# Data Modeling

# Data Modeling Overview

# Purpose of Conceptual Model

- Understand your data
- Essential objects
- Constraints



Users    Comments    Videos    Ratings

# Collaboration = Key

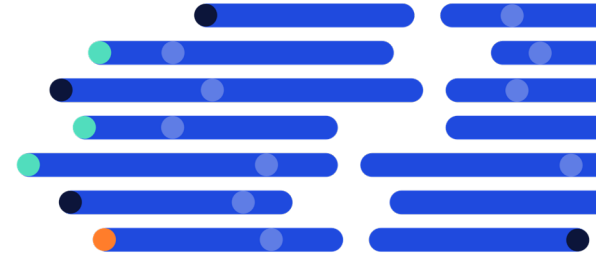# So What Does a Conceptual Model Look Like?

- Abstract view of your domain
- Technology independent
- Not specific to any database system

# Entity-Relationship (ER) Model

- Entity Types - Relationship Types - Attribute Types

**Entity Types**    **Relationship Types**

# Cardinality

- Number of times an entity can/must participate in the relationship
- Other possibilities:
  - 1-n
  - 1-1

# Attribute Types

- Fields to store data about an entity or relationship

# Key Attributes

- Identifies an Object

# Composite Attributes

- Groups related attributes together

# Multi Valued Attributes

- Stored multiple values per attribute

# Relationship Keys

# Weak Entity Types

- Entities that cannot exist without a String entity type

# Exercise 03.01

## Finish the Conceptual Model

# Exercise 03.01 – Finish the Conceptual Model

- Model the KillrVideo entities
- Identify the entity types, relationship types, and attribute types

DATASTAX

# Workflow and Access Patterns

# Where Are We Now?

# Workflow

- Each application has a workflow—Tasks/causal dependencies form a graph

- Access patterns help determine how data is accessed—Know what queries you will run first

- Example Task:

  - Have a user login to a site

DATASTAX

# Workflow & Access Patterns

Q1

User logs into site

**ACCESS PATTERNS**

Q1: Find a user with a **specified email**

DATASTAX

# Workflow & Access Patterns



**ACCESS PATTERNS**

Q1: Find a user with a **specified email**
Q2: Find most recently uploaded **videos**

# Workflow & Access Patterns



**ACCESS PATTERNS**
Q1: Find a user with a **specified email**
Q2: Find most recently uploaded **videos**
Q3: Find a **user with a specified id**

# Workflow & Access Patterns



Q1 — User logs into site

Q2 — Show latest videos added to the site

Q3 — Show basic information about a user

Q4 — Show videos added by a user

**ACCESS PATTERNS**

Q1: Find a user with a **specified email**

Q2: Find most recently uploaded **videos**

Q3: Find a **user with a specified id**

Q4: Find videos uploaded by a **user with a known id**(show most recently uploaded videos first)

# Workflow & Access Patterns



**ACCESS PATTERNS**

Q1: Find a user with a **specified email**

Q2: Find most recently uploaded **videos**

Q3: Find a **user with a specified id**

Q4: Find videos uploaded by a **user with a known id**(show most recently uploaded videos first)

Q5: Find a video with a **specified video id**

# Exercise 03.02

## Finish the Workflow & Access Patterns

# Exercise 03.02 – Finish the Workflow & Access Patterns

- Finish the application workflow based KillrVideo query requirements

- Add the remaining access patterns based on KillrVideo query requirements

DATASTAX

# Mapping Conceptual to Logical

# Where are we now?

# Query Driven Data Modeling

DATASTAX

# Chebotko Diagrams

- Graphical representation of Apache Cassandra™ database schema design
- Documents the logical and physical data model



Application Workflow

UDT Diagrams

Table Diagrams

Query List

**videos**

| | |
|---|---|
| video_id | K |
| uploaded_timestamp | |
| title | |
| description | |
| type | |
| url | |
| *encoding* | |
| {tags} | |
| <preview_thumbnails> | |

