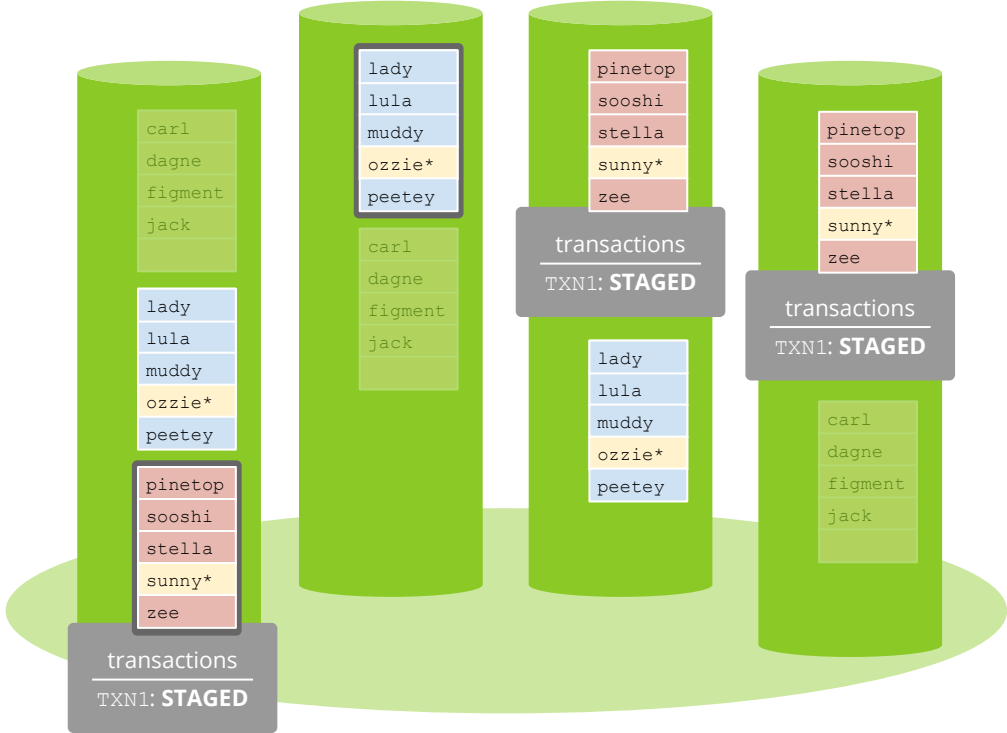
```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



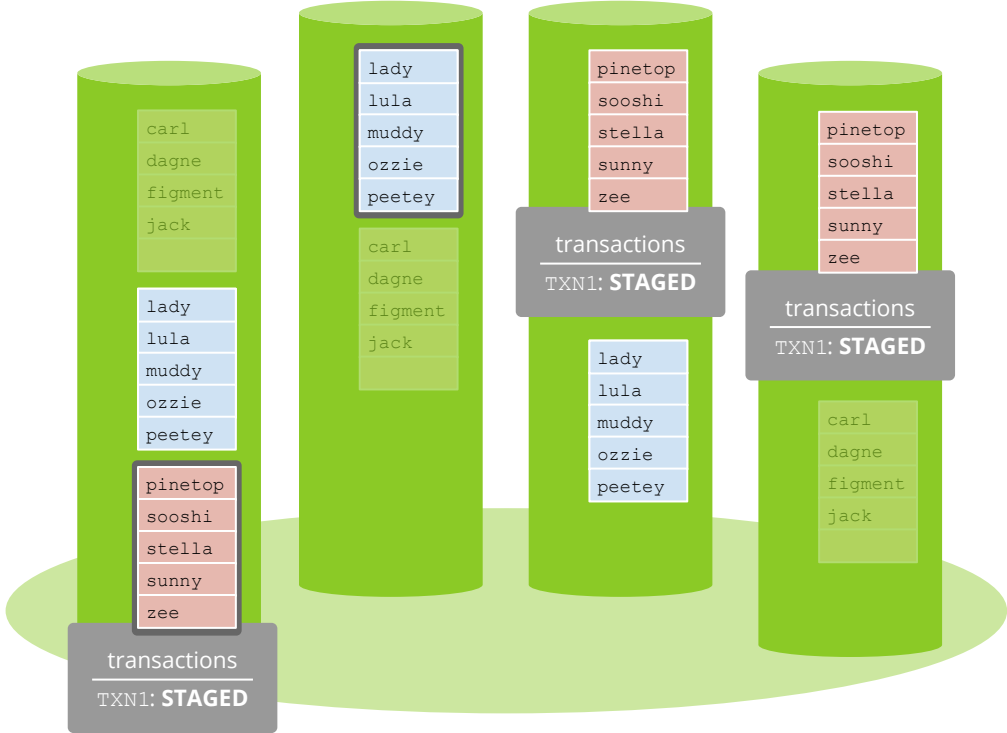
t = 1 RTT

# Parallel Commits

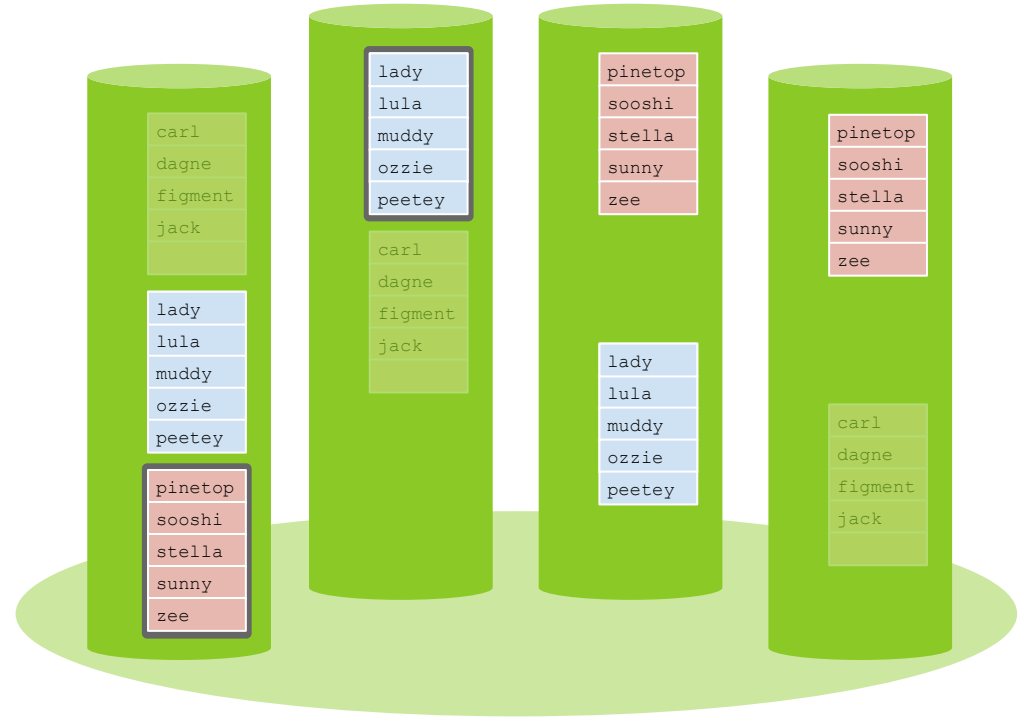


$t = 1 \text{ RTT}$

# Parallel Commits



# Parallel Commits



$t = 1 \text{ RTT}$





# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

- I. Spanner/Pipelined Transactions
- II. Parallel Commits
- III. Replicated Commit
- IV. Carousel
- V. MDCC
- VI. SLOG/OceanVista
- VII. TAPIR

# Replicated Commit



From UC Santa Barbara, 2013. Distinguishes between intra-DC (10ms) and inter-DC (100ms) latencies. Each DC contains replicas of all ranges.

2013-03.pdf  
Page 1 of 7

## Low-Latency Multi-Datacenter Databases using Replicated Commits

Hatem A. Mahmoud, Alexander Pucher, Faisal Nawab,  
Divyakant Agrawal, Amr El Abbadi  
University of California  
Santa Barbara, CA, USA  
{hatem.pucher,nawab,agrawal,amr}@cs.ucsb.edu

**ABSTRACT**

Web service providers have been using NoSQL databases to provide scalability and availability for globally distributed data at the cost of sacrificing transactional guarantees. Recently, major web service providers like Google have moved towards building storage systems that provide ACID transactional guarantees for globally distributed data. For example, the newly published system, Spanner, uses Two-Phase Commit and Two-Phase Locking to provide atomicity and isolation for globally distributed data, running on top of Paxos to provide fault-tolerant log replication. We show in this paper that it is possible to provide the same ACID transactional guarantees for multi-datacenter databases with fewer cross-datacenter communication trips, compared to replicated logging, by using a more efficient architecture. Instead of replicating the transactional log, we replicate the commit operation itself, by running Two-Phase Commit multiple times in different datacenters, and use Paxos to reach consensus among datacenters as to whether the transaction should commit. Doing so not only replaces several inter-datacenter communication trips with intra-datacenter communication trips, but also allows us to integrate atomic commitment and isolation protocols with consistent replication protocols so as to further reduce the number of cross-datacenter communication trips needed for consistent replication; for example, by eliminating the need for an election phase in Paxos.

**1. INTRODUCTION**

The rapid increase in the amount of data that is handled by web services as well as the globally-distributed client base of those web services has driven many web service providers towards adopting

cently, however, major web service providers have moved towards building storage systems that provide unrestricted ACID transactional guarantees. Google's Spanner [8] is a prominent example of such new trend. Spanner uses Two-Phase Commit and Two-Phase Locking to provide atomicity and isolation, running on top of a Paxos-replicated log to provide fault-tolerant synchronous replication across datacenters. The same architecture is also used in Scatter [11], a distributed hashtable datastore that provides ACID transactional guarantees for sharded, globally replicated data, through a key-value interface. Such layered architecture, in which the protocols that guarantee transactional atomicity and isolation are separated from the protocol that guarantees fault-tolerant replication, has many advantages from an engineering perspective, such as modularity, and clarity of semantics.

We show in this paper that it is possible to provide ACID transactional guarantees for cross-datacenter databases with a smaller number of cross-datacenter roundtrips, compared to a system that uses log replication, such as Spanner, by using a more efficient architecture. Instead of running Two-Phase Commit and Two-Phase Locking on top of Paxos to replicate the transactional log, we run Paxos on top of Two-Phase Commit and Two-Phase Locking to replicate the commit operation itself. That is, we execute the Two-Phase commit multiple time, once per datacenter, with each datacenter executing Two-Phase Commit and Two-Phase Locking internally, and we use Paxos to reach a consensus among datacenters as to whether the transaction should eventually commit. We refer to this approach as Replicated Commit, in contrast to the replicated log approach.

Replicated Commit has the advantage of replacing several inter-datacenter communication trips with intra-datacenter communi-

# Replicated Commit



Suggests running 2PC multiple times in parallel in each DC, and using Paxos across DCs to determine if txn should commit.

2013-03.pdf  
Page 1 of 7

## Low-Latency Multi-Datacenter Databases using Replicated Commits

Hatem A. Mahmoud, Alexander Pucher, Faisal Nawab,  
Divyakant Agrawal, Amr El Abbadi  
University of California  
Santa Barbara, CA, USA  
{hatem.pucher,nawab,agrawal,amr}@cs.ucsb.edu

**ABSTRACT**

Web service providers have been using NoSQL databases to provide scalability and availability for globally distributed data at the cost of sacrificing transactional guarantees. Recently, major web service providers like Google have moved towards building storage systems that provide ACID transactional guarantees for globally distributed data. For example, the newly published system, Spanner, uses Two-Phase Commit and Two-Phase Locking to provide atomicity and isolation for globally distributed data, running on top of Paxos to provide fault-tolerant log replication. We show in this paper that it is possible to provide the same ACID transactional guarantees for multi-datacenter databases with fewer cross-datacenter communication trips, compared to replicated logging, by using a more efficient architecture. Instead of replicating the transactional log, we replicate the commit operation itself, by running Two-Phase Commit multiple times in different datacenters, and use Paxos to reach consensus among datacenters as to whether the transaction should commit. Doing so not only replaces several inter-datacenter communication trips with intra-datacenter communication trips, but also allows us to integrate atomic commitment and isolation protocols with consistent replication protocols so as to further reduce the number of cross-datacenter communication trips needed for consistent replication; for example, by eliminating the need for an election phase in Paxos.

**1. INTRODUCTION**

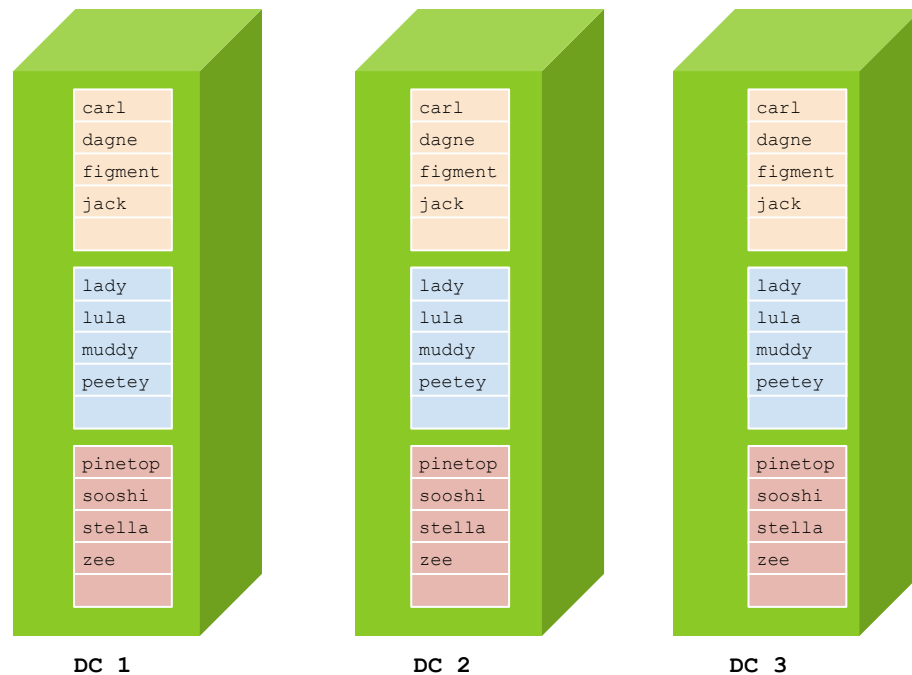
The rapid increase in the amount of data that is handled by web services as well as the globally-distributed client base of those web services has driven many web service providers towards adopting

cently, however, major web service providers have moved towards building storage systems that provide unrestricted ACID transactional guarantees. Google's Spanner [8] is a prominent example of such new trend. Spanner uses Two-Phase Commit and Two-Phase Locking to provide atomicity and isolation, running on top of a Paxos-replicated log to provide fault-tolerant synchronous replication across datacenters. The same architecture is also used in Scatter [11], a distributed hashtable datastore that provides ACID transactional guarantees for sharded, globally replicated data, through a key-value interface. Such layered architecture, in which the protocols that guarantee transactional atomicity and isolation are separated from the protocol that guarantees fault-tolerant replication, has many advantages from an engineering perspective, such as modularity, and clarity of semantics.

We show in this paper that it is possible to provide ACID transactional guarantees for cross-datacenter databases with a smaller number of cross-datacenter roundtrips, compared to a system that uses log replication, such as Spanner, by using a more efficient architecture. Instead of running Two-Phase Commit and Two-Phase Locking on top of Paxos to replicate the transactional log, we run Paxos on top of Two-Phase Commit and Two-Phase Locking to replicate the commit operation itself. That is, we execute the Two-Phase commit multiple time, once per datacenter, with each datacenter executing Two-Phase Commit and Two-Phase Locking internally, and we use Paxos to reach a consensus among datacenters as to whether the transaction should eventually commit. We refer to this approach as Replicated Commit, in contrast to the replicated log approach.

Replicated Commit has the advantage of replacing several inter-datacenter communication trips with intra-datacenter communi-

# Replicated Commit



$t = 0$

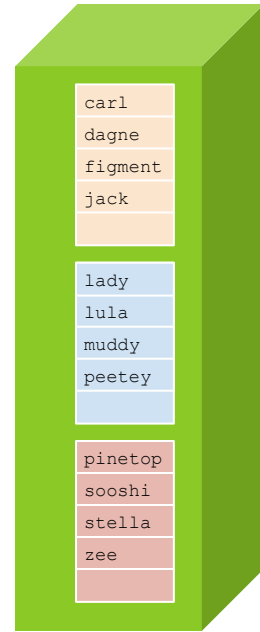
# Replicated Commit



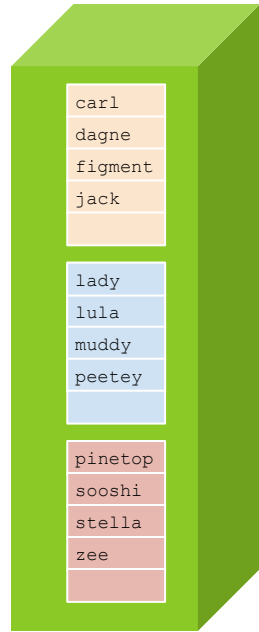
```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]
```

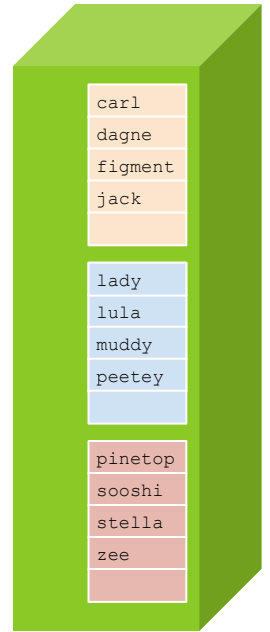
GATEWAY



DC 1



DC 2



DC 3

t = 0

# Replicated Commit



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

**BEGIN TXN1**  
**WRITE[sunny]**

GATEWAY

carl
dagne
figment
jack

lady
lula
muddy
peetey

pinetop
sooshi
stella
sunny*
zee

transactions  
TXN1: **PENDING**

carl
dagne
figment
jack

lady
lula
muddy
peetey

pinetop
sooshi
stella
sunny*
zee

transactions  
TXN1: **PENDING**

carl
dagne
figment
jack

lady
lula
muddy
peetey

pinetop
sooshi
stella
sunny*
zee

transactions  
TXN1: **PENDING**

t = 0

# Replicated Commit



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]
```

