

# DATA ACQUISITION FOR AI

Fatemeh Nargesian, University of Rochester

# OPEN DATA MOVEMENT

- AI has become ubiquitous.
- Data-centric AI: focus has shifted from big data to good data.
- Open data repositories and data markets have become prevalent.

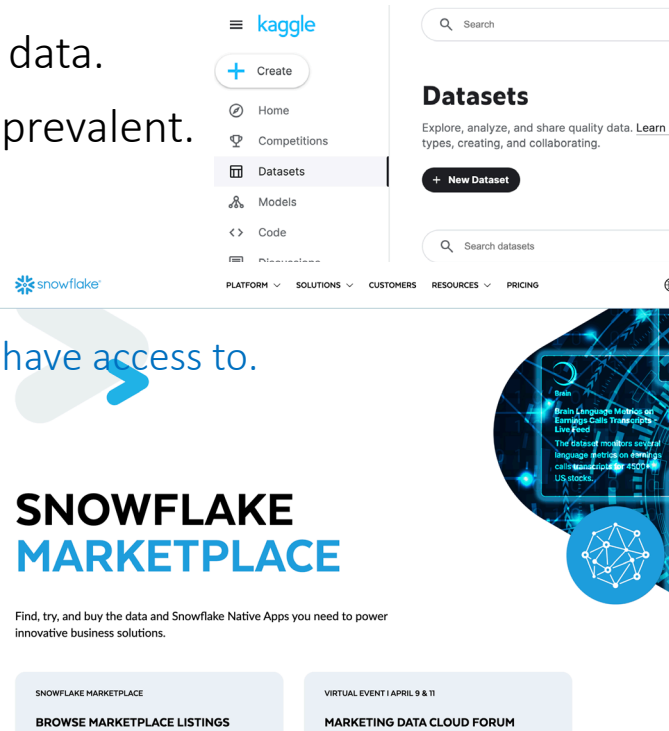
- Improving governments

- Empowering citizens

- Solving big public problem

- After Web, open data is the biggest database that citizens have access to.

- Focus of this talk: tabular and structured data



# DATA LAKES VS. TRADITIONAL DATABASES

- Data is stored in raw files (csv, xls, xml, ...) and must be extracted
- Large number of (medium-sized) datasets
- No centralized data design or data quality control
- Sparse and non-standardized metadata for datasets
- Skewed data distribution
- Vast number of datasets

35,675  
Canada

46,626  
France

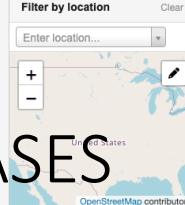
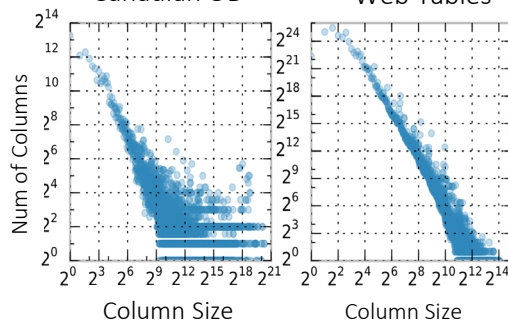


DATA REPORTS OPEN GOVERNMENT CONTACT

DATA CATALOG

Canadian OD

Web Tables



Topics

Local Government - 21214

Climate - 528

Older Adults... - 88

Energy - 21

Topic Categories

Arctic - 134

Ecosystem Vulnerability - 92

Water - 89

Human Health - 70

Arctic Ocean, Sea... - 66

Transportation - 61

Energy Infrastructure - 60

Atmospheric, Earth... - 63

Food Resilience - 52

Coastal Flooding - 42

Show More Topic Categories

Dataset Type

geospatial - 231609

Tags

earth science - 77182

county or equivalent entity - 71495

oceans - 66104

noaa - 54349

ocean - 54313

u.s. department of commerce - 47280

state fips code - 42137

county fips code - 42022

nessis - 40763

Show More Tags

Formats

XML - 143849

292,134 datasets found

Electric Vehicle Population Data 4135

State of Washington — This dataset shows the Hybrid Electric Vehicles (PHEVs) that are currently registered in the State of Washington.