**actors_by_video**

| | |
|---|---|
| video_id | K |
| actor_name | C↑ |
| character_name | C↑ |

**UDTs**

*encoding*

encoding
height
width
{bit_rates}

**ACCESS PATTERNS**

Q1: Find a video with a **specified video id**
Q2: Find actors for a **video with a known id**(show actor names in ascending order)

# Chebotko Diagram Notations

- Logical-level shows column names and properties
- Physical-level also shows the column data type

# Logical UDT Diagram

- Represents user defined types and tuples

# Physical UDT Diagram
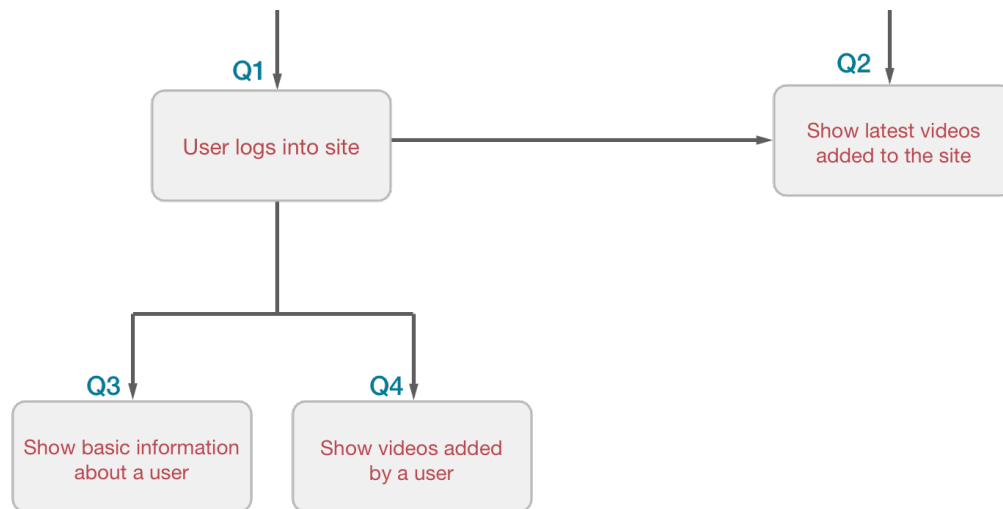
- Represents user defined types and tuples
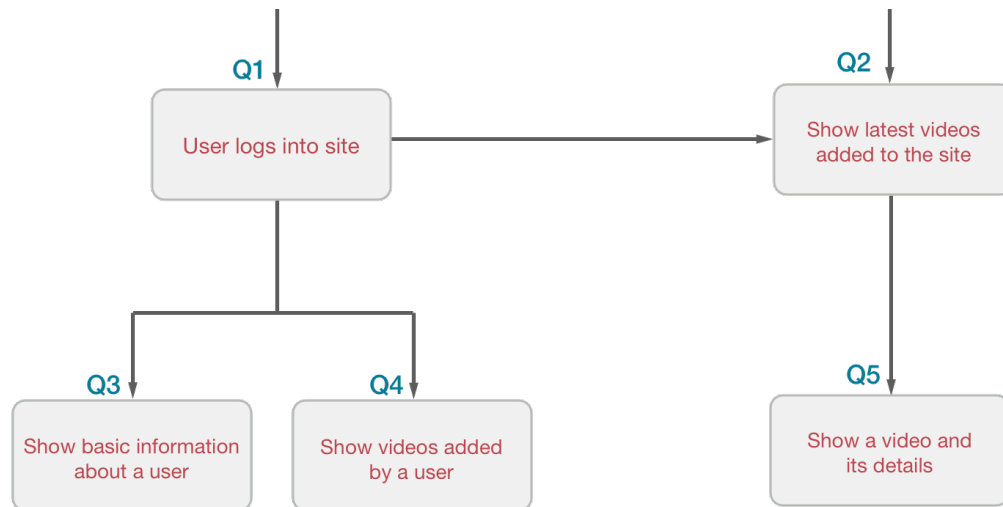
DATASTAX

