# AWS re:Invent

**DECEMBER 2 – 6, 2024 | LAS VEGAS, NV**

DAT427-NEW

# Deep dive into Amazon Aurora DSQL and its architecture

**Marc Brooker**

VP/Distinguished Engineer
Amazon Web Services

# You don't need to know this stuff!

# Agenda

- Quick reminder: What is Amazon Aurora DSQL?

- Writes and concurrency control

- Reads and SQL execution

- Cross-Region and scalability

# Amazon Aurora DSQL is…

A relational SQL database
optimized for transactional workloads.

aws

Scalable,
up and down.

Serverless.

# Active-active, and multi-Region.

# PostgreSQL compatible.

# Built on our experience.

# Rethinking transactional databases

```
BEGIN;
INSERT INTO dogs VALUES ('snuffles', 4);
INSERT INTO dogs VALUES ('sophie', 8);
COMMIT;
```

**A**tomic

**C**onsistent

**I**solated

**D**urable

# **A**tomic and **D**urable



Journal

Awesome internal primitive
providing atomicity and durability.

# **A**tomic and **D**urable and **I**solated



Adjudicator → Journal

Looks for conflicts between this transaction and other recent transactions.

# **A**tomic and **D**urable and **I**solated and **S**calable



Adjudicator

Distributed
commit protocol

Adjudicator

Journal

```
BEGIN;
UPDATE dogs SET age = age + 1
    WHERE name IN ('snuffles', 'max');
COMMIT;
```