GATEWAY

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

transactions  
TXN1: PENDING

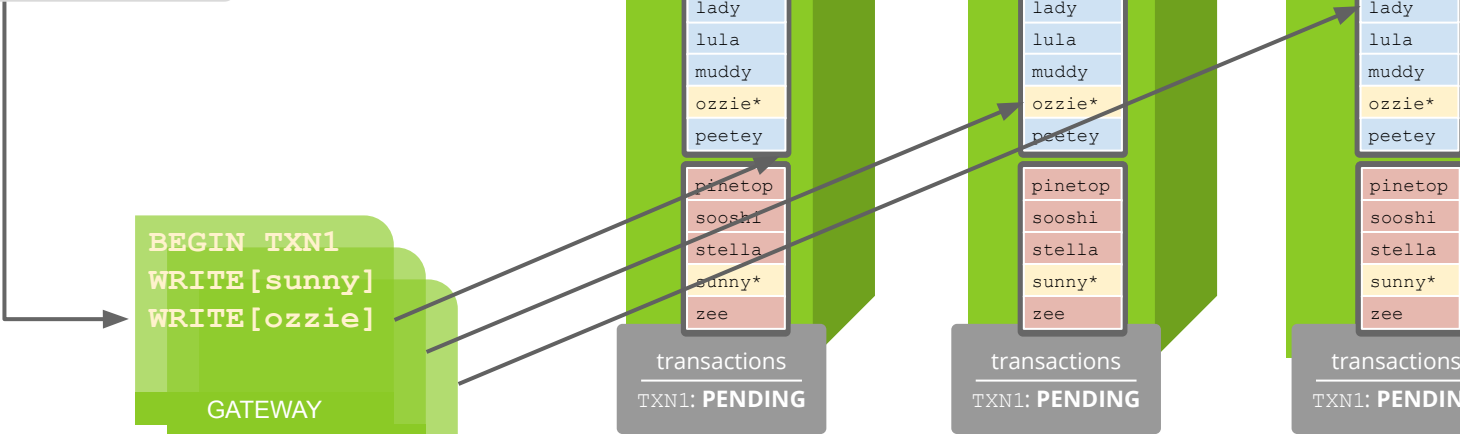
carl
dagne
figment
jack
lady
lula
muddy
ozzie*
peetey
pinetop
sooshi
stella
sunny*
zee

transactions  
TXN1: PENDING

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

transactions  
TXN1: PENDING

t = 1 LAN



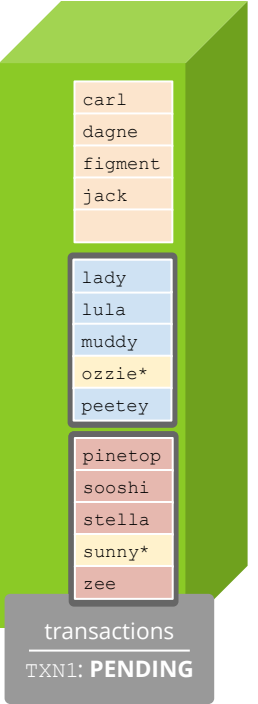
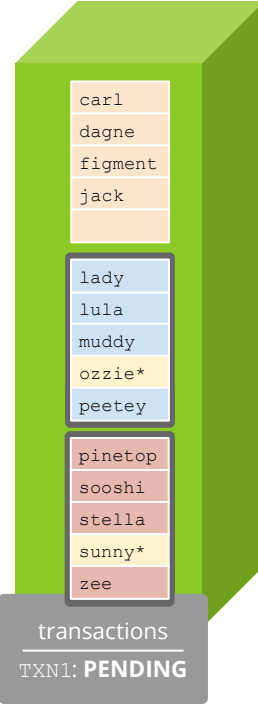
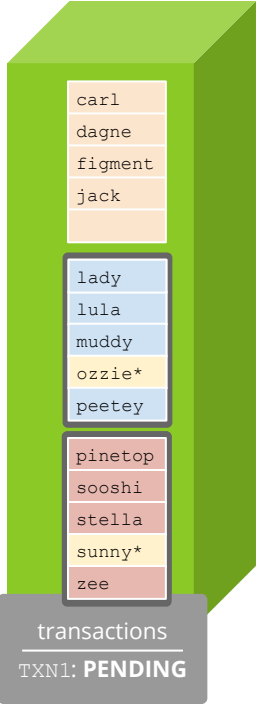
# Replicated Commit



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```

t = 1 LAN



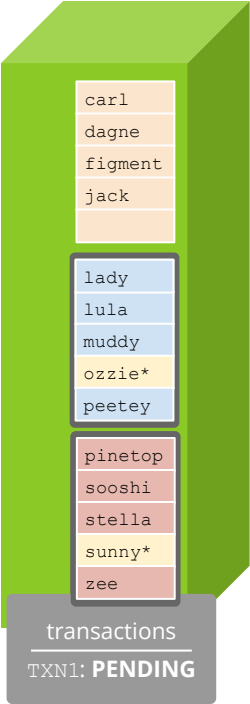
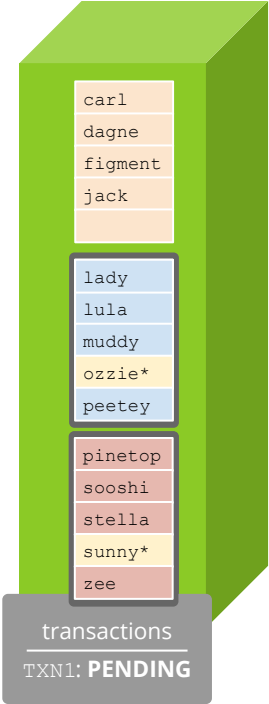
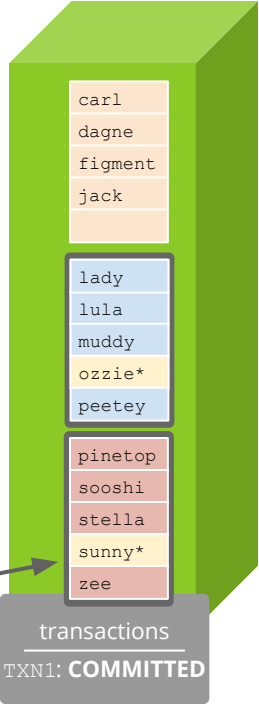


# Replicated Commit



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 1 LAN

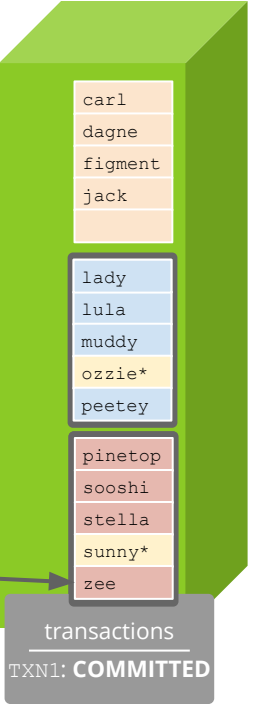
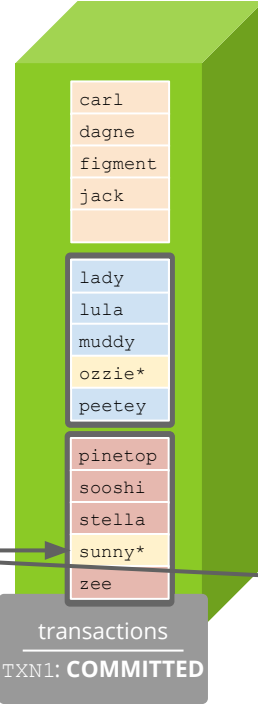
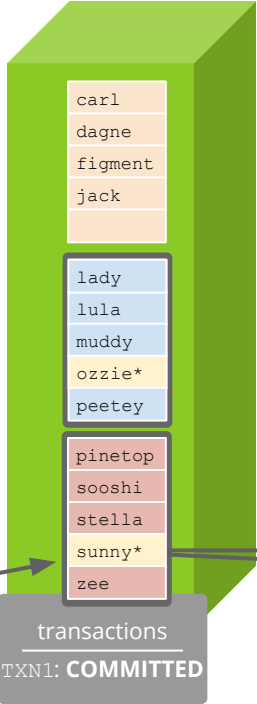
# Replicated Commit



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT
```

GATEWAY



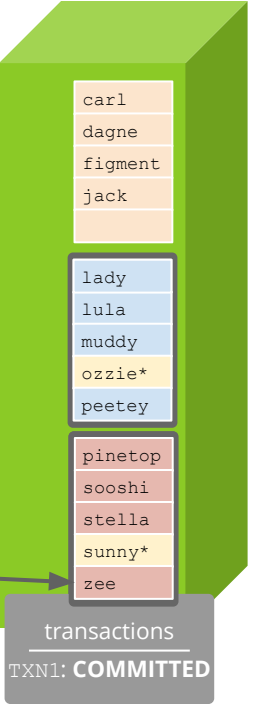
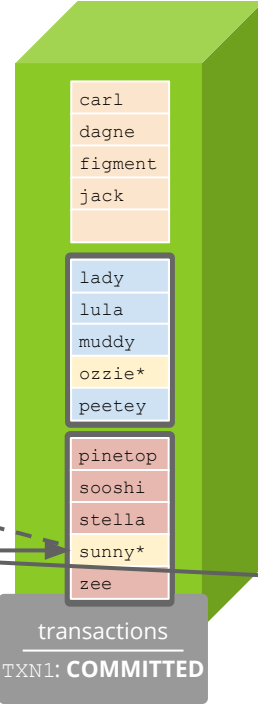
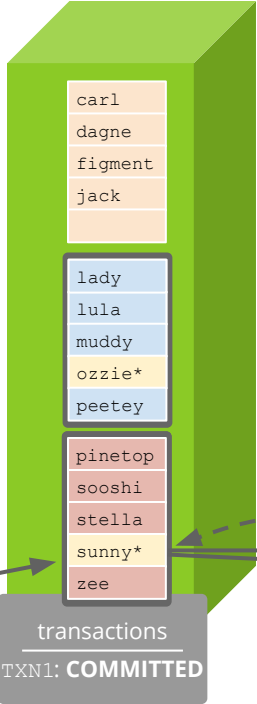
$$t = 1 \text{ LAN} + \frac{1}{2} \text{ WAN}$$

# Replicated Commit



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE[sunny]  
WRITE[ozzie]  
COMMIT  
GATEWAY
```



$$t = 1 \text{ LAN} + 1 \text{ WAN}$$

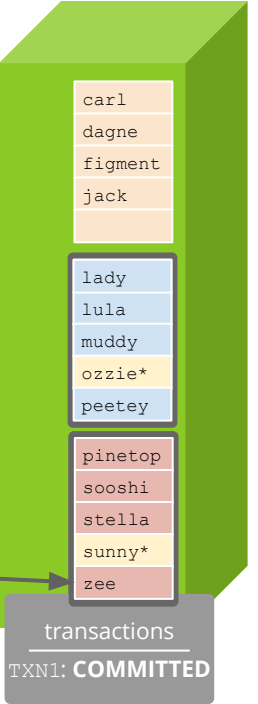
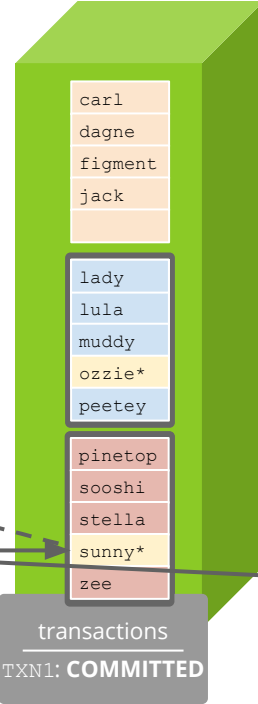
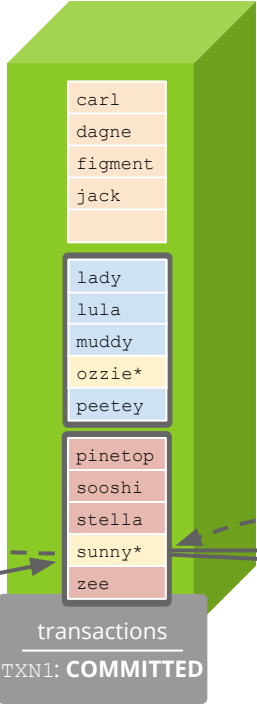
# Replicated Commit



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE[sunny]  
WRITE[ozzie]  
COMMIT
```

GATEWAY



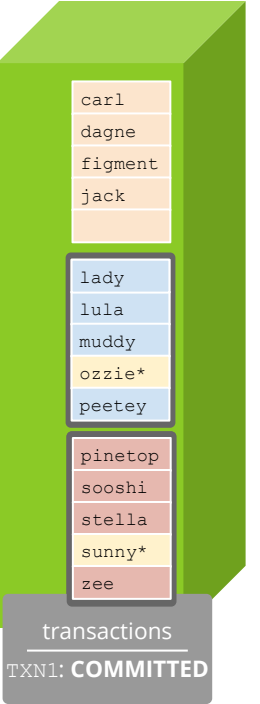
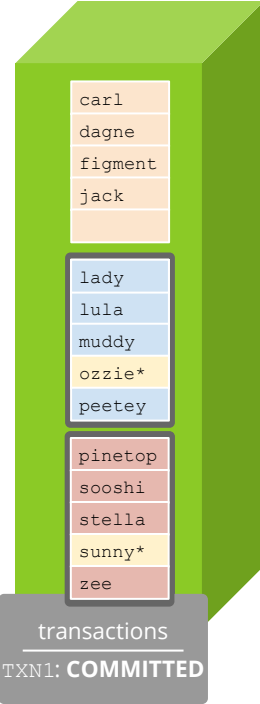
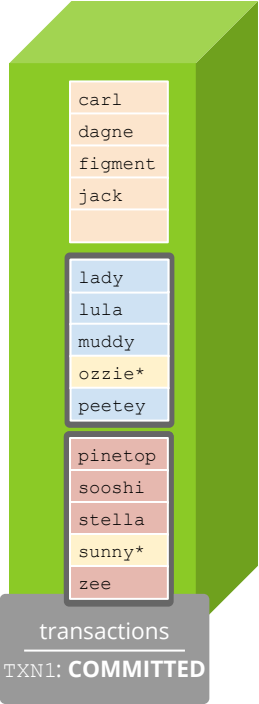
$t = 1 \text{ LAN} + 1 \text{ WAN}$

# Replicated Commit



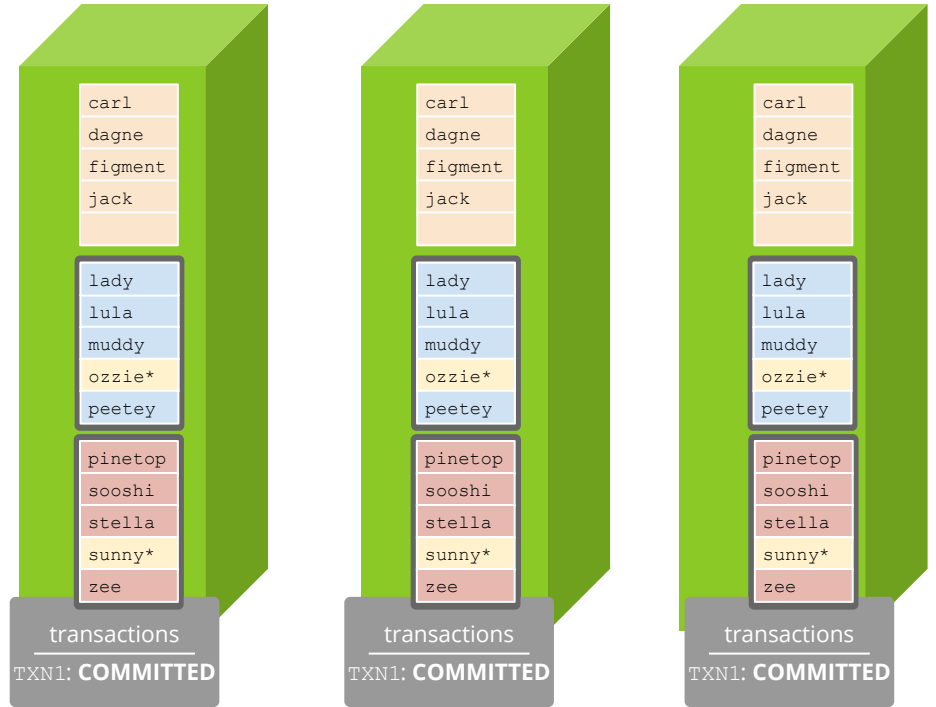
```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN TXN1  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



