# CASSANDRA SUMMIT 2016

# Scalable data modelling by example

Carlos Alonso (@calonso)

# Carlos Alonso

- Ex-Londoner

- MSc Salamanca University, Spain

- Software Engineer @ Jobandtalent

- Cassandra certified developer

- Datastax Cassandra MVP 2015 & 2016

- @calonso /

CASSANDRA SUMMIT 2016

# Jobandtalent

- Revolutionising how people find jobs and how businesses hire employees.

- Leveraging data to produce a unique job matching technology.

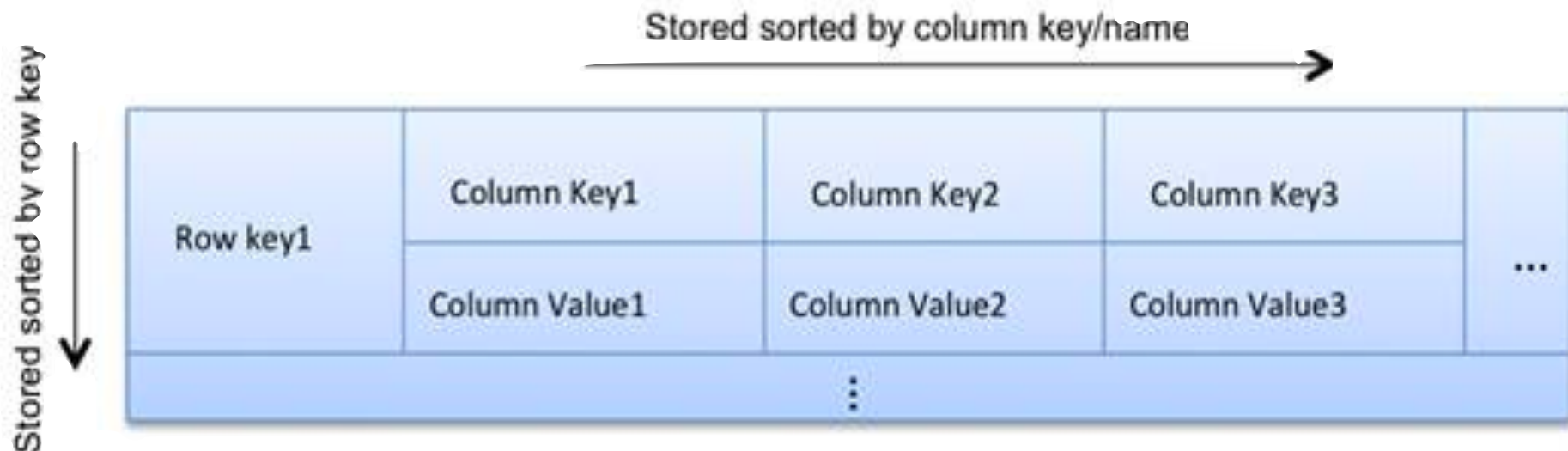- 10M+ users and 150K+ companies worldwide

- @jobandtalentEng / http://jobandtalent.com

- We are hiring!!

CASSANDRA SUMMIT 2016

CASSANDRA
SUMMIT 2016

Cassandra Concepts

The data model is the only thing you can't change once in production.

# Data organisation

Token

CASSANDRA
SUMMIT 2016

# Physical Data Layout

# Consistent Hashing

"Carlos"          185664

Hash function

17734567388476665283 49

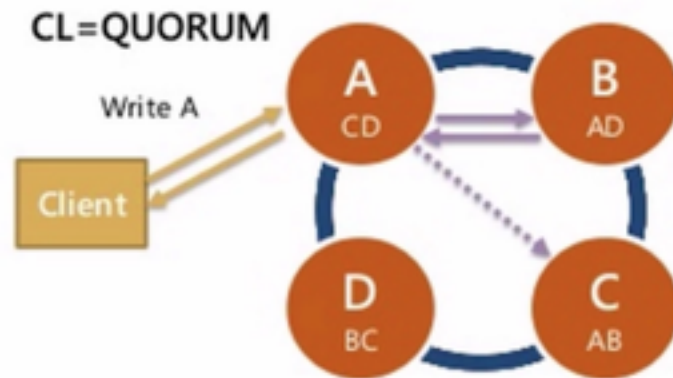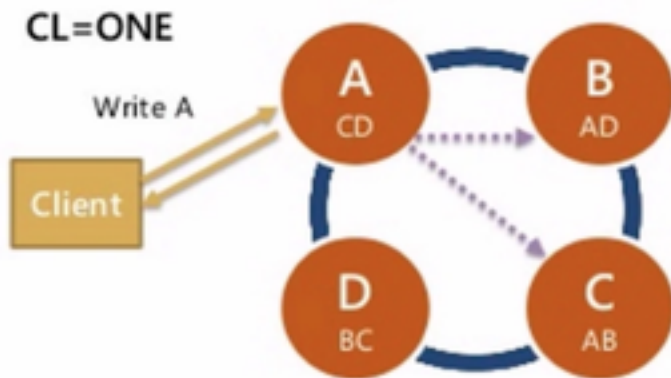-894763734895827651234