BEGIN;
Find the current age of max and snuffles.
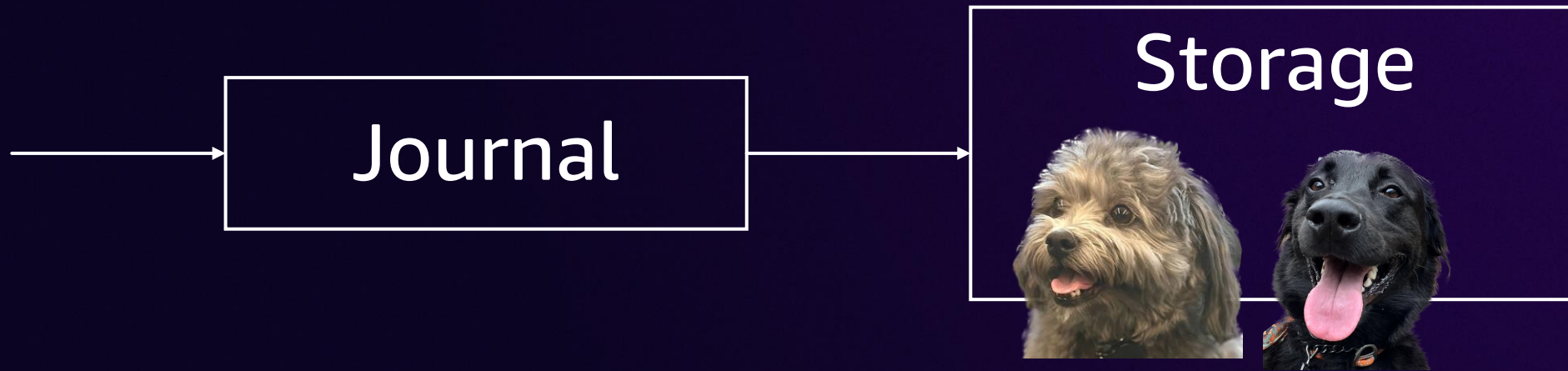Add one.
Overwrite the old values with the new ones.
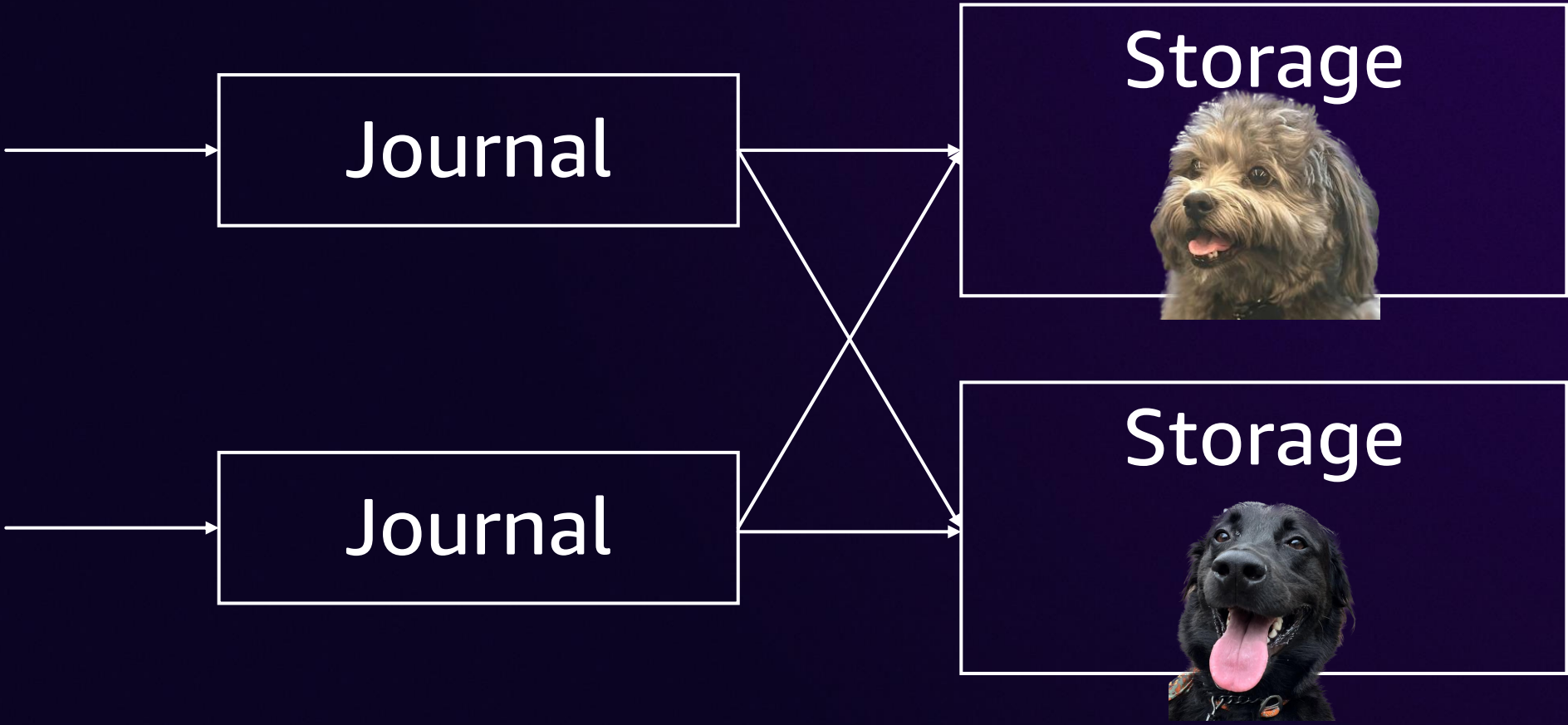COMMIT;

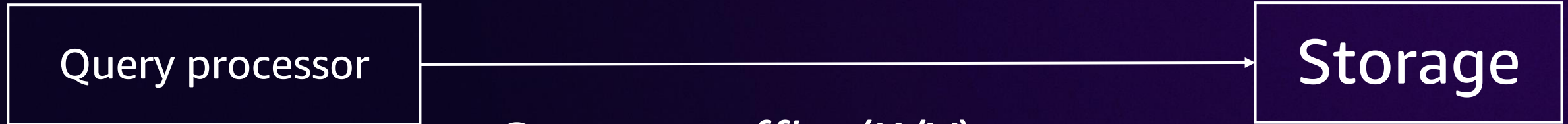but querying from a log isn't fun or efficient.