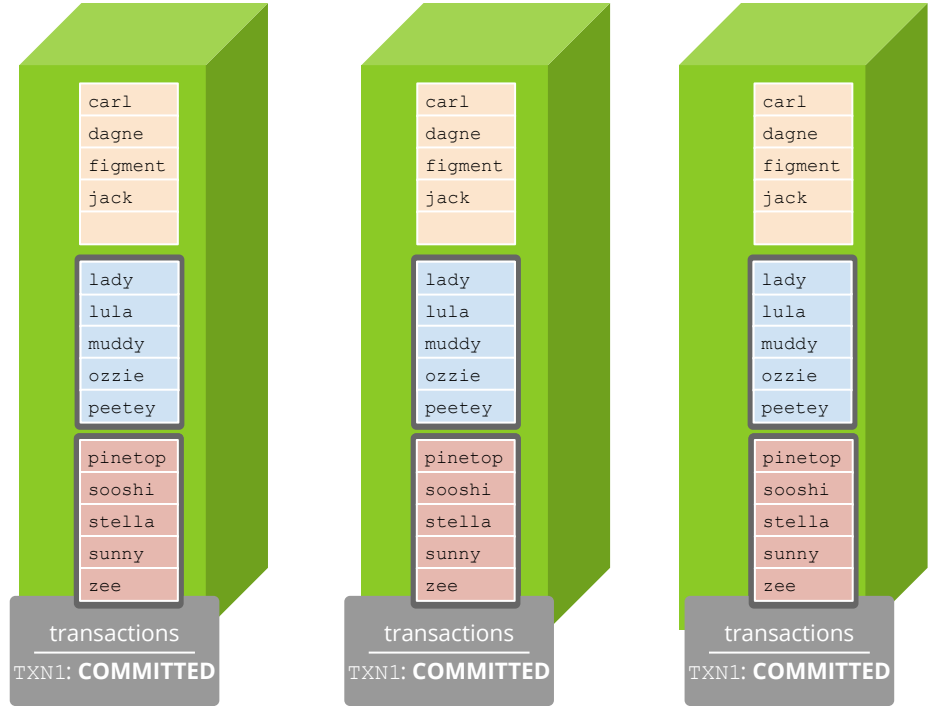
$$t = 1 \text{ LAN} + 1 \text{ WAN}$$

# Replicated Commit



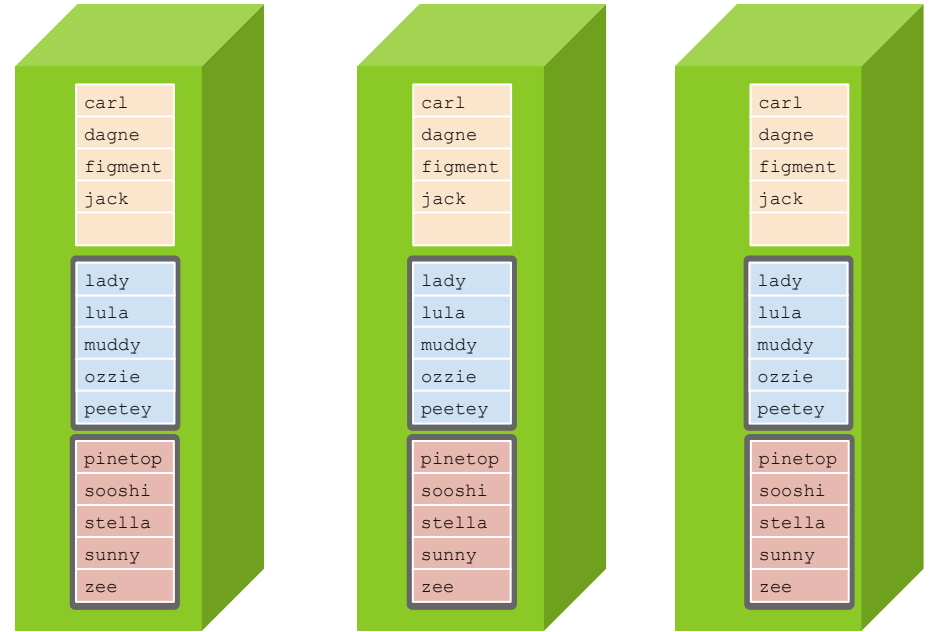
$$t = 1 \text{ LAN} + 1 \text{ WAN}$$

# Replicated Commit



$$t = 1 \text{ LAN} + 1 \text{ WAN}$$

# Replicated Commit



$$t = 1 \text{ LAN} + 1 \text{ WAN}$$



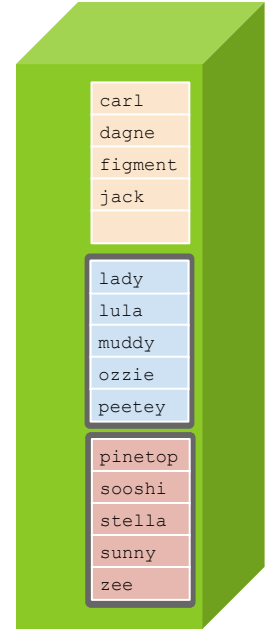
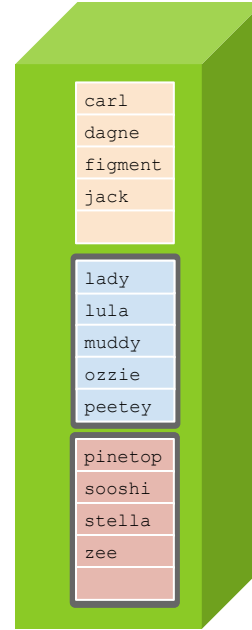
# Replicated Commit (reads)



```
BEGIN;  
READ (sunny) FROM dogs;  
COMMIT;
```

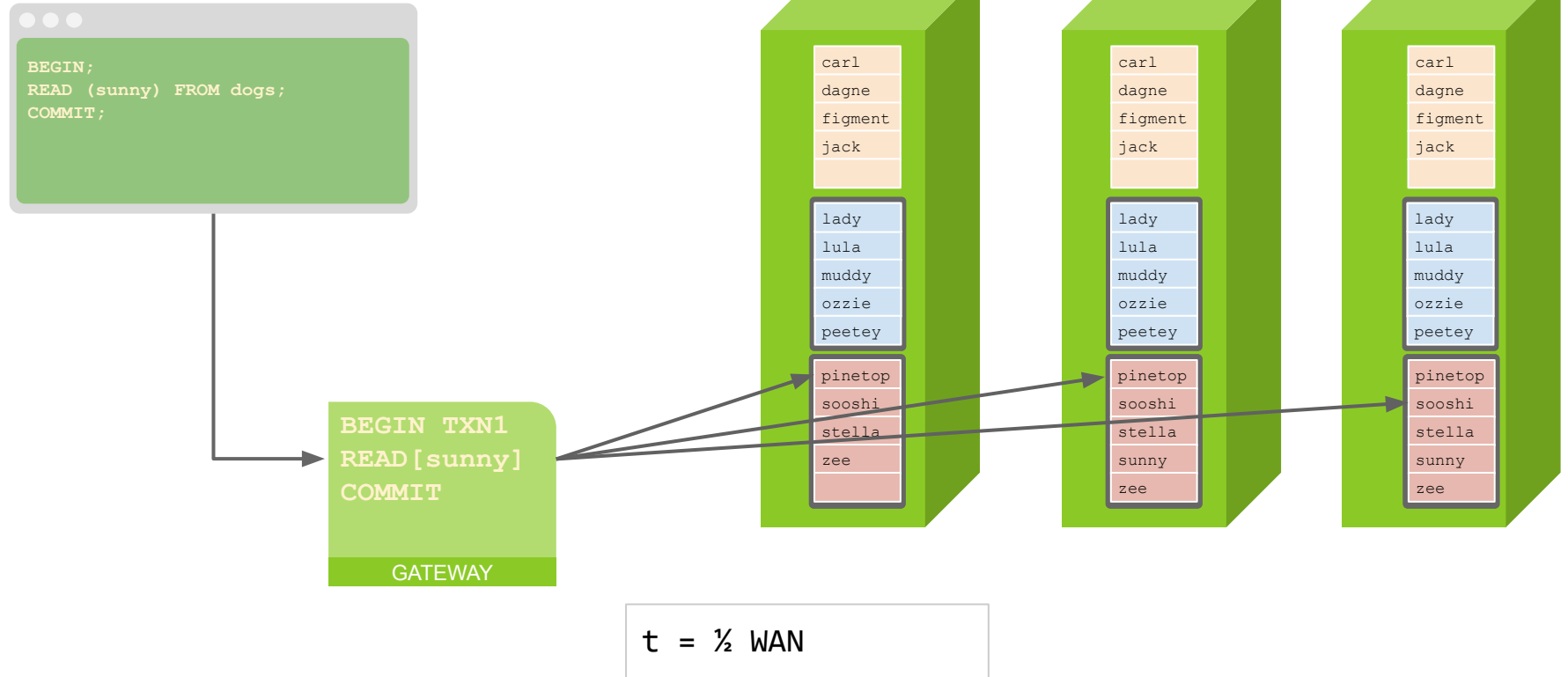
```
BEGIN TXN1  
READ [sunny]  
COMMIT
```

GATEWAY

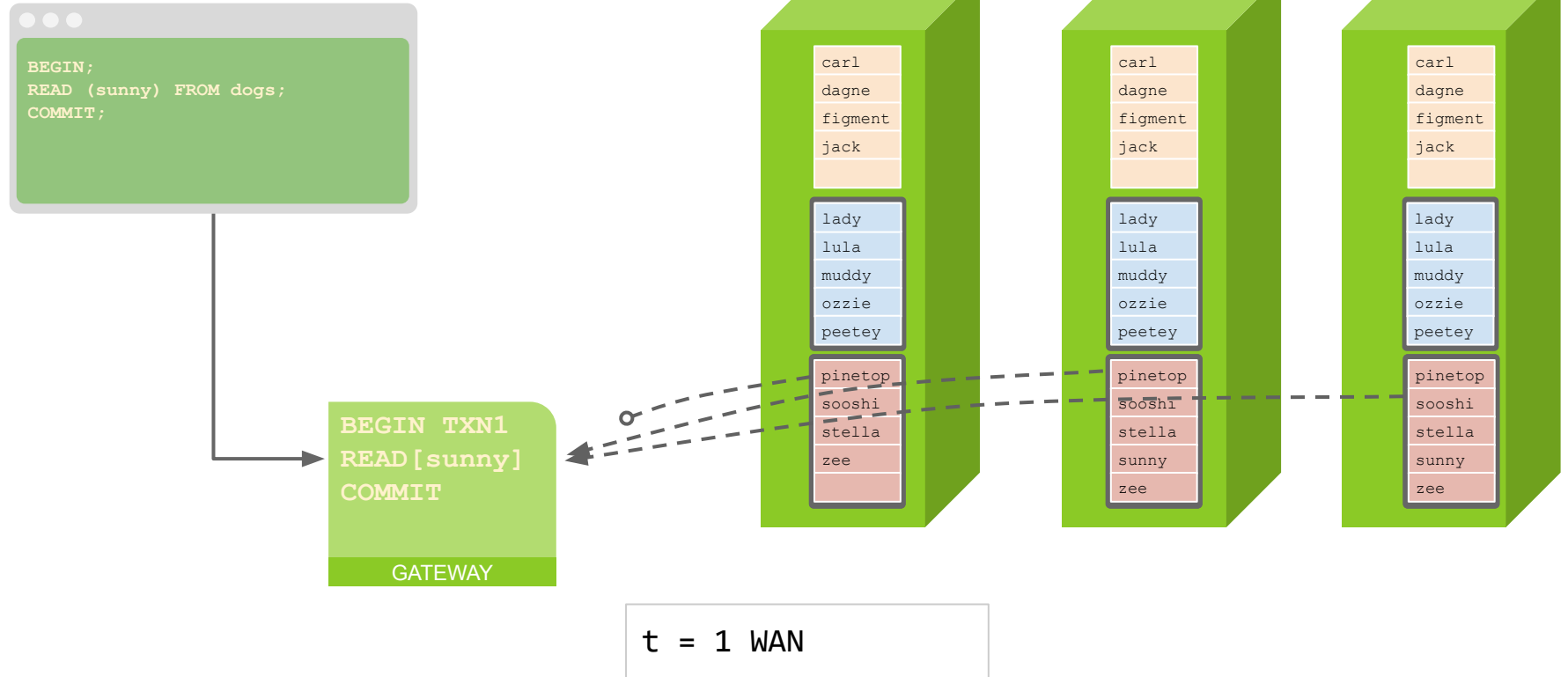


t = 0

# Replicated Commit (reads)



# Replicated Commit (reads)





# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

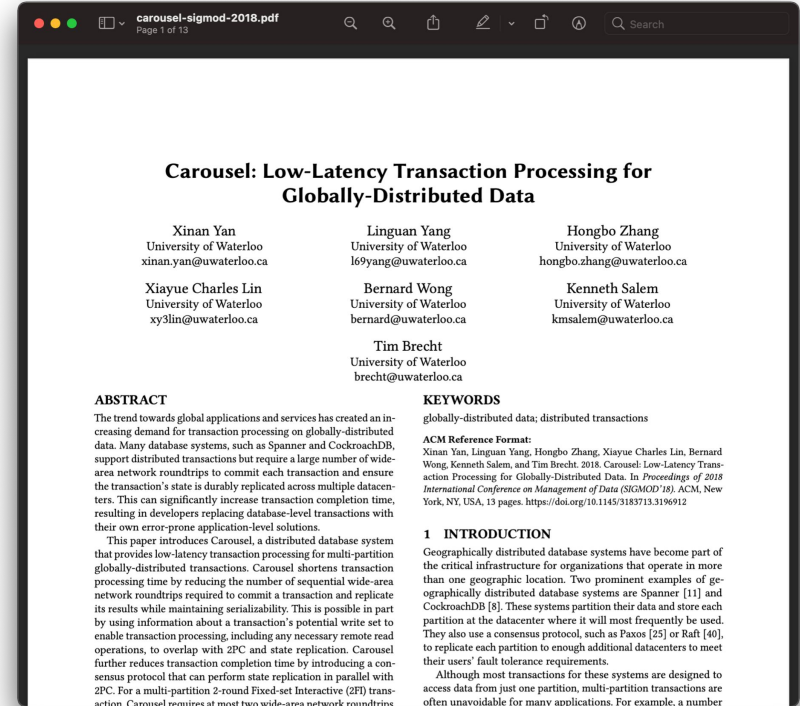
- I. Spanner/Pipelined Transactions
- II. Parallel Commits
- III. Replicated Commit
- IV. Carousel
- V. MDCC
- VI. SLOG/OceanVista
- VII. TAPIR



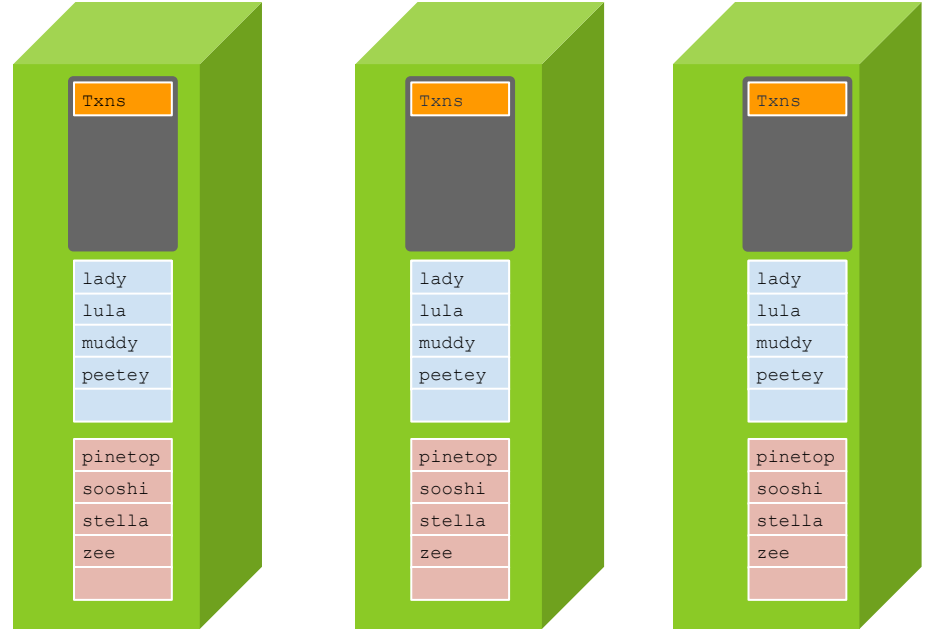
# Carousel

From Waterloo, 2018. Limits transaction model to 2FI (2-fixed set interactive): a round of reads followed by round of writes, with all keys declared in advance. Write values can depend on reads, but keys written to/read from cannot.

Uses Fast Paxos to replicate to all replicas directly, avoiding the leader hop. Retries on conflict, given it's optimistic.