# Example Chebotko Diagram



**ACCESS PATTERNS**

Q1: Find a user with a **specified email**
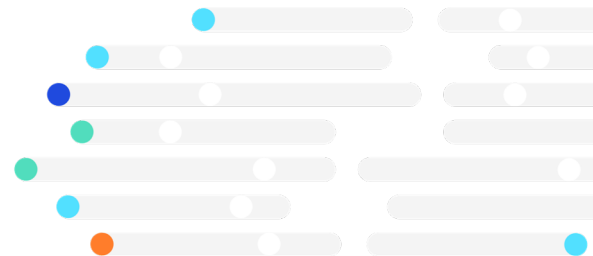Q2: Find most recently uploaded videos
Q3: Find a **user with a specified id**
Q4: Find videos uploaded by a **user with a known id**(show most recently uploaded videos first)
Q5: Find comments posted by a **user with a known id**(show most recently posted comments first)
Q6: Find comments posted for a **video with a known id**(show most recent comments first)
Q7: Find a video with a **specified video id**

**UDTs**

*encoding*

encoding
height
width
{bit_rates}

DATASTAX

# Data Modeling Principles

- Know your data

- Know your queries

- Nest data

- Duplicate data

DATASTAX

# Know Your Data

- Data captured by conceptual data model
- Define what is stored in database
- Preserve properties so that data is organized correctly

DATASTAX

# Know your Data

- Entity and relationship keys affect the table primary keys
- Primary key uniquely identifies a row / entity / relationship
- Composed of a key and possibly additional columns

| videos | |
|---|---|
| video_id | K |
| uploaded_timestamp | |
| user_id | |
| title | |
| description | |
| type | |
| *encoding* | |
| {tags} | |
| <preview_thumbnails> | |
| {genres} | |

| videos_by_user | |
|---|---|
| user_id | K |
| uploaded_timestamp | C↓ |
| video_id | C↑ |
| title | |
| type | |
| {tags} | |
| <preview_thumbnails> | |

DATASTAX

# Cardinality Constraints Affect the Key for Relationships

# Know your Queries

- Queries captured by application workflow model
- Table schema design changes if queries change



ACCESS PATTERNS

Q1: Find a user with a **specified email**
Q2: Find most recently uploaded **videos**

# Single Partition per Query = Ideal

- Most efficient access pattern
- Query accesses only one partition to retrieve results
- Partition can be single-row or multi-row

DATASTAX

# Partition+ Per Query = Acceptable

- Less efficient access pattern but not necessarily bad
- Query needs to access multiple partitions to retrieve results

# Table Scan/Multi-Table Scan = Anti-Pattern

- Least efficient type of query but may be needed in some cases
- Query needs to access all partitions in a table(s) to retrieve results

# Logical Model

# Where Are We Now?



Conceptual data model

Application workflow

Mapping conceptual to logical

Logical data model

Physical data model

Optimization tuning

DATASTAX

# Nest Data

- Data nesting is the main data modeling technique

- Nesting organizes multiple entities into a single partition

- Supports partition per query data access

- Three data nesting mechanisms:

  - Clustering columns – multi-row partitions

  - Collection columns

  - User-defined type columns

DATASTAX

# Nest Data - Clustering Columns

- Clustering column—primary data nesting mechanism
- Partition key identifies an entity that other entities will nest into
- Values in a clustering column identify the nested entities
- Multiple clustering columns implement multi-level nesting

| videos | |
|---|---|
| video_id | K |
| uploaded_timestamp | |
| user_id | |
| title | |
| description | |
| type | |
| {tags} | |
| <preview_thumbnails> | |
| {genres} | |

| actors_by_video | |
|---|---|
| video_id | K |
| actor_name | C↑ |
| character_name | C↑ |

