



# Table Discovery in Data Lakes: State-of-the-Art and Future Directions

Grace Fan, Jin Wang, Yuliang Li, Renée J. Miller

<https://northeastern-datalab.github.io/table-discovery-tutorial/>

# Outline

- Introduction
- Table Understanding
- Table Search Engine
- Table Navigation and Exploration
- Data Science and Application Support
- Conclusion / Future work

# Data Analytics Landscape

Major shift from **Data Warehouses** → **Data Lakes**

***Thousands of tables*** → ***Millions of files***

***Well-designed Metadata*** → ***Unreliable, inconsistent metadata***

COMPUTING

## A new paradigm for managing data

Open data lakehouse architectures speed insights and deliver self-service analytics capabilities.

By MIT Technology Review Insights

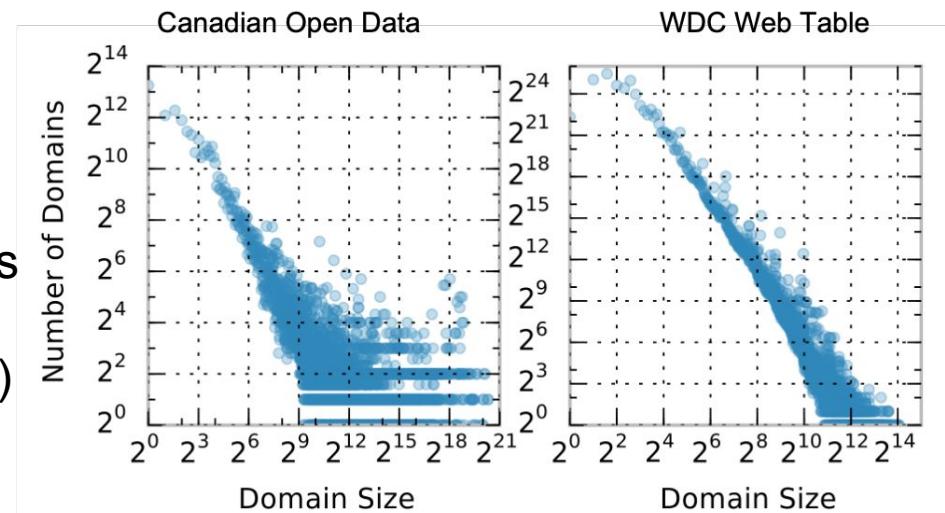
March 17, 2023

# Example Data Lakes

- Open government data
  - [Miller VLDB 18]
- Web tables
  - [Eberius et al. BDC 15], [Lehmer et al. WWW 16] and others
- Enterprise data lakes
- Domain-specific lakes

Real data lakes are skewed & large

- Large differences in attribute sizes
- Large differences in table sizes
- Massive vocabularies (no. values)

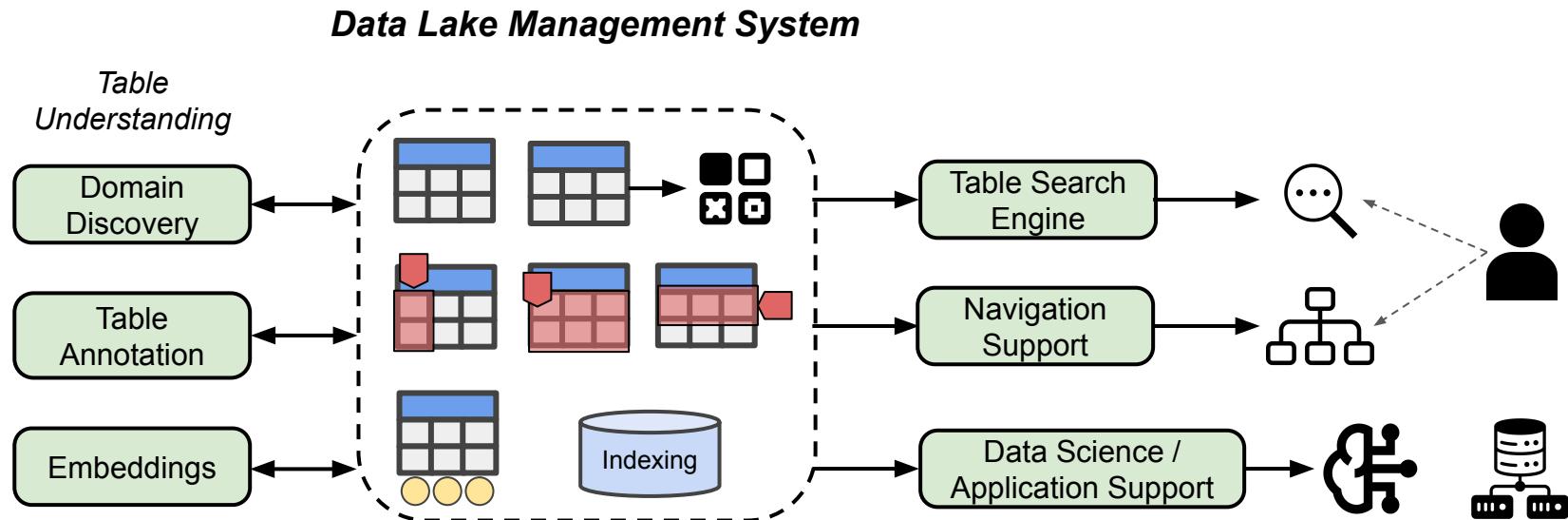


# Data Science Over Data Lakes

In data science, it is increasingly the case that the main challenge is not in *integrating known data*, rather it is in ***finding the right data to solve a given data science problem.***

How can we facilitate data science over data lakes?

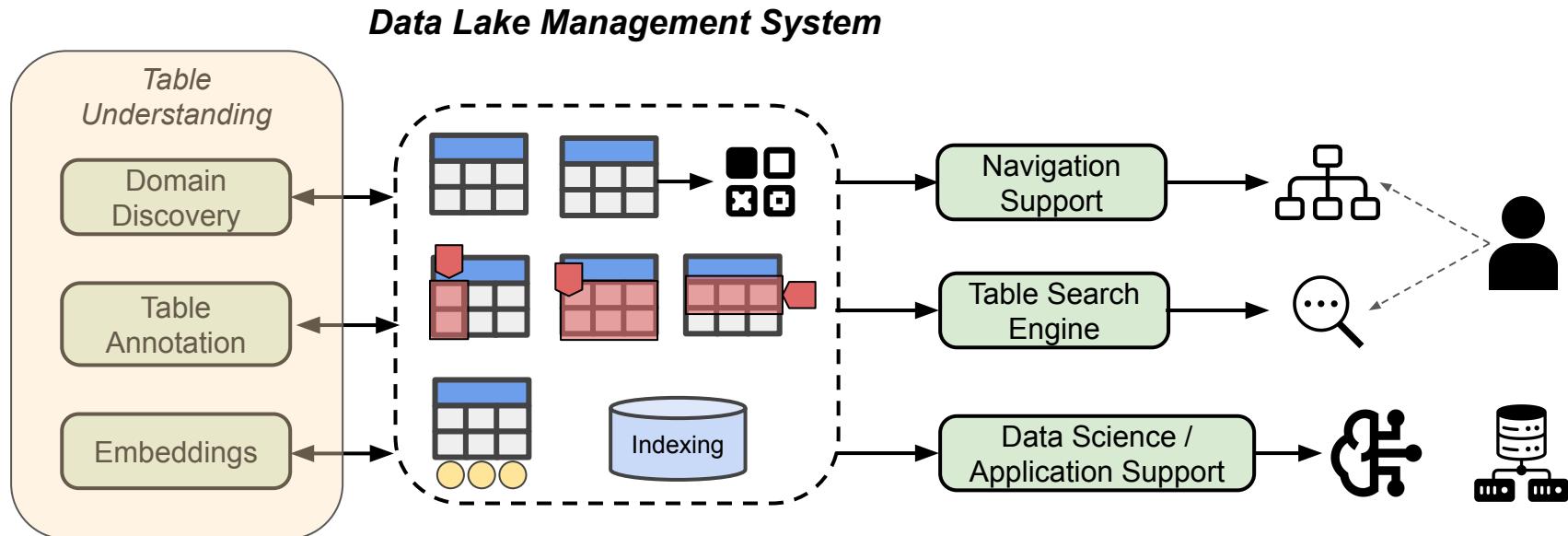
# Architecture of a Table Discovery System



# Outline

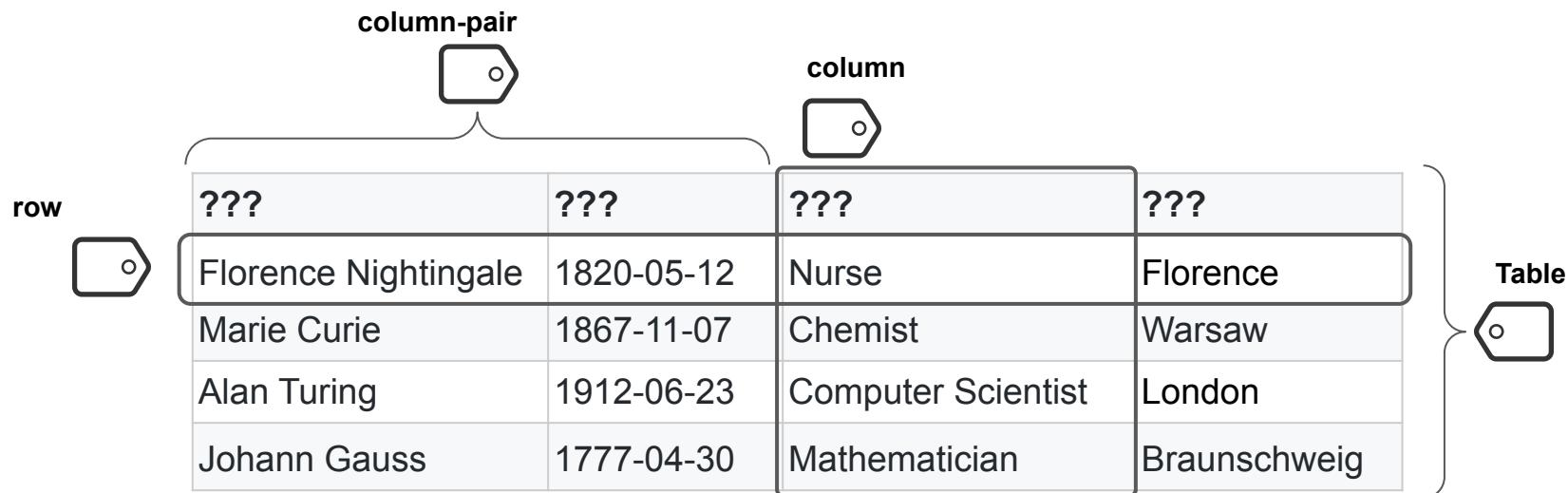
- Introduction
- **Table Understanding**
- Table Search Engine
- Table Navigation and Exploration
- Data Science and Application Support
- Conclusion / Future work

# Table Understanding



# Table Understanding tasks on different levels

- Column: Semantic Type Detection, Domain Discovery
- Column-pair: Relation Extraction
- Row: Row-level entity linking
- Table: Table-level entity linking, topic modeling