# Carousel

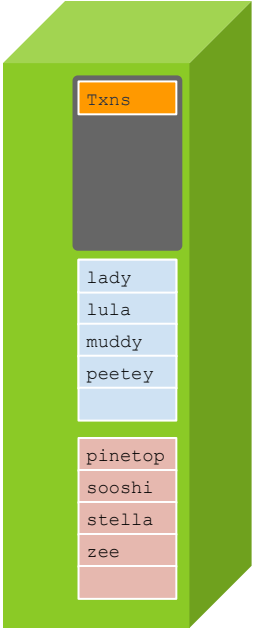
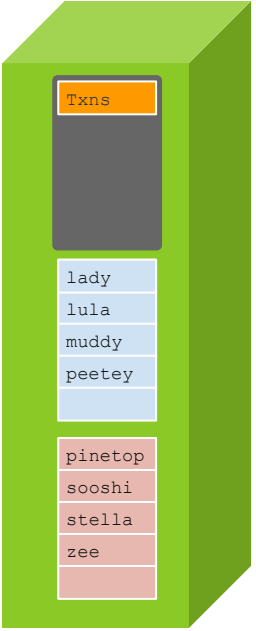
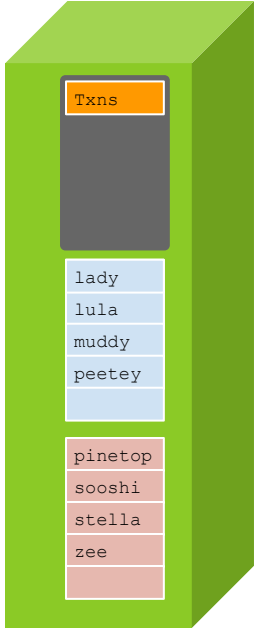
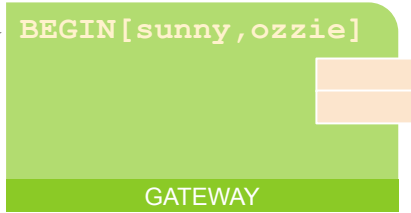


t = 0

# Carousel



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```



t = 0

# Carousel



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

BEGIN [sunny, ozzie]

