Sharding e replicação com Citus

# Relational Model

# Document Store

```
{
    first name: 'Paul',
    surname: 'Miller',
    cell: 447557505611,
    city: 'London',
    location: [45.123,47.232],
    Profession: ['banking', 'finance', 'trader'],
    cars: [
        { model: 'Bentley',
          year: 1973,
          value: 100000, … },
        { model: 'Rolls Royce',
          year: 1965,
          value: 330000, … }
    ]
}
```

String — Number — Geo-Coordinates → Typed field values

Fields

Fields can contain arrays

Fields can contain an array of sub-documents

# "RDBMS não escalam"

➔ Single-server

➔ Escala apenas verticalmente

➔ Flexibilidade X Performance

➔ Disponibilidade

# Scaling a Relational Database



Expensive single servers

Time to buy and configure hardware

**X** NO ELASTICITY

MBs

GBs

# NoSQL cresceu!



DB-Engines Ranking

© August 2018, DB-Engines.com

Legend:
- Oracle
- MySQL
- Microsoft SQL Server
- PostgreSQL
- MongoDB
- DB2
- Redis
- Elasticsearch
- Microsoft Access
- Cassandra
- SQLite
- Teradata
- Splunk
- MariaDB
- Solr
- SAP Adaptive Server
- HBase
- Hive
- FileMaker
- SAP HANA
- Amazon DynamoDB
- Neo4j
- Couchbase
- Memcached
- Microsoft Azure SQL Database
- Informix
- Firebird
- Vertica
- Microsoft Azure Cosmos DB
- CouchDB
- Netezza
- Amazon Redshift
- Google BigQuery
- Impala
- Spark SQL

1/11

# ACID

```
1    START TRANSACTION;
2    SELECT balance FROM checking WHERE customer_id = 10233276;
3    UPDATE checking SET balance = balance - 200.00 WHERE customer_id = 10233276;
4    UPDATE savings  SET balance = balance + 200.00 WHERE customer_id = 10233276;
5    COMMIT;
```

# PostgreSQL: Melhorias

→Particionamento nativo

→FDW

→FTS para JSON e JSONB

→Queries paralelas

→Quorum Commit

# Como aliar PostgreSQL à necessidade de escalar horizontalmente?

# Citus é uma extensão!
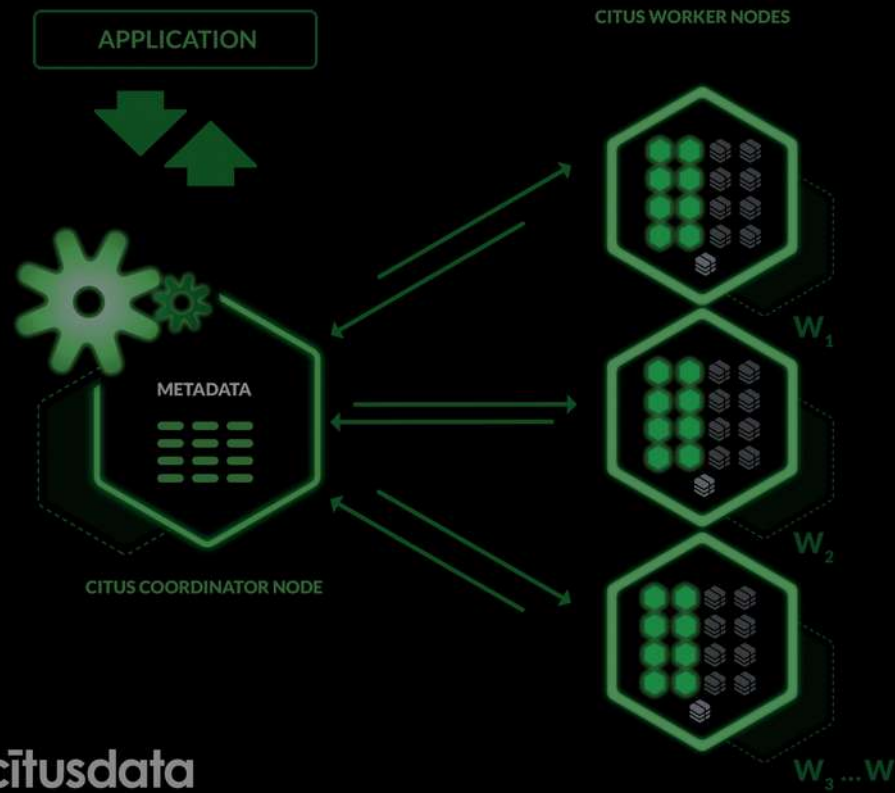


Distributed Postgres Database from Citus Data

# Problemas

→ Multi-tenancy

→ Real-time analytics

# Multi-tenancy

# Ad Analytics

```sql
CREATE TABLE companies (
  id bigserial PRIMARY KEY,
  name text NOT NULL,
  image_url text,
  created_at timestamp without time zone NOT NULL,
  updated_at timestamp without time zone NOT NULL
);

CREATE TABLE campaigns (
  id bigserial PRIMARY KEY,
  company_id bigint REFERENCES companies (id),
  name text NOT NULL,
  cost_model text NOT NULL,
  state text NOT NULL,
  monthly_budget bigint,
  blacklisted_site_urls text[],
  created_at timestamp without time zone NOT NULL,
  updated_at timestamp without time zone NOT NULL
);
```

# Ad Analytics

```sql
CREATE TABLE ads (
  id bigserial PRIMARY KEY,
  campaign_id bigint REFERENCES campaigns (id),
  name text NOT NULL,
  image_url text,
  target_url text,
  impressions_count bigint DEFAULT 0,
  clicks_count bigint DEFAULT 0,
  created_at timestamp without time zone NOT NULL,
  updated_at timestamp without time zone NOT NULL
);

CREATE TABLE clicks (
  id bigserial PRIMARY KEY,
  ad_id bigint REFERENCES ads (id),
  clicked_at timestamp without time zone NOT NULL,
  site_url text NOT NULL,
  cost_per_click_usd numeric(20,10),
  user_ip inet NOT NULL,
  user_data jsonb NOT NULL
);
```

