# Motivations

## Queries Featurization

Queries featurization is crucial for query-driven estimators.

- The query is represented as a collection of four sets:
  - $< Tables >, < Joins >, < Columns >, < Values >$
  - e.g.,
- Query

```sql
SELECT COUNT(*)
FROM Title t, Company c
WHERE t.t_id = c.t_id
      AND t.Year >= 2000
      AND c.c_id <= 3
      AND c.Zip = 125
```
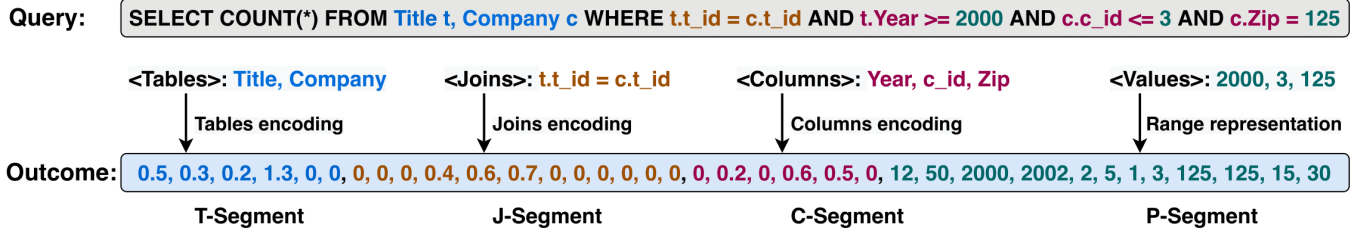
- Query representation
  - `<Tables>: Title, Company`
  - `<Joins>: t.t_id = c.t_id`
  - `<Columns>: Year, c_id, Zip`
  - `<Values>: 2000, 3, 125`
- Existing queries featurization methods cannot capture the fine-grained correlations among (Tables), (Joins), (Columns)
- Existing estimators do not give any quantification of uncertainty of the estimation.
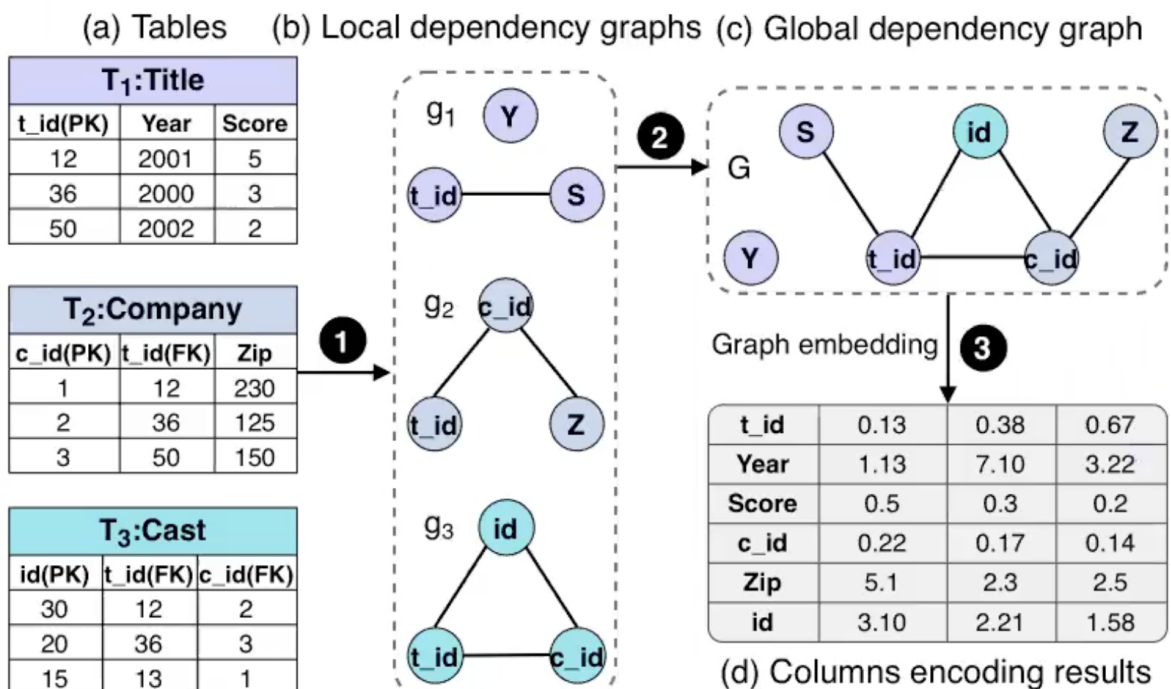
## Fauce Overview

