

# CRATEDB: A SEARCH ENGINE OR A DATABASE? BOTH!

HOW WE BUILT A SQL DATABASE ON TOP OF ELASTICSEARCH AND LUCENE



Maximilian Michels

[@stadtlegende](https://twitter.com/stadtlegende)

[max@crate.io](mailto:max@crate.io)

[mxm@apache.org](mailto:mxm@apache.org)

# WHY ARE WE TALKING ABOUT THIS?

- Traditional databases are well-researched and there are plenty of them (Postgres, MySQL, Oracle...)
  - Scalable search using these can be tricky
- Search engines are databases optimized for search and scale (Lucene, Solr, Elasticsearch)
  - You can't typically use SQL with Search Engines
- Why not stick with an mature query language standard which everybody knows?

“A scalable SQL database optimized for search without the NoSQL bullshit.”



# CRATEDB IN A NUTSHELL

- Since 2014: <https://github.com/crate/crate>
- Apache 2.0 licensed (community edition)
- Built using Elasticsearch, Lucene, Netty, Antlr, ...
- SQL-99 compatible
- REST / Postgres Wire Protocol / JDBC / Python ...

# WHAT TO EXPECT

- What is great about CrateDB
  - Easy to setup
  - No funny APIs / SQL
  - Great scale out - Massive reads / writes
  - Container aware
- Not so great
  - Transactions
  - Foreign keys

# USING CRATEDB

 CrateDB

# CRATEDB IS JUST LIKE A SQL DB

- SQL is the only query API
- `CREATE TABLE fosdem.speakers (id int PRIMARY KEY, name string)`
- `CREATE TABLE fosdem.talks (id INT PRIMARY KEY, title STRING, abstract STRING, speaker INT);`
- `INSERT INTO fosdem.speakers (id, name) VALUES (1, 'max')`
- `INSERT INTO fosdem.talks (id, title, abstract, speaker) VALUES (1, 'Talk about CrateDB', 'bla', 1)`
- `SELECT * FROM fosdem.talks t1 LEFT JOIN fosdem.speakers t2 ON t1.id = t2.id`

# BUT THERE IS MORE

- CrateDB denormalized (no joins necessary)
- `CREATE TABLE fosdem.speakers (name STRING, talk OBJECT AS (title STRING, abstract STRING))`
- `INSERT INTO fosdem.speakers (name, talk) VALUES ('max', {title = 'CrateDB', abstract = 'Lorem ipsum'})`
- `SELECT talk['title'] as title FROM fosdem.speakers ORDER BY title`

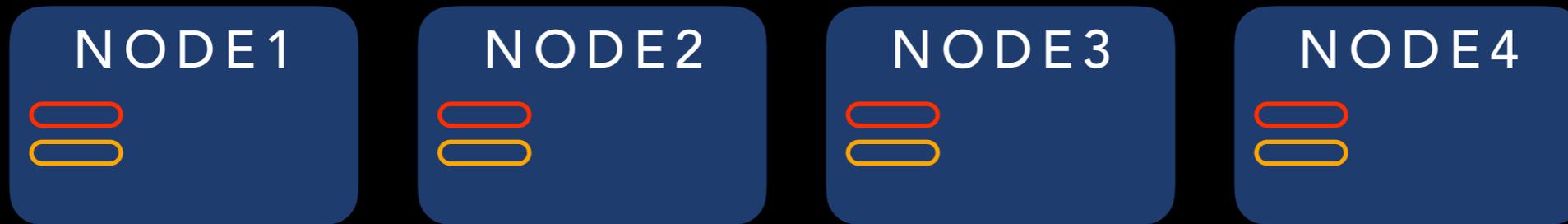
# CLUSTERING / REPLICATION



- `CREATE TABLE fosdem.speakers (name STRING, talk OBJECT AS (title STRING, abstract STRING))`
- `CLUSTERED BY` name into 4 shards