# **Q**ueryable: Pushdown
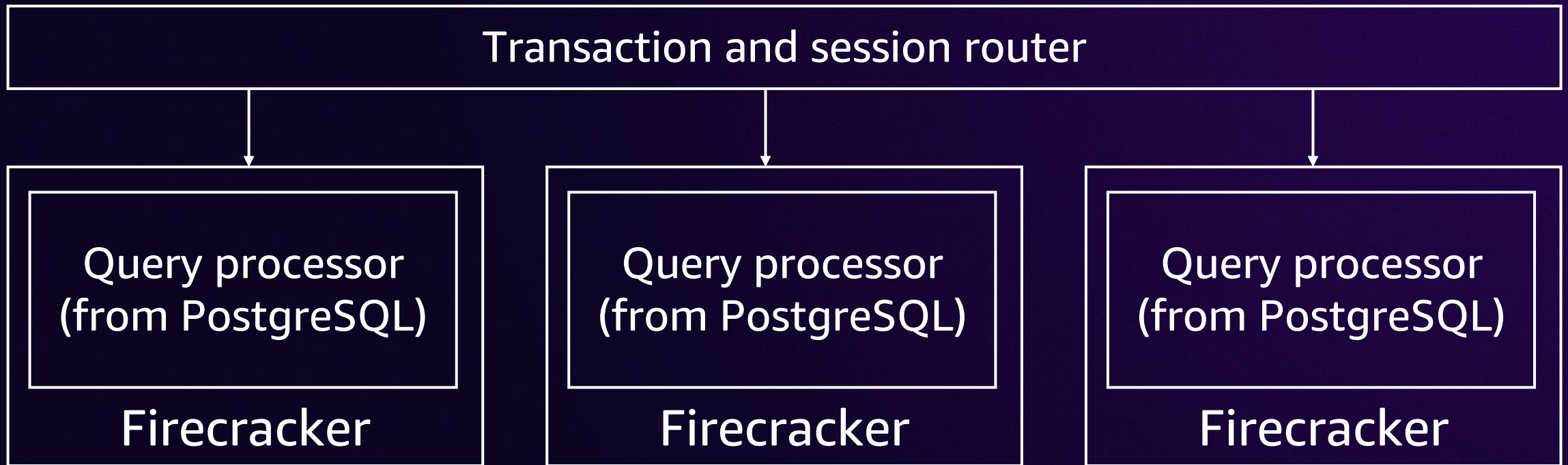
Query processor ───────────────▶ Storage

- Get me *snuffles* (K/V)
- Get me all good dogs (scan)
- Count all bad dogs (aggregate)
- Get me their ages and favorite treat (project)
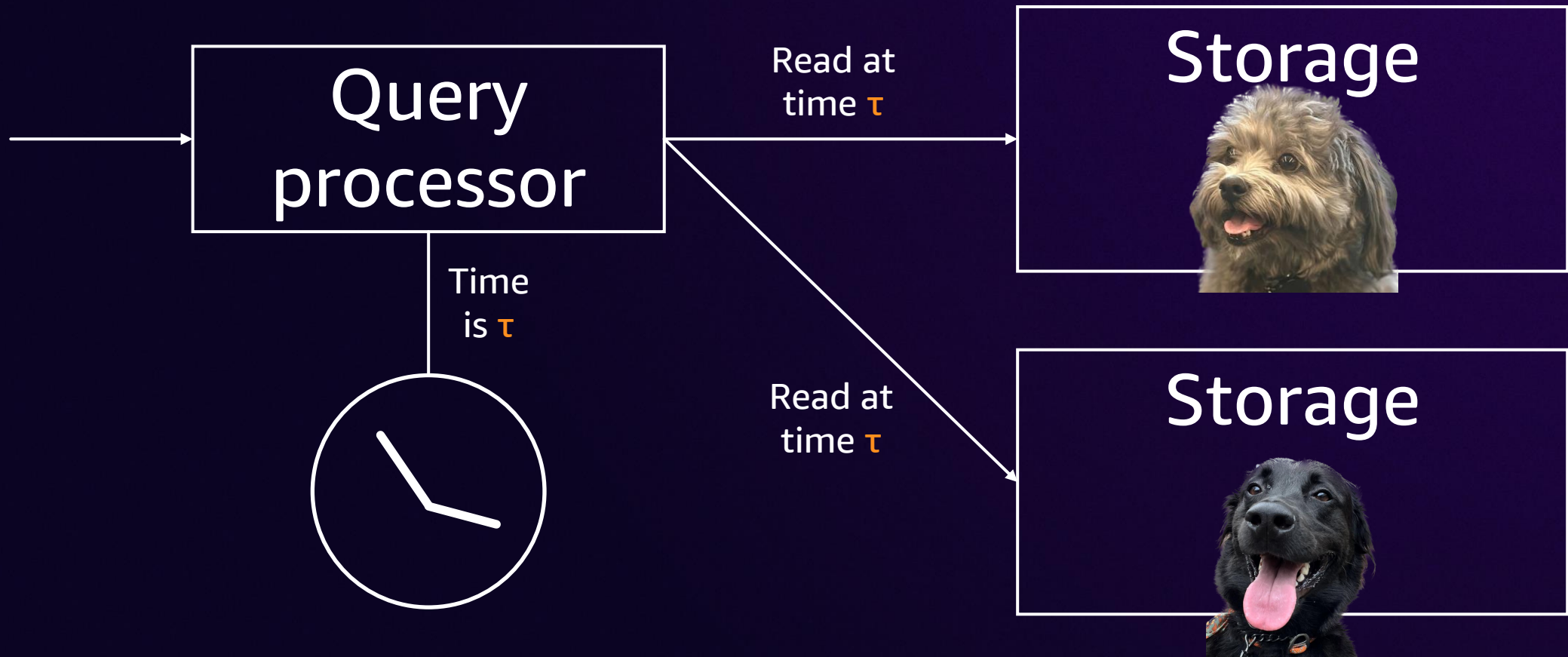- etc.

Saves a lot of round trips!

```sql
BEGIN;
SELECT count(1) FROM dogs
    WHERE state = 'hungry';
UPDATE food SET quantity = quantity – 2
    WHERE type = 3;
UPDATE dogs SET state = 'well fed'
    WHERE name IN ('fido', 'max');
COMMIT;
```

Transaction and session router

Query processor (from PostgreSQL) — Firecracker
Query processor (from PostgreSQL) — Firecracker
Query processor (from PostgreSQL) — Firecracker

Each database can have any number of these. Just keep scaling.

# Isolation of reads



Query processor

Read at time $\tau$

Time is $\tau$

Read at time $\tau$

Storage

Storage

# Isolation of reads: Multiversioning

Storage

Read at time τ = 3

t = 1        t = 2        t = 3        t = 5

Pick this one.