GATEWAY

Txns

lady  
lula  
muddy  
peetey

pinetop  
sooshi  
stella  
sunny\*  
zee

Txns

lady  
lula  
muddy  
peetey

pinetop  
sooshi  
stella  
sunny\*  
zee

Txns

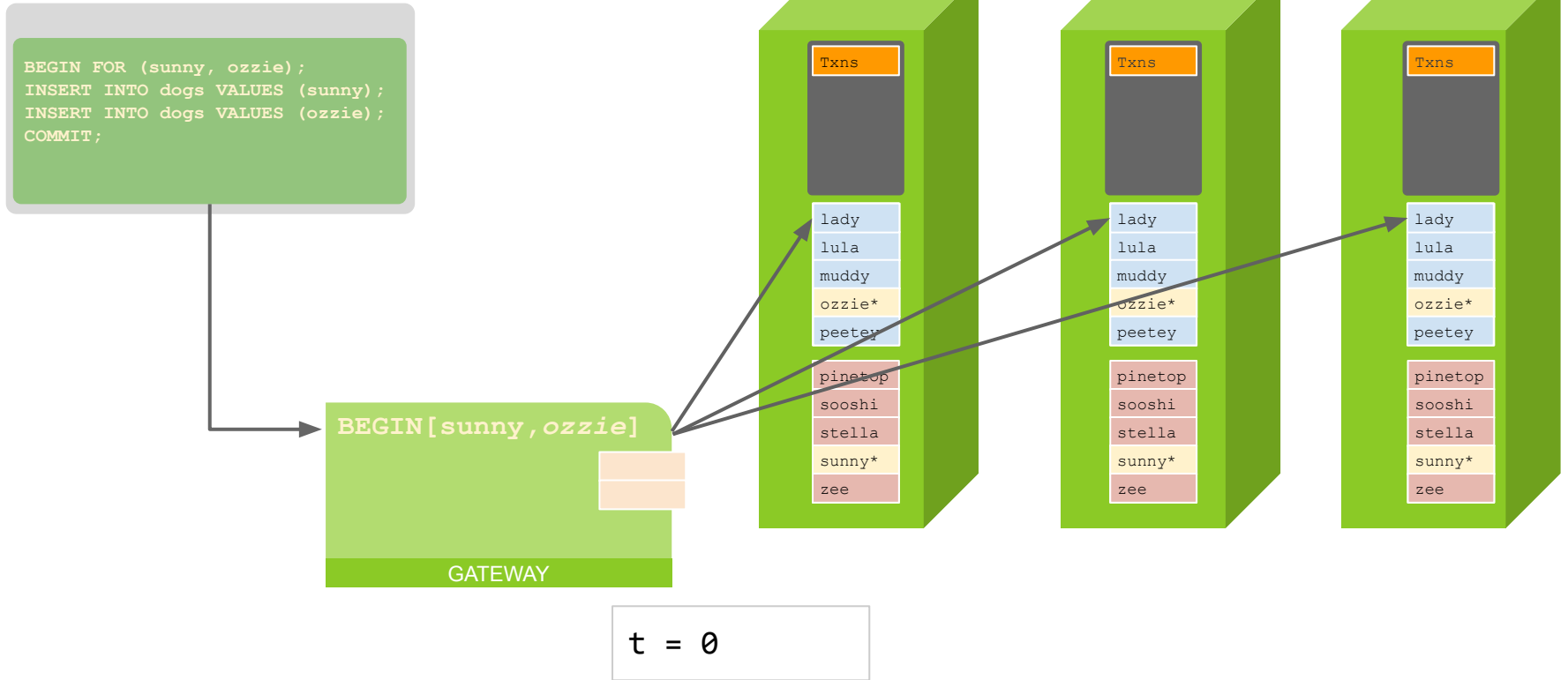
lady  
lula  
muddy  
peetey

pinetop  
sooshi  
stella  
sunny\*  
zee

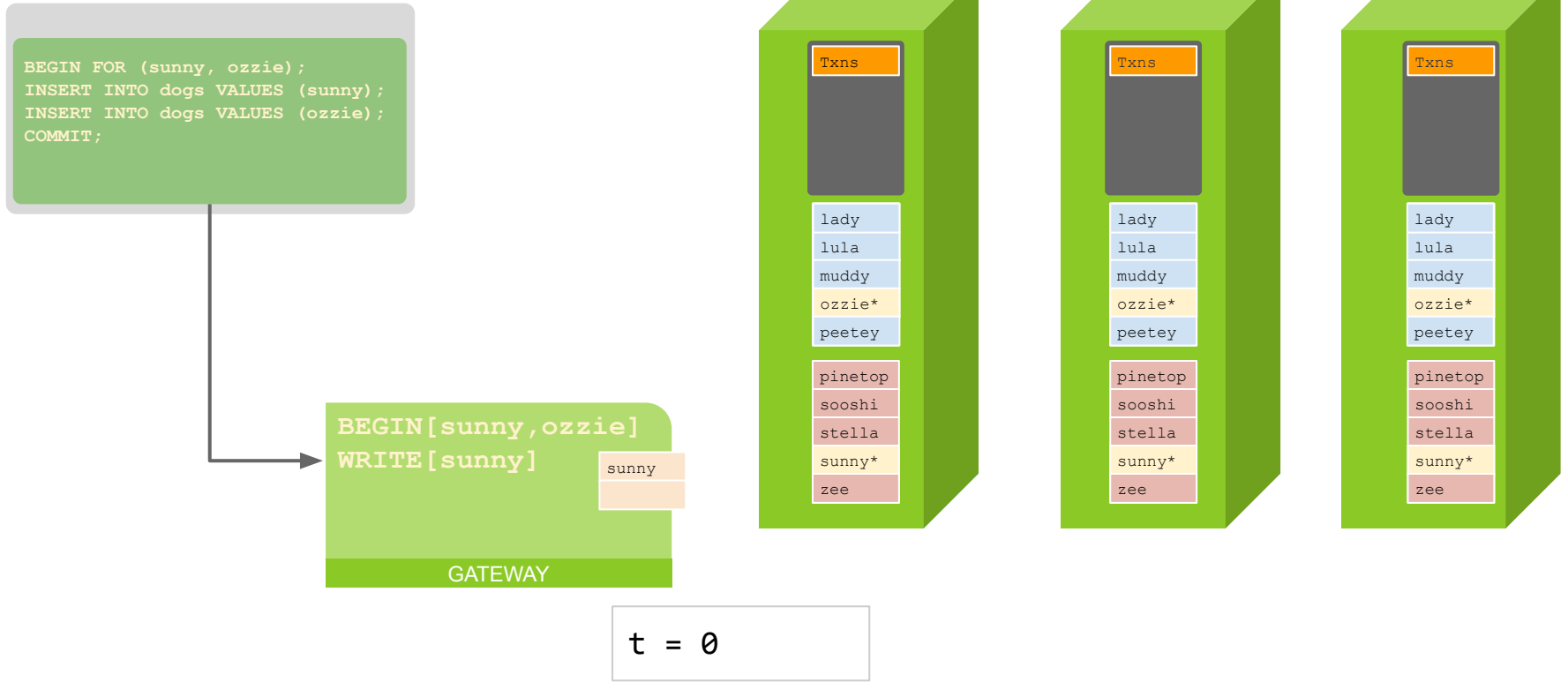
t = 0



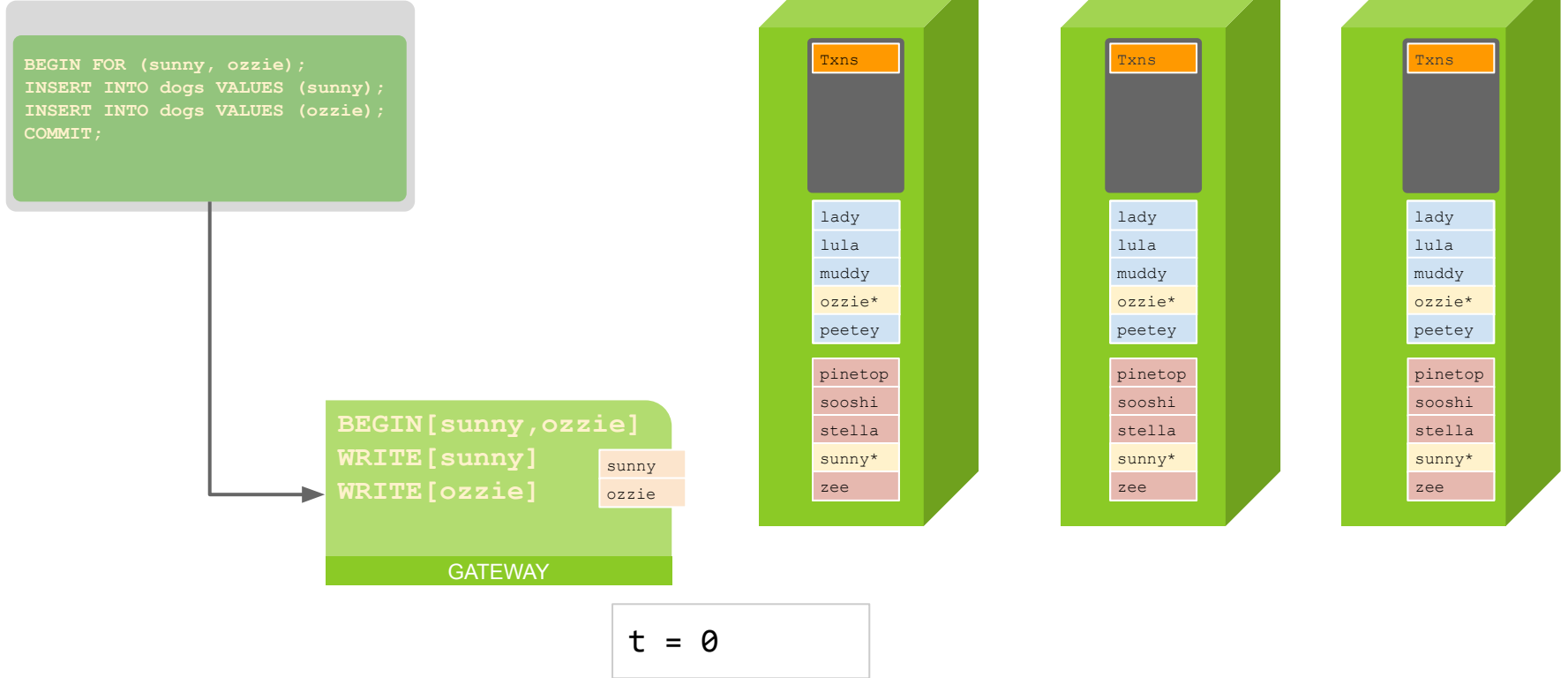
# Carousel



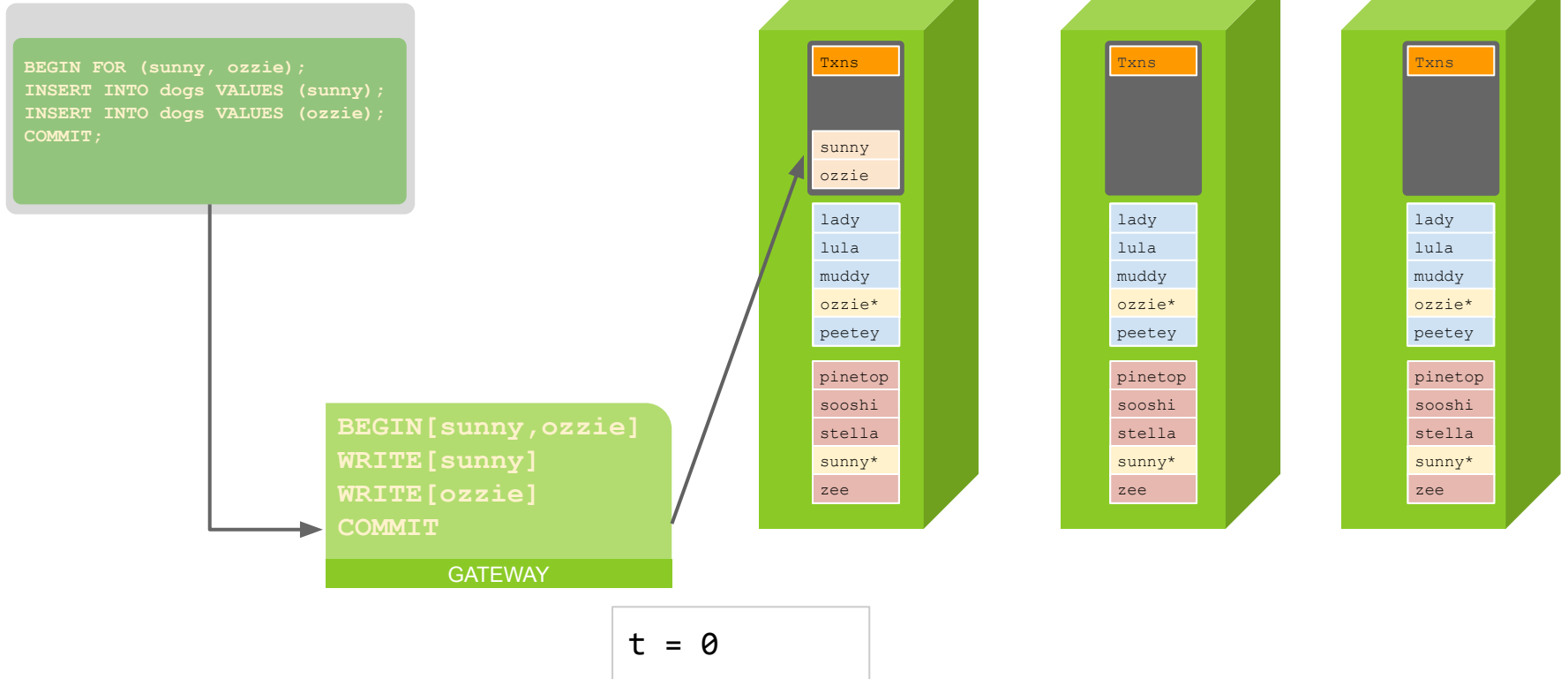
# Carousel



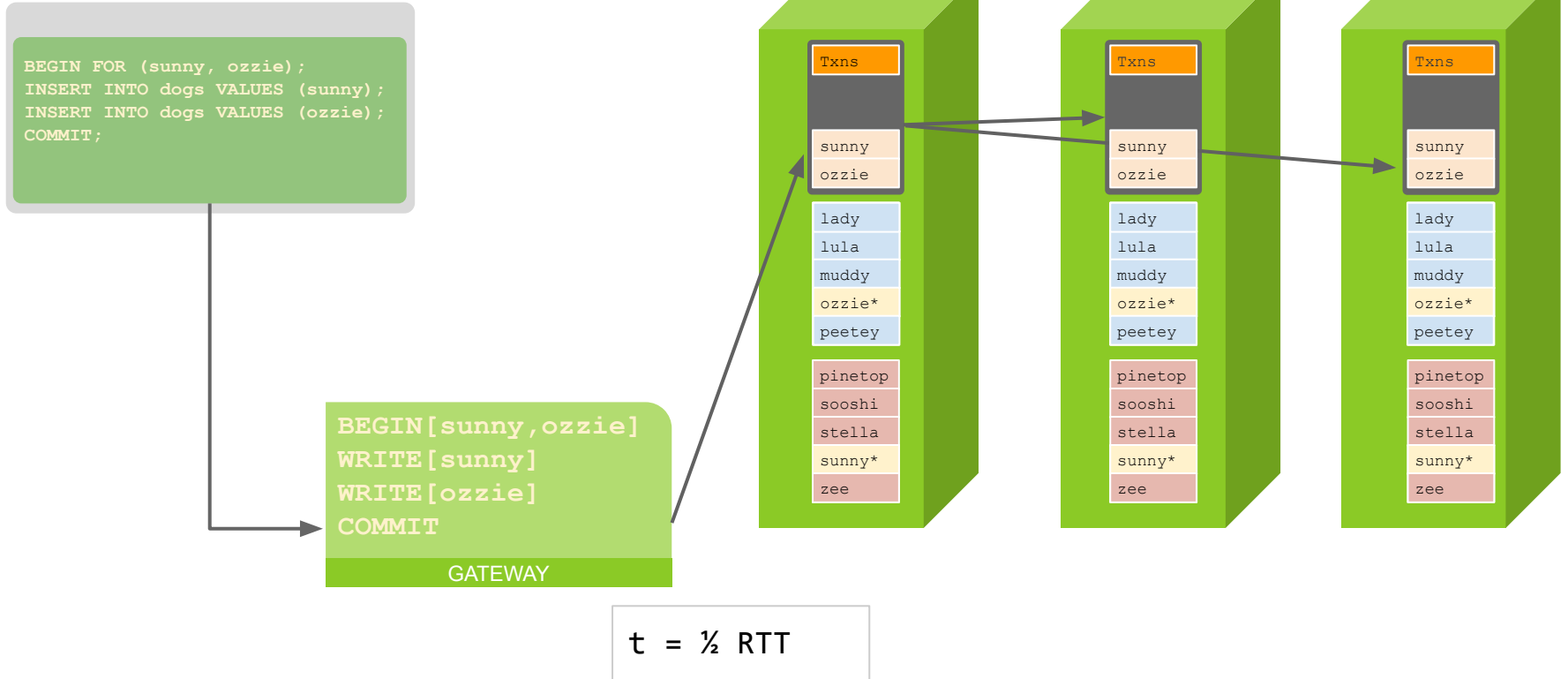
# Carousel



# Carousel

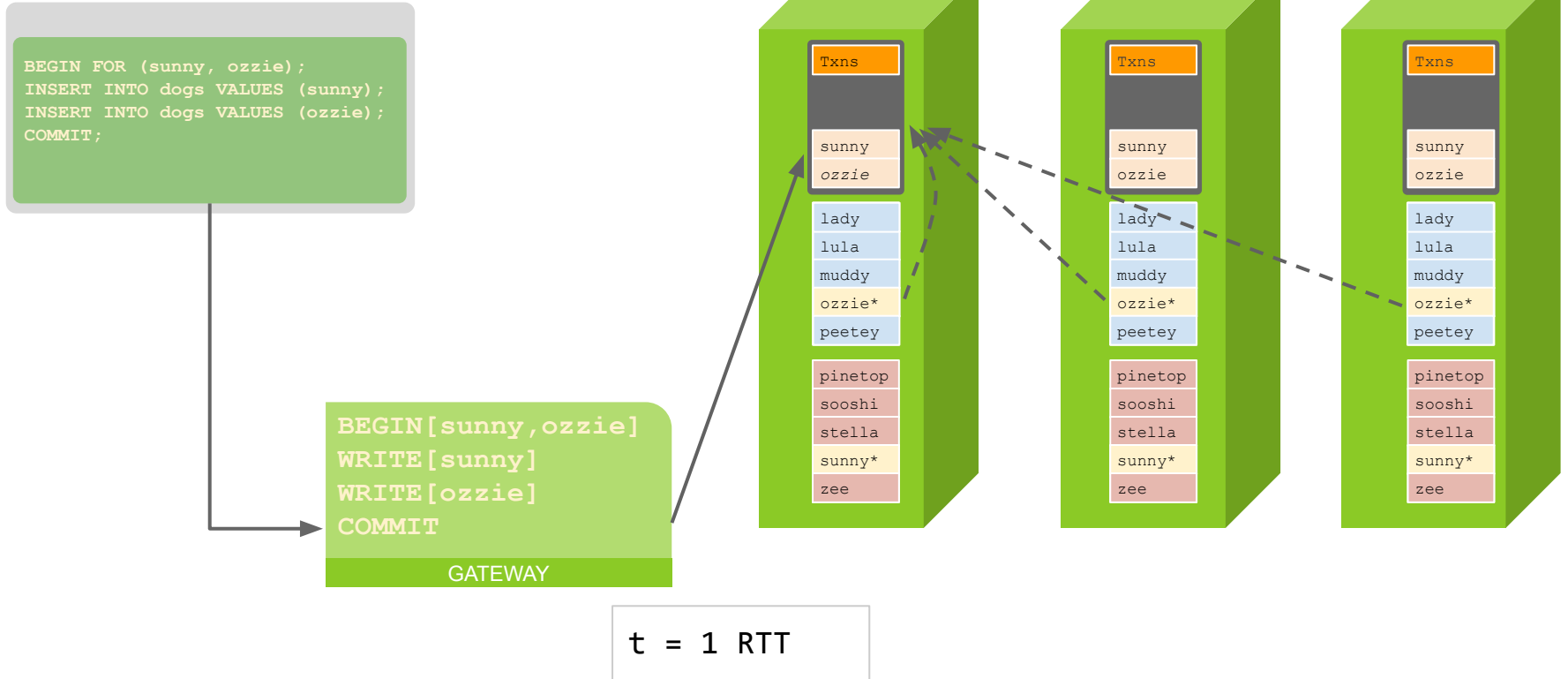


# Carousel

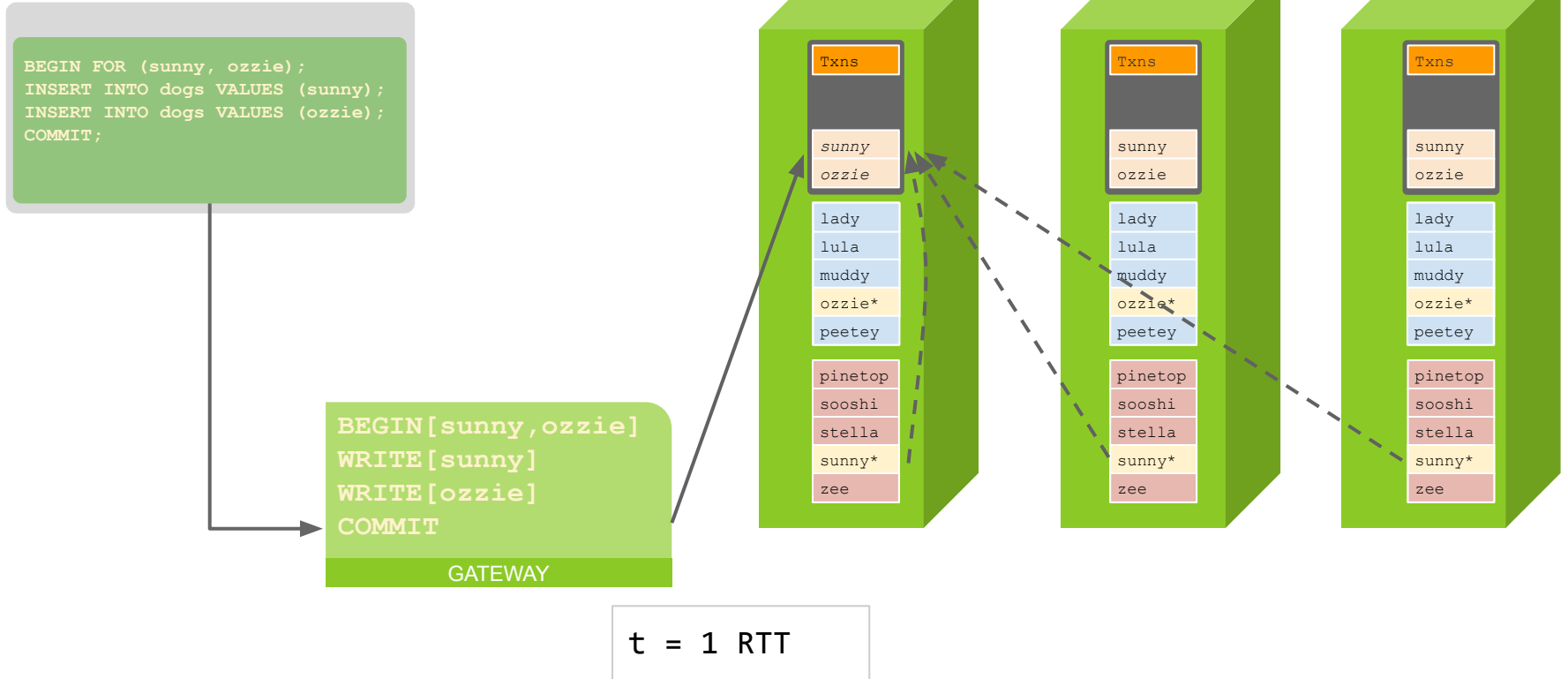




# Carousel

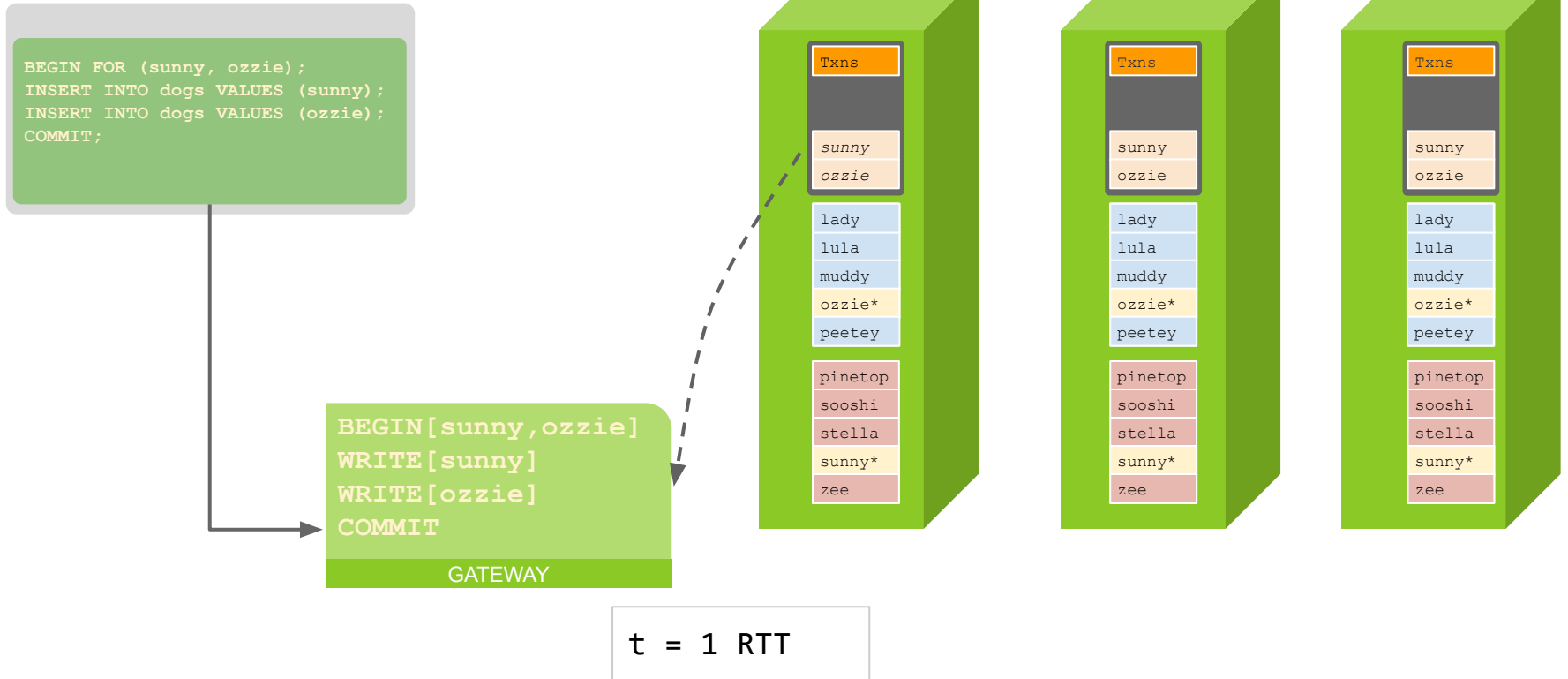


# Carousel

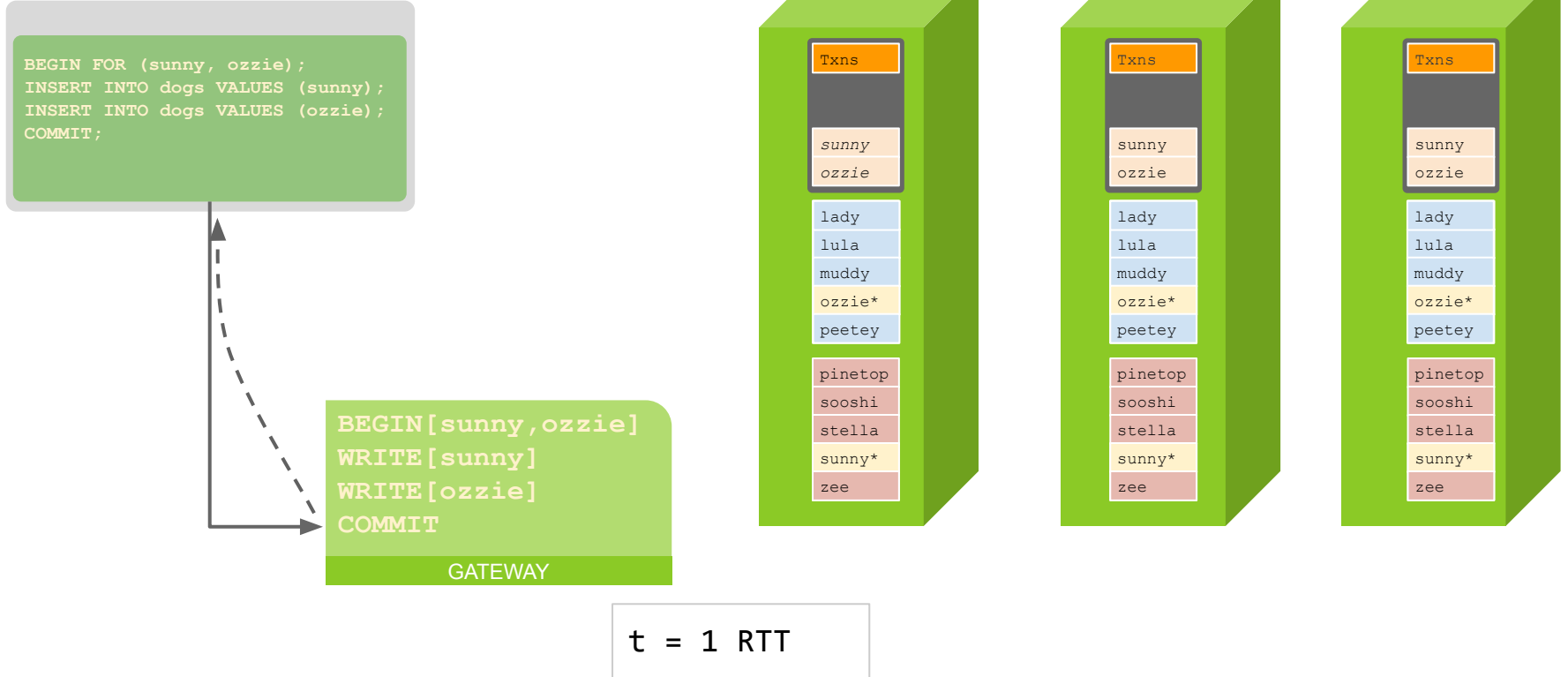




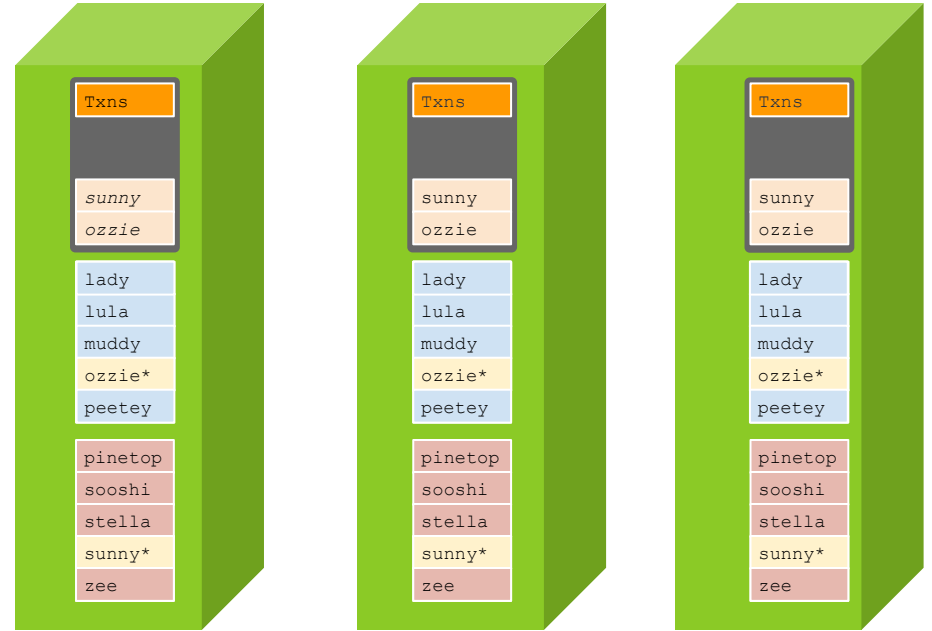
# Carousel



# Carousel

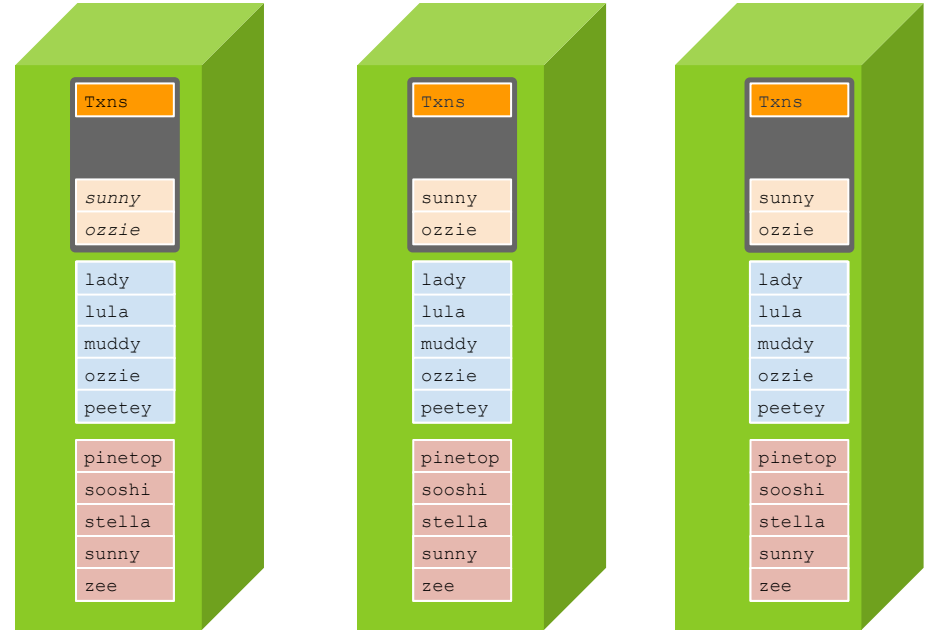


# Carousel



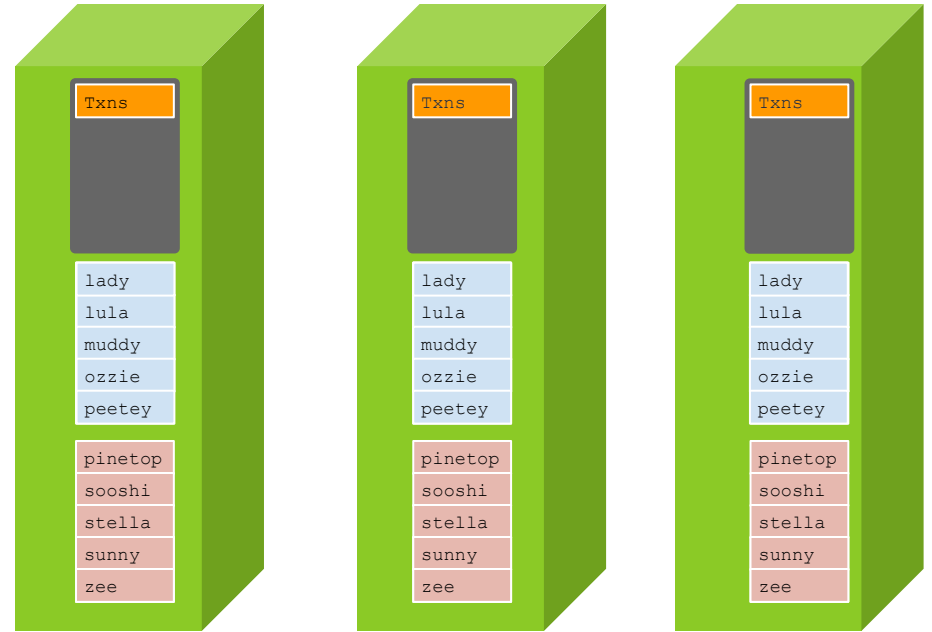
$t = 1 \text{ RTT}$

# Carousel



$t = 1 \text{ RTT}$

# Carousel



t = 1 RTT



# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

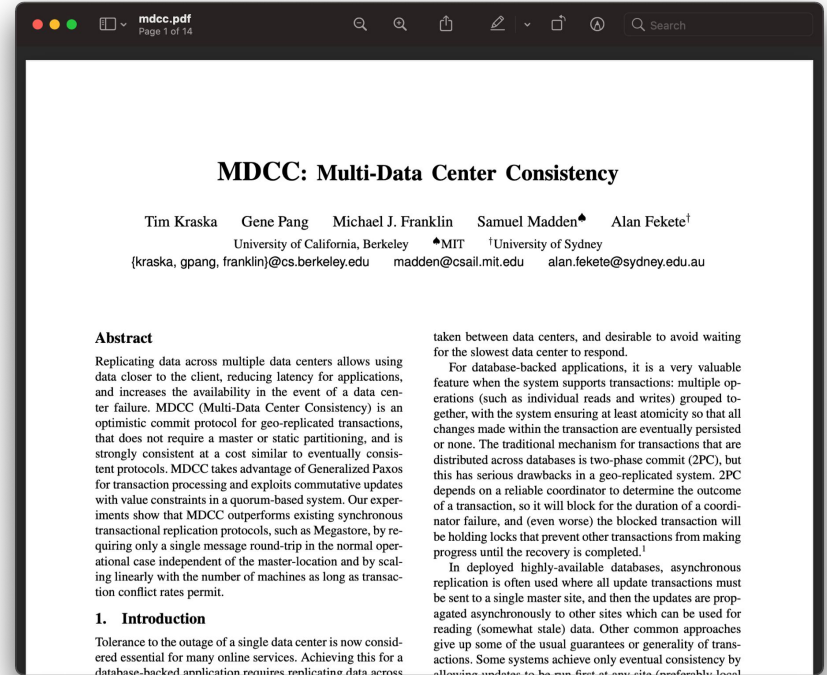
- I. Spanner/Pipelined Transactions
- II. Parallel Commits
- III. Replicated Commit
- IV. Carousel
- V. MDCC
- VI. SLOG/OceanVista
- VII. TAPIR