SHARD

# CLUSTERING / REPLICATION

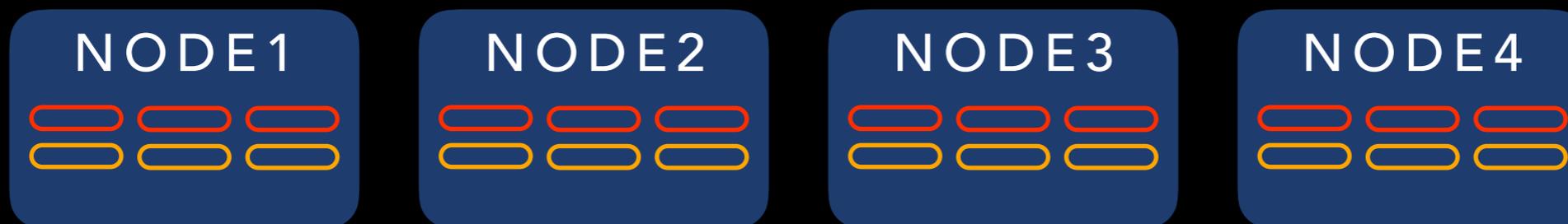


- **CREATE TABLE** fosdem.speakers (name **STRING**, talk **OBJECT AS** (title **STRING**, abstract **STRING**))
- **CLUSTERED BY** name into 4 shards
- **WITH** (number\_of\_replicas = 1)

PRIMARY

REPLICA

# PARTITIONED TABLES



- **CREATE TABLE** fosdem.speakers (name **STRING**, talk **OBJECT** as (title = **STRING**, abstract = **STRING**), year **INT**)
- **CLUSTERED BY** name into 4 shards
- **PARTITIONED BY** (year, ...)
- **WITH** (number\_of\_replicas = 1)

PRIMARY

REPLICA

# MORE FEATURES

- Aggregations
- Geo search
- Text Analyzers
- UDFs
- Snapshots
- User management
- Schema / Table privileges
- SSL encryption
- MQTT Ingestion

# ARCHITECTURE

 CrateDB

# ON THE SHOULDERS OF GIANTS



- **CrateDB:** Distributed SQL Execution Engine
  - **Antlr:** Parsing of SQL statements
  - **Netty:** REST, Postgres Wire Protocol, Web interface
  - **Lucene:** Storage, Indexing, Queries
  - **Elasticsearch:** Transport, Routing, Replication



# INTRODUCTION TO

- Lucene stores documents which are CrateDB's rows
- Documents have fields
- ```
{ _id      : '123',  
  name    : 'Bob',  
  title   : 'How I Learned to Stop Worrying  
            and Love the Bomb',  
  text    : 'Lorem ipsum...'  
}
```
- Fields are indexed for efficient lookup
- Fields have column store for efficient aggregation

# INTRODUCTION TO ELASTICSEARCH

- Elasticsearch core concepts revolve around indices, shards, and replicas
- An index is a document store with n parts, called shards
- Each shard has 0 or more replicas which hold copies of the shard data
- Replicas are not only useful for fault tolerance but also increase the search performance