# Isolation of reads: The locking alternative



Storage

Read at time $\tau = 3$

t = 1          t = 2          t = 3          t = 5

I read this one, so you can't add any new Snuffles till I'm done.
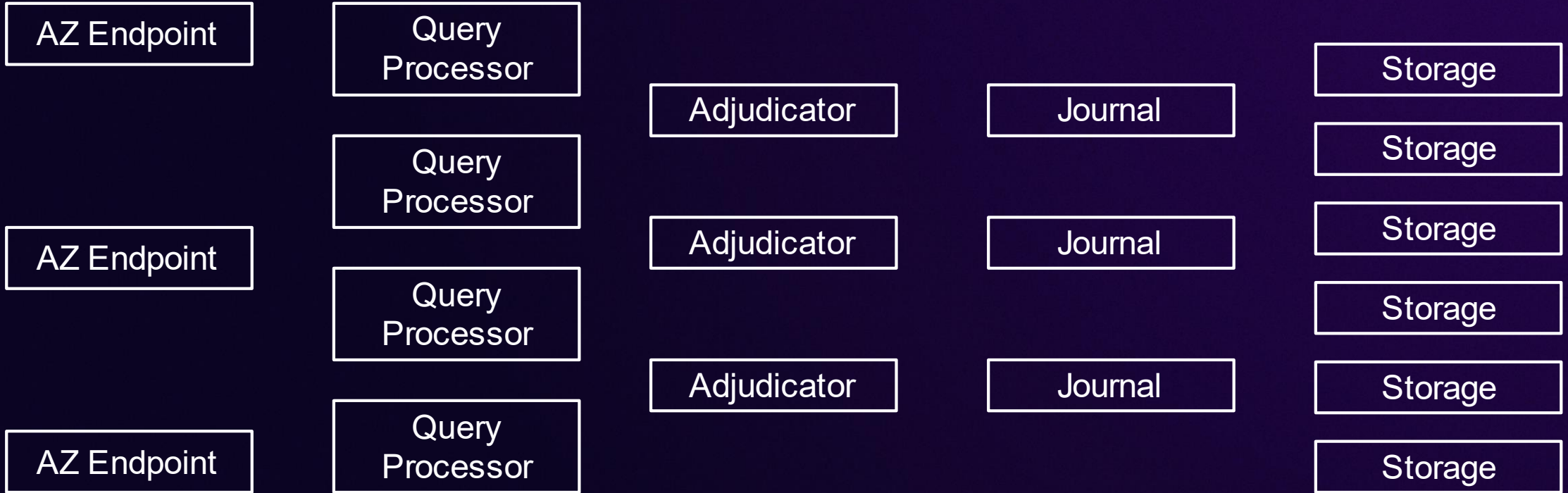
Transaction and session router

Query processor

Adjudicator

Journal

Storage

Each layer scales:

- Horizontally
- Independently
- Dynamically

Based on the demands of your workload.

# Deep dive: Isolation

BEGIN;        Query processor chooses $\tau_{start}$

SELECT ...   Reads done at $\tau_{start}$

INSERT ...   Writes spooled in QP

UPDATE ...  Reads at $\tau_{start}$, writes in QP

COMMIT;   Check isolation rules.

# No coordination needed before COMMIT!

# Optimistic Concurrency Control (aka OCC)

No locks, no coordination before commit.

# Strong Snapshot Isolation

(equivalent to PostgreSQL's REPEATABLE READ level).

# What is Snapshot Isolation?

- Never see uncommitted data.

- Reads are repeatable.

- Reads all come from a single point in (logical) time.

- Conflicting writes are rejected (writes are not lost).

But it's not serializable.

START TRANSACTION;

SELECT n FROM t WHERE id IN (1, 2);

UPDATE t SET n = 2 WHERE id = 1;

COMMIT; ✅

START TRANSACTION;

SELECT n FROM t WHERE id IN (1, 2);

UPDATE t SET n = 2 WHERE id = 2;

COMMIT; ✅

# Snapshot Isolation recipe

- Perform all reads at $\tau_{start}$

- At commit time, choose $\tau_{commit}$

- The transaction can commit if (and only if) no other transaction has written to the same keys between $\tau_{start}$ and $\tau_{commit}$

- Perform the writes at $\tau_{commit}$

QP ──────────────────────────────────────────→ Adjudicator

Here are the keys I intend to write, and my $\tau_{start}$

If no other transaction has written these keys since $\tau_{start}$, pick a $\tau_{commit}$ and write these changes to the Journal.

Never allow another conflicting transaction to pick a lower $\tau_{commit}$.

# Snapshot Isolation is a sweet spot.

# Deep dive: Cross-Region

# Optimize for round trips.

Data travels at 200km per ms,
or 123 miles per ms.

Can happen entirely locally!