CSV RDF JSON XML

Crime Data from 2020 to Present 311

City of Los Angeles — Starting on March 7th, 2021, the City of Los Angeles (LAPD) will adopt a new Records Management system is...

CSV RDF JSON XML

FDIC Failed Bank List 2457 recent views

Federal Deposit Insurance Corporation — The FDIC is a government corporation that insures the deposits of banks and bank branches. This list includes banks which have failed since 1999.

CSV HTML

Dynamic Small Business Search (DSBS)

Small Business Administration — The Small Business Search (DSBS) database. As a result of the Small Business Administration's (SBA) Award Management, there...

HTML

Fruit and Vegetable Prices 1999 recent views

Department of Agriculture — How much do fruit and vegetable prices for 153 commonly consumed fresh and processed fruits and vegetables cost?

XLS

Motor Vehicle Collisions - Crashes 1

City of New York — The Motor Vehicle Collision Database (MVCDB) is a database of motor vehicle collisions that occur in the City of New York. Each row represents a crash event. The data is updated daily.

CSV RDF JSON XML

Walkability Index 1913 recent views

U.S. Environmental Protection Agency — The U.S. Environmental Protection Agency (EPA) has developed a Walkability Index (WI) based on the 2010 Census 2010 block group in the U.S. based on characteristics of...

ZIP CSV Esri REST HTML

Lottery Powerball Winning Numbers: E

State of New York — Go to <http://on.ny.gov/1Gp> for the latest Powerball results and payouts.

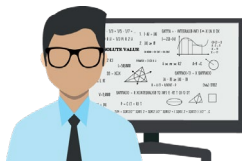
CSV RDF JSON XML

Supply Chain Greenhouse Gas Emission

1568 recent views

U.S. Environmental Protection Agency — The U.S. Environmental Protection Agency (EPA) has developed a Supply Chain Greenhouse Gas (GHG) emission factors (Factors) for 1,016 U.S. manufacturing sectors and the North American industry...

CSV CSV



# DATA SEARCH AND DATA ENRICHMENT

Analyzing the driving factors of GHG emission!

Enrich data scientist's work in progress  
with right data!

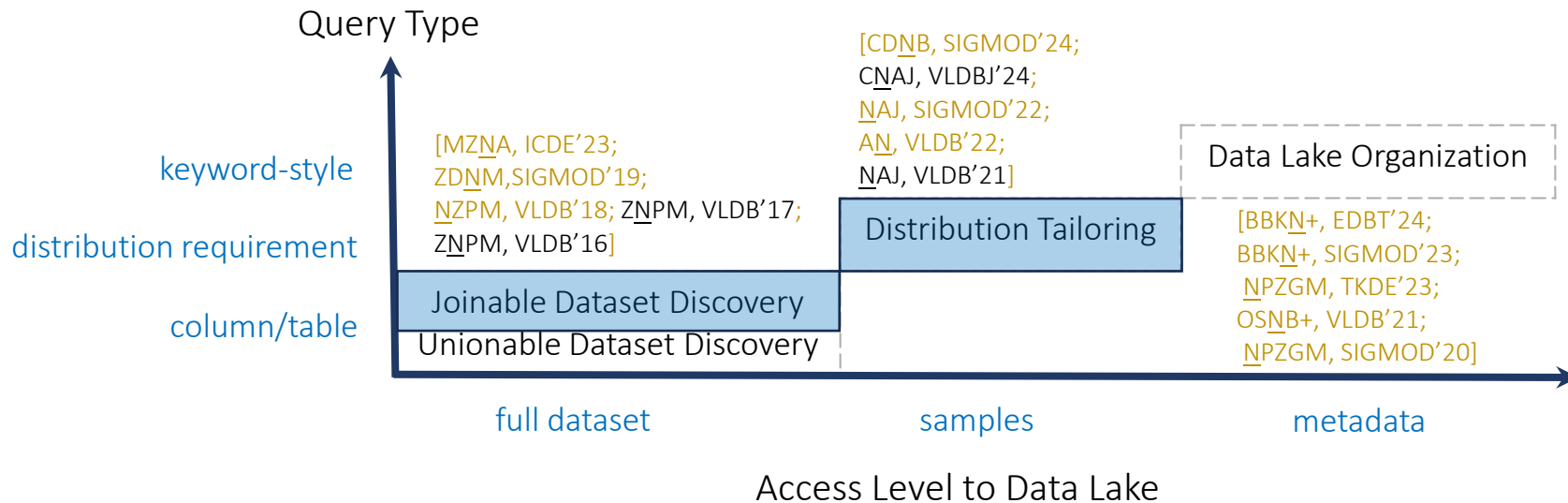
Geo	Date	Fuel	t CO2	Sector	...
Cambridge	2015	electricity	2	Waste	...
Worcester	2021	diesel	20	Metal	...
Camden	2014	coal	12	Oil&Gas	...
NYC	2019	electricity	11	Oil&Gas	...
Boston	2023	diesel	8	Metal	...
Rochester	2021	coal	9	Metal	...
...	...	...	...	...	...