# HOW TABLES RELATE TO INDICES AND SHARDS

- Each table in CrateDB is represented by an ES index with a mapping
- Each partition in a partitioned table is represented by an ES index
- Partition indices are created by encoding the partition value in the index name

```

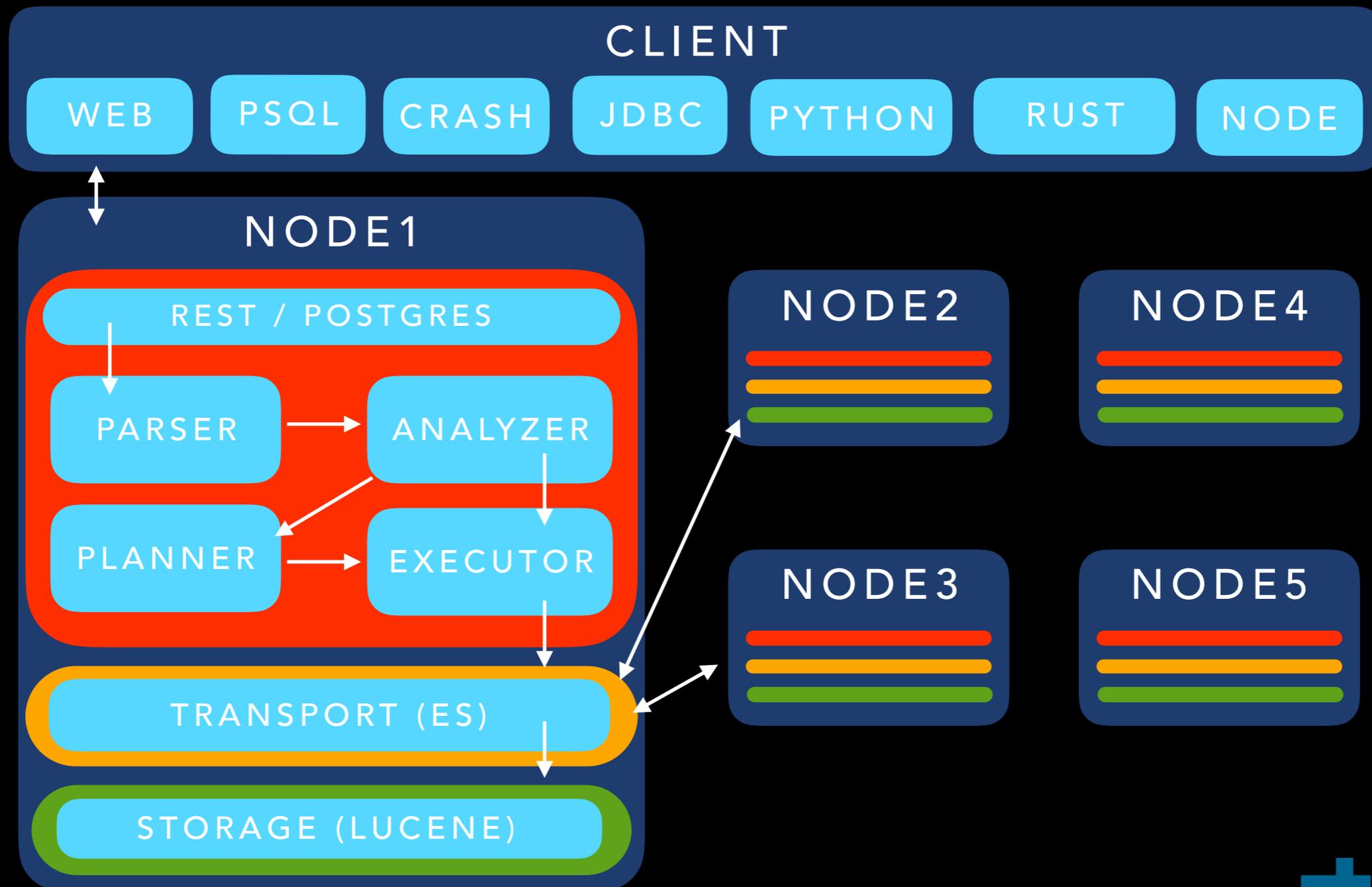
"properties":{
  "name":{"type":"keyword"},
  "talks":{"dynamic":"true",
    "properties":{
      "abstract":{"type":"keyword"},
      "title":{"type":"keyword"}
    }
  }
}

```

| TABLE  | T1 | T2      |         |     | T3 | ... |
|--------|----|---------|---------|-----|----|-----|
| INDEX  | t1 | t2.day1 | t2.day2 | ... | t3 | ... |
| SHARD1 | X  | X       | X       | ... | X  | ... |
| SHARD2 | X  | X       | X       | ... | X  | ... |
| SHARD3 | X  |         |         |     | X  | ... |
| SHARD4 | X  |         |         |     |    | ... |
| ...    |    |         |         |     |    | ... |