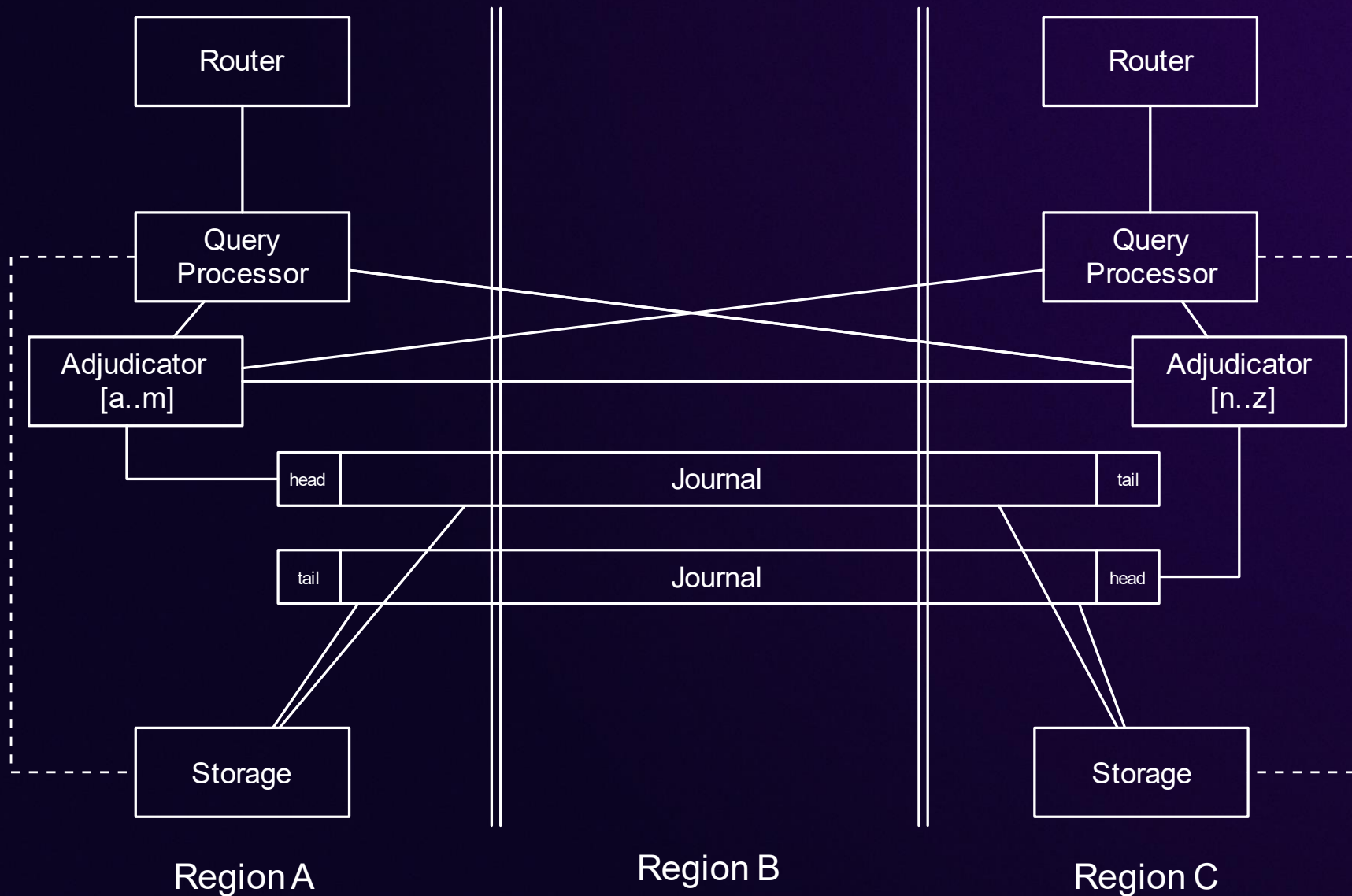| | |
|---|---|
| BEGIN; | Query processor chooses $\tau_{start}$ |
| SELECT ... | Reads done at $\tau_{start}$ |
| INSERT ... | Writes spooled in QP |
| UPDATE ... | Reads at $\tau_{start}$, writes in QP |
| COMMIT; | Check isolation rules. |
| | Make data durable. |