# Ad Analytics

```sql
CREATE TABLE impressions (
  id bigserial PRIMARY KEY,
  ad_id bigint REFERENCES ads (id),
  seen_at timestamp without time zone NOT NULL,
  site_url text NOT NULL,
  cost_per_impression_usd numeric(20,10),
  user_ip inet NOT NULL,
  user_data jsonb NOT NULL
);
```

# Ad Analytics

# Ad Analytics

# Agrupe os dados!

```sql
CREATE TABLE companies (
  id bigserial PRIMARY KEY,
  name text NOT NULL,
  image_url text,
  created_at timestamp without time zone NOT NULL,
  updated_at timestamp without time zone NOT NULL
);

CREATE TABLE campaigns (
  id bigserial,           -- was: PRIMARY KEY
  company_id bigint REFERENCES companies (id),
  name text NOT NULL,
  cost_model text NOT NULL,
  state text NOT NULL,
  monthly_budget bigint,
  blacklisted_site_urls text[],
  created_at timestamp without time zone NOT NULL,
  updated_at timestamp without time zone NOT NULL,
  PRIMARY KEY (company_id, id) -- added
);
```
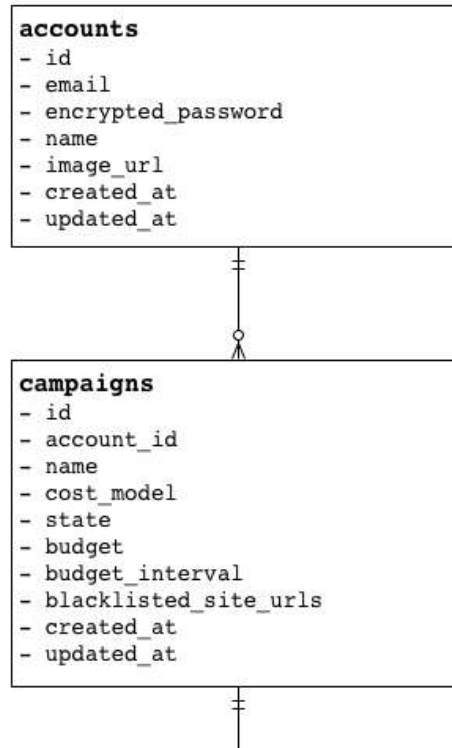
```sql
CREATE TABLE ads (
  id bigserial,         -- was: PRIMARY KEY
  company_id bigint,  -- added
  campaign_id bigint, -- was: REFERENCES campaigns (id)
  name text NOT NULL,
  image_url text,
  target_url text,
  impressions_count bigint DEFAULT 0,
  clicks_count bigint DEFAULT 0,
  created_at timestamp without time zone NOT NULL,
  updated_at timestamp without time zone NOT NULL,
  PRIMARY KEY (company_id, id),            -- added
  FOREIGN KEY (company_id, campaign_id) -- added