# Replication factor

How many copies (replicas) for your data

# Consistency Level

How many replicas of your data must acknowledge?

# A complete read/write example



Partitioner

834

f81d4fae-…

- RF = 3
- CL = QUORUM
- SELECT * … WHERE id = f81d4fae-…

750

749

Client    Driver

999    0

A

B    249
250

D

C

500    499

CASSANDRA
SUMMIT 2016

# Coordinator

Write <3, Betty, Blue, 63>

Acknowledge

Periodically …

## Memtable (corresponds to a CQL table)

| partition key1 | first:Oscar | last:Orange | level:42 |
|---|---|---|---|
| partition key2 | first:Ricky | last:Red | |
| partition key3 | first:Betty | last:Blue | level:63 |

Flush current state to SSTable

Node memory

Node file system

Append Only

**CommitLog**

**SSTables**

**Row and Key miss** · Off Heap · On Heap

Coordinator

Read <pk7>

**Row Cache** (optional)
pk1, pk2, **pk7**

Miss

| pk7 | first:**Elizabeth** | last:**Blue** | level:**42** |

**MemTable** (e.g., player)

| ... | ... | ... | ... |
|-----|-----|-----|-----|
| **pk7** | ... | ... | level:**42** timestamp 1114 |

Node memory
Node file system

Bloom Filter **?**

Key Cache **pk7**

Miss → Partition Summary → Partition Index

| pk1 | ... | ... | ... |
|-----|-----|-----|-----|
| **pk7** | first:Betty timestamp 541 | last:*Blue* timestamp 541 | level:63 timestamp 541 |

Bloom Filter **?**

Miss → Partition Summary → Partition Index

| pk2 | ... | ... | ... |
|-----|-----|-----|-----|
| **pk7** | first:*Elizabeth* timestamp 994 | | |