Needs cross-Region coordination.
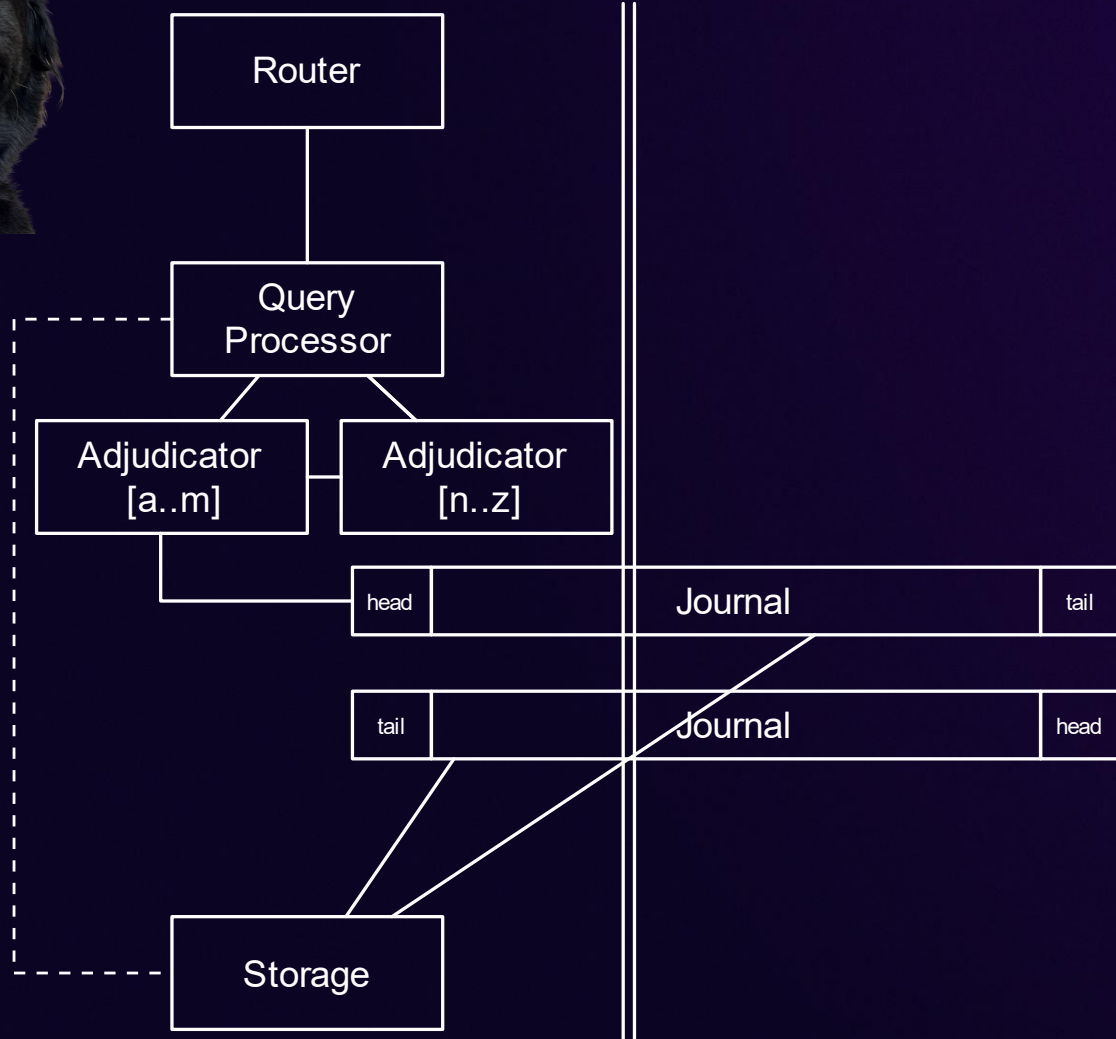
# Coordinate once, at commit time.