    REFERENCES campaigns (company_id, id)
);
```

```sql
CREATE TABLE clicks (
  id bigserial,            -- was: PRIMARY KEY
  company_id bigint,       -- added
  ad_id bigint,            -- was: REFERENCES ads (id),
  clicked_at timestamp without time zone NOT NULL,
  site_url text NOT NULL,
  cost_per_click_usd numeric(20,10),
  user_ip inet NOT NULL,
  user_data jsonb NOT NULL,
  PRIMARY KEY (company_id, id),        -- added
  FOREIGN KEY (company_id, ad_id)      -- added
    REFERENCES ads (company_id, id)
);

CREATE TABLE impressions (
  id bigserial,            -- was: PRIMARY KEY
  company_id bigint,       -- added
  ad_id bigint,            -- was: REFERENCES ads (id),
  seen_at timestamp without time zone NOT NULL,
  site_url text NOT NULL,
  cost_per_impression_usd numeric(20,10),
  user_ip inet NOT NULL,
  user_data jsonb NOT NULL,
  PRIMARY KEY (company_id, id),        -- added
  FOREIGN KEY (company_id, ad_id)      -- added
    REFERENCES ads (company_id, id)
);
```
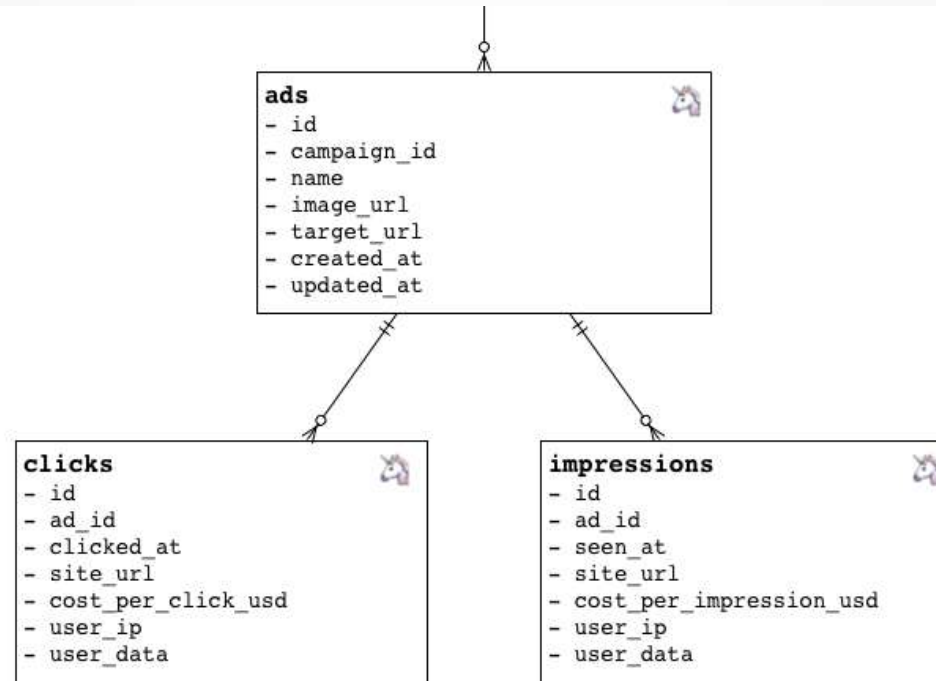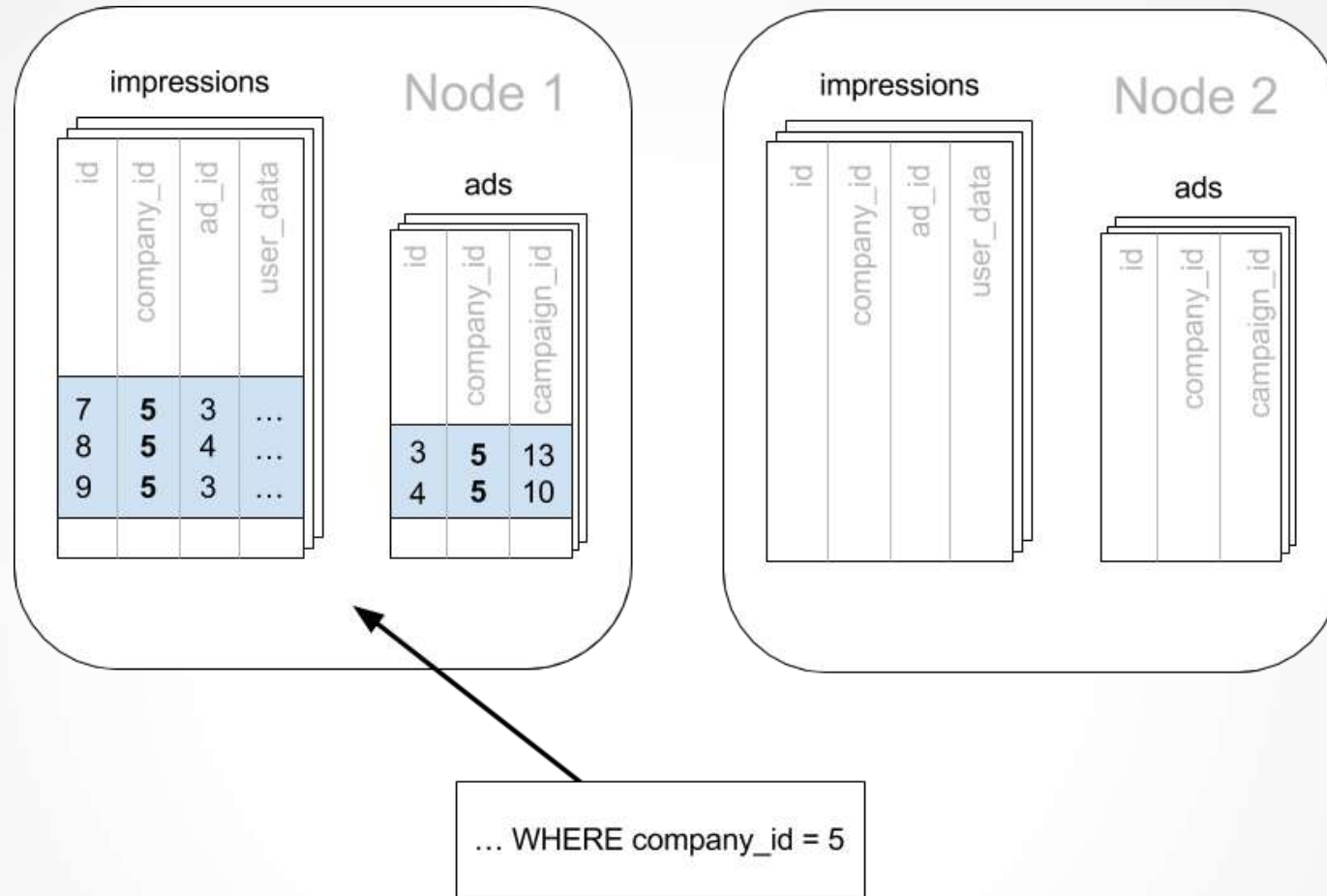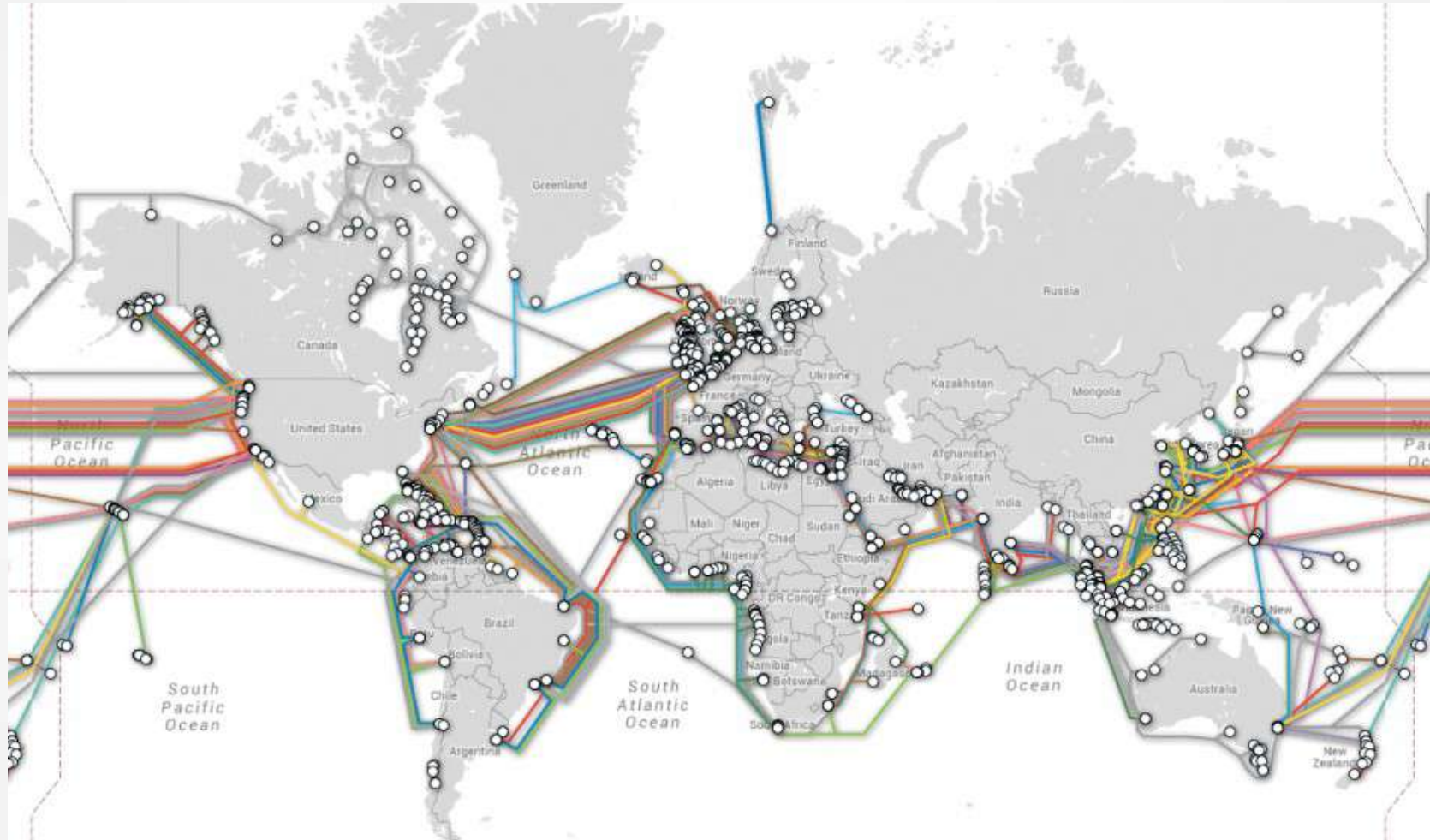
# Distributed tables

```sql
SELECT create_distributed_table('companies',   'id');
SELECT create_distributed_table('campaigns',   'company_id');
SELECT create_distributed_table('ads',         'company_id');
SELECT create_distributed_table('clicks',      'company_id');
SELECT create_distributed_table('impressions', 'company_id');
```