- Overview
  - Join Schema -(Basic Information Parsing)->
  - Query Featurization -(Training Data Generation)->
    - Query Encoding
    - Predicates Representation
  - Model Design
- What Fauce includes:
  1. A new query featurization method

# Query Featurization of Fauce

**Query:**  SELECT COUNT(*) FROM **Title t, Company c** WHERE **t.t_id = c.t_id** AND **t.Year >= 2000** AND **c.c_id <= 3** AND **c.Zip = 125**

&lt;Tables&gt;: Title, Company          &lt;Joins&gt;: t.t_id = c.t_id          &lt;Columns&gt;: Year, c_id, Zip          &lt;Values&gt;: 2000, 3, 125

Tables encoding          Joins encoding          Columns encoding          Range representation

**Outcome:**  0.5, 0.3, 0.2, 1.3, 0, 0, 0, 0, 0.4, 0.6, 0.7, 0, 0, 0, 0, 0, 0, 0, 0.2, 0, 0.6, 0.5, 0, 12, 50, 2000, 2002, 2, 5, 1, 3, 125, 125, 15, 30

T-Segment          J-Segment          C-Segment          P-Segment

- How to encode `<Joins>` of a query into vectors
  - Leverage the semantic information contained in the **Join Schema**.
  - Each sub-graph of Join Schema represent a join relationship among tables
- How to encode `<Columns>` into vectors
  - a method called Columns2Vec has been proposed for the `<Columns>` encoding.
  - Tables -> Local Dependency graphs -> Global dependency graph -> Columns encoding results



(a) Tables    (b) Local dependency graphs  (c) Global dependency graph
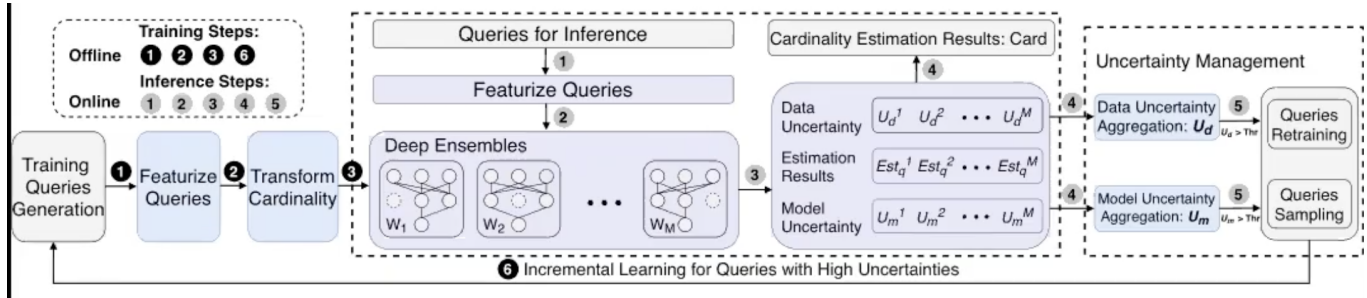
(d) Columns encoding results

- 

# Uncertainty Quantification

- The uncertainty consists of **model uncertainty** and **data uncertainty**
  - Model uncertainty: describes how confident the learned model is