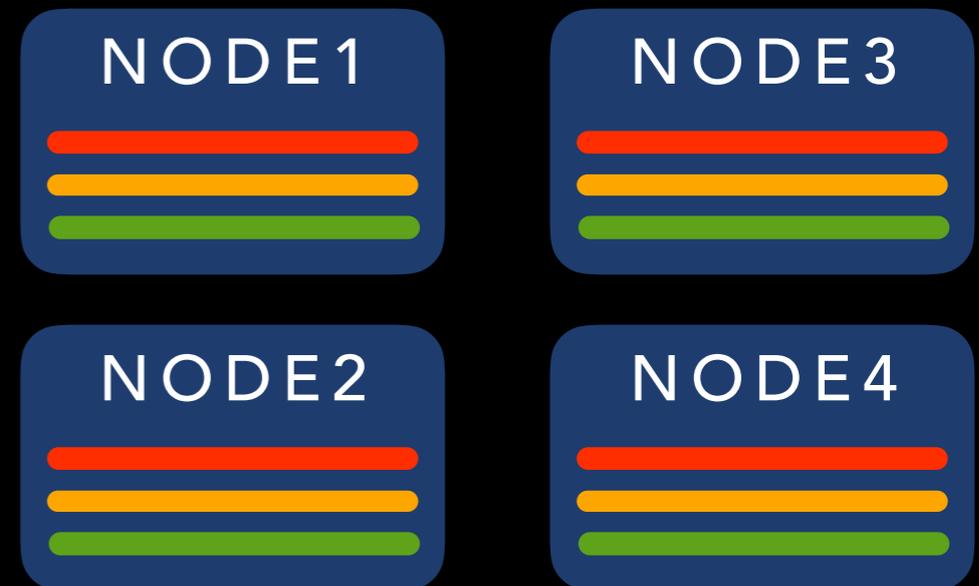
# FROM QUERY TO EXECUTION

- `SELECT name, count(*) as talks FROM fosdem.speakers  
WHERE room = 'hpc' AND year = 2018 GROUP BY name ORDER BY name`



# ARCHITECTURE HIGHLIGHTS

- Distributed storage / Distributed query execution
- Masterless
- Replication
- Only ephemeral storage needed (Container aware)
- Optimized for search: Indexing of all fields with Lucene (tuneable)



HANDS-ON

 CrateDB

# WHAT CAN YOU DO WITH CRATEDB?

- Monitoring (IoT, Industry 4.0, Cyber Security)
- Stream Analysis
- Text Analysis
- Time Series Analysis
- Geospatial Queries