# Reference tables

```sql
CREATE TABLE geo_ips (
  addrs cidr NOT NULL PRIMARY KEY,
  latlon point NOT NULL
    CHECK (-90  <= latlon[0] AND latlon[0] <= 90 AND
           -180 <= latlon[1] AND latlon[1] <= 180)
);
CREATE INDEX ON geo_ips USING gist (addrs inet_ops);
```

```sql
SELECT create_reference_table('geo_ips');
```

# DDLs

```
ALTER TABLE ads
   ADD COLUMN caption text;
```

# Real-time analytics

```
$ http PUT httpbin.org/put hello=world
HTTP/1.1 200 OK
Access-Control-Allow-Credentials: true
Access-Control-Allow-Origin: *
Connection: keep-alive
Content-Length: 434
Content-Type: application/json
Date: Sun, 08 Feb 2015 00:39:38 GMT
Server: nginx

{
    "args": {},
    "data": "{\"hello\": \"world\"}",
    "files": {},
    "form": {},
    "headers": {
        "Accept": "application/json",
        "Accept-Encoding": "gzip, deflate",
        "Content-Length": "18",
        "Content-Type": "application/json; charset=utf-8",
        "Host": "httpbin.org",
        "User-Agent": "HTTPie/0.9.1"
    },
    "json": {
        "hello": "world"
    },
    "origin": "109.81.210.175",
    "url": "http://httpbin.org/put"
}
```

# Data model

```sql
CREATE TABLE http_request (
  site_id INT,
  ingest_time TIMESTAMPTZ DEFAULT now(),

  url TEXT,
  request_country TEXT,
  ip_address TEXT,

  status_code INT,
  response_time_msec INT
);

SELECT create_distributed_table('http_request', 'site_id');
```