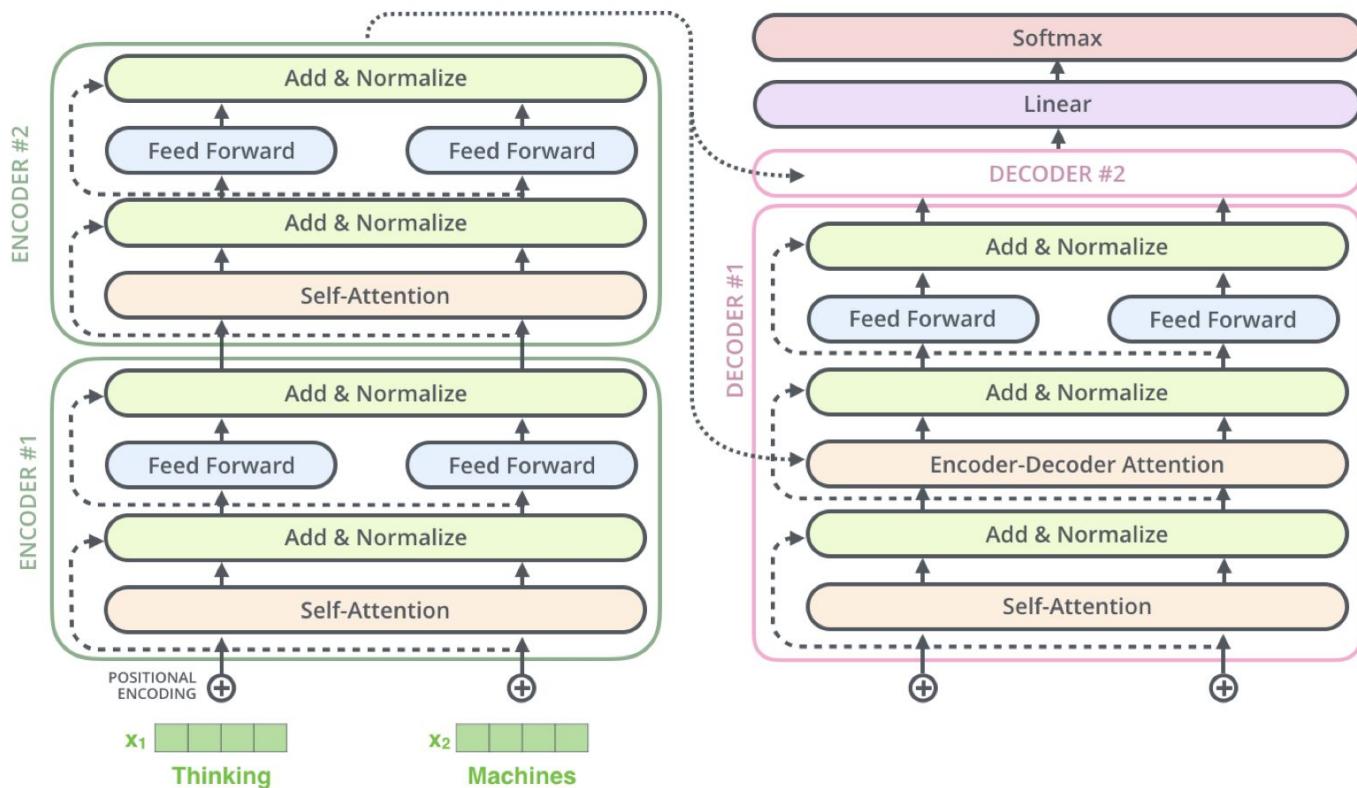


# Techniques for Table Understanding

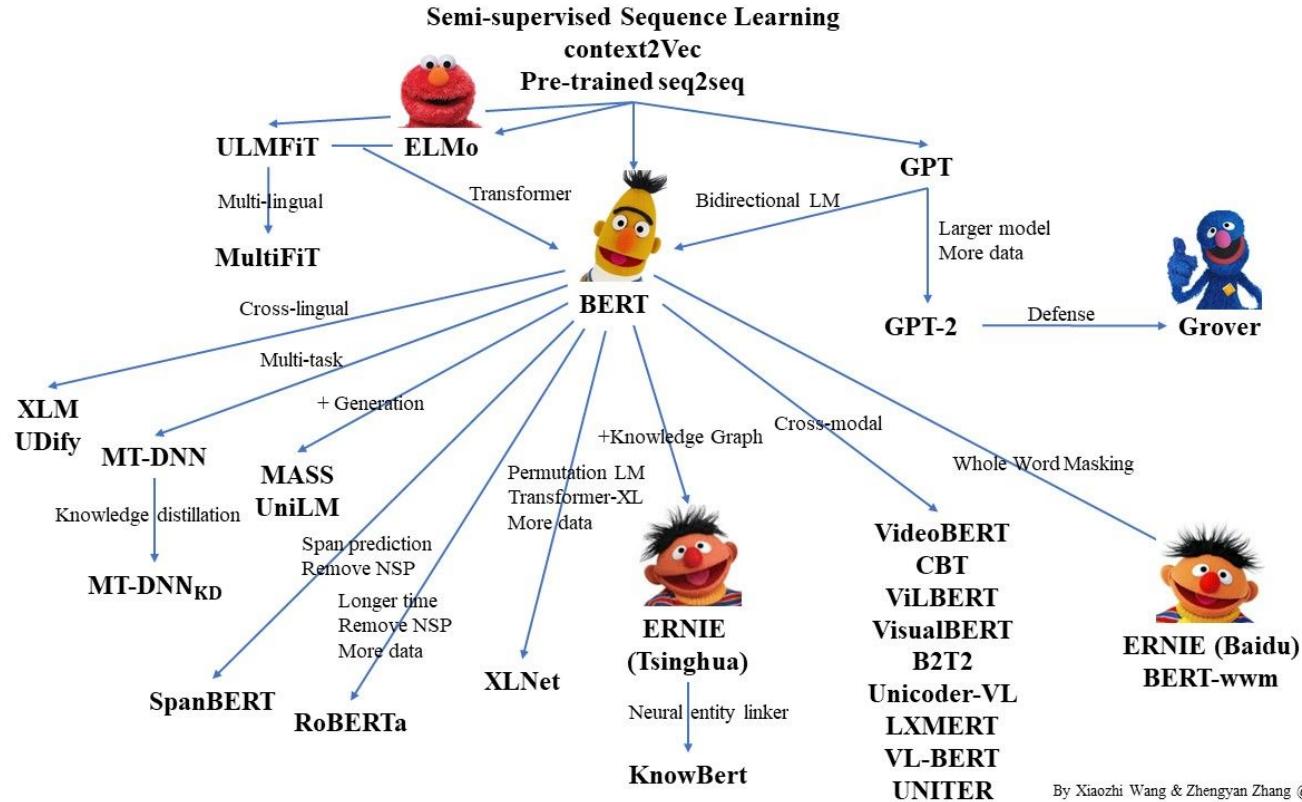
- Traditional Techniques, such as
  - Probabilistic models
  - Ontology matching
  - Feature engineering
- Nowadays Transformer-based techniques have been widely adopted in Table Understanding problems!

[Badaro et al. VLDB 22]

# Background: Transformer



# Many Transformer-based pre-trained models



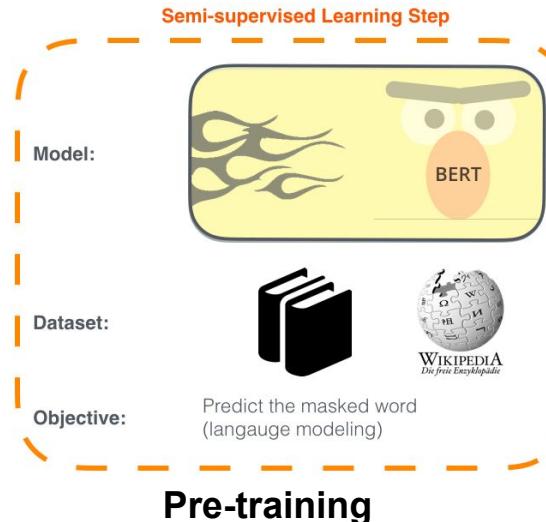
By Xiaozhi Wang & Zhengyan Zhang @THUNLP

# Background: Pre-trained language models

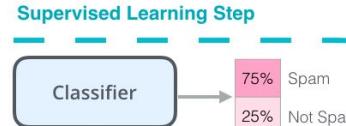
- Pre-training + fine-tuning framework

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



2 - **Supervised** training on a specific task with a labeled dataset.



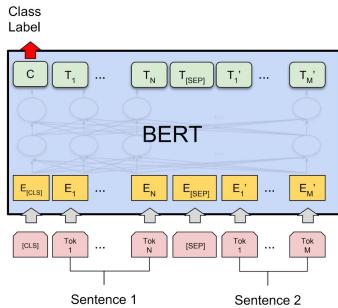
Email message      Class

Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

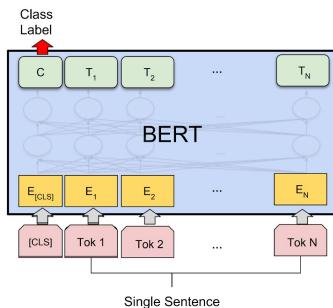
**Fine-tuning**

# Why are BERT-style models so popular?

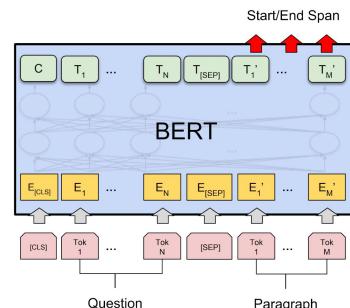
- The model can be used for many many NLP tasks
  - What we need is just to fine-tune on training data of the target task



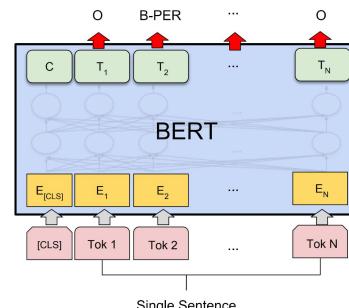
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



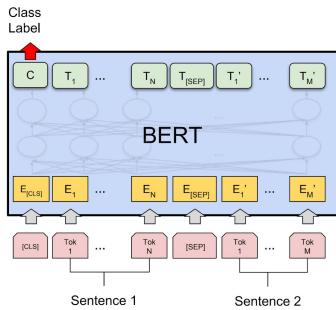
(c) Question Answering Tasks:  
SQuAD v1.1



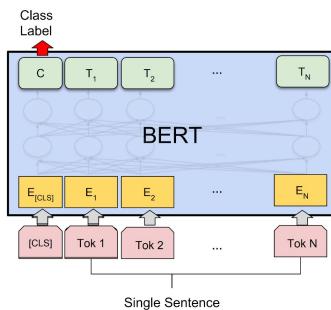
(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

# Why are BERT-style models so popular?

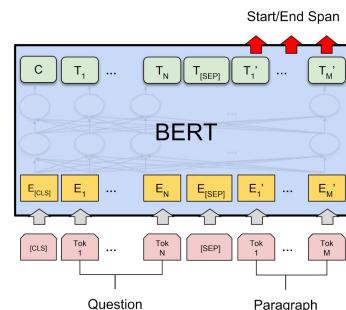
- The model can be used for many many NLP tasks
  - What we need is just to fine-tune on training data of the target task



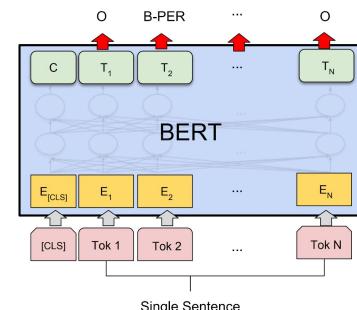
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



(c) Question Answering Tasks:  
SQuAD v1.1



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

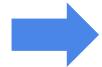


A Swiss army knife for NLP!

# Table Serialization

- Flattening a table into a token sequence

Year	Champion	Motorcycle
1949	 Nello Pagani	
1950	 Bruno Ruffo	
1951	 Carlo Ubbiali	
1957	 Tarquinio Provini	



[COL] Year [VAL] 1949 [COL] Champion [VAL] Nello Pagani  
Year | 1949 | Champion | Nello Pagani  
Year is 1949 and Champion is Nello Pagani  
...

There are several options for serialization

# Domain Discovery: Overview

- Goal: given a dataset of columns, find a set of domains representing all semantic types in the dataset

dvv_color	dvt_color	color	category2	dvt_make	bor	neighborhood1	city2
BEIGE	BEIGE	BLUE	ASIAN	4WD	BRONX	AVERNE	ASTORIA
BLACK	BLACK	BROWN	BLACK	BLACK	BROOKLYN	ASTORIA	BAY RIDGE
BLUE	BLUE	GOLD	ELL	BLUE	MANHATTAN	BAYSIDE	BRONX
BROWN	BROWN	GREEN	EP	BMW	QUEENS	BRIARWOOD	BROOKLYN
GOLD	GRAY	MAROON	FEMALE	BUICK	STATEN IS	CORONA	BUSHWICK
GRAY	GREEN	OLIVE	FORMER ELL	CHEVY	UNKNOWN	ELMHURST	CHELSEA
GREEN	KHAKI	ORANGE	HISPANIC	GRAY		UNKNOWN	CORONA
KHAKI	MAROON	PINK	MALE	GREEN		... (41 more)	MANHATTAN
ORANGE	OLIVE	RED	NOT SWD	HONDA			N/A
PINK	PINK	WHITE	SWD	TOYOTA			STATEN IS
RED	RED	YELLOW	WHITE	WHITE			... (121 more)
WHITE	YELLOW			... (87 more)			
category1	ethnicity	demographic	boroname	borough_name	borough	neighborhood2	rental_nbh
ASIAN	ASIAN	ASIAN	BRONX	BRONX	BRONX	ASTORIA	AVERNE
BLACK	BLACK	BLACK	BROOKLYN	BROOKLYN	BROOKLYN	BAYSIDE	ASTORIA
HISPANIC	HISPANIC	HISPANIC	MANHATTAN	MANHATTAN	BELLEROSE	BELLEROSE	BAYSIDE
WHITE	WHITE	N/A	QUEENS	QUEENS	CORONA	CORONA	JAMAICA
		WHITE	STATEN IS	STATEN IS	ELMHURST	ELMHURST	JAMAICA
					FAR ROCKAWAY	FAR ROCKAWAY	FLUSHING
					QUEENS	QUEENS	FLUSHING
					STATEN IS	STATEN IS	... (42 more)
							... (48 more)

Domain: Color = {BEIGE, BLACK, BLUE, BROWN, GOLD, GRAY, GREEN, KHAKI, MAROON, OLIVE, ORANGE, PINK, RED, WHITE, YELLOW}

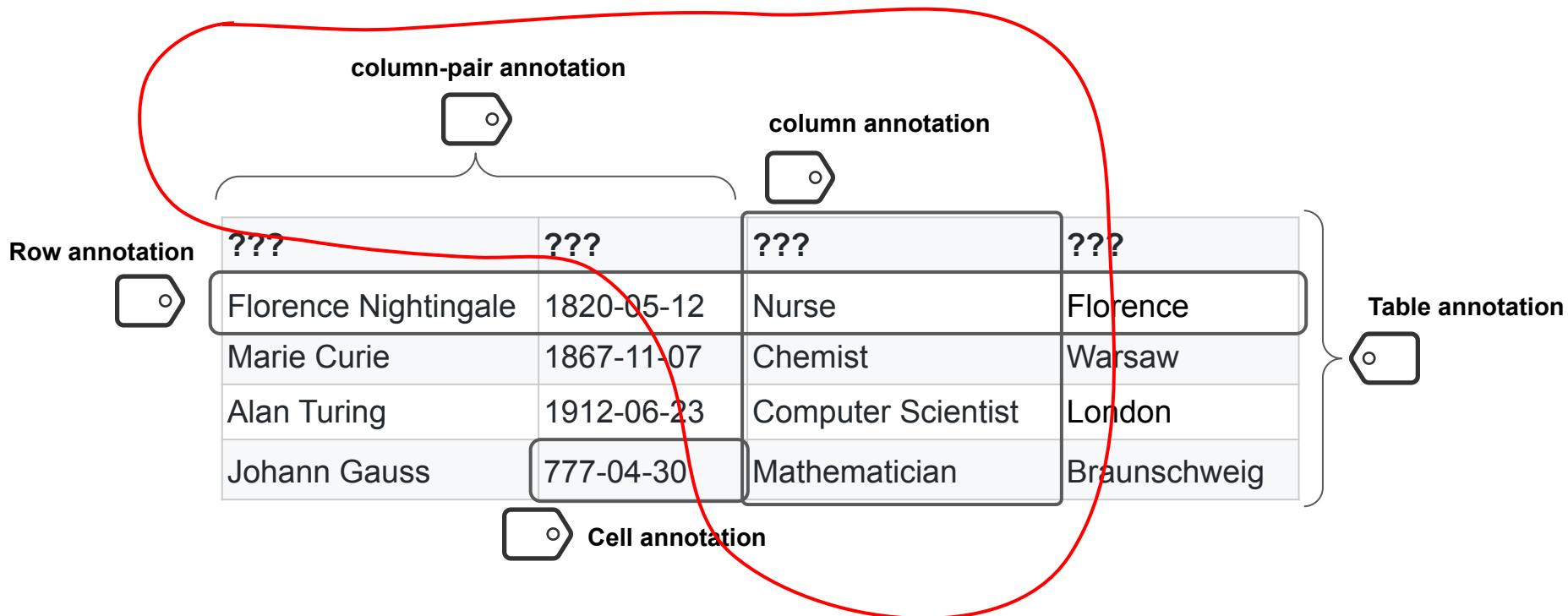
# Domain Discovery: solutions

- C<sup>4</sup>: Discovering Enterprise Concepts Using Spreadsheet Tables
  - Use co-occurrence information to build a concept hierarchy
  - Select nodes from the hierarchy that represent important concept
- D4: Data-Driven Domain Discovery for Structured Datasets
  - An end-to-end, data-driven approach
  - Robust context based signatures
  - Scalable column-based clustering algorithm

[Li et al. KDD 17]

[Ota et al. VLDB 20]

# Table annotation: Toward Making Tables Make More Sense



[Limaye et al. VLDB 10]