# Nest Data – UDT

- User-defined type—secondary data nesting mechanism
- Represents 1-1 relationship, but can use in conjunction with collections
- Easier than working with multiple collection columns

| videos_by_user | | |
| --- | --- | --- |
| user_id | UUID | K |
| [videos] | video_type | |

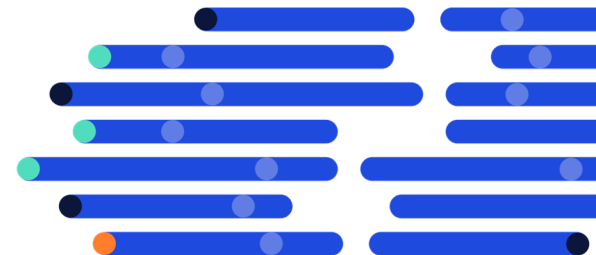| UDTs | |
| --- | --- |
| **video_type** | |
| id | TIMEUUID |
| title | TEXT |
| description | TEXT |
| type | TEXT |
| url | TEXT |
| release_date | TIMESTAMP |
| mpaa_rating | TEXT |

DATASTAX

# Duplicate Data

- Partition per query and data nesting may result in data duplication
- Query results are pre-computed and materialized
- Data can be duplicated across tables, partitions, and / or rows

| videos_by_actor | |
| --- | --- |
| actor | K |
| release_date | C↓ |
| video_id | C↑ |
| title | |
| type | |
| {tags} | |
| <preview_thumbnails> | |

| videos_by_genre | |
| --- | --- |
| genre | K |
| release_date | C↓ |
| video_id | C↑ |
| title | |
| type | |
| {tags} | |
| <preview_thumbnails> | |

| videos_by_tag | |
| --- | --- |
| tag | K |
| release_date | C↓ |
| video_id | C↑ |
| title | |
| type | |
| {tags} | |
| <preview_thumbnails> | |

DATASTAX

# Mapping Rules
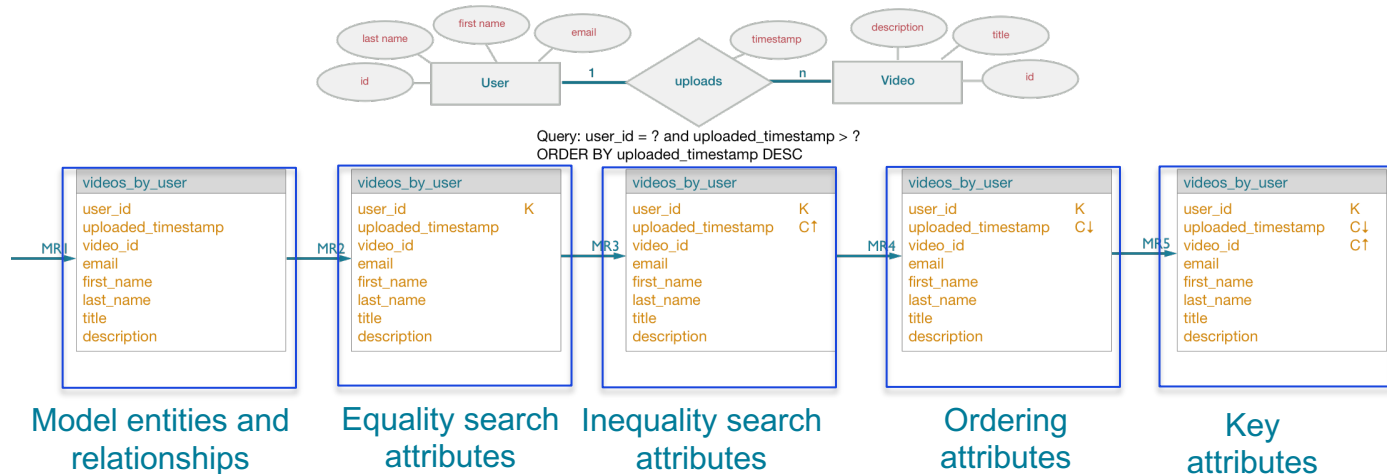
# Mapping Rules for Query Driven Methodology

- Mapping rules ensure that a logical data model is correct
- Each query has a corresponding table
- Tables are designed to allow queries to execute properly
- Tables return data in the correct order