Bloom Filter ✖

| pk1 | ... | ... | ... |
|-----|-----|-----|-----|
| pk2 | ... | ... | ... |

**SSTables** (e.g., player)

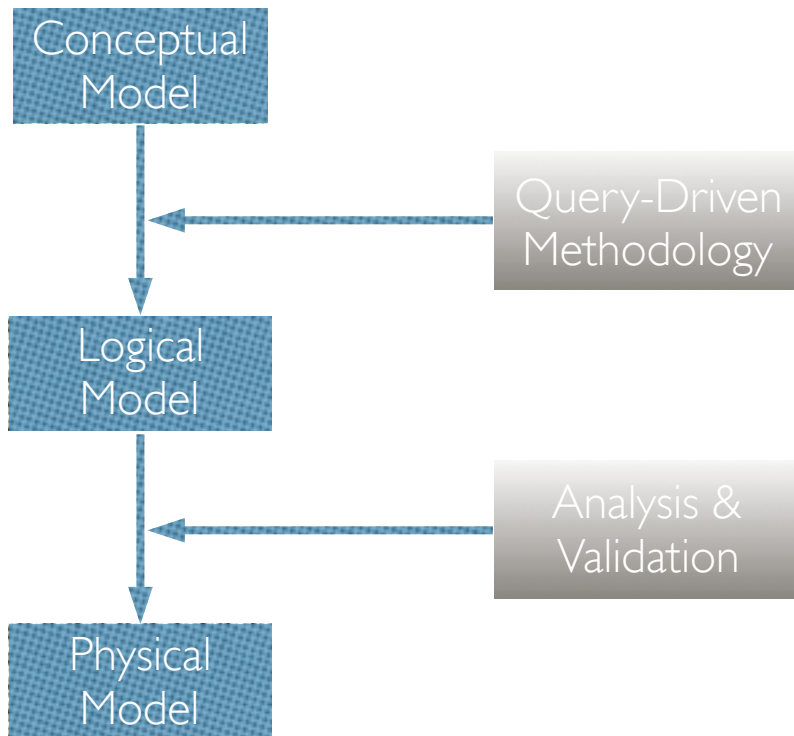CASSANDRA SUMMIT 2016

# CASSANDRA SUMMIT 2016

Data Modelling

# Data Modelling

- Understand your data

- Decide (know) how you'll query the data

- Define column families to satisfy those queries

- Implement and optimise

CASSANDRA
SUMMIT **2016**

# Data Modelling

```
┌──────────────┐
│  Conceptual  │
│    Model     │
└──────────────┘
        │
        ▼                    ┌──────────────────┐
        ◄────────────────────│   Query-Driven   │
        │                    │   Methodology    │
        ▼                    └──────────────────┘
┌──────────────┐
│   Logical    │
│    Model     │
└──────────────┘
        │
        ▼                    ┌──────────────────┐
        ◄────────────────────│    Analysis &    │
        │                    │    Validation    │
        ▼                    └──────────────────┘
┌──────────────┐
│   Physical   │
│    Model     │
└──────────────┘
```

CASSANDRA
SUMMIT **2016**

# Query Driven Methodology: goals

- Spread data evenly around the cluster
- Minimise the number of partitions read
- Keep partitions manageable

CASSANDRA
SUMMIT 2016

# Query Driven Methodology: process

- Entities and relationships: map to tables
- Key attributes: map to primary key columns
- Equality search attributes: must be at the beginning of the primary key
- Inequality search attributes: become clustering columns
- Ordering attributes: become clustering columns

# The Primary Key

PARTITION KEY        +        CLUSTERING COLUMN(S)

```
CREATE TABLE . . .(
  fields . . .
  PRIMARY KEY (part_key, clust1, . . .)
);
```
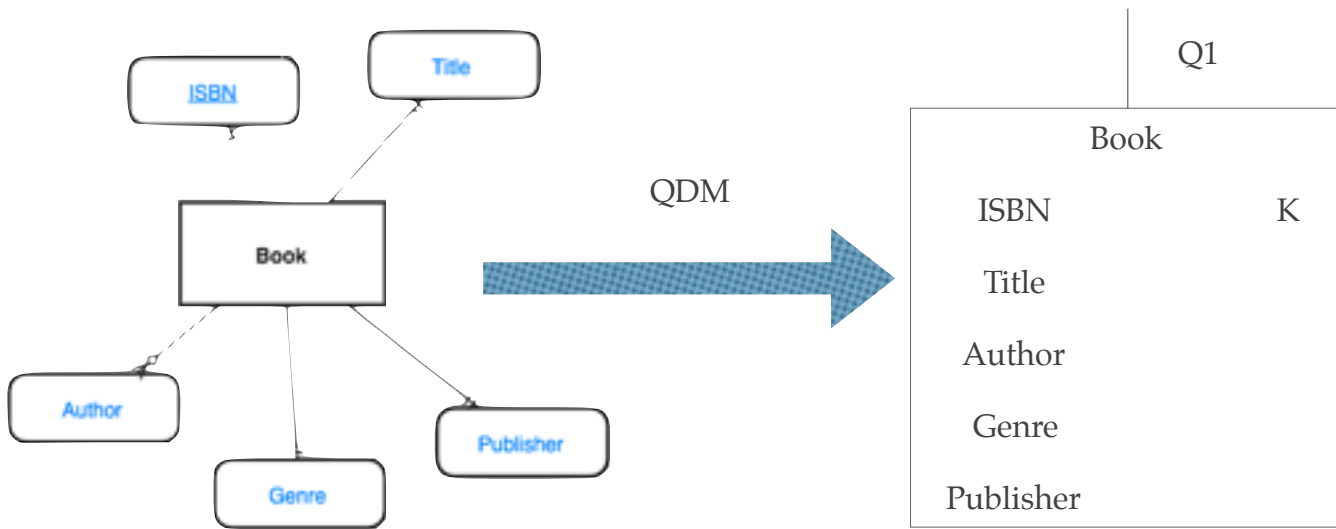
CASSANDRA
SUMMIT 2016

# Analysis & Validation

- Data evenly spread?

- 1 Partition per read?

- Are write conflicts (overwrites) possible?

- How large are partitions?