# Data model

```
DO $$
  BEGIN LOOP
    INSERT INTO http_request (
      site_id, ingest_time, url, request_country,
      ip_address, status_code, response_time_msec
    ) VALUES (
      trunc(random()*32), clock_timestamp(),
      concat('http://example.com/', md5(random()::text)),
      ('{China,India,USA,Indonesia}'::text[])[ceil(random()*4)],
      concat(
        trunc(random()*250 + 2), '.',
        trunc(random()*250 + 2), '.',
        trunc(random()*250 + 2), '.',
        trunc(random()*250 + 2)
      )::inet,
      ('{200,404}'::int[])[ceil(random()*2)],
      5+trunc(random()*150)
    );
    PERFORM pg_sleep(random() * 0.25);
  END LOOP;
END $$;
```
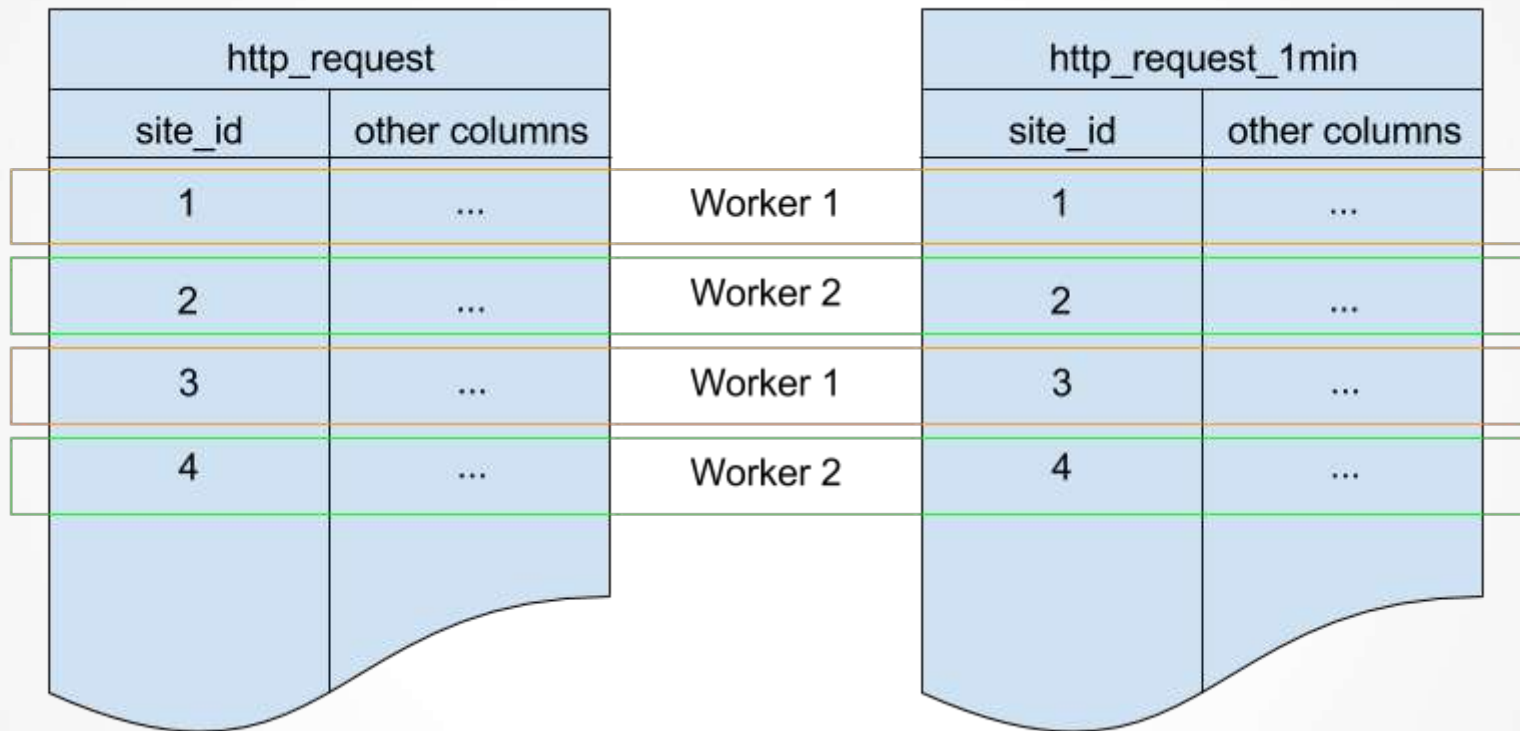
# Dashboard query

```sql
SELECT
    site_id,
    date_trunc('minute', ingest_time) as minute,
    COUNT(1) AS request_count,
    SUM(CASE WHEN (status_code between 200 and 299) THEN 1 ELSE 0 END) as success_count,
    SUM(CASE WHEN (status_code between 200 and 299) THEN 0 ELSE 1 END) as error_count,
    SUM(response_time_msec) / COUNT(1) AS average_response_time_msec
FROM http_request
WHERE date_trunc('minute', ingest_time) > now() - '5 minutes'::interval
GROUP BY site_id, minute
ORDER BY minute ASC;
```

# ROLLUPS

```sql
CREATE TABLE http_request_1min (
  site_id INT,
  ingest_time TIMESTAMPTZ, -- which minute this row represents

  error_count INT,
  success_count INT,
  request_count INT,
  average_response_time_msec INT,
  CHECK (request_count = error_count + success_count),
  CHECK (ingest_time = date_trunc('minute', ingest_time))
);

SELECT create_distributed_table('http_request_1min', 'site_id');

CREATE INDEX http_request_1min_idx ON http_request_1min (site_id, ingest_time);
```