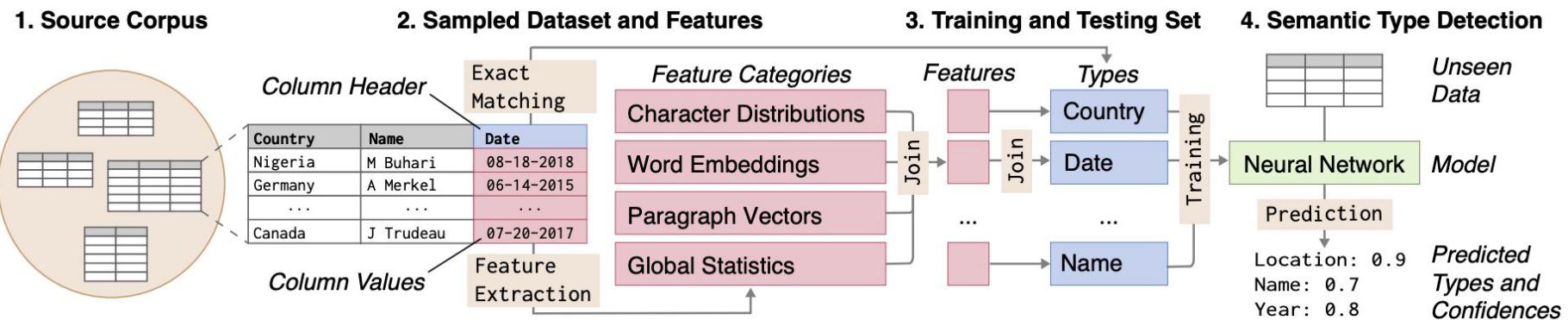
# Sherlock: Feature based approach

**Column type detection:** formulated as multi-class classification (w/ supervised learning)

Input: Raw column values

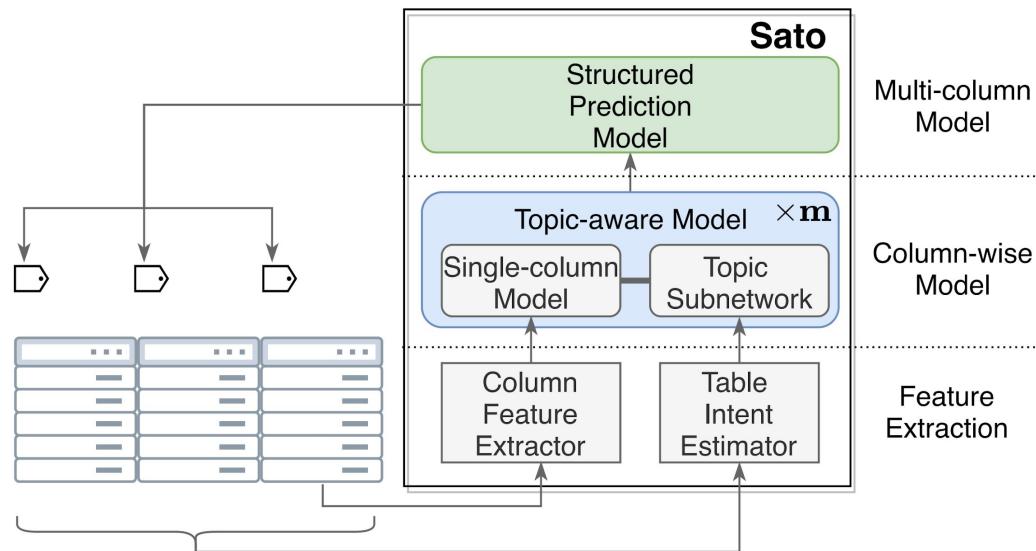
Model: Multi-input Neural Network

Output: Semantic type of column (from 78 candidate classes)

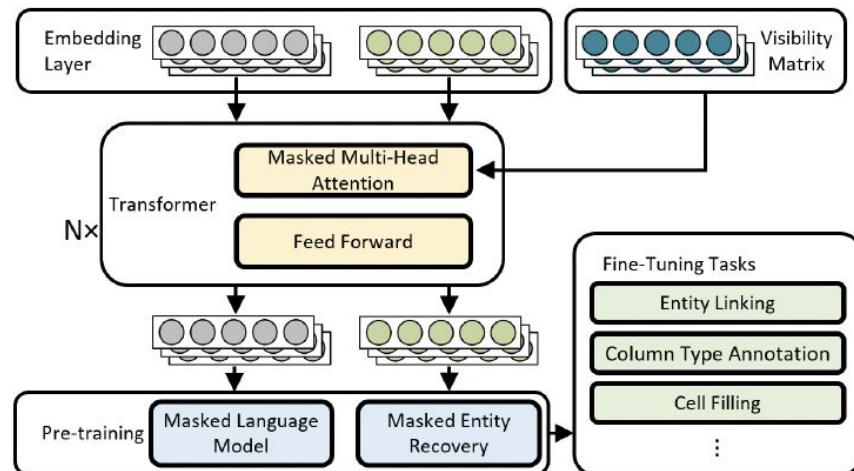


# Sato: Adding column context

- Column-wise + Table-wise features + Structured prediction
- Enhance with Topic models
- Manually crafted features



# TURL: Pre-trained model for tabular data



Transformer-based Architecture



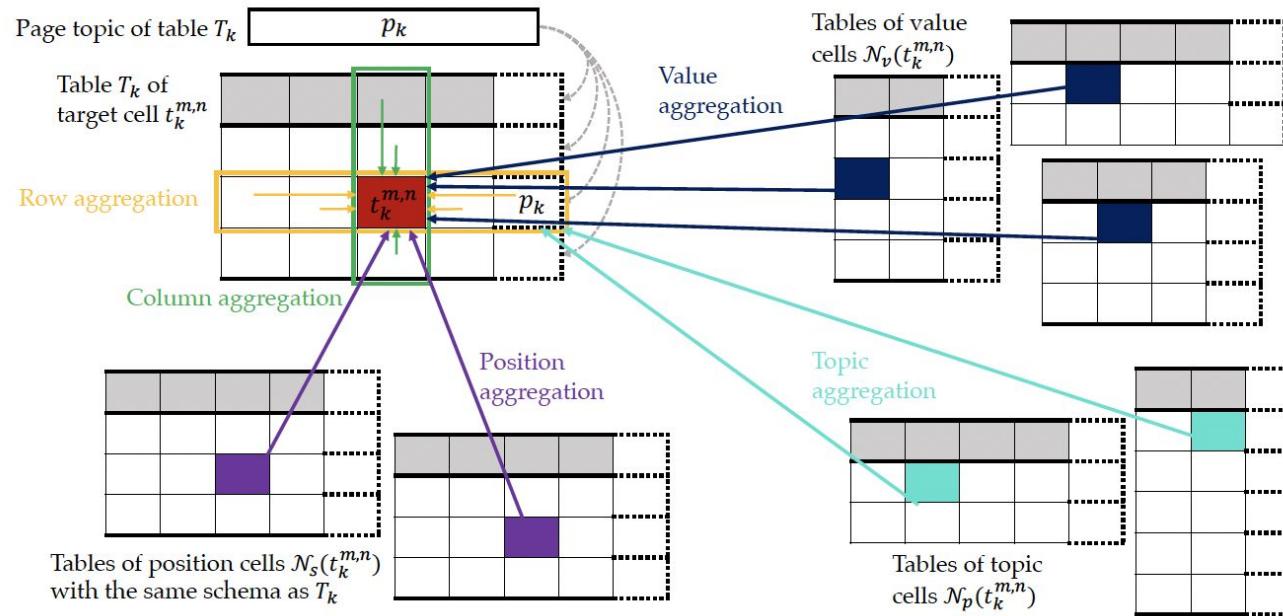
Pre-training with masked self-attention: Each token/entity in a table can only attend to its directly connected neighbors

# TURL: Fine-tuning

- Perform Fine-tuning with task-specific datasets
- Strategies for different categories of tasks
  - Table interpretation
  - Table augmentation

Task	Finetune Strategy
Table Interpretation	<p>Entity Linking</p>
	<p>Column Type Annotation</p>
	<p>Relation Extraction</p>
Table Augmentation	<p>Row Population</p>
	<p>Cell Filling</p>
	<p>Schema Augmentation</p>

# TCN: Table convolution network

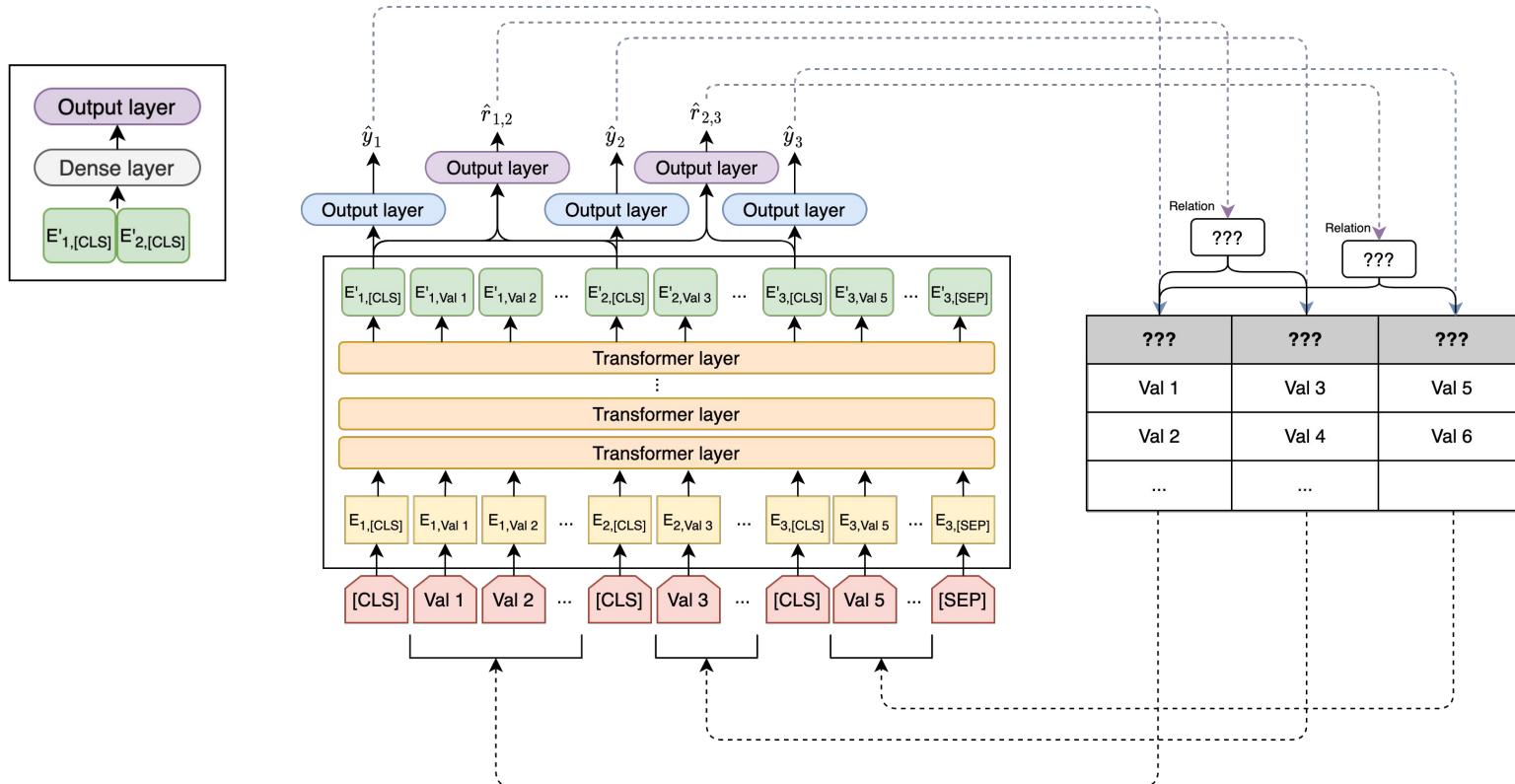


Learn Table representation by taking both intra- and inter-table contexts

# Doduo: Multi-task fine-tuning

[Suhara et al. SIGMOD 22]

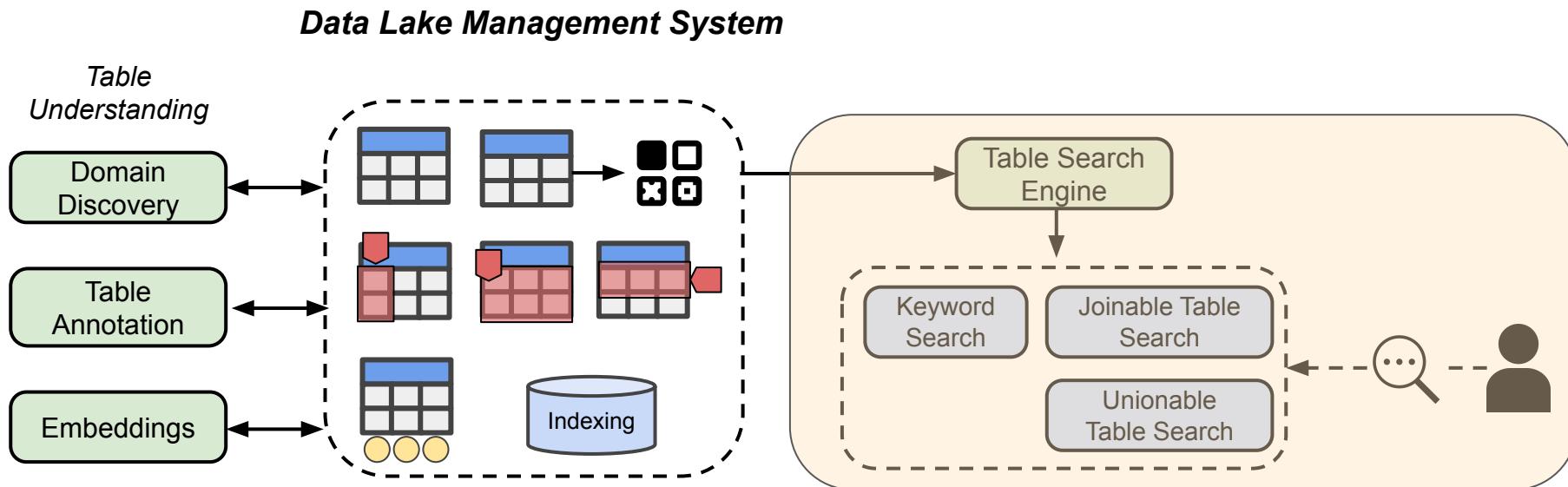
Support both column type detection and relation extraction with pre-trained LMs



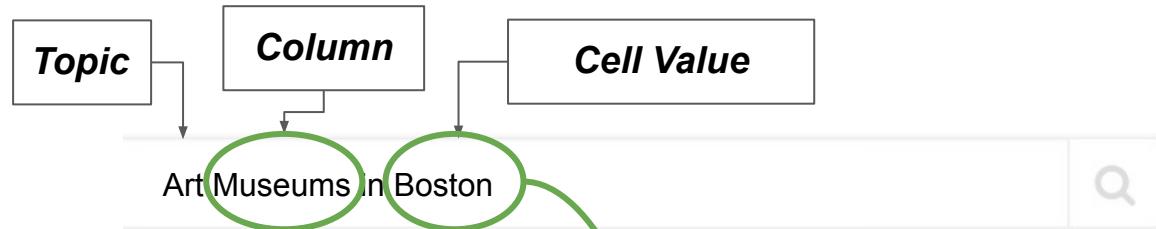
# Outline

- Introduction
- Table Understanding
- **Table Search Engine**
- Table Navigation and Exploration
- Data Science and Application Support
- Conclusion / Future work

# Table Search Engine



# Table Discovery using Keyword Search



[Pimplikar, Sarawagi PVLDB 12]

Octopus  
[Cafarella et al. PVLDB 09]

Goods  
[Brickley et al. WWW 19]

Museum	City	Zip Code
MFA	Boston	02115
MOMA	New York City	10019
Art Institute	Chicago	60603
...	...	...

# Tables-as-Queries for Table Discovery

- In addition to searching for tables using keyword queries, there are systems that allow users to input **tables** as queries to find relevant tables to their input tables

E.g.