  - $N_{cells} = N_{row} \times (N_{cols} - N_{pk} - N_{static}) + N_{static} < 1M$

- How much data duplication? (batches)

CASSANDRA SUMMIT 2016

An E-Library project.

# Requirement: 1

Books can be uniquely identified and accessed by ISBN, we also need a title, genre, author and publisher.



QDM

Q1

**Book**

| | |
|---|---|
| ISBN | K |
| Title | |
| Author | |
| Genre | |
| Publisher | |

Q1: Find books by ISBN

# Analysis & Validation

Q1: Find books by ISBN

- Data evenly spread? ✓
- 1 Partition per read? ✓
- Are write conflicts (overwrites) possible? ✓
- How large are partitions? ✓
  - Ncells = Nrow X ( Ncols – Npk – Nstatic ) + Nstatic < 1M
  - 1 X (5 - 1 - 0) + 0 < 1M
- How much data duplication? 0 ✓

Q1

| Book | |
|------|------|
| ISBN | K |
| Title | |
| Author | |
| Genre | |
| Publisher | |

CASSANDRA SUMMIT 2016

# Physical data model

```
CREATE TABLE books (
  ISBN VARCHAR PRIMARY KEY,
  title VARCHAR,
  author VARCHAR,
  genre VARCHAR,
  publisher VARCHAR
);



SELECT * FROM books WHERE ISBN = '…';
```
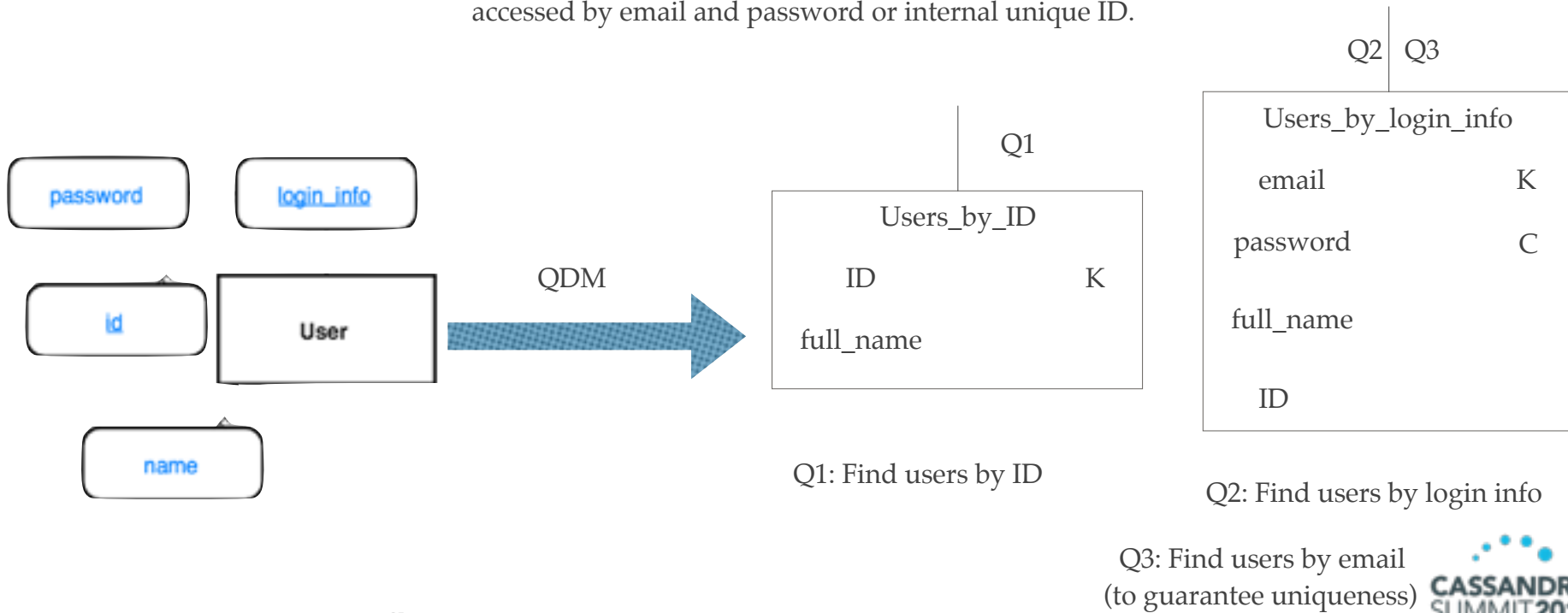
Q1: Find books by ISBN