# Co-location

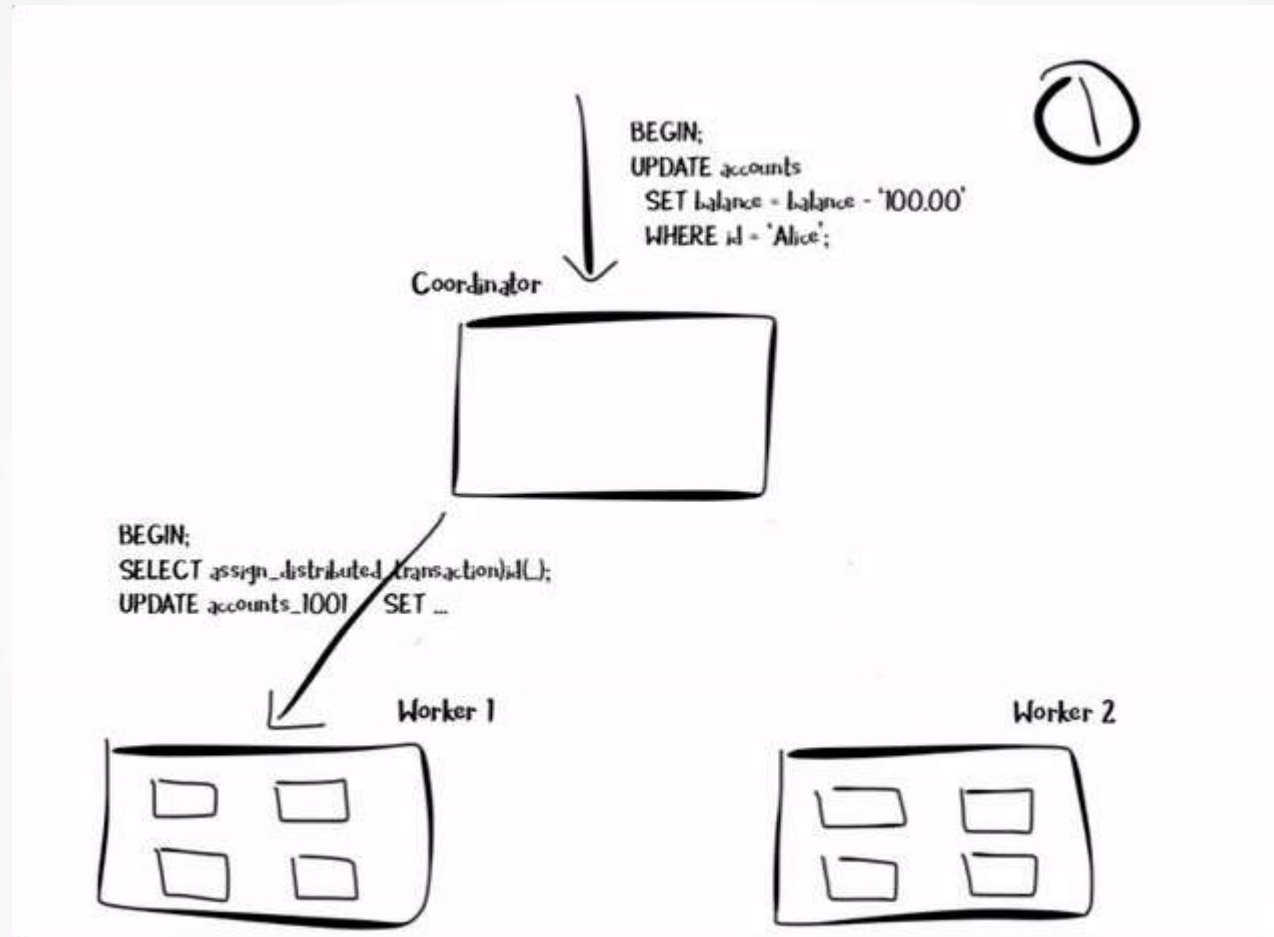| http_request | | | http_request_1min | |
|---|---|---|---|---|
| site_id | other columns | | site_id | other columns |
| 1 | ... | Worker 1 | 1 | ... |
| 2 | ... | Worker 2 | 2 | ... |
| 3 | ... | Worker 1 | 3 | ... |
| 4 | ... | Worker 2 | 4 | ... |

# INSERT INTO SELECT

```sql
-- single-row table to store when we rolled up last
CREATE TABLE latest_rollup (
  minute timestamptz PRIMARY KEY,
  -- "minute" should be no more precise than a minute
  CHECK (minute = date_trunc('minute', minute))
);
-- initialize to a time long ago
INSERT INTO latest_rollup VALUES ('10-10-1901');
-- function to do the rollup
CREATE OR REPLACE FUNCTION rollup_http_request() RETURNS void AS $$
DECLARE
  current_time     timestamptz := date_trunc('minute', now());
  last_rollup_time timestamptz := minute from latest_rollup;
BEGIN
  INSERT INTO http_request_1min (
    site_id, ingest_time, request_count,
    success_count, error_count, average_response_time_msec
  ) SELECT
    site_id,
    date_trunc('minute', ingest_time),
    COUNT(1) as request_count,
    SUM(CASE WHEN (status_code between 200 and 299) THEN 1 ELSE 0 END) as success_count,
    SUM(CASE WHEN (status_code between 200 and 299) THEN 0 ELSE 1 END) as error_count,
    SUM(response_time_msec) / COUNT(1) AS average_response_time_msec
  FROM http_request
  -- roll up only data new since last_rollup_time
  WHERE date_trunc('minute', ingest_time) <@
        tstzrange(last_rollup_time, current_time, '(]')
  GROUP BY 1, 2;
  -- update the value in latest_rollup so that next time we run the
  -- rollup it will operate on data newer than current_time
  UPDATE latest_rollup SET minute = current_time;
END;
$$ LANGUAGE plpgsql;
```