# MDCC

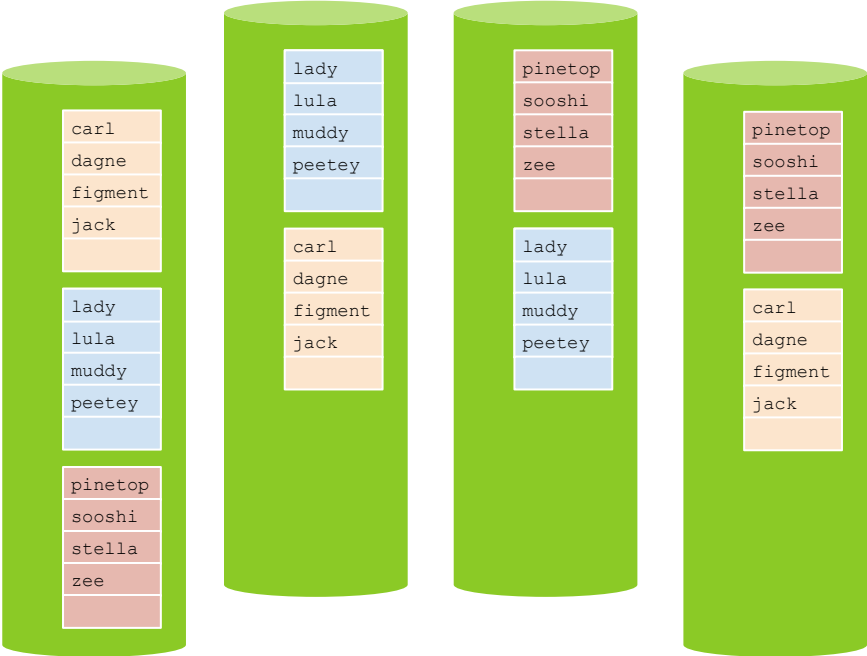


2013, from Berkeley. The closest one to Parallel Commits. Replicates “options” instead of update values directly, options record the full set of keys written to in the txn. Client can be notified of commit once all options are replicated.

Uses fast paxos to avoid leader hop, generalized paxos to reason about commutative operations. Supports only upto read-committed isolation.



# MDCC



$t = 0$

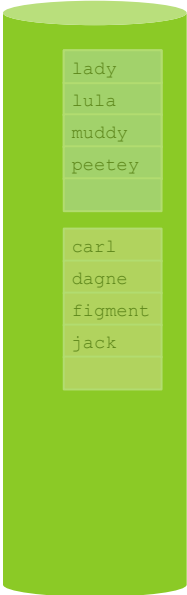
NB: with fast paxos we need at least 4 replicas, not shown here



# MDCC



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```



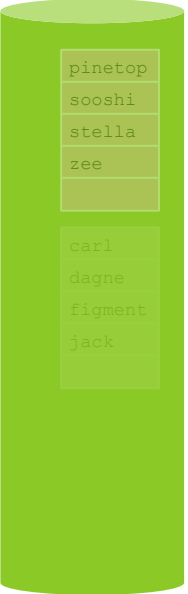
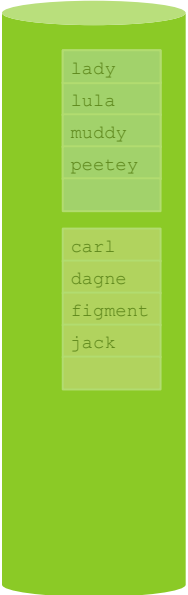
t = 0

# MDCC



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

```
BEGIN  
WRITE [sunny] sunny  
GATEWAY
```



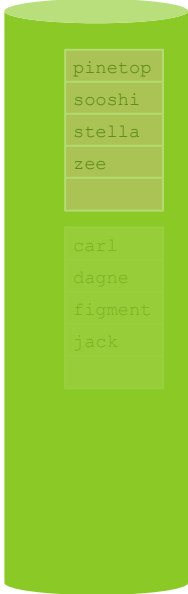
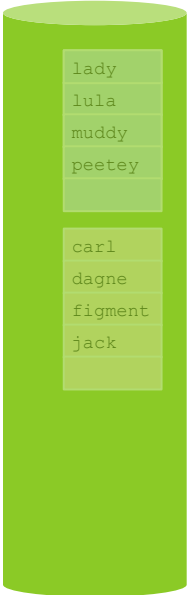
t = 0

# MDCC



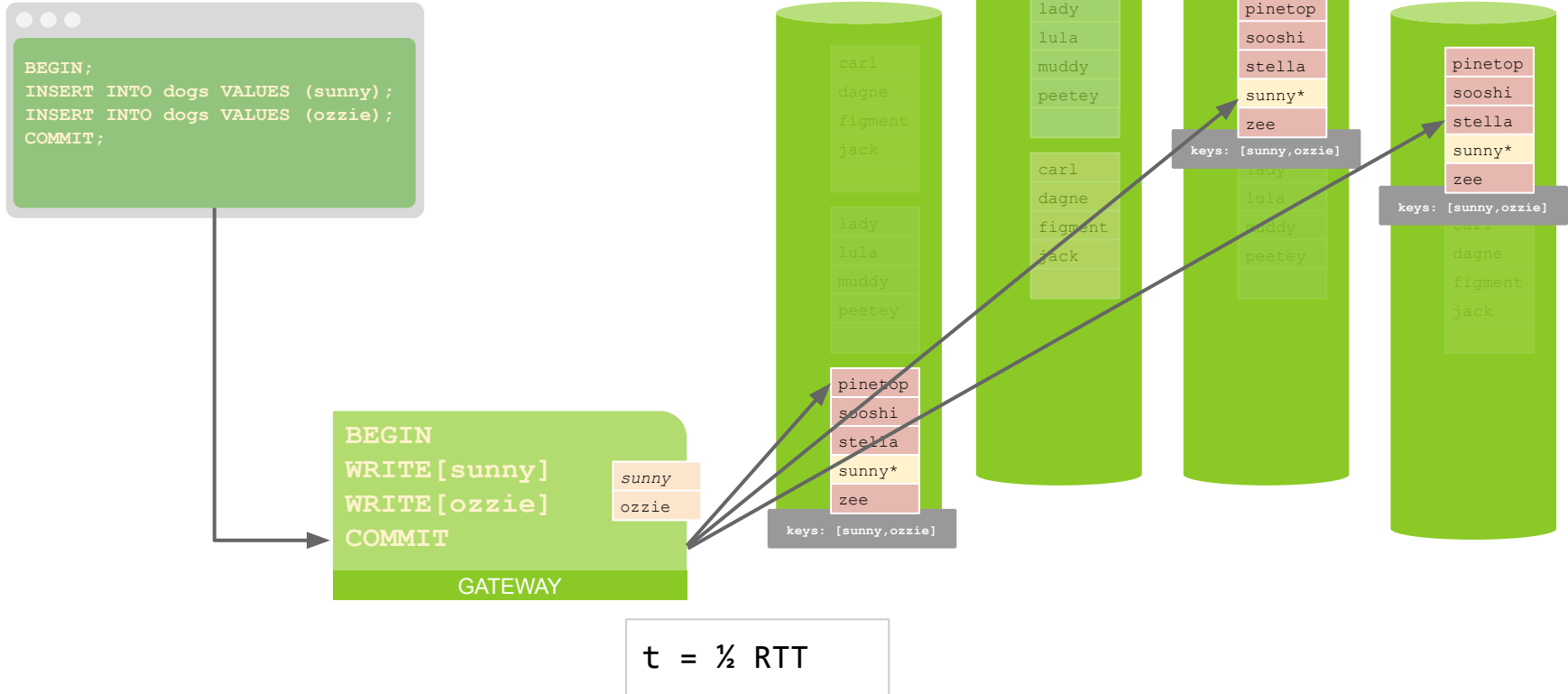
```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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

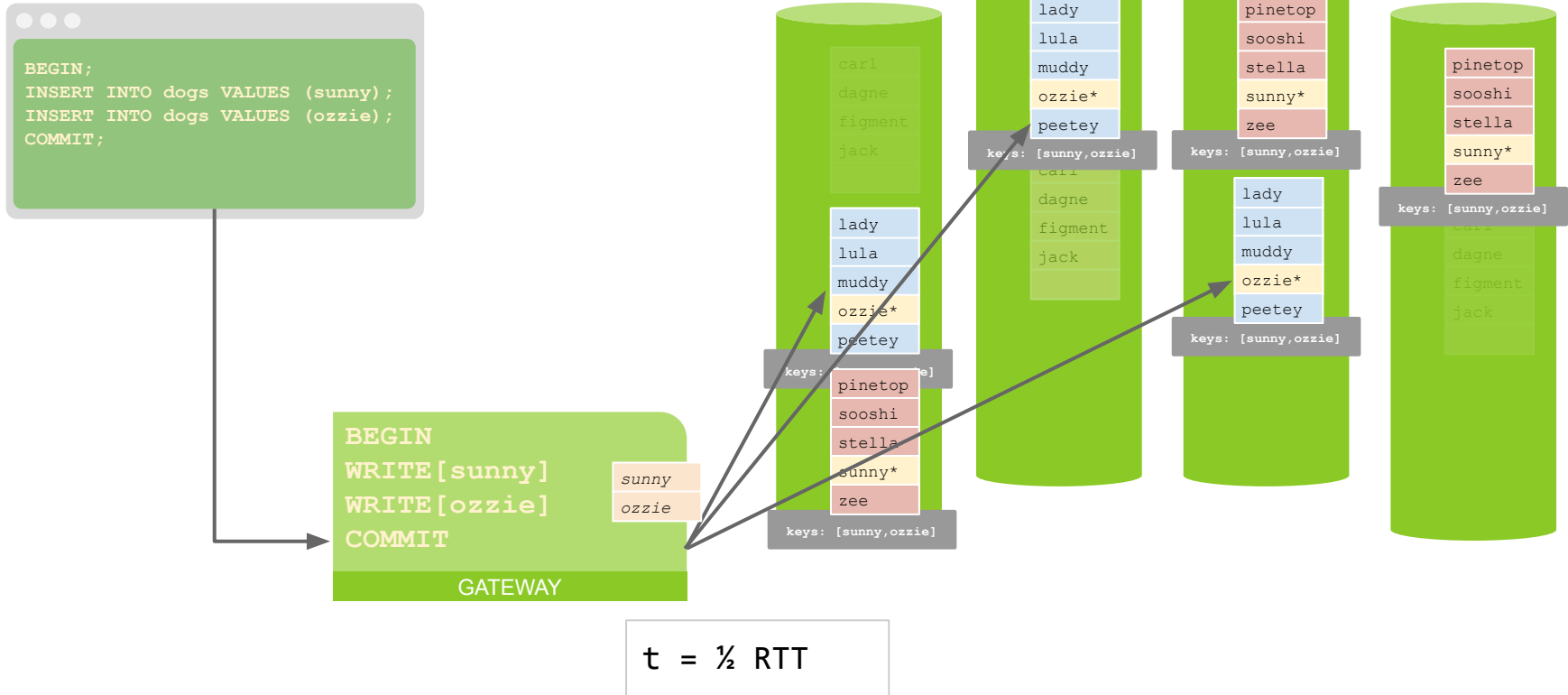


t = 0

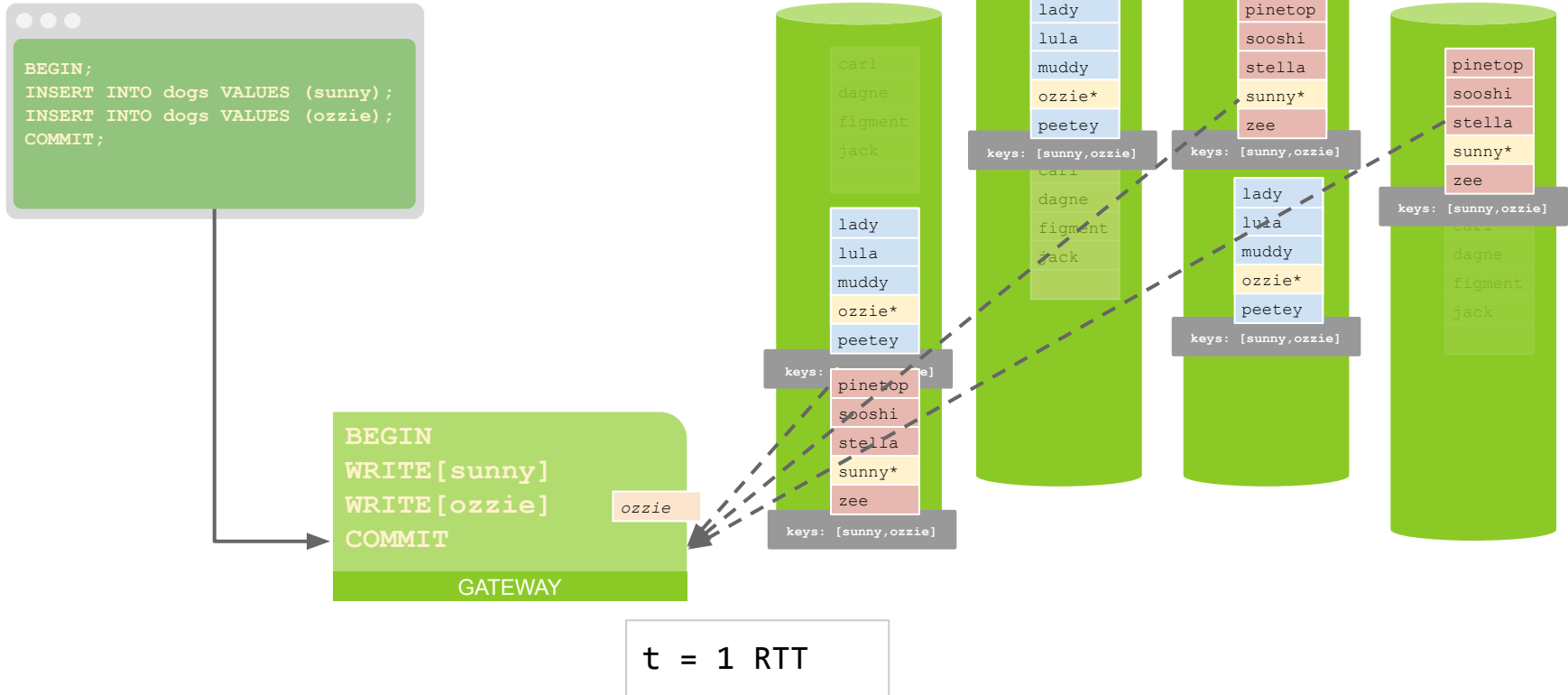
# MDCC



# MDCC



# MDCC



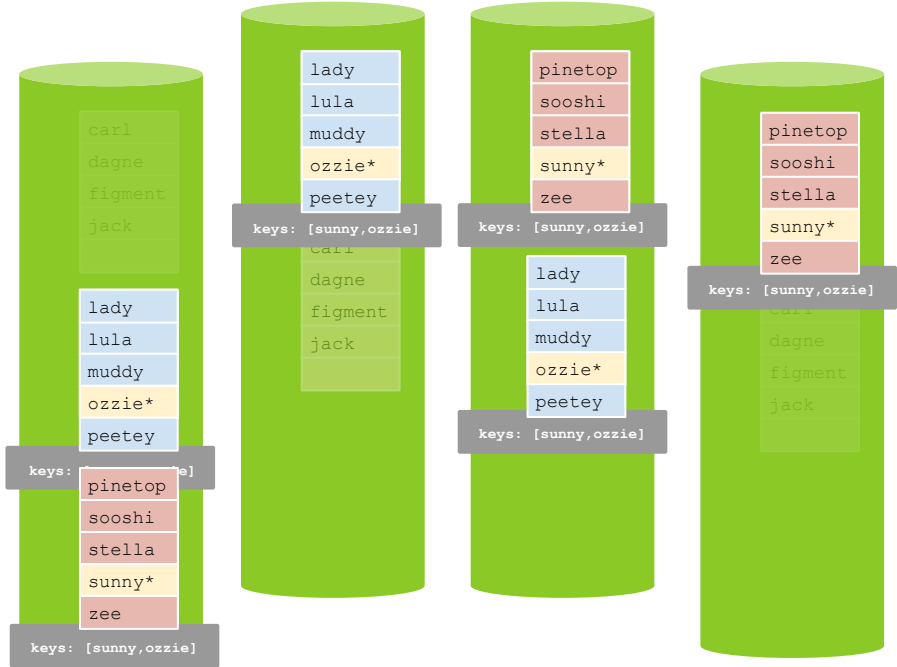


# MDCC