# Coordinate once, at commit time.

Read-only transactions don't need to coordinate at all!
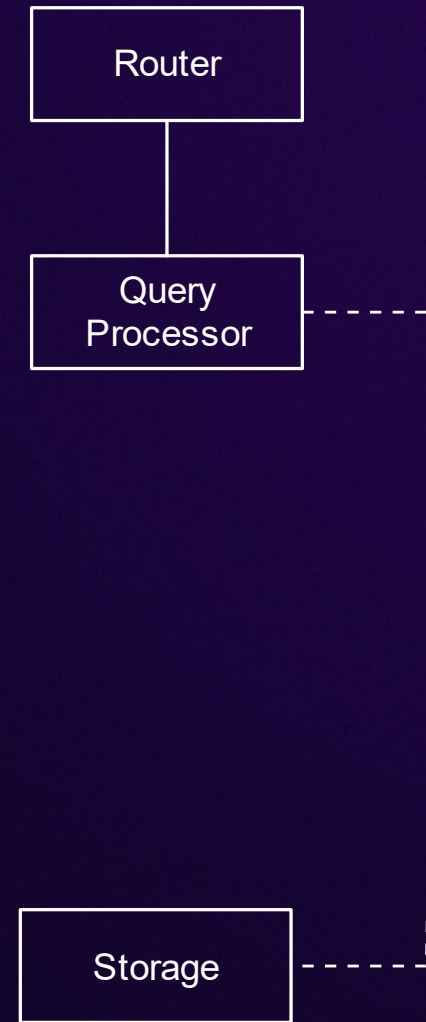
# Optimize for fast failover.

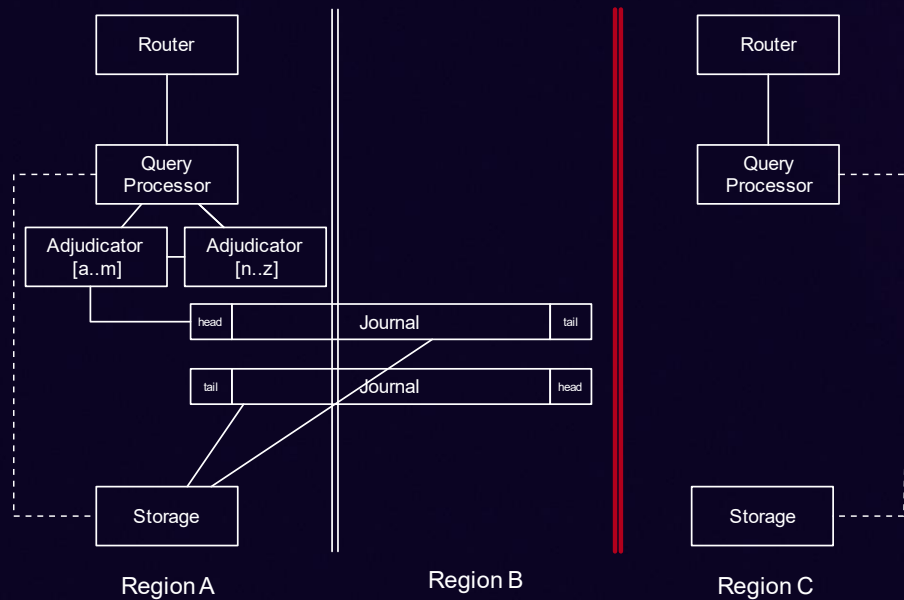Changing leadership doesn't require moving data (or lock state).

Router

Query
Processor

Adjudicator
[a..m]

Adjudicator
[n..z]

| head | Journal | tail |

| tail | Journal | head |

Storage

Region A

Router

Query
Processor

Storage

Region B

Region C

- Read path: no change
- Adjudicators:
  move into healthy regions
- Journals:
  majority already in healthy regions
- No data loss, no availability loss

# Implementation quality

# Rust.

# Deterministic simulation testing.

# Fuzzing.

# Formal methods.

# Runtime monitoring.

# Thank you!

**Marc Brooker**

mbrooker@amazon.com

@marcjbrooker


Please complete the session survey in the mobile app