# What are the Rules?

- Mapping Rule 1: Entities and relationships

- Mapping Rule 2: Equality search attributes

- Mapping Rule 3: Inequality search attributes

- Mapping Rule 4: Ordering attributes

- Mapping Rule 5: Key attributes
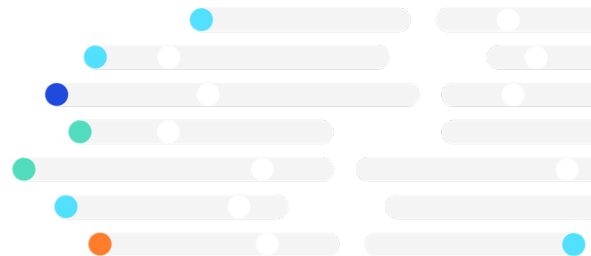
# Mapping Rules in Action

- Create a table schema from the conceptual data model and for each query
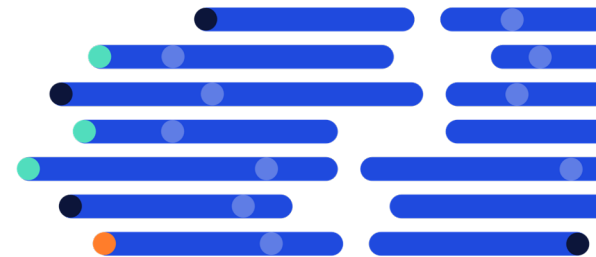- Apply the mapping rules in order