Joinable Table Search

Unionable Table Search

# Finding Related Tables

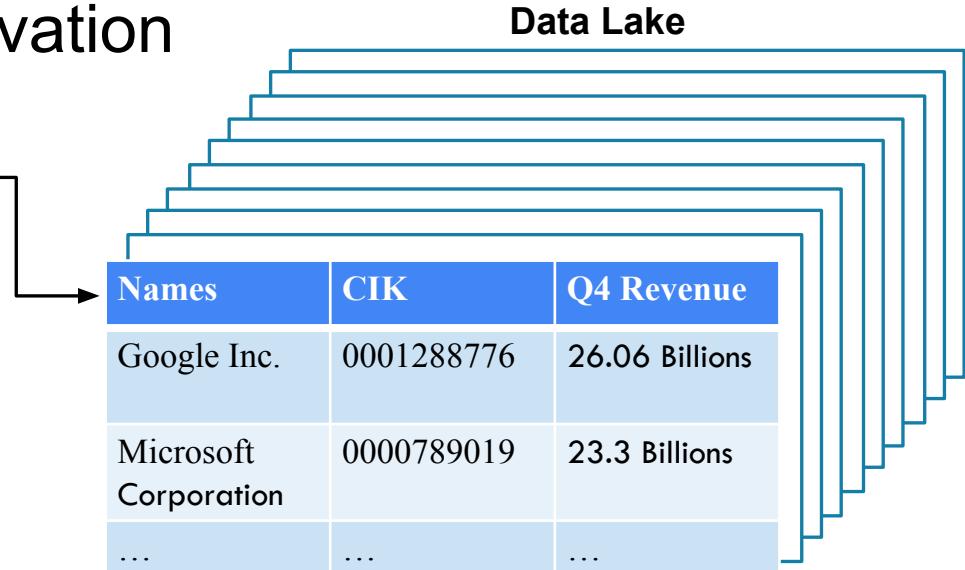
- Given a table corpus and a query table, return a ranked list of relevant tables
- Schema Complements are found using schema matching and are returned as joinable tables with additional attributes
  - Joinable Table Search
- Entity Complements share a subject attribute and similar schemas
  - Unionable Table Search

# Joinable Table Search: Motivation

↓

**Query Table**

Company	Headquarter	CEO
Google Inc.	Mountain View	Sundar Pichai
Databricks	San Francisco	Ali Ghodsi
Microsoft Corporation	Redmond	Satya Nadella
...	...	...



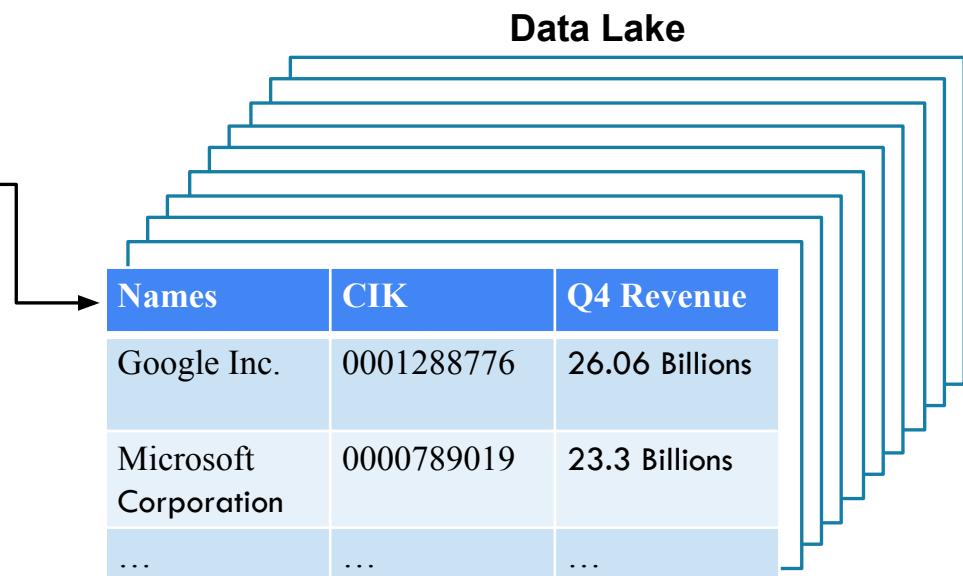
**Data Science Question:** How can I find more features for my model of corporate health?

**Data Management Task:** Find tables that can be joined with a query table.

# Infogather

**Query Table**

Company	Headquarter	CEO
Google Inc.	Mountain View	Sundar Pichai
Databricks	San Francisco	Ali Ghodsi
Microsoft Corporation	Redmond	Satya Nadella
...	...	...



*Uses schema matching to find attributes that join  
Company matches Names*

# Set Containment

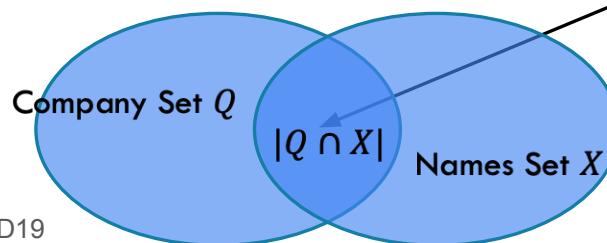
**Query Table**

Company	Headquarter	CEO
Google Inc.	Mountain View	Sundar Pichai
Databricks	San Francisco	Ali Ghodsi
Microsoft Corporation	Redmond	Satya Nadella
...	...	...

**Data Lake**

Names	CIK	Q4 Revenue
Google Inc.	0001288776	26.06 Billions
Microsoft Corporation	0000789019	23.3 Billions
...	...	...

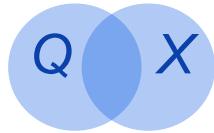
**Query Column**



Number of distinct values in the query that are joinable – **Overlap**

# LSH Ensemble: Approximate Containment Join Search

$$Jaccard(Q, X) = \frac{|Q \cap X|}{|Q \cup X|}$$



**Jaccard( $Q, X$ ) >> Jaccard( $Q, X'$ )**

Same intersection size,  
but the Jaccard similarity is much smaller on the right

*Jaccard is biased to small sets*  
[Fernandez et al. ICDE18], and others

$$Containment(Q, X) = \frac{|Q \cap X|}{|Q|}$$



**Containment( $Q, X$ ) = Containment( $Q, X'$ )**

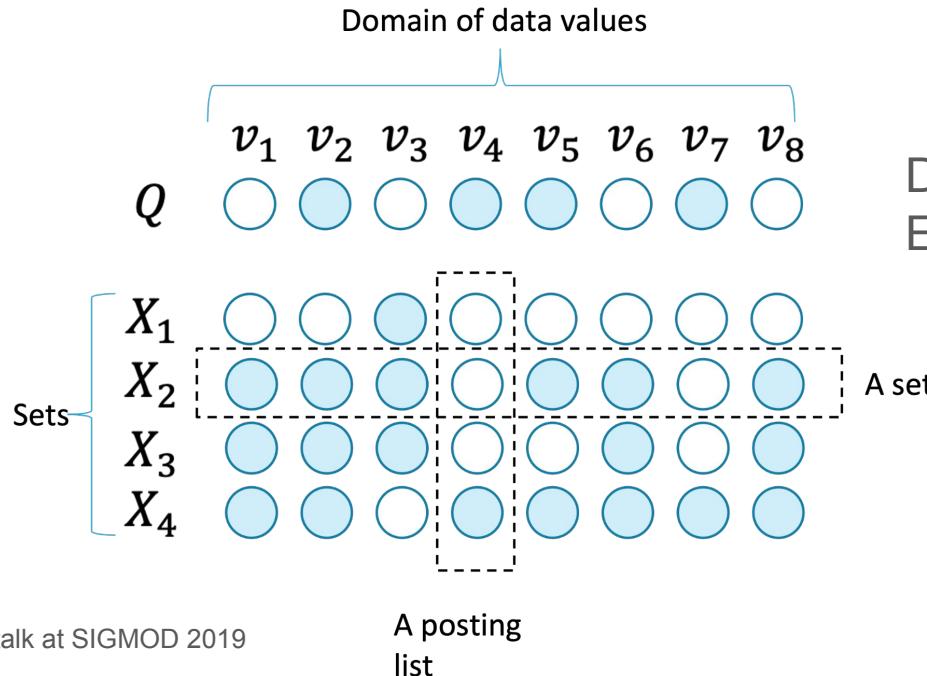
Containment is the same for both,  
independent of the size of  $X$  and  $X'$



**LSH Ensemble: Uses MinHash LSH for  
approximate, scalable set containment search**

# JOSIE: Exact Containment Join Search

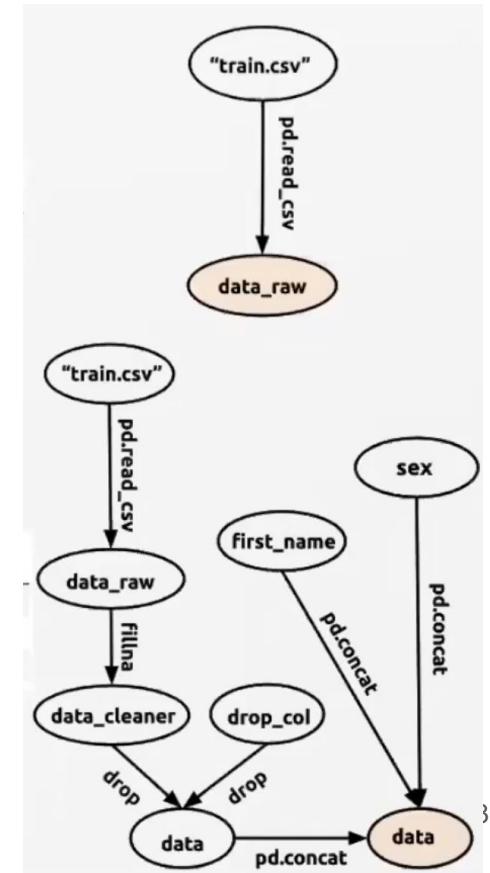
- Performs exact top-k joinable table search using set containment
- Scales inverted index to massive queries and massive number of posting lists



Data-driven cost model  
 Experiments on data lakes with  
 Posting lists: .5 Billion  
 Queries: upto 10K in size

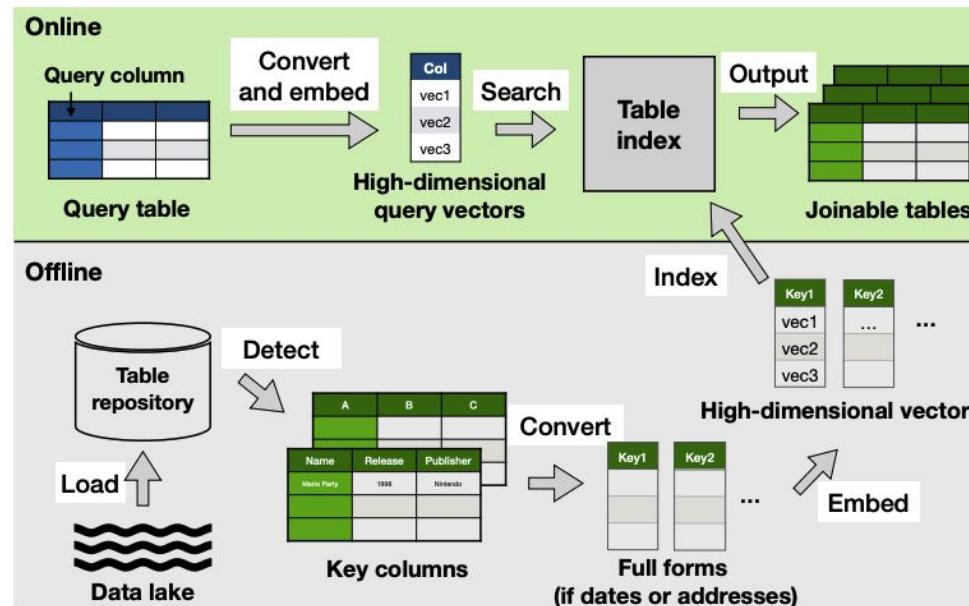
# Juneau: Joinable Search in Jupyter Notebooks

- Joinable Table Search in Jupyter Notebooks
- Create data profiles using data values and domains



# PEXESO: Fuzzy-Join Search

Use word embeddings to retrieve tables that can be fuzzy-joined



# A Sketch-based Index for Correlated Dataset Search

- Find top- $k$  joinable tables that also contain numerical columns correlated with numerical column(s) in the Query Table
- Employs a QCR hashing scheme to estimate correlations between numerical values

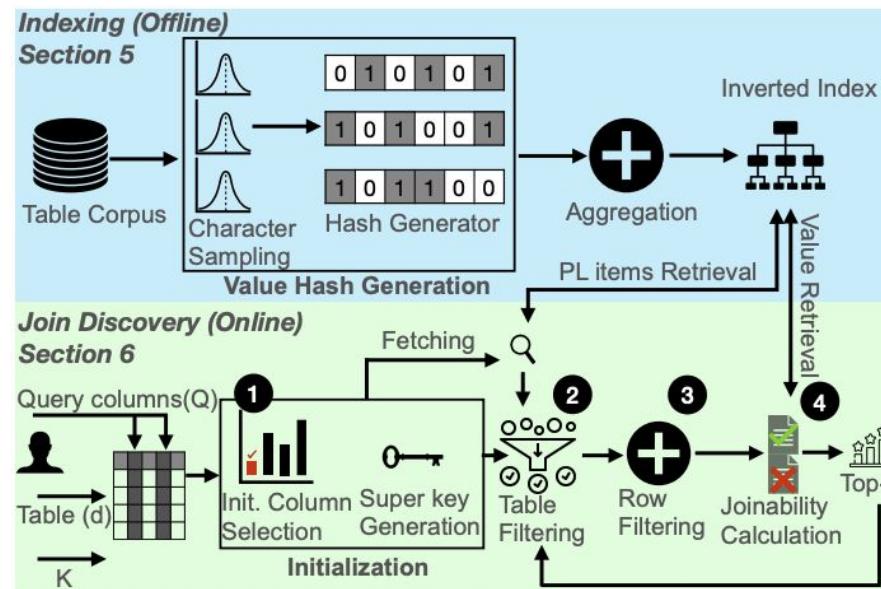
## WarpGate

[Cong et al. CIDR 23]