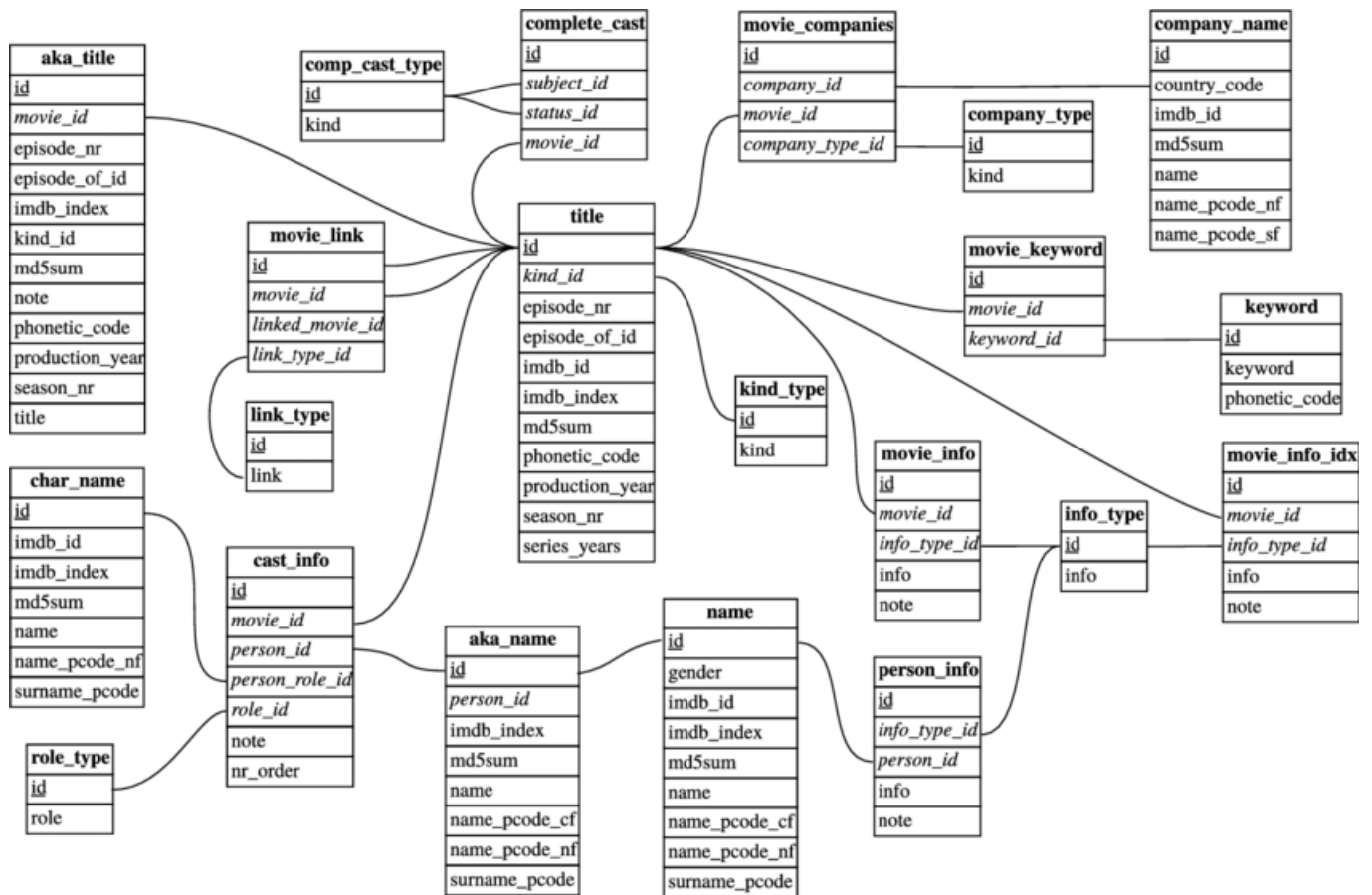- Data uncertainty: measures how noisy the collected query data are
- $Var(y) = Var(E[y|x]) + E[Var(y|x)]$
- Model uncertainty: $Var(E[y|x])$
- Data uncertainty: $E[Var(y|x)]$
- x denotes the feature vector of a query after featurization. y denotes the query's estimated cardinality

# Training and Inference of Fauce



- Data preparation
- Deep ensembles training
- Uncertainty Management

# Dataset

**aka_title**
id
movie_id
episode_nr
episode_of_id
imdb_index
kind_id
md5sum
note
phonetic_code
production_year
season_nr
title

**comp_cast_type**
id
kind

**complete_cast**
id
subject_id
status_id
movie_id

**movie_companies**
id
company_id
movie_id
company_type_id

**company_name**
id
country_code
imdb_id
md5sum
name
name_pcode_nf
name_pcode_sf

**company_type**
id
kind

**movie_link**
id
movie_id
linked_movie_id
link_type_id

**link_type**
id
link

**title**
id
kind_id
episode_nr
episode_of_id
imdb_id
imdb_index
md5sum
phonetic_code
production_year
season_nr
series_years

**kind_type**
id
kind

**movie_keyword**
id
movie_id
keyword_id

**keyword**
id
keyword
phonetic_code

**char_name**
id
imdb_id
imdb_index
md5sum
name
name_pcode_nf
surname_pcode

**cast_info**
id
movie_id
person_id
person_role_id
role_id
note
nr_order

**aka_name**
id
person_id
imdb_index
md5sum
name
name_pcode_cf
name_pcode_nf
surname_pcode

**name**
id
gender
imdb_id
imdb_index
md5sum
name
name_pcode_cf
name_pcode_nf
surname_pcode

**movie_info**
id
movie_id
info_type_id
info
note

**info_type**
id
info

**movie_info_idx**
id
movie_id
info_type_id
info
note

**person_info**
id
info_type_id
person_id
info
note

**role_type**
id
role

t-mi

t-mi_idx

t-mi

t-mc

t-mk