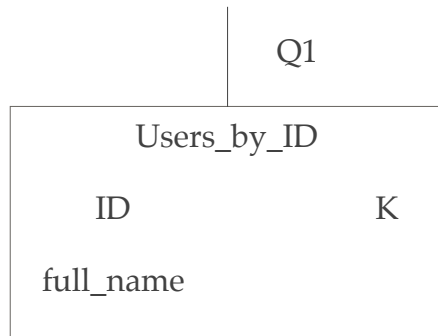| Book | |
|---|---|
| ISBN | K |
| Title | |
| Author | |
| Genre | |
| Publisher | |

Q1

CASSANDRA
SUMMIT 2016

# Requirement 2

Users register into the system uniquely identified by an email and a password. We also want their full name. They will be accessed by email and password or internal unique ID.

password

login_info

id

**User**

QDM →

Q1

**Users_by_ID**

| ID | K |
|---|---|
| full_name | |

Q1: Find users by ID

name

Q2 | Q3

**Users_by_login_info**

| email | K |
|---|---|
| password | C |
| full_name | |
| ID | |

Q2: Find users by login info

Q3: Find users by email
(to guarantee uniqueness)

CASSANDRA SUMMIT 2016

# Analysis & Validation

- Data evenly spread? ✓
- 1 Partition per read? ✓
- Are write conflicts (overwrites) possible? ✓
- How large are partitions? ✓
  - $N_{cells} = N_{row} \times (N_{cols} - N_{pk} - N_{static}) + N_{static} < 1M$
  - $1 \times (2 - 1 - 0) + 0 < 1M$
- How much data duplication? 0 ✓

Q1: Find users by ID

Q1

Users_by_ID

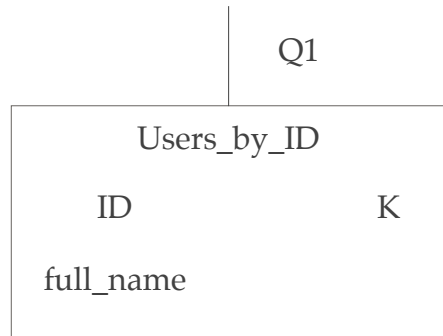ID                    K

full_name

CASSANDRA SUMMIT 2016

# Physical Data Model

```
CREATE TABLE users_by_id (
    ID TIMEUUID PRIMARY KEY,
    full_name VARCHAR
);




SELECT * FROM users_by_id WHERE ID = …;
```

Q1: Find users by ID

Q1

| Users_by_ID | |
|---|---|
| ID | K |
| full_name | |

CASSANDRA SUMMIT 2016

# Analysis & Validation

- Data evenly spread? ✓
- 1 Partition per read? ✓
- Are write conflicts (overwrites) possible? ✓
- How large are partitions? ✓
  - Ncells = Nrow X ( Ncols – Npk – Nstatic ) + Nstatic < 1M
  - 1 X (4 - 1 - 0) + 0 < 1M
- How much data duplication? 1 ✓

Q2: Find users by login info

Q3: Find users by email
(to guarantee uniqueness)

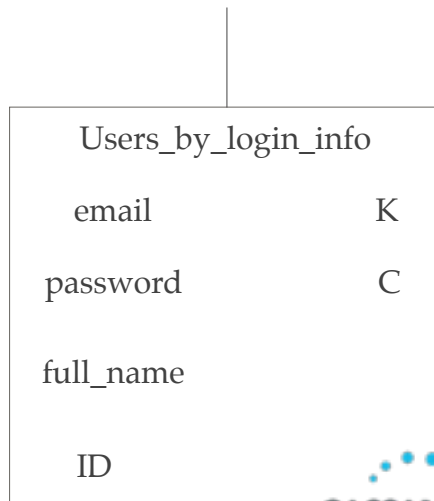| Users_by_login_info | |
| --- | --- |
| email | K |
| password | C |
| full_name | |
| ID | |

CASSANDRA SUMMIT 2016

# Physical Data Model

```
CREATE TABLE users_by_login_info (
  email VARCHAR,
  password VARCHAR,
  full_name VARCHAR,
  ID TIMEUUID,
  PRIMARY KEY (email, password)
);



SELECT * FROM users_by_login_info
WHERE email = '…' [AND password = '…'];
```