```
BEGIN;  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

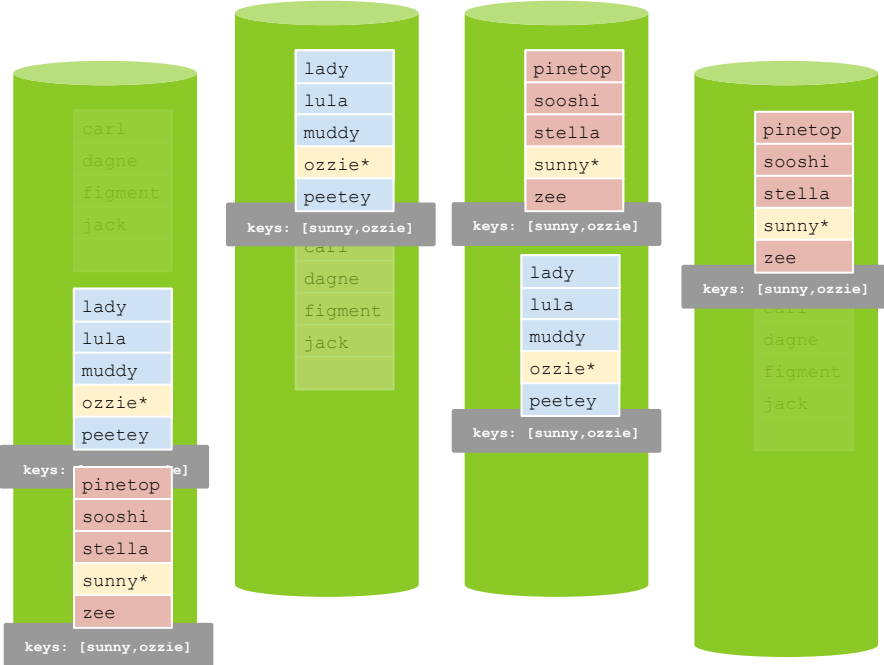
```
BEGIN  
WRITE [sunny]  
WRITE [ozzie]  
COMMIT  
GATEWAY
```



t = 1 RTT

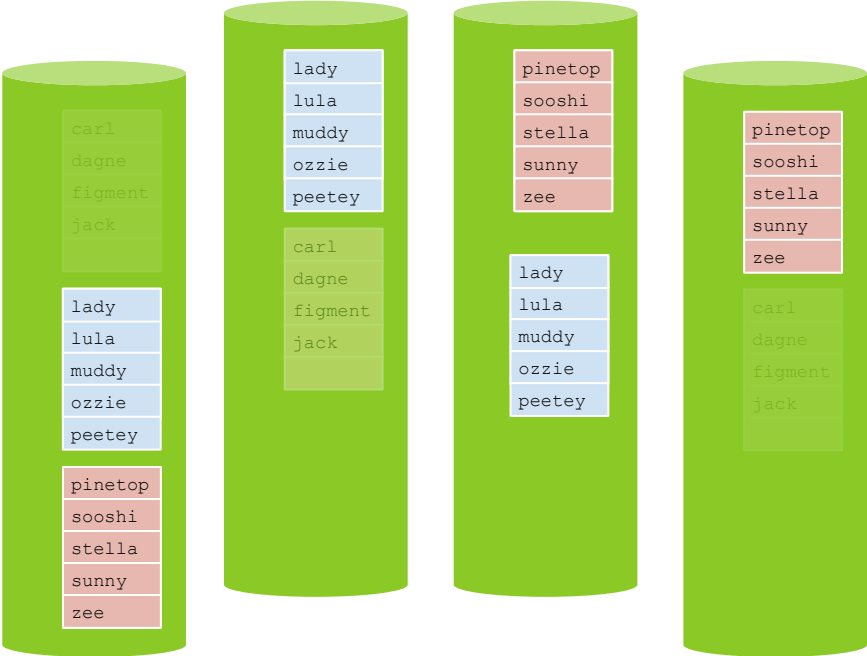


# MDCC



t = 1 RTT

# MDCC



t = 1 RTT



# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

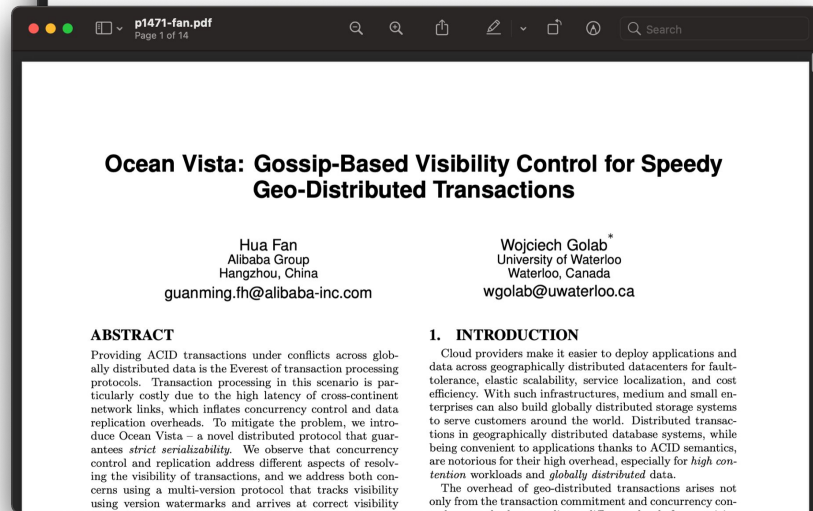
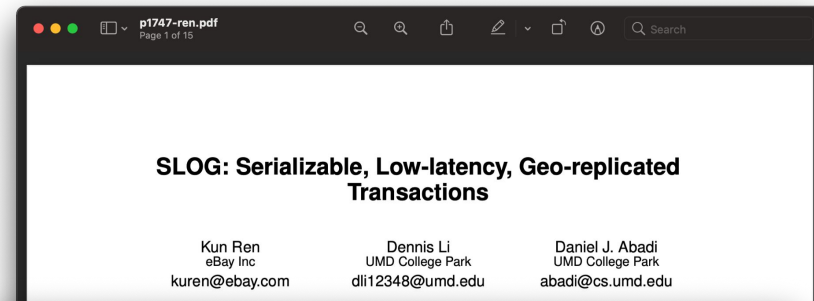
- I. Spanner/Pipelined Transactions
- II. Parallel Commits
- III. Replicated Commit
- IV. Carousel
- V. MDCC
- VI. SLOG/OceanVista
- VII. TAPIR

# SLOG/OceanVista

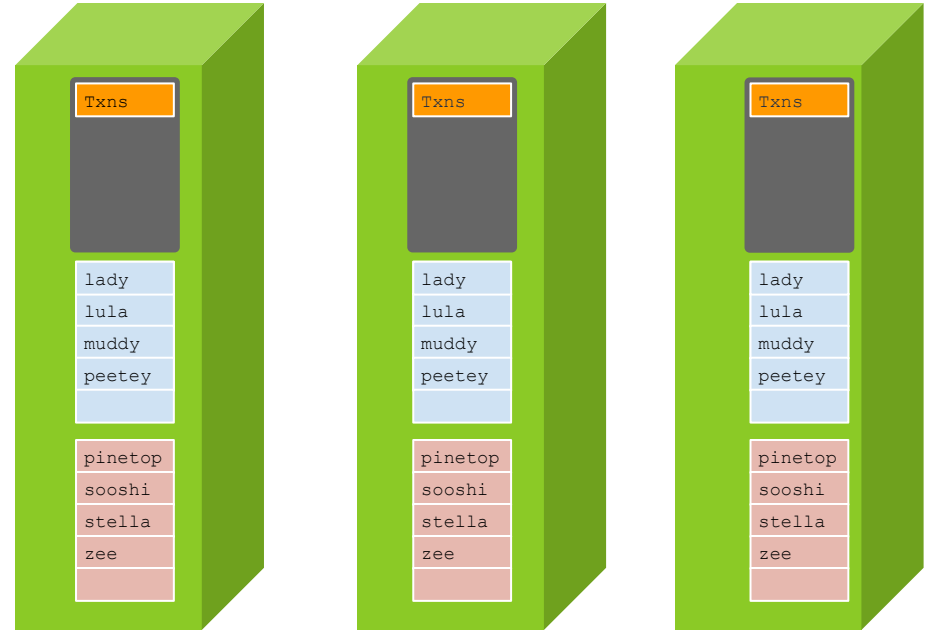


2019, from University of Maryland. Fully deterministic transactions, with read and write sets pre-declared. Lets it avoid 2PC; replicates the txn code itself.

Ditto for OceanVista. 2019, from University of Waterloo. Though SLOG provides local latencies for single-region transactions.



# OceanVista



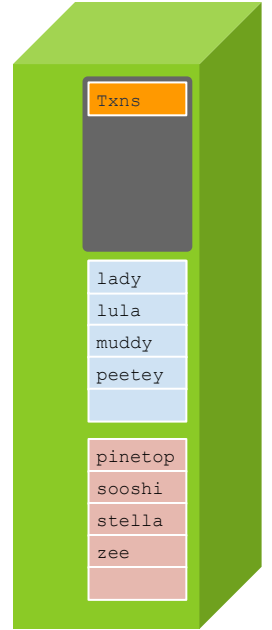
t = 0

# OceanVista



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



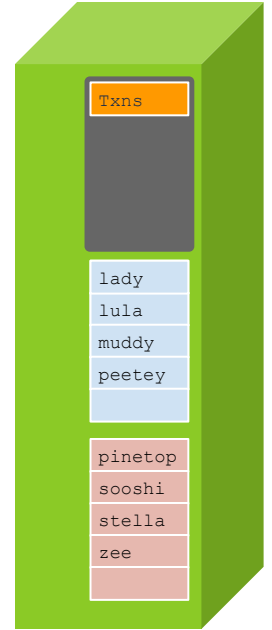
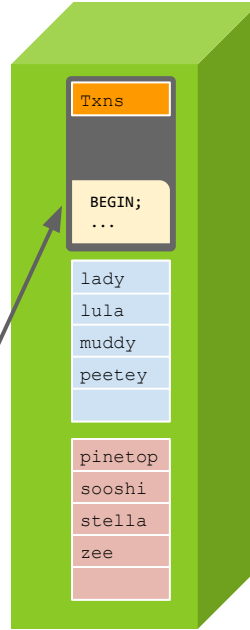
t = 0

