# Search Engine in PySpark
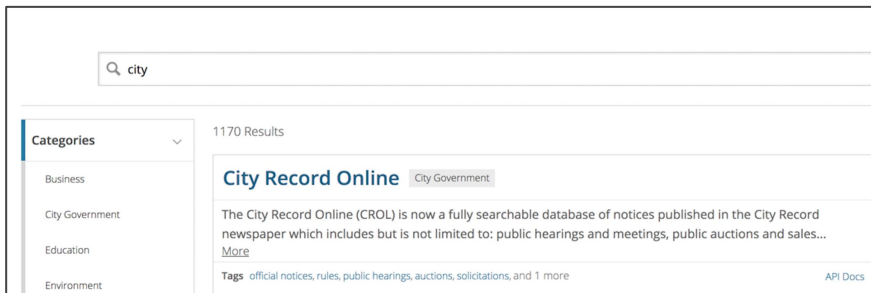
Big Data Term Project
Spring 2018

# Project Goal: search engine for structured data



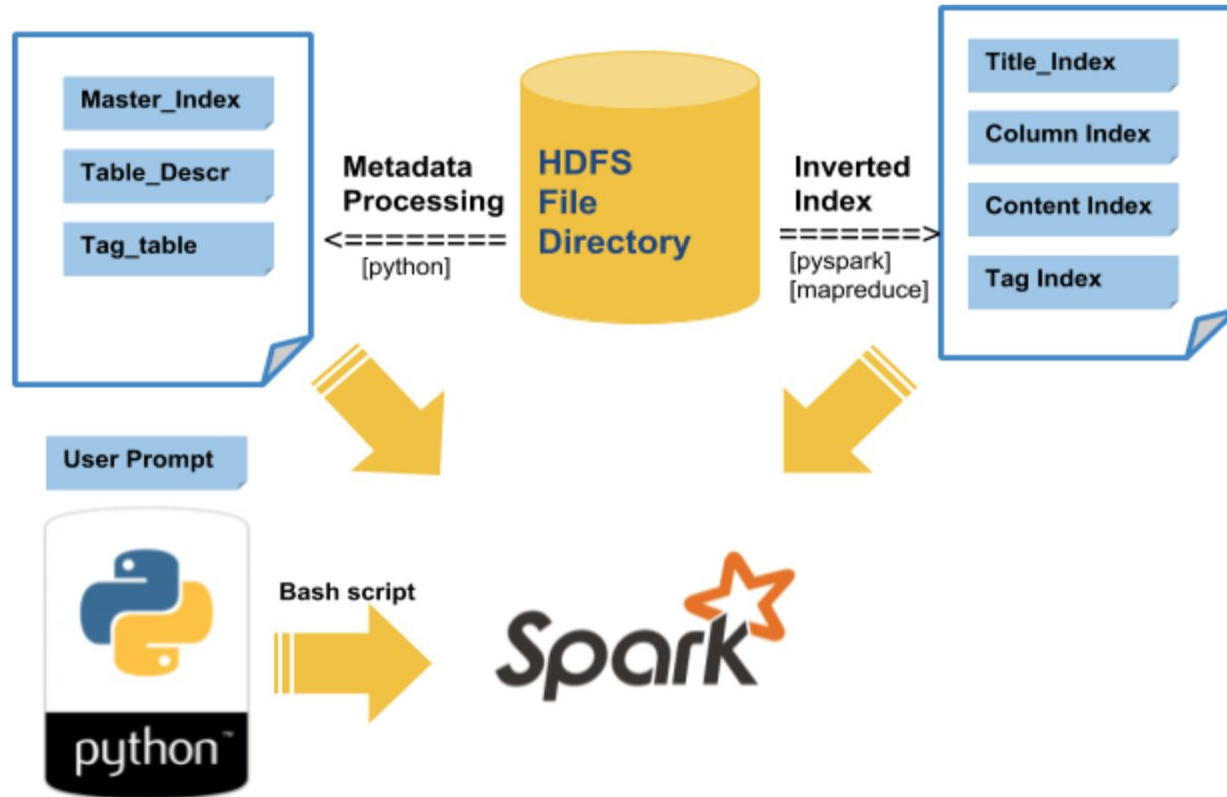## We will offer search by:

- **Title**
- **Column**
- **Content**
- **Topic**
- **Filter by table length**

**Only Title Search Available!**

# Architecture & Design: Method Design

# Architecture & Design: Application Design

# Json Processing and Other Tables

metadata

**Tag_table** (Doc_ID, tags) (generate tag_index)

**Table_Desc** (Doc_ID, Category, Description)

**Table_Schema**(used to generate column_index) (Doc_ID, Schema)

**HDFS Directory**

**Master_Index**(used to generate title_index) (Doc_ID, Table_Name, Table_Length)

**All tables in .csv format for content index**

data

5

# Inverted Indices: title, column, context, tags

1.  **Collect** the data
2.  Parse collection of tables
    a.  **Lowercase** all words: "SPARK" -> "spark"
    b.  **Tokenize**: "search+engine" -> "search", "engine"
    c.  Filter out **stop words**: "the", "a", "an", etc.
    d.  Stem each token using **Porter Stemmer Algorithm**: "fisher", "fishing", "fishes", "fished" -> "fish"
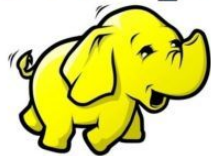
# Inverted Index