## Cluster: demo

Health

good

Replicated Data

100.0%

Available Data

100.0%

Total Records

11.5 Billion

Underrepl. Records

0

Unavail. Records

0

## Cluster Load

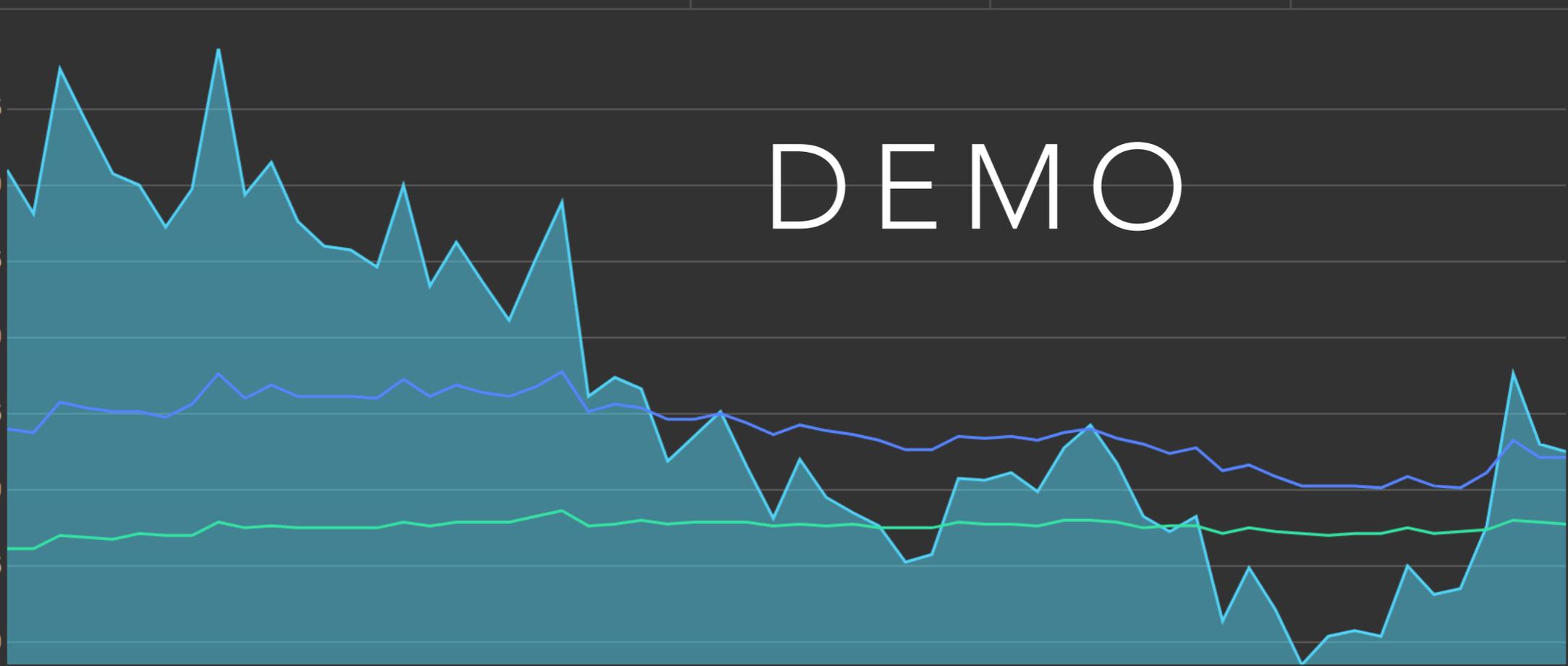
● Load 1

● Load 5

● Load 15

2.15  
2.10  
2.05  
2.00  
1.95  
1.90  
1.85  
1.80

DEMO





Filter tables ...



Doc Tables

**metrics**

3.0 Billion Records (565.7 GB)  
48 Shards / 0-1 Replicas



**pg\_settings**

0 Records (162.0 B)  
1 Shards / 0 Replicas



**statkraft**

984,344 Records (80.8 MB)  
16 Shards / 0-1 Replicas



**uk\_power**

0 Records (3.3 KB)  
20 Shards / 0-1 Replicas



jodok Tables

uk\_dale Tables

**energy\_meta**

1,380 Records (520.6 KB)  
8 Shards / 0-1 Replicas



gantner Tables

**home**

8.3 Billion Records (1.0 TB)  
72 Shards / 0 Replicas



**home\_calc**

2.8 Million Records (201.4 MB)  
4 Shards / 0-1 Replicas



**metric**

15.5 Million Records (1.3 GB)



## Tables

Name

metrics (partitioned)

Health

good

Configured Replicas

0-1

Configured Shards

48

Started Shards

96

Missing Shards

0

Underrepl. Shards

0

Total Records

3.0 Billion

Unavailable Records

0

Underrepl. Records

0

Size

565.9 GB

Recovery

100.0%

QUERY TABLE

## Partitions

Partition Columns: day\_generated

| Health | Ident ^                  | Partition Values | Conf. Replicas | Conf. Shards | Started Shards | Missing Shards | Under Shards |
|--------|--------------------------|------------------|----------------|--------------|----------------|----------------|--------------|
|        |                          | day_generated    |                |              |                |                |              |
| good   | 04732d9h6so3idpm60o30c1g | 1517097600000    | 0-1            | 8            | 16             | 0              | 0            |