# DISTRIBUTED TRANSACTIONS
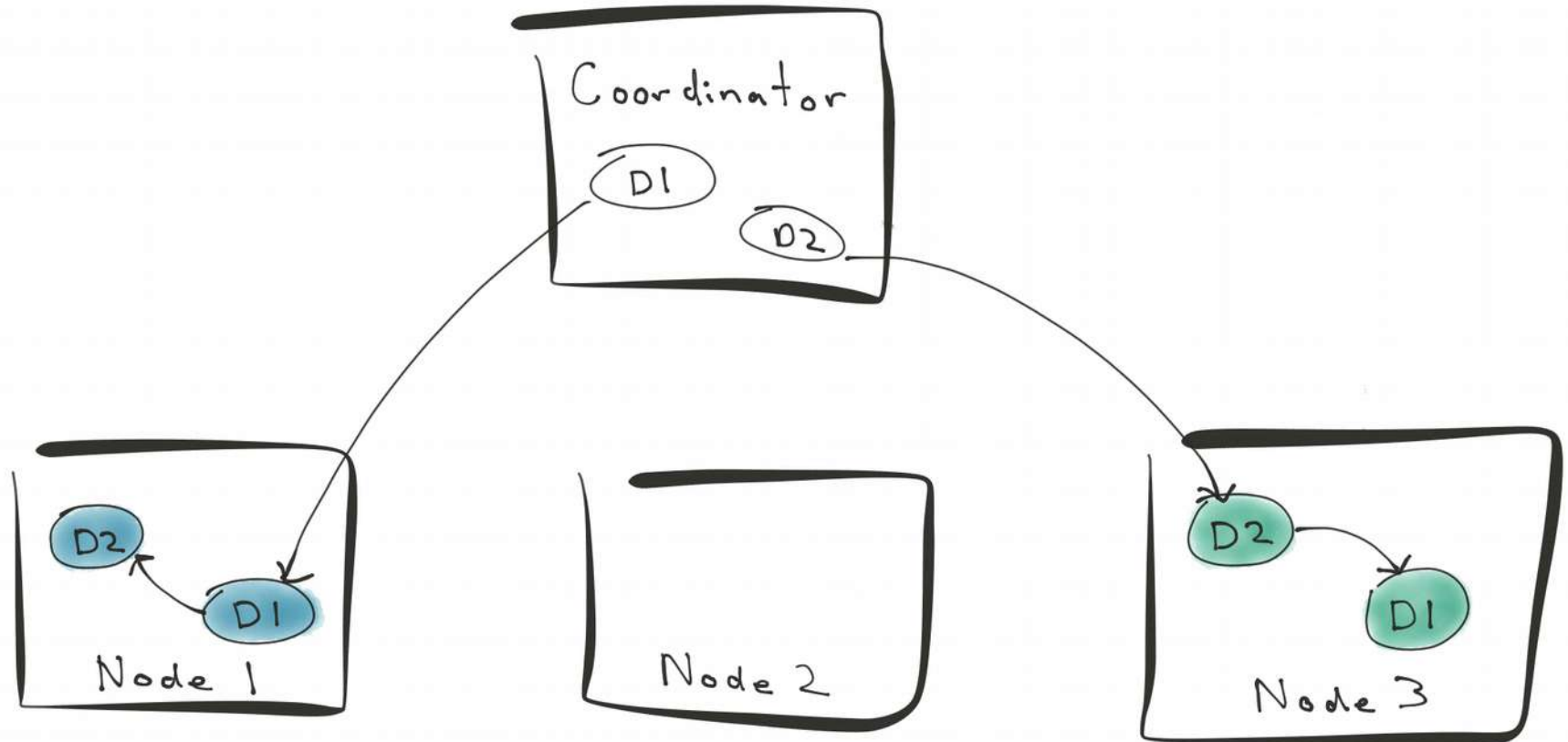
# DEADLOCK DETECTION

```
S1:   BEGIN;  // session 1 starts transaction block
S1:   UPDATE accounts
      SET balance = balance - '100.00'
      WHERE id = 'Alice';  // S1 takes 'Alice' lock

S2:   BEGIN;  // session 2 starts transaction block
S2:   UPDATE accounts
      SET balance = balance + '100.00'
      WHERE id = 'Bob';  // S2 takes 'Bob' lock

S1:   UPDATE accounts
      SET balance = balance + '100.00'
      WHERE key = 'Bob';  // waits for 'Bob' lock held by S2

S2:   UPDATE accounts
      SET balance = balance - '100.00'
      WHERE key = 'Alice';  // deadlocks on 'Alice' lock held by S1
```
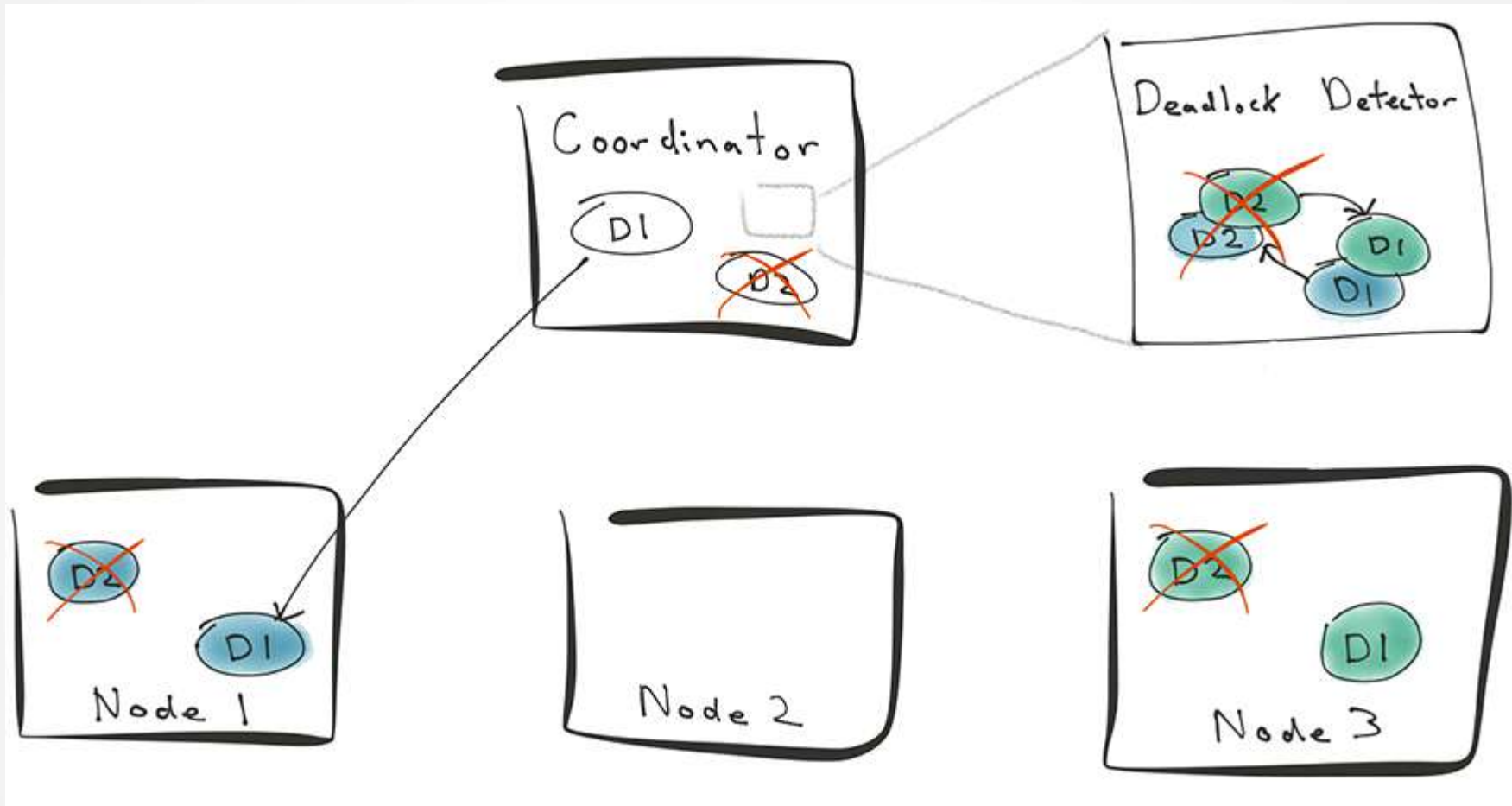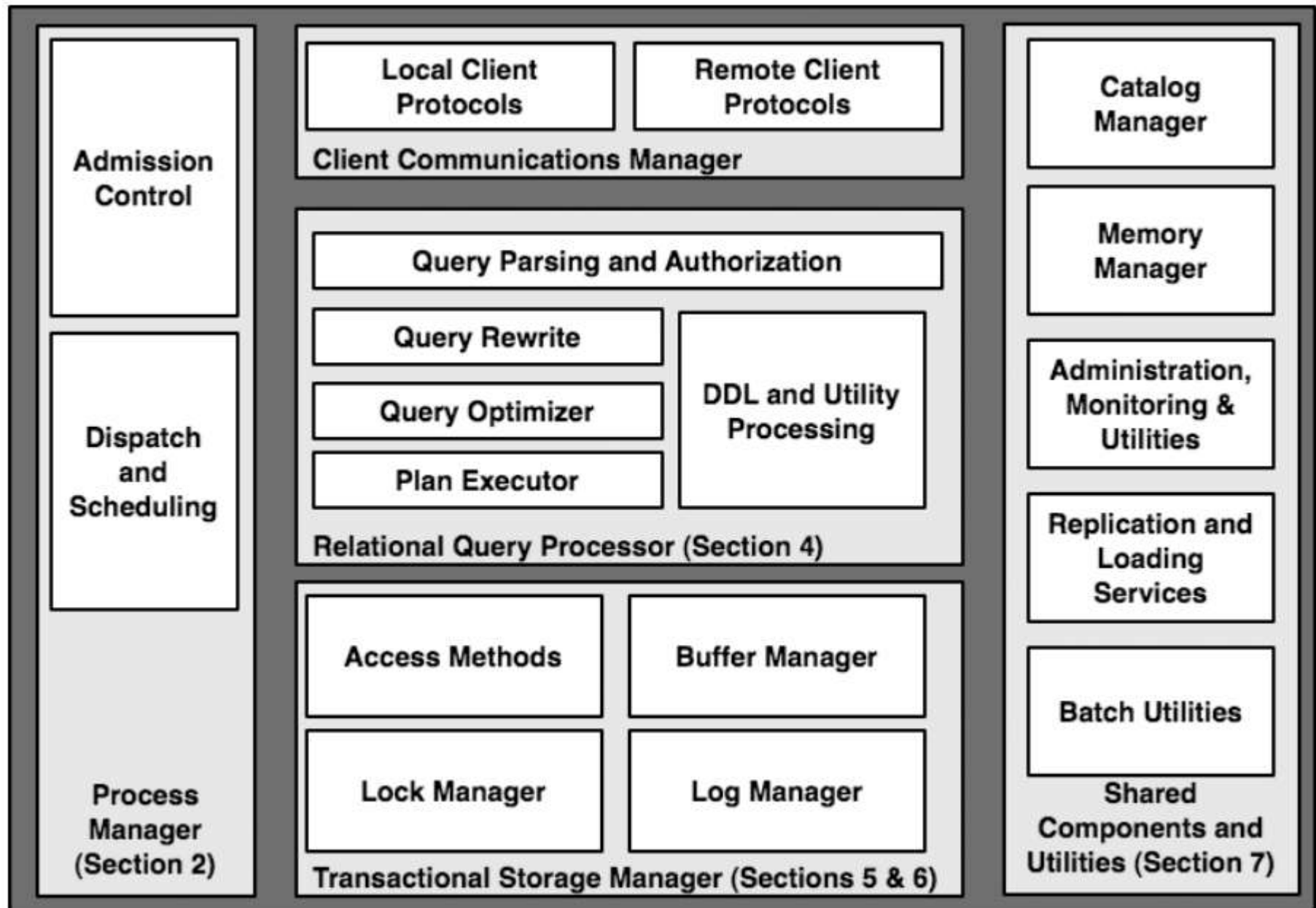
# Por que extender?

# Casos de uso

# Casos de uso

*„Hoje uma de nossas tabelas do PostgreSQL ultrapassou meio trilhão de registros [mais de 54TB de dados]. E ainda assim a maior parte de nossas queries rodam em menos de 600ms. Somente graças a @citusdata e o seu apoio incrível" - Pex CEO*

# F.A.Q.

➜ Posso criar PKs?

➜ Posso adicionar mais nós?

➜ Falha dos workers?

➜ Falha dos coordinators?

➜ Posso misturar tabelas distribuídas e locais?

➜ Como criar roles, functions, extensions em workers?

# F.A.Q.

➜ E se um só mudar de endereço?

➜ Posso distribuir por múltiplas colunas?

➜ Falha dos coordinators?

# Conclusão

*„...Quando escolhemos extender Postgres, nos disseram que SQL não escala. E acontece que é muito fácil dispensar um problema aparentemente intratável ao afirmar algo que banaliza o problema. E melhor forma de se resolver um problema complexo não é ao descartá-lo e sim quebrá-lo em pedaços menores e resolvê-los um por um...“ - Ozgun Erdogan, Citusdata CTO*

# Obrigado!!!