- Building the index using **MapReduce**
  - **Key**: vocabulary term
  - **Value**: postings list (tables where the term appears & in-table term frequency)

# We compared two methods for constructing indices

- Built **two versions** of inverted indices
  - PySpark
  - Hadoop
- Decided to go with **PySpark**
  - Built-in functionality for querying
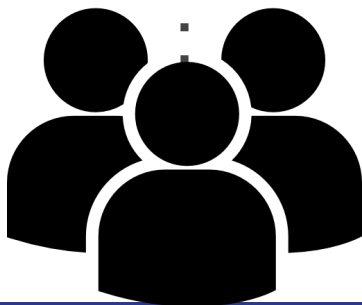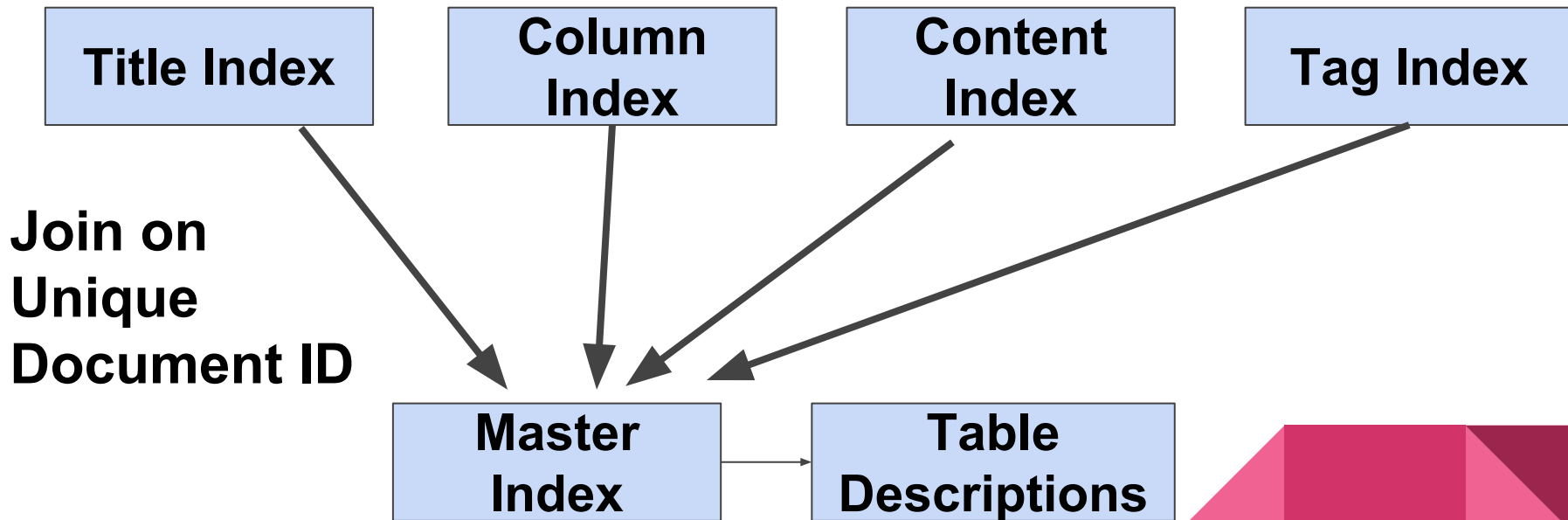
# Querying

**Spark.SQL Queries** →

## Python

- Search_Type:
  ---1,2
- Keywords:
  ---cab,taxi,city
- Row_Filter:
  ---100000

## Spark

title_search([cab,taxi,city],100000)

column_search([cab,taxi,city],100000)

# Processing Queries

# Search Ranking: TF-IDF

- Term Frequency & Document Frequency

$$TF(t, d) \qquad\qquad DF(t, D)$$

- TF.IDF Scoring

$$IDF(t, D) = \log \frac{|D| + 1}{DF(t, D) + 1}$$

$$TFIDF(t, d, D) = TF(t, d) \cdot IDF(t, D).$$

# Search Ranking: Vector Space Model

- Represent each document as a vector, with each entry being the TF.IDF weight of the corresponding term *t* in document *d*
- Represent each query as a vector in the same vector space as the documents
- Cosine similarity: the relevance between each document and the query

$$sim(q,d) \quad = \quad \frac{\vec{V}(q) \cdot \vec{V}(d)}{|\vec{V}(q)||\vec{V}(d)|}$$

# Search Ranking

## Preprocessing the tables

- Calculate TF.IDF scores from the inverted index files
- 4 inverted index files -> 4 TF.IDF models
- Save the models into files

## Searching in run-time

- Get the query & Parse
- Convert into vector
- Read the TF.IDF models
- Calculate the cosine similarities
- Rank & Output

# Demo

# Results

- Improved **Search Capabilities**
  - Title, Column, Content, Tags, Row Filter
- **Efficient Ranking**
  - TF-IDF using MLlib
- **Scalable** Implementation
  - Can built additional applications on top