# OceanVista



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



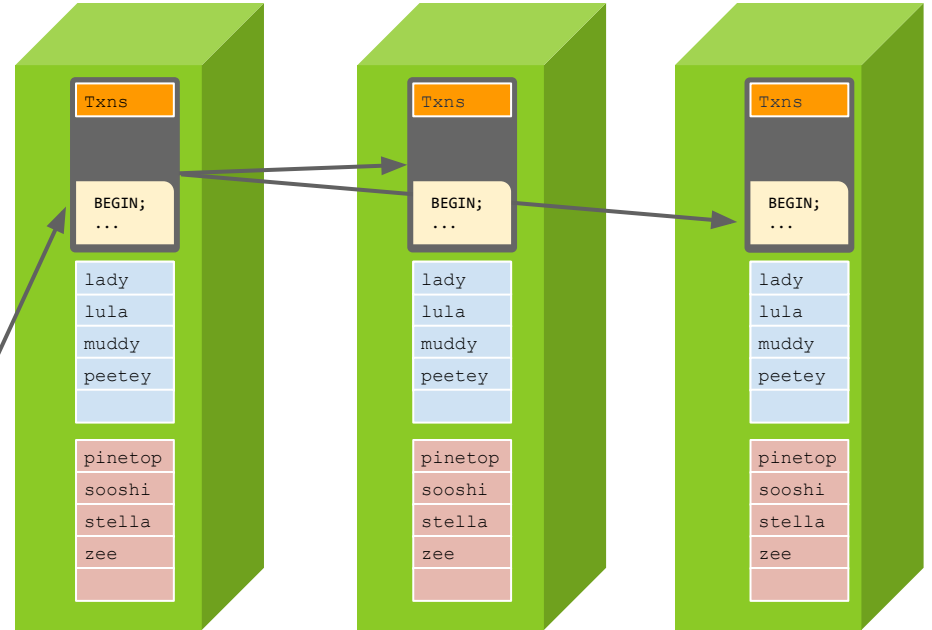
t = 0

# OceanVista



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



$$t = \frac{1}{2} \text{WAN}$$

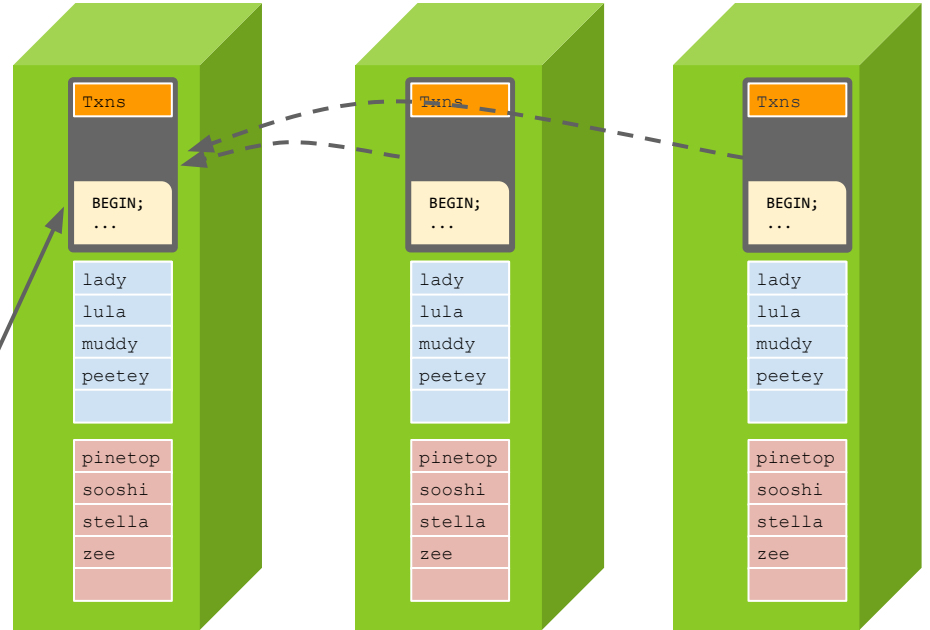


# OceanVista



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



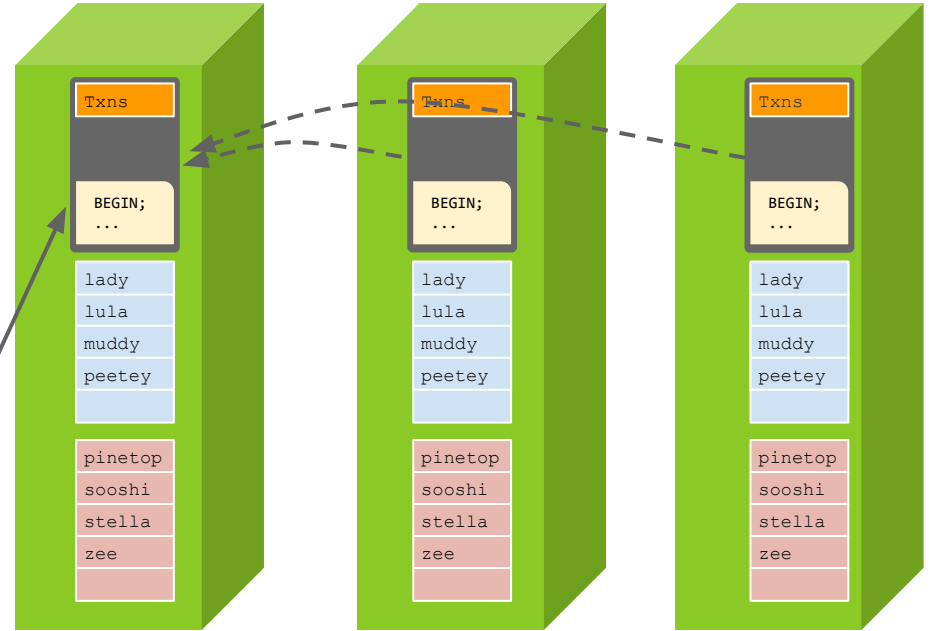
t = 1 WAN

# OceanVista



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



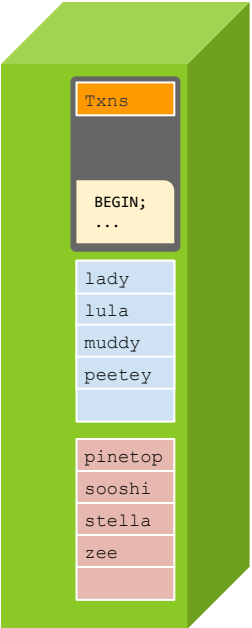
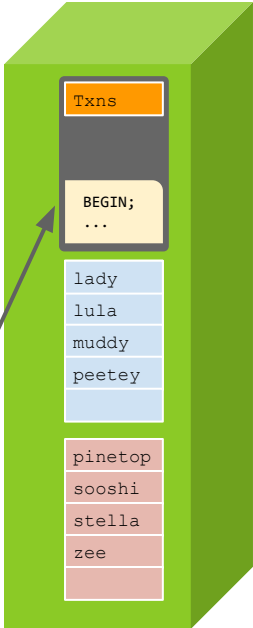
$t = 1$  WAN

# OceanVista



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



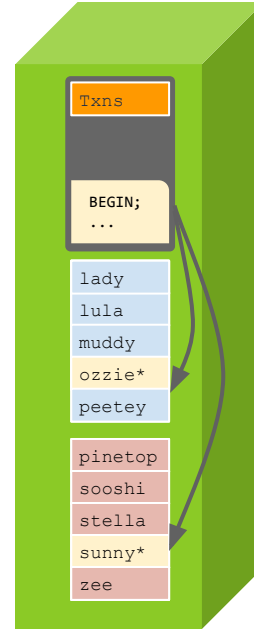
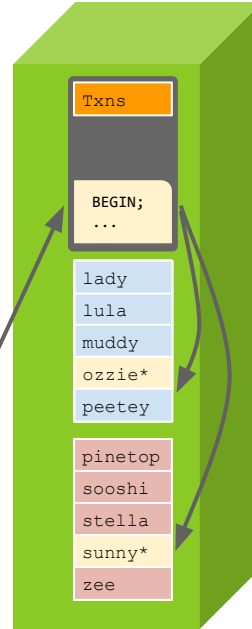
t = 1 WAN

# OceanVista



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



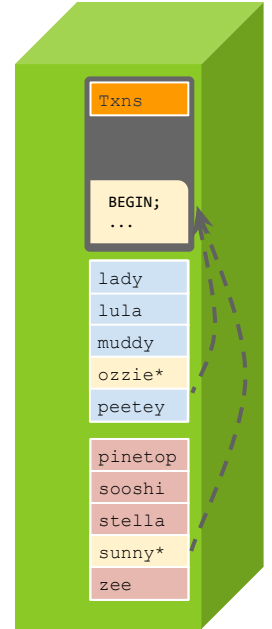
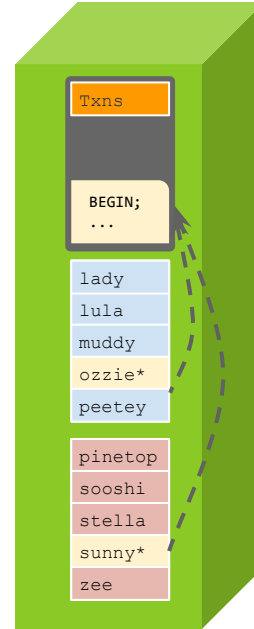
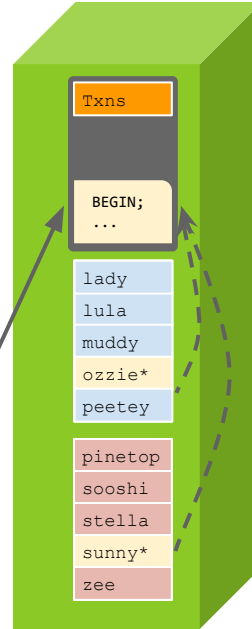
$$t = 1 \text{ WAN} + \frac{1}{2} \text{ LAN}$$

# OceanVista



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



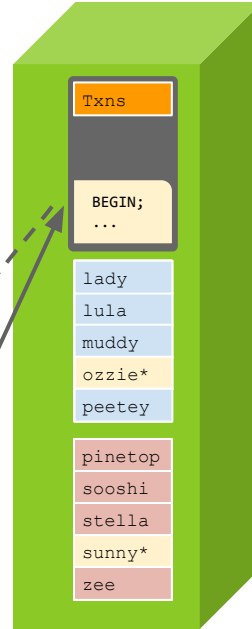
$$t = 1 \text{ WAN} + 1 \text{ LAN}$$

# OceanVista



```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



$$t = 1 \text{ WAN} + 1 \text{ LAN}$$

# OceanVista



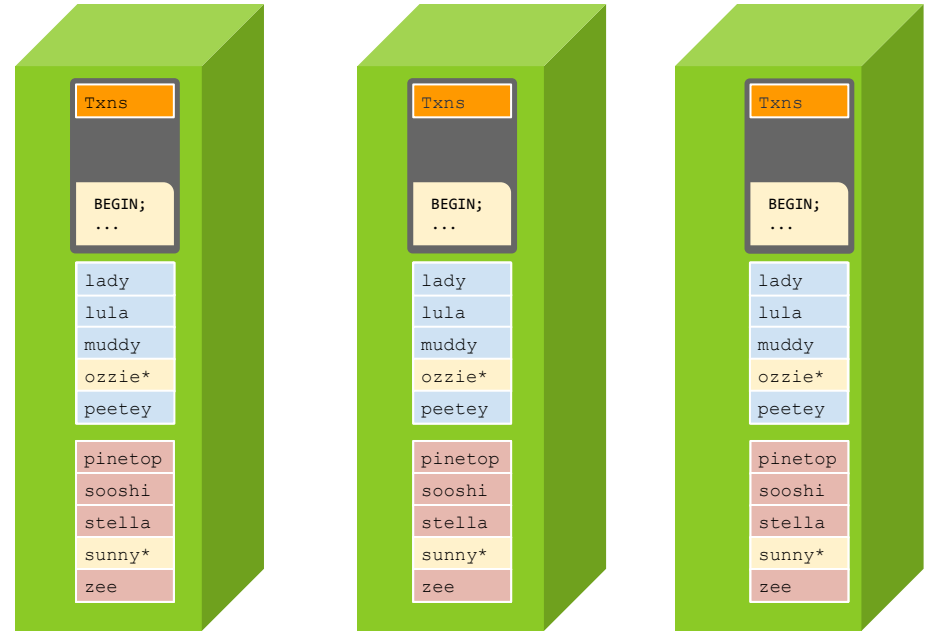
```
BEGIN FOR (sunny, ozzie);  
INSERT INTO dogs VALUES (sunny);  
INSERT INTO dogs VALUES (ozzie);  
COMMIT;
```

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



$$t = 1 \text{ WAN} + 1 \text{ LAN}$$

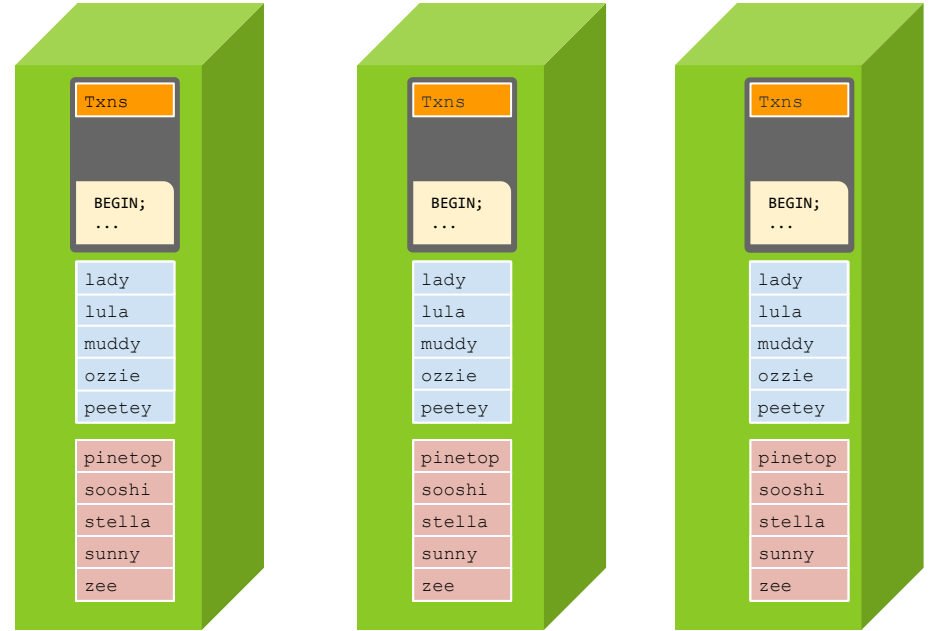
# OceanVista



$$t = 1 \text{ WAN} + 1 \text{ LAN}$$

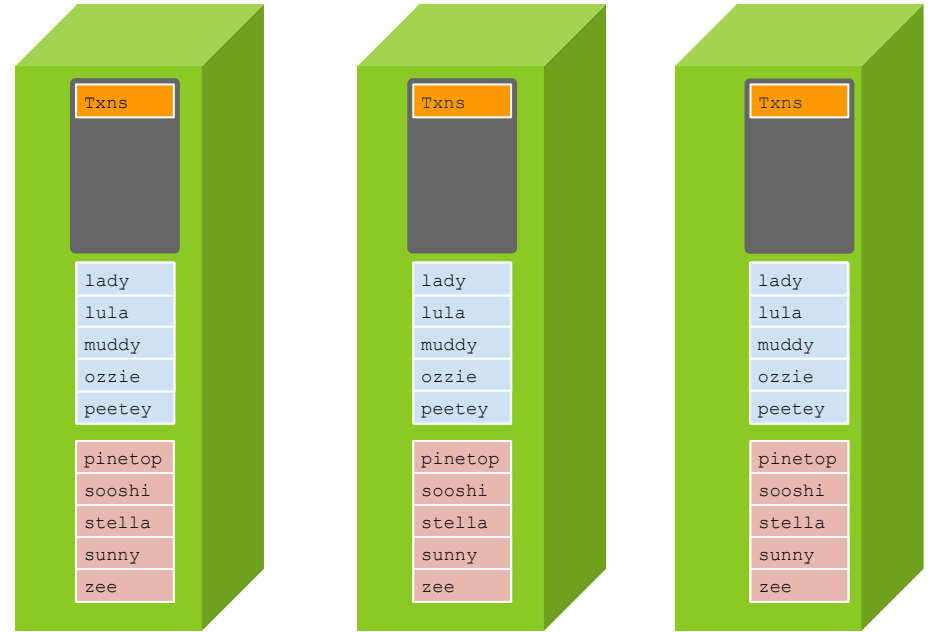


# OceanVista



$$t = 1 \text{ WAN} + 1 \text{ LAN}$$

# OceanVista



$$t = 1 \text{ WAN} + 1 \text{ LAN}$$



# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

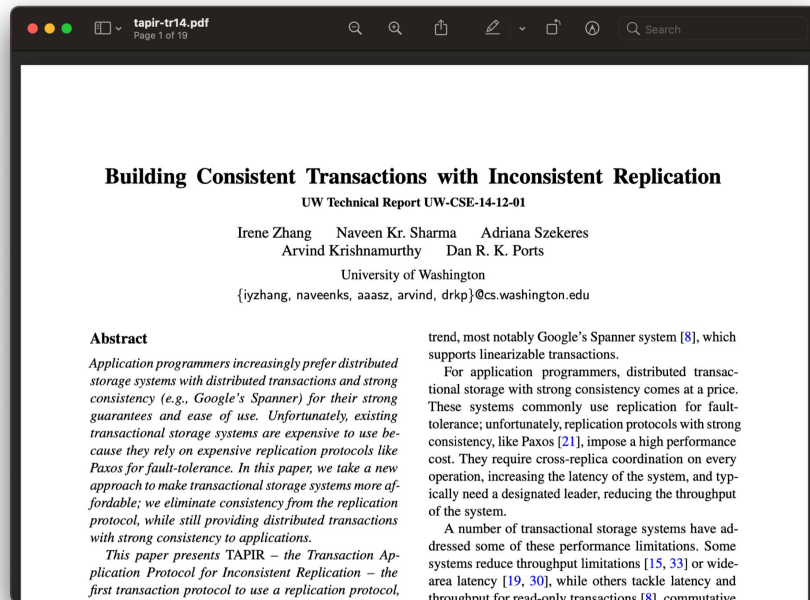
- I. Spanner/Pipelined Transactions
- II. Parallel Commits
- III. Replicated Commit
- IV. Carousel
- V. MDCC
- VI. SLOG/OceanVista
- VII. TAPIR

# TAPIR



2014, from University of Washington. Observes that 2PC and Consensus-based replication both order operations (each contributing 1RTT). Proposes unordered replication, and builds transaction protocol on top of it.

Primitive: replicas all contain a set of operations, but in no particular order. Two sets are guaranteed to overlap.





# Agenda

1. Foundations
2. Transactions
- 3. Implementations**

- I. Spanner/Pipelined Transactions **2 WAN RTT**
- II. Parallel Commits **1 WAN RTT**
- III. Replicated Commit **1 WAN + LAN RTT**
- IV. Carousel **1 WAN RTT**
- V. MDCC **1 WAN RTT**
- VI. SLOG/OceanVista **1 WAN + LAN RTT**
- VII. TAPIR **1 WAN RTT**



**Questions?**

`irfansharif.io`

`@irfansharifm`