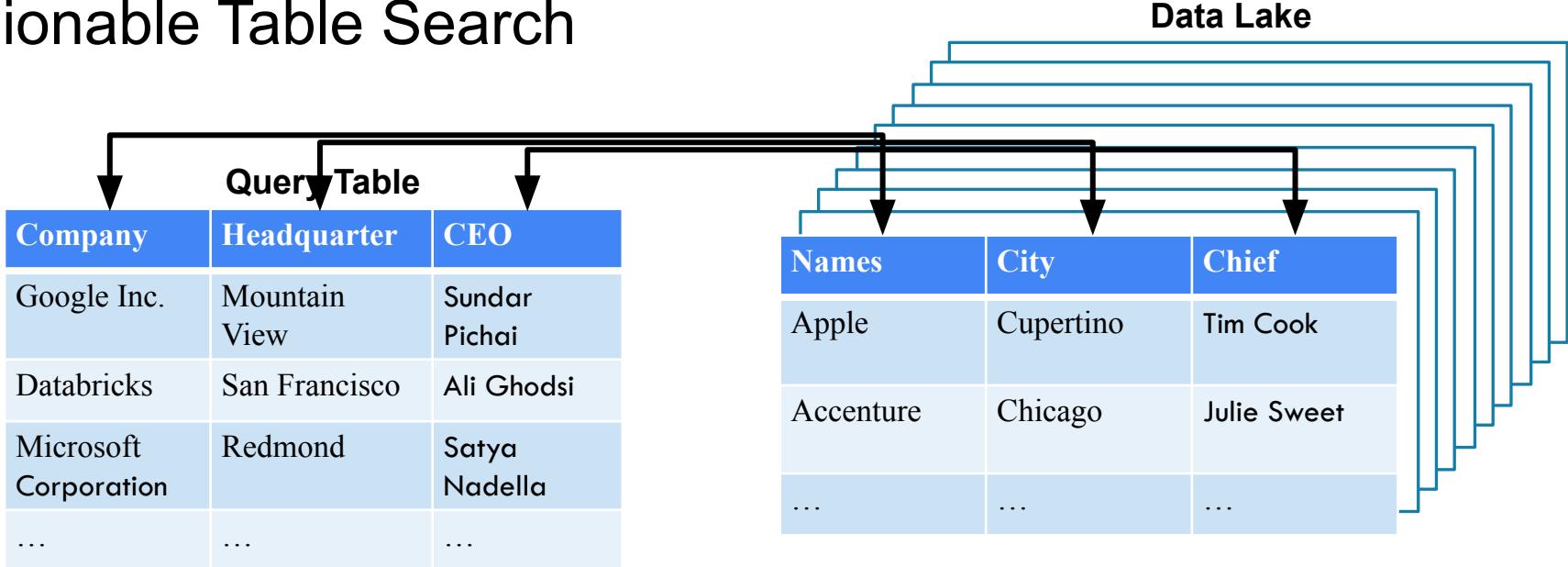
- Retrieves top- $k$  semantically joinable tables
- Use embeddings to capture semantic relationships

# MATE: Multi-Attribute Join

- Supports joins over multiple attributes
- Hash each row value into a super key



# Unionable Table Search



**Data Science Question:** Does my analysis generalize? To new Companies or new regions?

**Data Management Task:** Find tables that can be unioned with a query table.

# Table Union Search

[Nargesian et al. VLDB 18]

- Given a set of tables and a Query Table, return the top- $k$  tables with the highest table unionability
  - Some **attributes** may overlap
  - Some may refer to **entities** of a common type
  - Some may use **semantically similar words**

The diagram illustrates the search process for a query table and candidate tables. The Query Table contains data for various energy sources and their values. The Candidate Table contains data for different locations and fuel types. Colored arrows (blue, yellow, and red) highlight semantic similarities between attributes across the tables.

Electricity	Barnett	Domestic	240.99	...
Gas	Brent	Transport	164.44	
Coal	Camden	Transport	134.90	
Railways diesel	City of London	Domestic	10.52	
Gas	Brent	Domestic	169.69	
Coal	Brent	Transport	120.01	

Query Table

Benton	Transport	Gasoline	64413	62.9
Kittitas	Hydro	Fuel oil (1,2,...)	12838	66.0
Grays Harbor	Domestic	Aviation fuels	117039 3	66.1
Skagit	Transport	Liquified petroleum	59516	60.1

Candidate Table

[Bogatu et al. ICDE 20]

- D3L: Extends TUS to also consider formatting similarity and schema similarity

# SANTOS

Data Scientist's table (Query Table)

Park Name	Supervisor	City	Country
River Park	Vera Onate	Fresno	USA
West Lawn Park	Paul Veliotis	Chicago	USA
-----	-----	-----	-----

(a) A table about parks

Data lake tables (Candidate unionable Tables)

Park Name	Film Title	Park Location	Park Phone	Park City	Film Director	Film Studio
Chippewa Park	Bee Movie	6748 N. Sacramento Ave.	773 731-0380	Cook	Simon J. Smith	Dreamworks
Lawler Park	Coco	5210 W. 64 <sup>th</sup> St.	773 284-7328	Riverside	Adrian Molina	Pixar
-----	-----	-----	-----	-----	-----	-----

(b) A table about films shown in different parks

Person	Occupation	Birthplace	Country	Park Name
James Taylor	Singer	Boston	USA	Central Park
Anthony Pelissier	Film Director	Barnet	UK	Cairngorms National Park
Akram Afif	Football Player	Doha	Qatar	Aspire Park
Ivan A. Getting	Physicist	NYC	USA	El Segundo Park
Abby May	Social Worker	Boston	USA	Fenway Park
Stevie Ray Vaughan	Singer	Texas	USA	Chastain Park

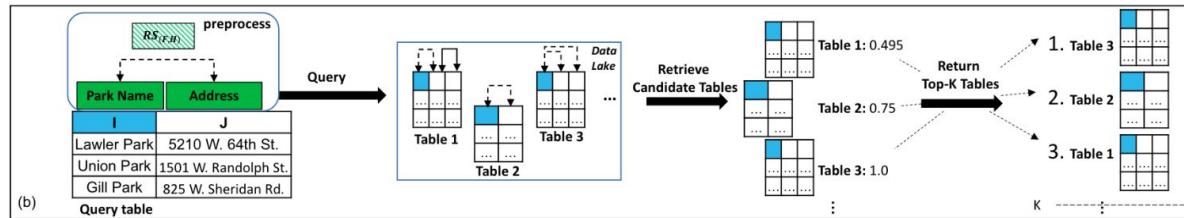
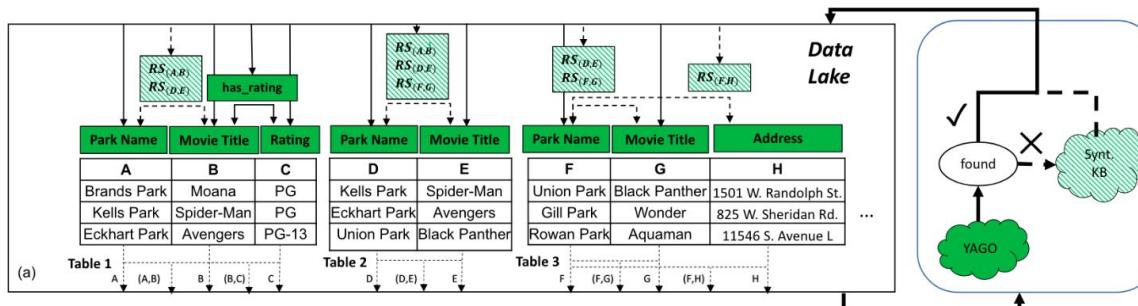
(c) A table about people

Along with column semantics,

consider the **binary relationships** between the column pairs.

# SANTOS

- Semantic Table Union Search: considers both attribute semantics *and binary Relationship semantics between attributes*
  - Existing Knowledge Base
  - KB has limited coverage → construct a data-driven synthesized Knowledge Base



# Starmie

Table A:

Name	Mode of Travel	Purpose	Destination	Day	Month	Year	Expense
Philip Duffy	Air	Regional Meeting	London	10	April	2019	189.06
Jeremy Oppenheim	Taxi	Exchange Visit	Ottawa	30	Jul	2019	8.08
Mark Sedwill	Air	Evening Meal	Bristol	02	September	2019	50

Table B:

Name	Date	Destination	Purpose
Clark	23/07	France	Discuss EU
Gyimah	03/09	Belgium	Build Relations
Harrington	05/08	China	Discuss Productivity

Table C:

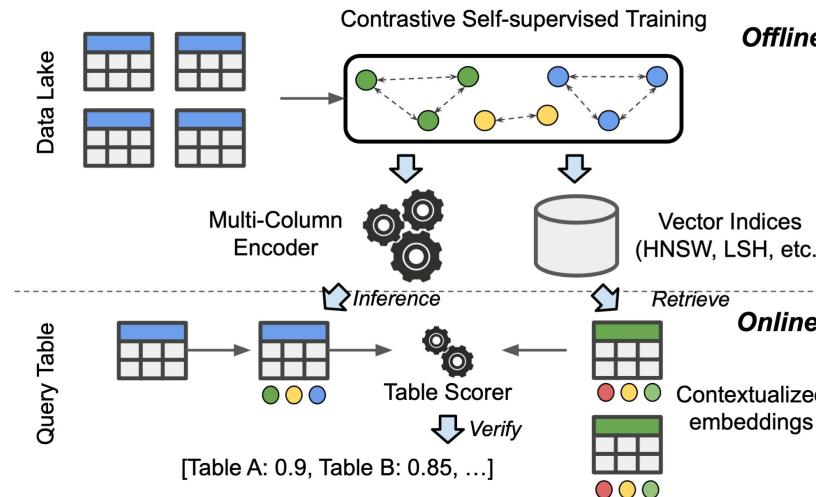
Bird Name	Scientific Name	Date	Location
Pine Siskin	Carduelis Pinus	2019	Ottawa
American Robin	Turdus migratorius	2019	Ottawa
Northern Flicker	Colaptes auratus	2019	London

- Only Column Unionability → wrongly aligned columns
- Also Using Binary relationship between columns → influenced by value overlap

Use the context from the entire table to capture the semantics

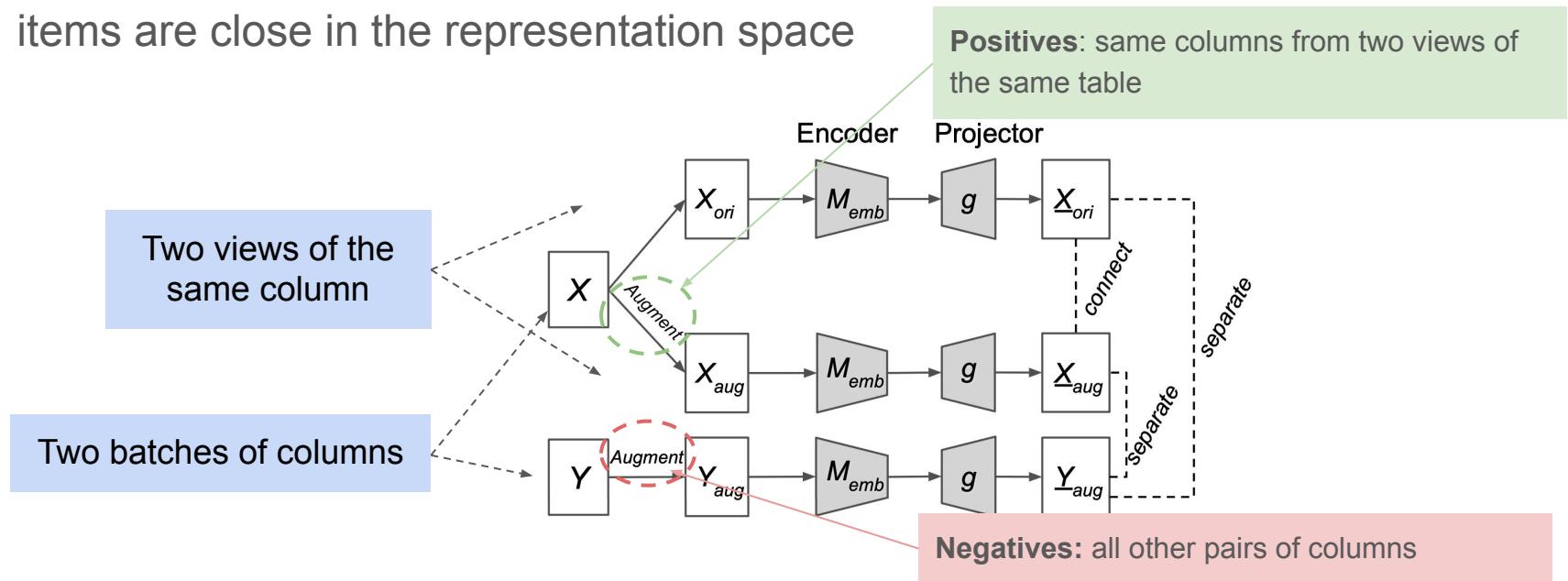
# Starmie

- Semantic Table Union Search: Extend column semantics and binary relationship semantics to *use the entire table context*
- Natural language approach: Language Model uses the table context to capture and encode the semantics of each column



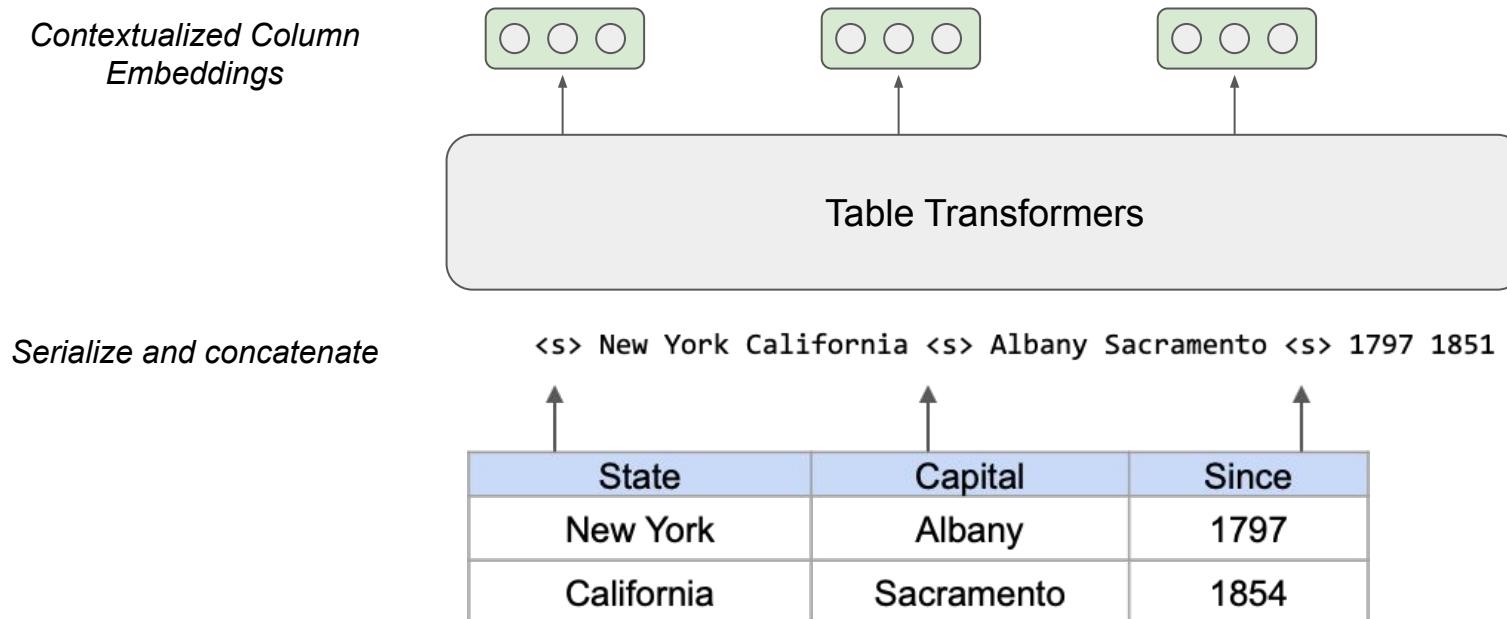
# Starmie

- Difficult to manually label training data → Self-Supervised Technique
- Contrastive Learning: Learn representations (w/o labels) such that similar items are close in the representation space