# Exercise 03.03
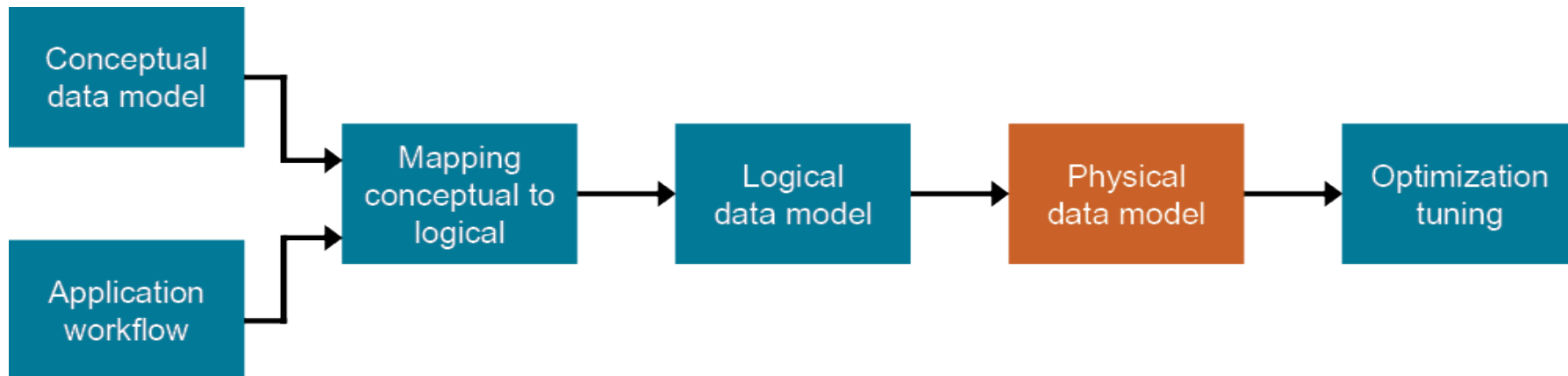
# Finalize your Logical Model

# Exercise 03.03 – Finalize your Logical Model

- Add tables to the logical model to support additional queries
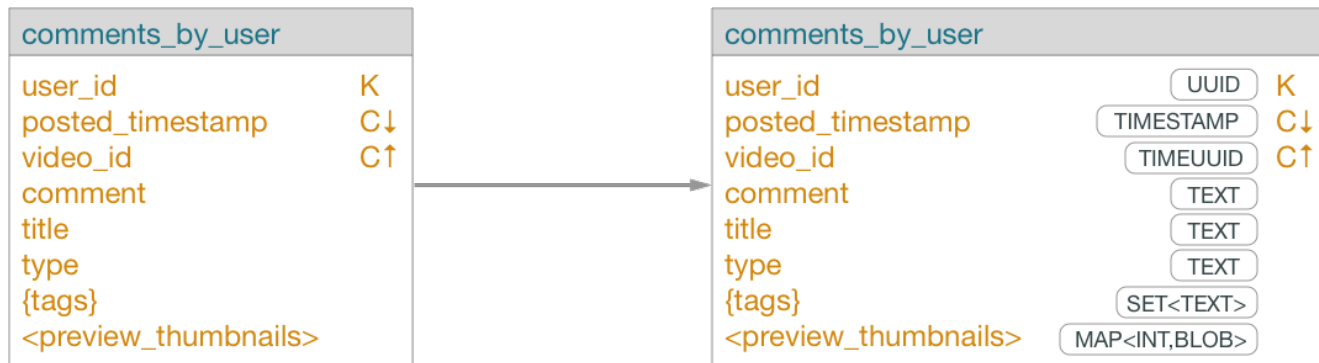- Ensure that the mapping rules are applied appropriately

DATASTAX®

# Physical Model

# Where Are We Now?

# Adding Data Types

# Creating Tables

```
CREATE TABLE comments_by_user (
    user_id UUID,
    posted_timestamp TIMESTAMP,
    video_id TIMEUUID,
    comment TEXT,
    title TEXT,
    type TEXT,
    tags SET<TEXT>,
    preview_thumbnails MAP<INT, BLOB>,
    PRIMARY KEY ((user_id), posted_timestamp, video_id)
) WITH CLUSTERING ORDER BY (posted_timestamp DESC, video_id ASC);
```

# Data Loading Methods

- COPY Command

- SSTable Loader

- Spark for Data Loading

- More about loading in Module 5!

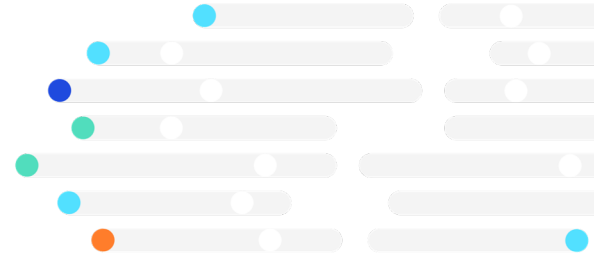DATASTAX

# Basic Data Loading with CQLCOPY

- COPY TO exports data from a table to a CSV file

- COPY FROM imports data to a table from a CSV file

- The process verifies the PRIMARY KEY and updates existing records

- If HEADER = false is specified the fields are imported in deterministic order

- When column names are specified, fields are imported in that order—missing and empty fields set to null

- Source cannot have more fields than the target table--can have fewer fields

```
COPY table1 (column1, column2, column3) FROM 'table1data.csv'
WITH HEADER=true;
```

DATASTAX

# Finalize your Physical Model

# Exercise 03.04 – Finalize your Physical Model

- Add data types to the physical data model
- Run the `CQL CREATE TABLE` statements for each table in physical model
- Load data and run some queries to test the physical data model

DATASTAX