# CrateDB Web Interface



## Nodes

Name ▲ Health

|                                                              |  |
|--------------------------------------------------------------|--|
| <b>c01</b><br>c01.demo.crate <span>db</span> .cloud<br>2.3.0 |  |
| <b>c02</b><br>c02.demo.crate <span>db</span> .cloud<br>2.3.0 |  |
| <b>c03</b><br>c03.demo.crate <span>db</span> .cloud<br>2.3.0 |  |
| <b>c04</b><br>c04.demo.crate <span>db</span> .cloud<br>2.3.0 |  |
| <b>c05</b><br>c05.demo.crate <span>db</span> .cloud<br>2.3.0 |  |
| <b>c06</b><br>c06.demo.crate <span>db</span> .cloud<br>2.3.0 |  |
| <b>c07</b><br>c07.demo.crate <span>db</span> .cloud<br>2.3.0 |  |
| <b>c08</b><br>c08.demo.crate <span>db</span> .cloud<br>2.3.0 |  |

## Nodes

Name

Hostname

**c01**

c01.demo.cratedb.cloud

CrateDB Version

REST URL

2.3.0

10.4.34.21:4200

CPU Usage



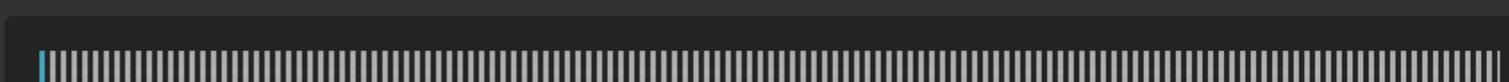
Heap Usage



Disk Usage



CrateDB CPU Usage





## Shards


 Show Shard IDs

**Started Primary**

**Started Replica**

**Initializing**

**Relocating**

**Unassigned**


|  | fhv                |                    |                      |                     |          | gantner             |      |   |   |
|--|--------------------|--------------------|----------------------|---------------------|----------|---------------------|------|---|---|
|  | parkplatz          | part               |                      | person              | replicas | test                | home |   |   |
|  | 2<br>4<br>10<br>14 | 7<br>9<br>10<br>14 | 11<br>12<br>14<br>15 | 1<br>5<br>10<br>14  | 0        | 2<br>4<br>5<br>10   | 3    | 3 | 3 |
|  | 1<br>3<br>5<br>7   | 0<br>1<br>2<br>15  | 2<br>3<br>4<br>10    | 1<br>11<br>12<br>13 | 0        | 3<br>7<br>11<br>12  | 4    | 4 | 4 |
|  | 1<br>5<br>7<br>13  | 0<br>3<br>4<br>6   | 1<br>3<br>6<br>8     | 4<br>6<br>8<br>15   |          | 1<br>3<br>5<br>8    | 1    | 1 | 1 |
|  | 0<br>8<br>11<br>15 | 5<br>12<br>15      | 5<br>7<br>9<br>13    | 7<br>8<br>11<br>12  | 0        | 0<br>11<br>12<br>15 | 0    | 0 | 0 |
|  | 0<br>2<br>3        | 1<br>2<br>5        | 0<br>4<br>5          | 0<br>3              | 0        |                     |      |   |   |

CONCLUSION

 CrateDB

# WHAT WE HAVE LEARNED

- Elasticsearch used Lucene and Netty to built a distributed search engine
- CrateDB used Elasticsearch, Lucene, and Netty to built a distributed SQL database
- CrateDB is perfect when you
  - want or have to use SQL
  - store large amounts of structured or unstructured data
  - have many thousands of queries per second

# SEE FOR YOURSELF!

- Try out CrateDB
  - Download from <https://crate.io/download/>
  - or `$ curl try.crate.io | bash`
  - or `$ docker run crate`
  - or build from source <https://github.com/crate/crate>
- Check out <https://crate.io/docs>
- Contributions welcome
  - Check out <https://github.com/crate/crate/blob/master/devs/docs/index.rst>
  - Check out the issues
  - Stackoverflow
  - Join our Slack channel



THANK YOU!

Maximilian Michels

@stadtlegende

[max@crate.io](mailto:max@crate.io)

[mxm@apache.org](mailto:mxm@apache.org)