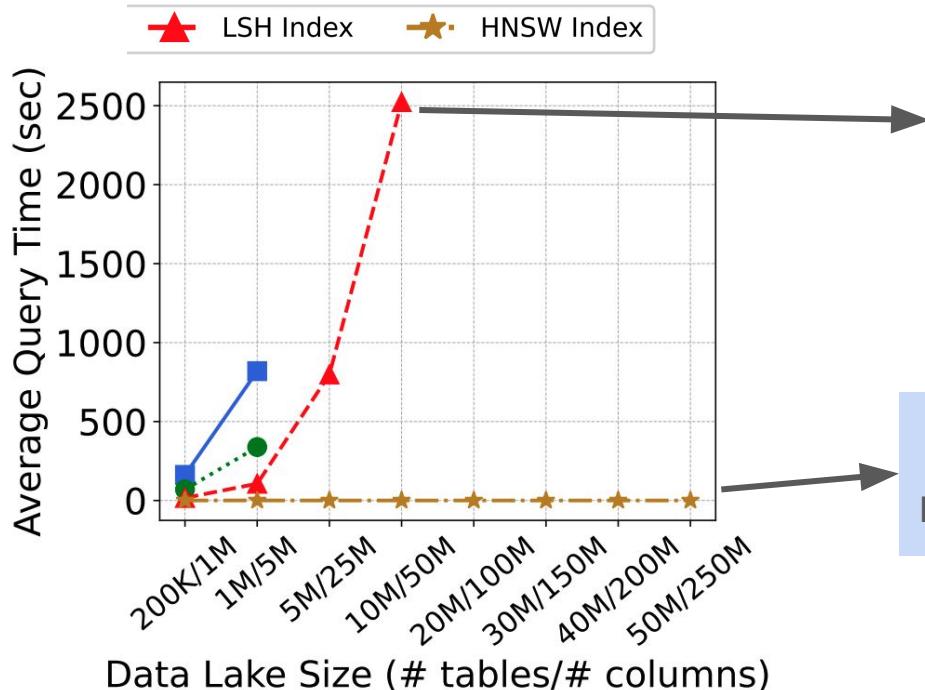
# Starmie

- Multi-column table models: Preserve the **context** of the column during encoding



# Starmie

- Hierarchical Navigable Small Worlds (**HNSW**) index: graph-based indexing
- LSH Index: the state-of-the-art hash-based indexing



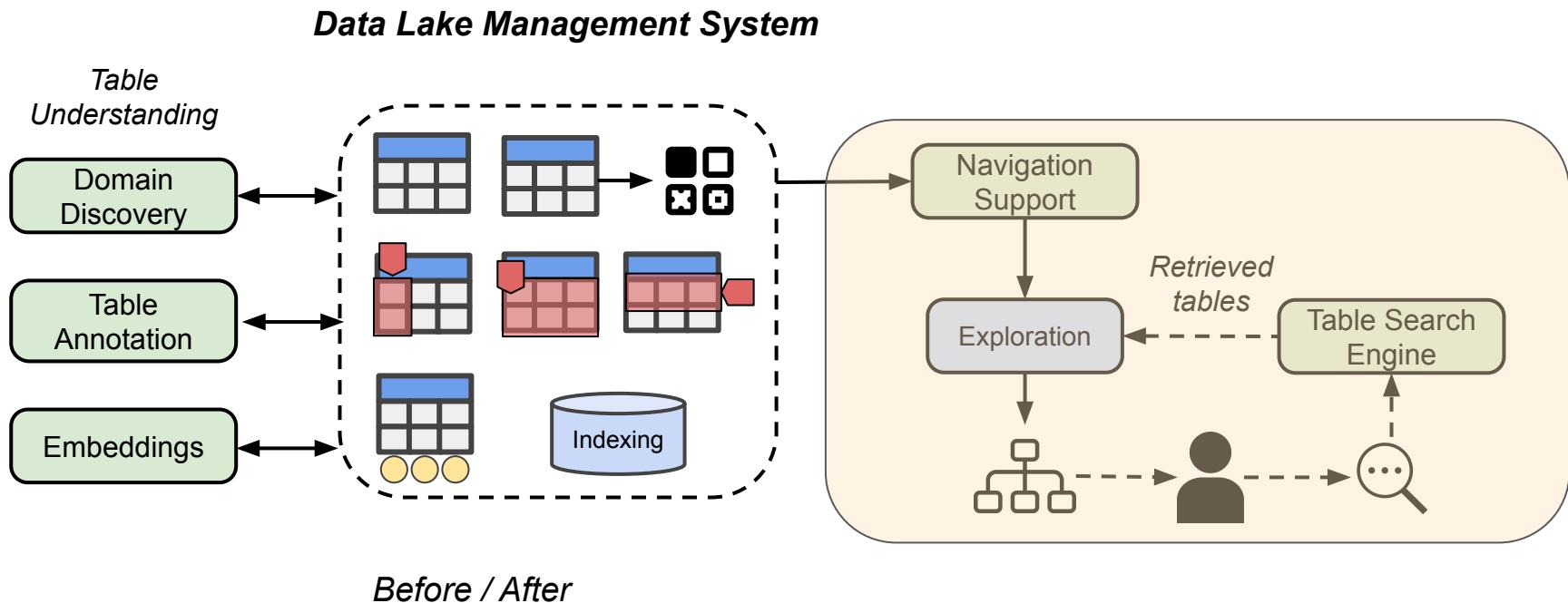
LSH times out after the data lake size exceeds 10M tables

HNSW index has a query time of 60 ms, even with a data lake of 50M tables

# Outline

- Introduction
- Table Understanding
- Table Search Engine
- **Table Navigation and Exploration**
- Data Science and Application Support
- Conclusion / Future work

# Table Navigation and Exploration



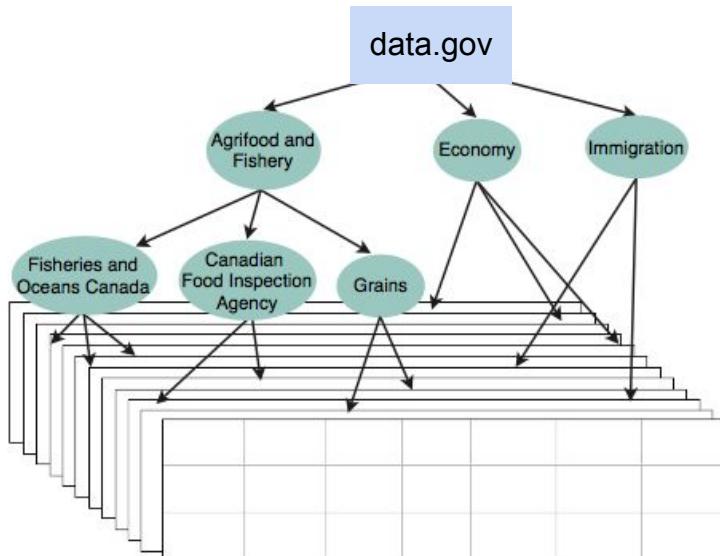
**Complementary to search:** present the data lake in a navigation-friendly structure to the user

# Table Navigation / Exploration

[Nargesian et al. SIGMOD 20]

## Problem Definition:

- Automatically building a directory structure that enables a user to most efficiently find tables of interest in a data lake.

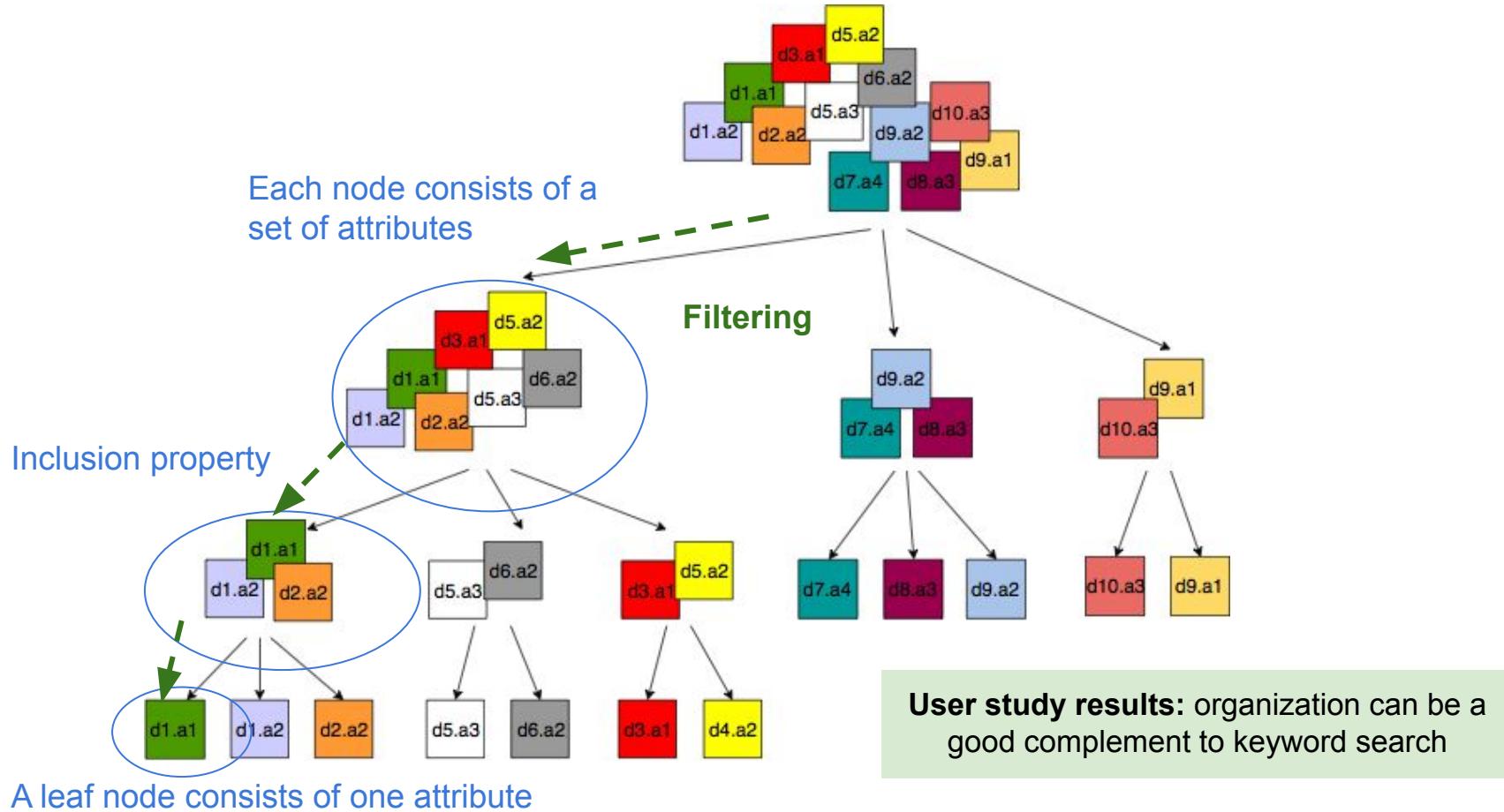


## Example properties to be optimized:

- reasonable number of choices at each step.
- reasonable number of steps to find a dataset.
- unambiguous choices at each step.

# Data Lake Organization: a DAG

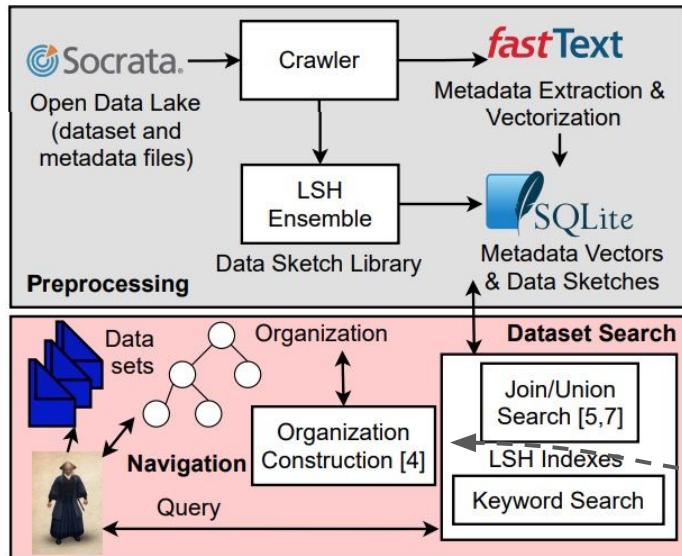
[Nargesian et al. SIGMOD 20]



# RONIN: Data Lake Exploration

[Ouellette et al. VLDB 21]

- A follow-up demo system that integrates navigation and search over data lakes
- Support different kinds of search: keyword, unionable and joinable



**Goal:** Enables the users to seamlessly switch between two operations

The navigation organization is constructed on-the-fly as the search result is generated

# An alternative organization: Enterprise KG

## Aurum: A Data Discover System

[Fernandez et al. ICDE 18a]

- **EKG:** a hypergraph that captures relationships between columns of datasets
- Discussed **building, maintaining, and querying** of EKG

Algorithms to decide when to re-index a column

a graph query language SRQL

[Fernandez et al. ICDE 18b]

## Sleeping Semantics: Linking Datasets using Embeddings

- Studied how to use aggregated embeddings (e.g., GloVe) to identify column/dataset relationships

We have better tools for this problem today

# Dataset Search Engine: Auctus

[Castelo et al. VLDB 21]

keyword search

Auctus

Advanced Search: Any Date, Any Location, Related File, Source, Data Type

Temporal   
Start: 01/01/2016 End: 03/17/2021 Resolution: Any Resolution

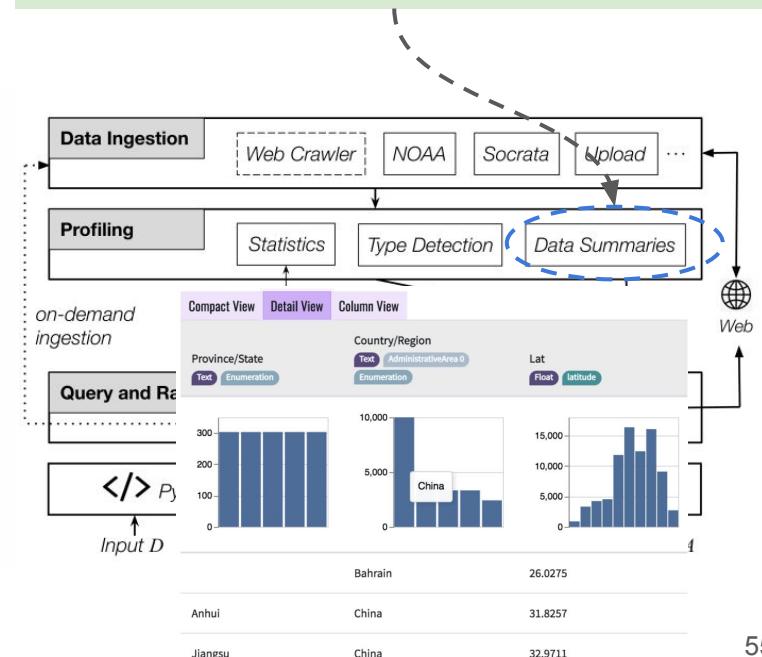
Geo-Spatial   
Search Map

Left-click to start selection. Right-click to clear selection.