Q2: Find users by login info

Q3: Find users by email
(to guarantee uniqueness)

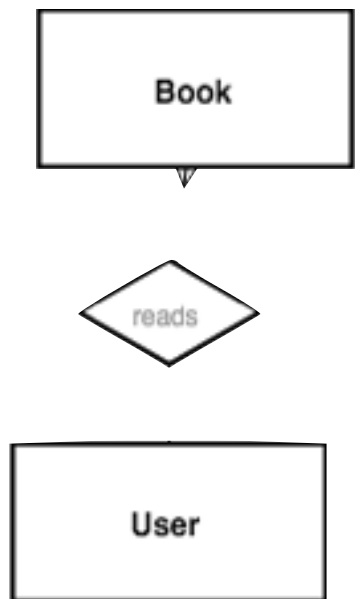| Users_by_login_info | |
| --- | --- |
| email | K |
| password | C |
| full_name | |
| ID | |

# Physical Data Model

```
BEGIN BATCH
  INSERT INTO users_by_id (ID, full_name) VALUES (…) IF NOT EXISTS;
  INSERT INTO users_by_login_info (email, password, full_name, ID) VALUES (…);
APPLY BATCH;
```

# Requirement 3

Users read books.
We want to know which books has a user read and
show them sorted by title and author

Q1: Find all books a logged
user has read

Q1

QDM

### Books_read_by_user

| | |
|---|---|
| user_ID | K |
| title | C |
| author | C |
| full_name | S |
| ISBN | |
| genre | |
| publisher | |

Book

reads

User

# Analysis & Validation

- Data evenly spread? ✔
- 1 Partition per read? ✔
- Are write conflicts (overwrites) possible? ✔
- How large are partitions? ✔
  - Ncells = Nrow X ( Ncols – Npk – Nstatic ) + Nstatic < 1M
  - Books X (7 - 1 - 1) + 1 < 1M => 200,000 books per user
- How much data duplication? 0 ✔

Q1: Find all books a logged user has read

|  | Q1 |
| --- | --- |
| Books_read_by_user | |
| user_ID | K |
| title | C |
| author | C |
| full_name | S |
| ISBN | |
| genre | |
| publisher | |

CASSANDRA SUMMIT 2016

# Physical Data Model

```
CREATE TABLE books_read_by_user (
    user_id TIMEUUID,
    title VARCHAR,
    author VARCHAR,
    full_name VARCHAR STATIC,
    ISBN VARCHAR,
    genre VARCHAR,
    publisher VARCHAR,
    PRIMARY KEY (user_id, title, author)
);



SELECT * FROM books_read_by_user
WHERE user_ID = …;
```
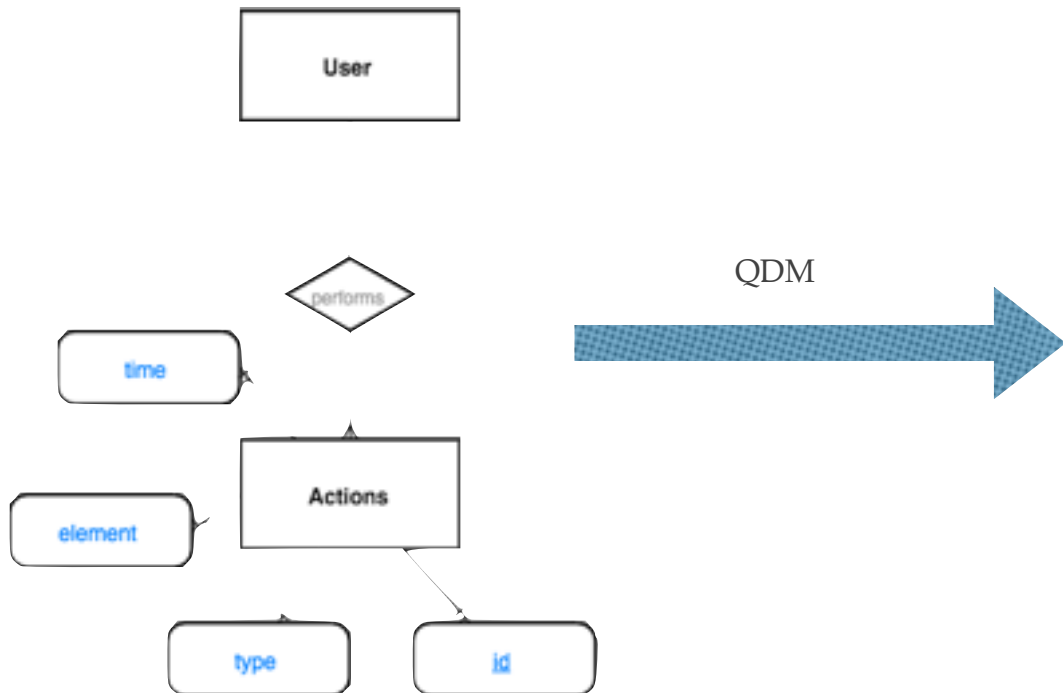
Q1: Find all books a logged
user has read