Bottom Left: 40° 25' 06" N 74° 21' 44" W Top Right: 41° 12' 18" N 73° 38' 32" W

© OpenStreetMap contributors.

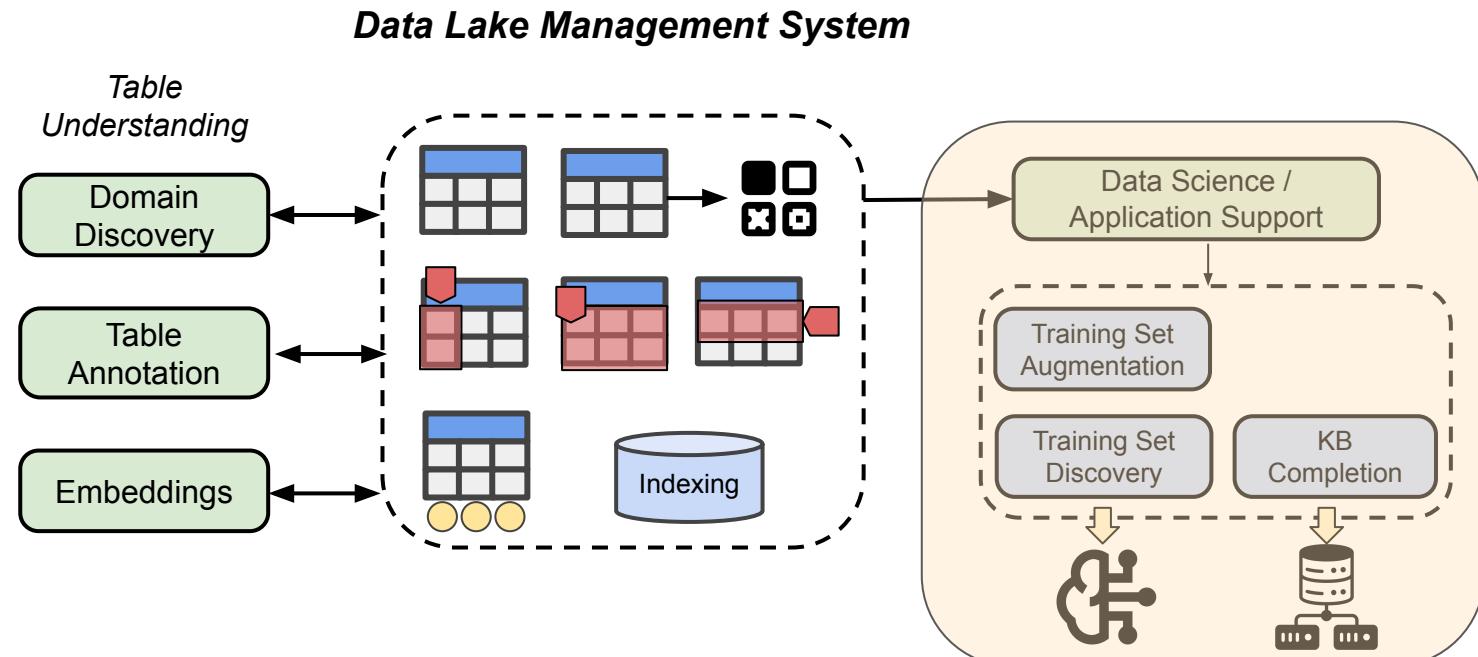
Compact dataset summaries allow users to quickly navigate through the search results



# Outline

- Introduction
- Table Understanding
- Table Search Engine
- Table Navigation and Exploration
- **Data Science and Application Support**
- Conclusion / Future work

# Table Discovery for Data Science / Application Support



Data/Table discovery is the key to success for many **data-hungry applications**, especially ML

# Table Discovery for Supporting Data Science Applications

Training Set  
Augmentation

**ML:** augment a target table (for regression/classification) with features learned or extracted from data lake tables

Training Set  
Discovery

**ML:** create novel regression/classification tasks using joins/unions of data lake tables

KB Completion

**KB:** combine webtables to form structured records for KB creation and/or completion

# Case study from Starmie: Data Discovery for ML tasks

- **Goal:** find joinable tables from WDC to improve downstream ML tasks
- **Query tables / ML tasks:** 25 tables of  $\geq 200$  rows, target column = “Rating”
- **Data lake tables:** 4,130 VizNet tables of  $\geq 50$  rows

State	Office	District	Name	Party	Rating
AZ	U.S. House	8	Trent Franks	Republican	95.0
TX	U.S. House	3	Sam Johnson	Republican	95.0
OH	U.S. House	4	Jim Jordan	Republican	95.0
CO	U.S. House	5	Doug Lamborn	Republican	95.0

**Goal:** find useful features for predicting the “Rating”

# Case study from Starmie: Data Discovery for ML tasks

- **Baselines:** NoJoin, Jaccard, Overlap
- **Starmie:** contextualized embeddings for join+target columns

	NoJoin	Jaccard	Overlap	Starmie
Avg. MSE	0.0820	0.0753	0.0748	<b>0.0699</b>
Improvement	-	8.23%	8.82%	<b>14.75%</b>
#improved	-	13 / 25	12 / 25	<b>15 / 25</b>
avg. Improve	-	14.74%	14.05%	<b>20.64%</b>

**Findings:** learned table representations help find more meaningful joinable tables for improving downstream ML tasks

# Example output (Starmie)

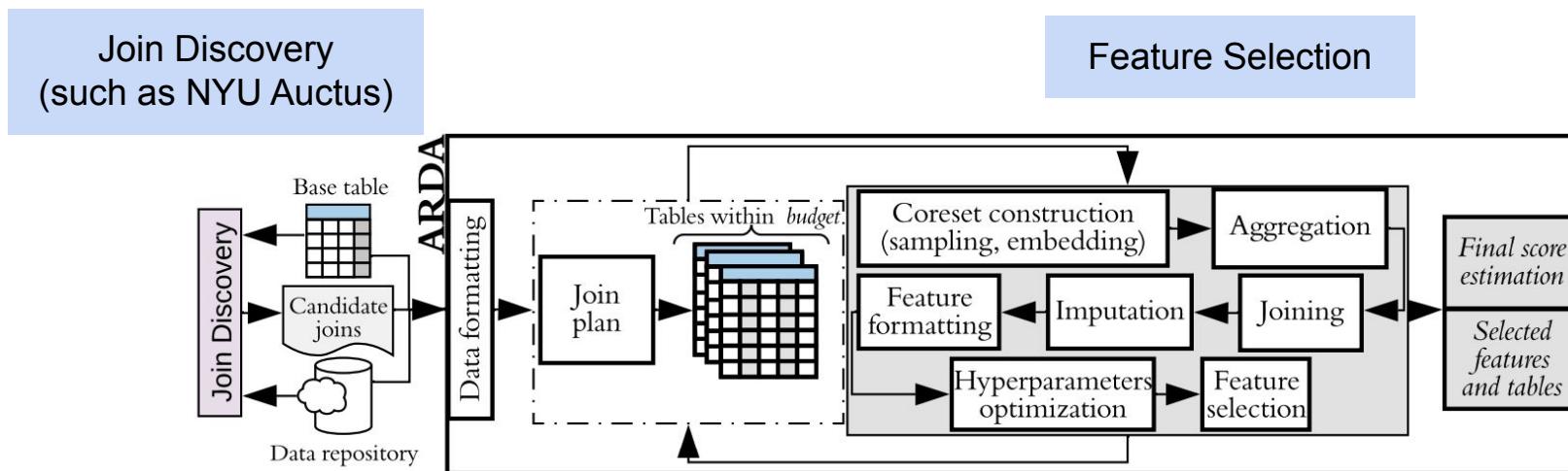
Joining with a table about money raised on the “name” column

State	Office	District	Name	Party	Rating	Party_r	State_r	\$ From Interest Groups That Supported	\$ From Interest Groups That Opposed	Vote
AZ	U.S. House	8	Trent Franks	Republican	95.0	R	AZ-2	\$12,000	\$0	Yes
TX	U.S. House	3	Sam Johnson	Republican	95.0	R	TX-3	\$4,000	\$0	Yes
OH	U.S. House	4	Jim Jordan	Republican	95.0	R	OH-4	\$22,000	\$0	Yes
CO	U.S. House	5	Doug Lamborn	Republican	95.0	R	CO-5	\$2,000	\$0	Yes

Model performance (MSE) NoJoin: **0.1598** → Jaccard: **0.1544** → CL: **0.1195**

# Training set augmentation: ARDA

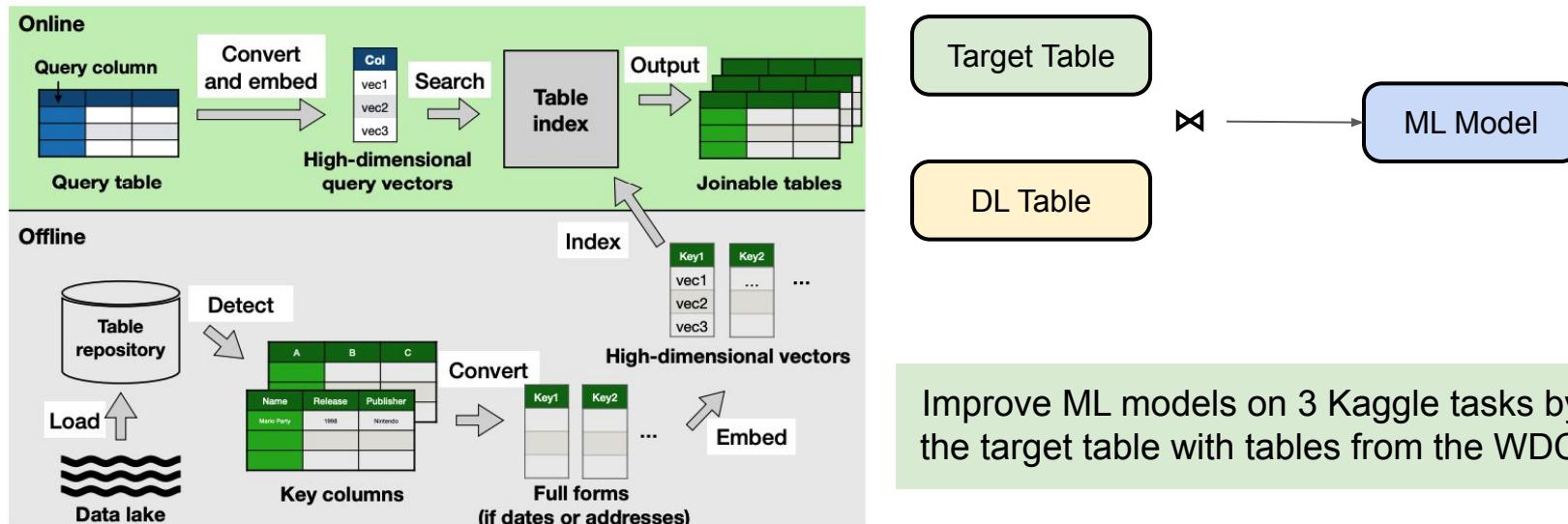
- **Goal:** input a dataset and a data repository, and outputs an augmented dataset => improved predictive performance.



Consistently improved regression models (e.g., RF) across 5 real-world tasks from DARPA D3M

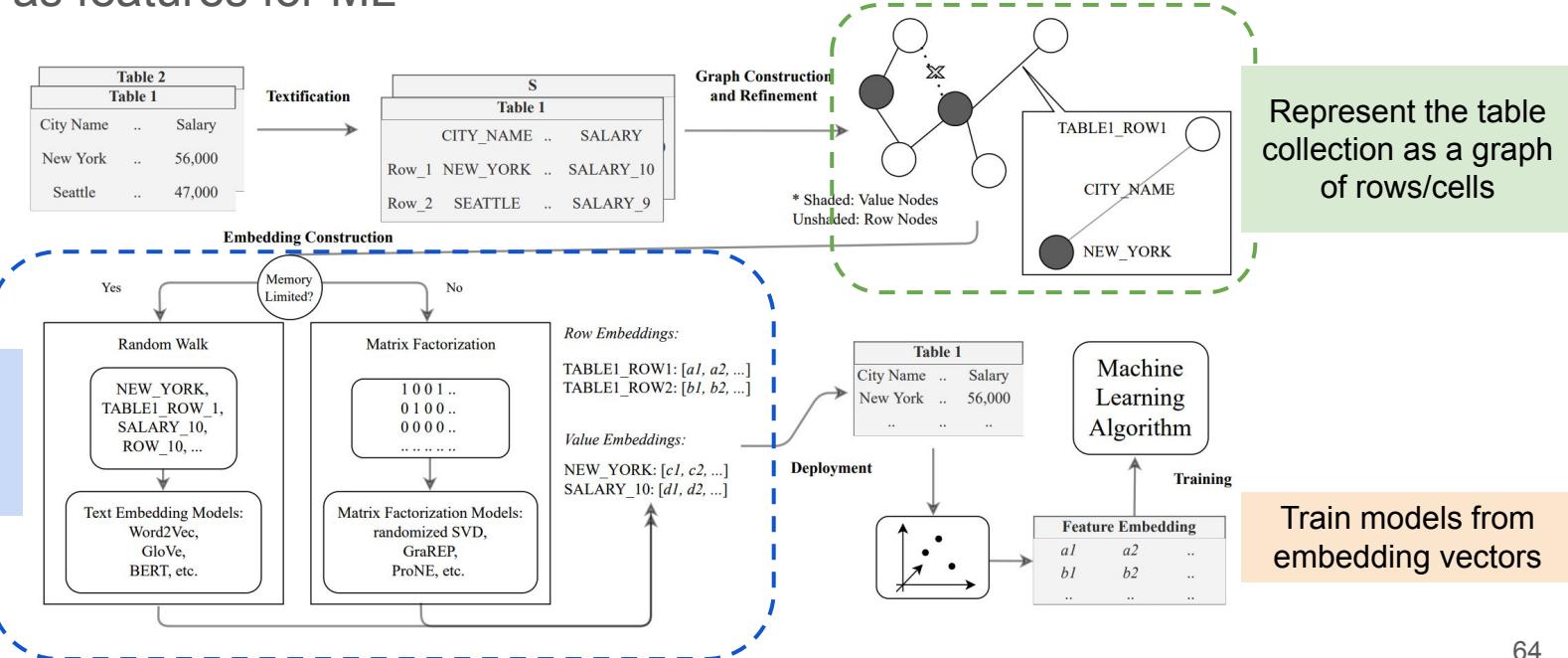
# Recap: Augment by finding joinable tables - PEXESO

- High-dimensional similarity search based on averaging GloVe embeddings



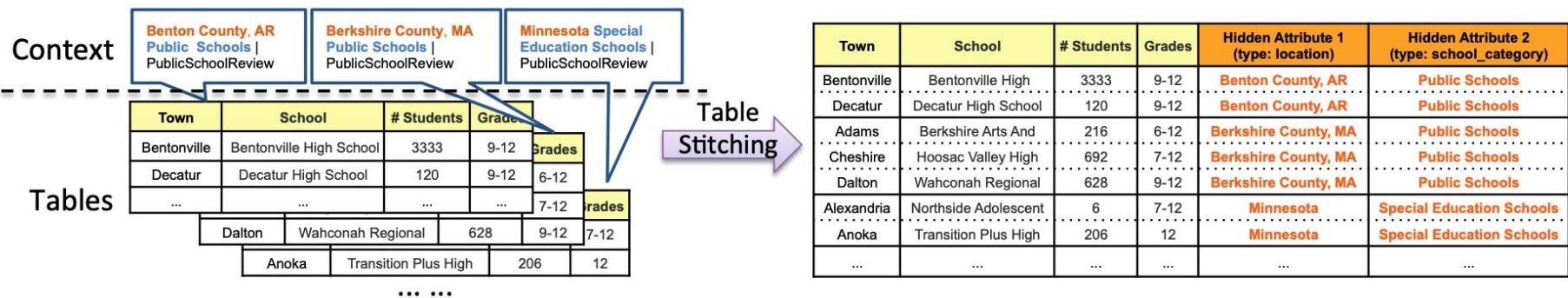
# Learning relational embeddings: Leva

- Instead of joining tables, use vector representations learned from the data corpus as features for ML



# KB Completion: Table Stitching

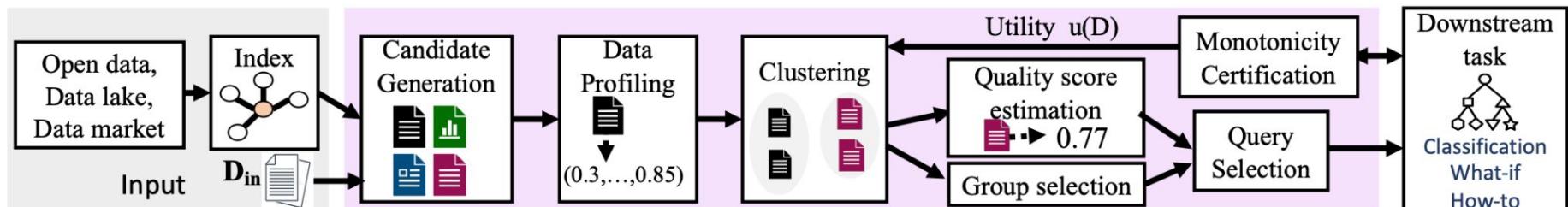
- Goal: discover unionable tables (with context) to match/augment with a KB



- **Table stitching [Ling et al. IJCAI 13]:** identify webtables of the same schema from the same sites + synthesizing new columns from page context
- **[Lehmberg and Bizer, VLDB 17]:** experimental study on stitching tables + KB alignment

# METAM: Goal-Oriented Data Discovery

- **Goal:** build a generic framework for querying the downstream task with a candidate dataset to automatically steer the discovery and augmentation process

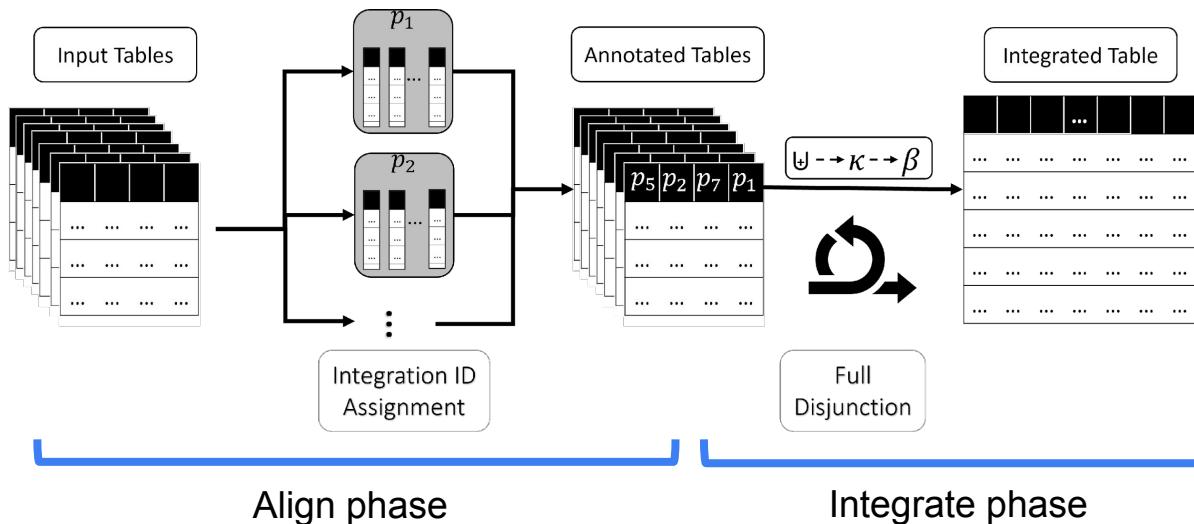


Profile the datasets using (1) **data stats**, (2) **utility function**, and (3) **solution set size**

Adapt to a range of downstream tasks:  
classification, regression, causal inference

# Integrating Data Lake Tables

- **ALITE** (Align and Integrate)



- **Align:** Identify the matching columns across the set of tables and annotate them with a dummy column header.
- **Integrate:** Apply a novel algorithm for Full Disjunction that scales better than prior work. [Galindo-Legaria SIGMOD 94]

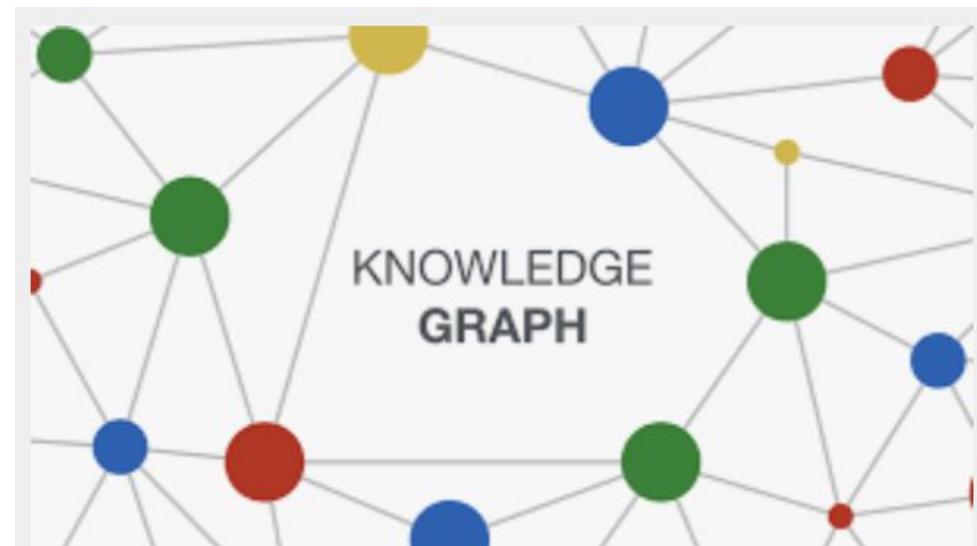
# Outline

- Introduction
- Table Understanding
- Table Search Engine
- Table Navigation and Exploration
- Data Science and Application Support
- Conclusion / Future work

# Theory of how to use LM and KB Symbiotically

**LLM** & **KB** invaluable tools for Table Discovery & Semantics Recovery

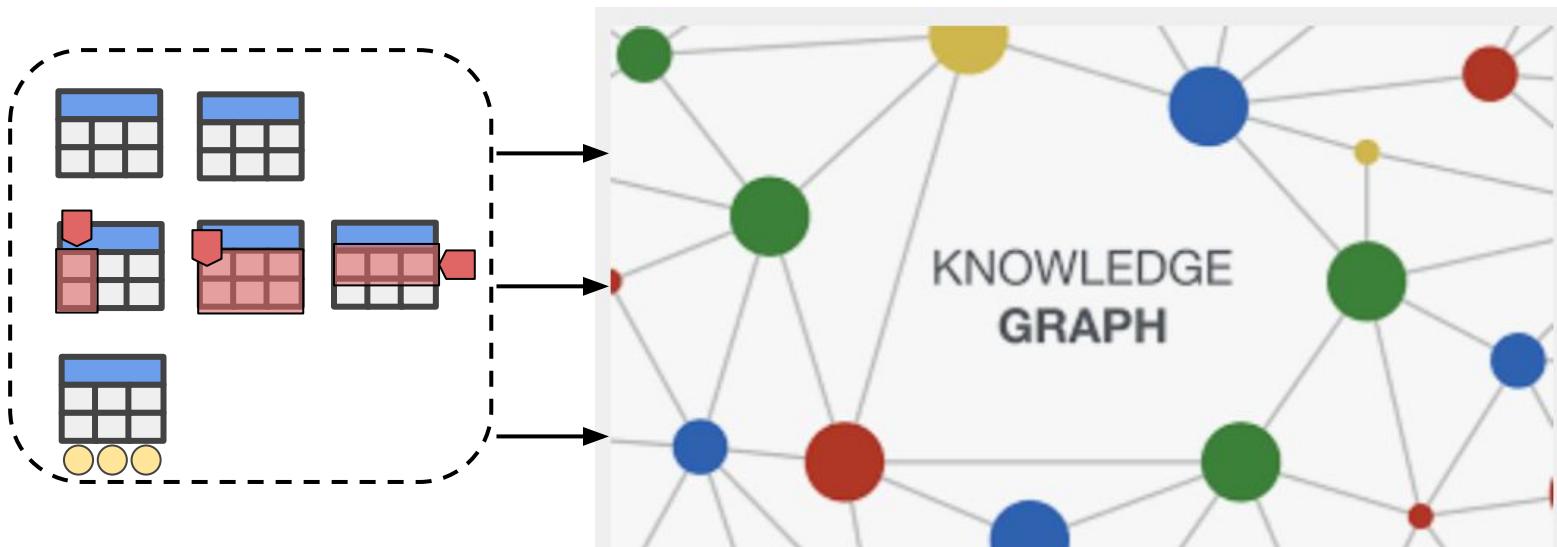
Predation or Symbiosis?



# Data Lake as KB

## KB augmentation or completion

*Judiciously* map tuples, columns, tables to KB to augment knowledge



[Zhu et al. PVLDB 16]  
[Zhu et al. SIGMOD 19]  
[Santos et al. ICDE 22]  
[Fan et al. PVLDB 23]

# Theory of Data Lake Indexing

For DBMS/DW have Periodic Table of Data Structures

- For one (or few) tables
- Goal speed query or update processing

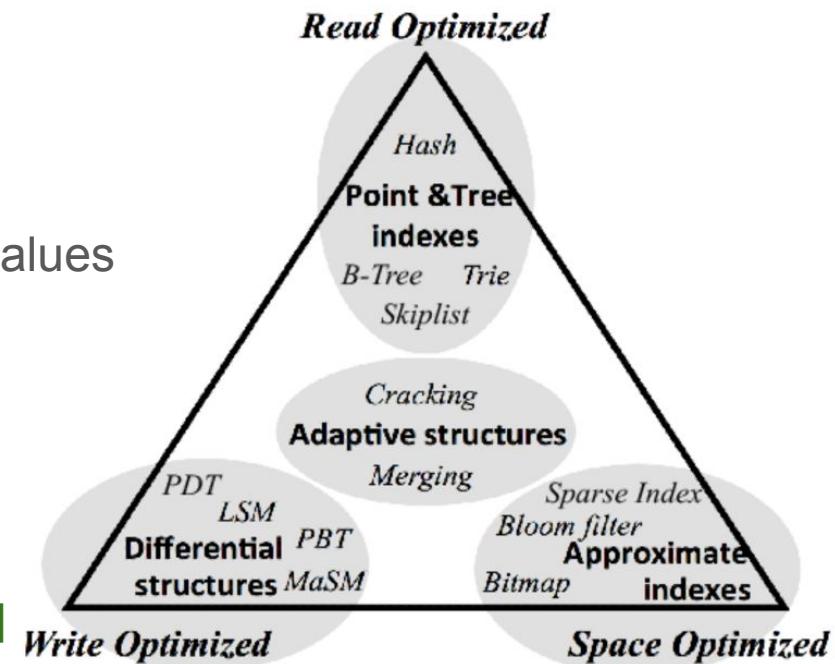
What is equivalent for Data Lakes?

- Index millions of tables, billions+ columns/values
- Goal is accurate (and fast) table discovery

Inverted indices, LSH, HNSW, Sketch-based...

- All required innovation for real data lakes

[Idreos et al. DEBull18]



# Batch vs. Query-Time

## Table Understanding

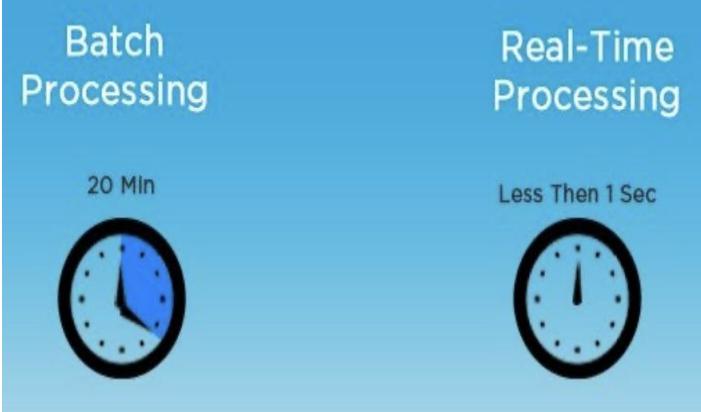
- Today requires expensive offline processing
- Can we ingest new tables in real-time?

## Table Discovery

- Incremental, real-time update of indices for lakes
- Data-driven cost models

## Table Navigation, Exploration

- Today requires offline organization construction
- Can we produce organizations at query time



[Nargesian et al. SIGMOD 20]  
[Oulette et al. PVLDB 21]

# Responsible Data Science

Table Discovery can be ***unfair***

E.g., Jaccard can lead to unfairness for larger datasets

***Even when all datasets are indexed***

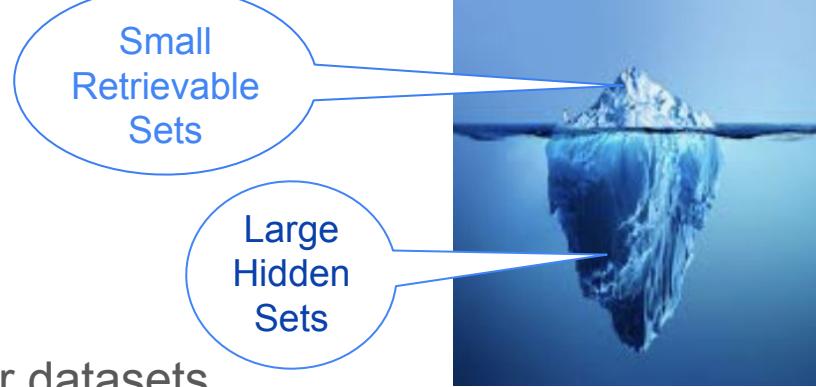
How to ensure all relevant answers are equally likely to be retrieved?

Beyond size ensuring fairness for other data characteristics

Table Discovery to satisfy ***Group Representation*** or ***Distribution Requirements***

Responsible Data Integration      [Nargesian et al. SIGMOD 22]

Table Understanding to identify bias





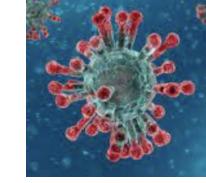
*Grace Fan*



*Jin Wang*



*Yuliang Li*



*Renée J. Miller*

## Table Discovery in Data Lakes: State-of-the-Art and Future Directions

<https://northeastern-datalab.github.io/table-discovery-tutorial/>