Q1

| Books_read_by_user | |
| --- | --- |
| user_ID | K |
| title | C |
| author | C |
| full_name | S |
| ISBN | |
| genre | |
| publisher | |

33

# Requirement 4

In order to improve our site's usability we need to understand how our users use it by tracking every interaction they have with our site.



User

performs

time

Actions

element

type                    id

QDM →

Q1

| Actions_by_user | |
| --- | --- |
| user_ID | K |
| time | C |
| element | |
| type | |

Q1: Find all actions a user does in a time range

CASSANDRA SUMMIT 2016

# Analysis & Validation

- Data evenly spread?  ✓

- 1 Partition per read?  ✓

- Are write conflicts (overwrites) possible?  ✓

- How large are partitions?  ✗
  - Ncells = Nrow X ( Ncols – Npk – Nstatic ) + Nstatic < 1M
  - Actions X (4 - 1 - 0) + 0 < 1M => 333.333

- How much data duplication? 0

Q1: Find all actions a user does in a time range

Q1

| Actions_by_user | |
| --- | --- |
| user_ID | K |
| time | C |
| element | |
| type | |

CASSANDRA SUMMIT 2016

# Requirement 4: Bucketing

– Ncells = Nrow X ( Ncols – Npk – Nstatic ) + Nstatic < 1M

– Actions X (5 - 2 - 0) + 0 < 1M => 333.333
per user every <bucket_size>

bucket_size = 1 year => 38 actions / h

bucket_size = 1 month => 462 actions / h

bucket_size = 1 week => 1984 actions / h

| Actions_by_user | |
| --- | --- |
| user_ID | K |
| month | K |
| time | C |
| element | |
| type | |

# Analysis & Validation

- Data evenly spread? ✓

- 1 Partition per read? ✓

- Are write conflicts (overwrites) possible? ✓

- How large are partitions? ✓
  - Ncells = Nrow X ( Ncols – Npk – Nstatic ) + Nstatic < 1M
  - Actions X (5 - 2 - 0) + 0 < 1M => 333.333 / month

- How much data duplication? 0 ✓

Q1: Find all actions a user does in a time range

Q1

| Actions_by_user | |
| --- | --- |
| user_ID | K |
| month | K |
| time | C |
| element | |
| type | |

CASSANDRA SUMMIT 2016

# Physical Data Model

```
CREATE TABLE actions_by_user (
  user_ID TIMEUUID,
  month INT,
  time TIMESTAMP,
  element VARCHAR,
  type VARCHAR,
  PRIMARY KEY ((user_ID, month), time)
);



SELECT * FROM actions_by_user
WHERE user_ID = … AND month = … AND time < … AND time > …;
```

Q1: Find all actions a user
does in a time range

| Q1 | |
|----|---|
| **Actions_by_user** | |
| user_ID | K |
| month | K |
| time | C |
| element | |
| type | |

CASSANDRA
SUMMIT 2016

# Further validation

$$\sum sizeOf(pk) + \sum sizeOf(sc) + Nr \times \sum(sizeOf(rc) + \sum sizeOf(clc)) + 8 \times Nv < 200 \text{ MB}$$

- – pk = Partition Key column
- – sc = Static column
- – Nr = Number of rows
- – rc = Regular column
- – clc = Clustering column
- – Nv = Number of values

CASSANDRA
SUMMIT **2016**

# Next Steps

- Test your models against your hardware setup
  - cassandra-stress
  - http://www.sestevez.com/sestevez/CassandraDataModeler/ (kudos Sebastian Estevez)
- Monitor everything
  - DataStax OpsCenter
  - Graphite
  - Datadog
  - . . .

CASSANDRA
SUMMIT 2016