Data enrichment requires **dataset discovery**

- Adding novel features: **joining** the query dataset with some datasets in the lake.
- Adding samples: **unioning** the query dataset with some datasets in the lake.

# DATASET DISCOVERY

TASK OF FINDING RELEVANT DATASETS TO A QUERY

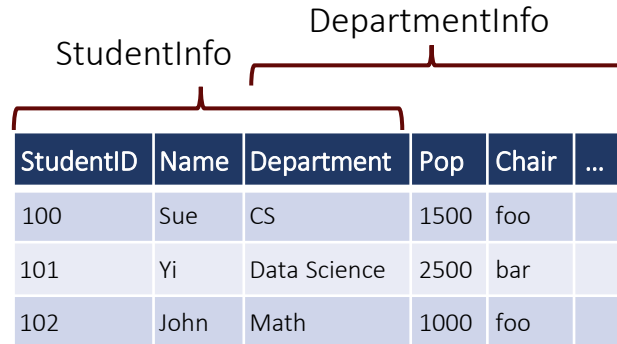
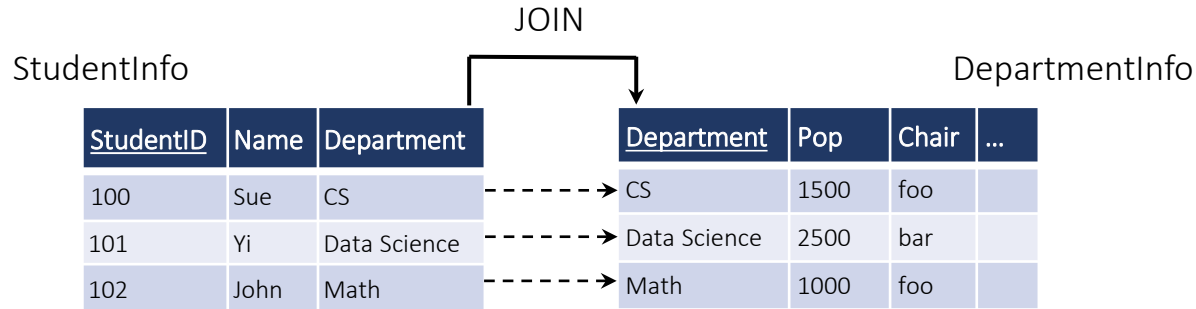


# JOINABLE DATASET DISCOVERY

How to enrich a query dataset with novel columns and features?

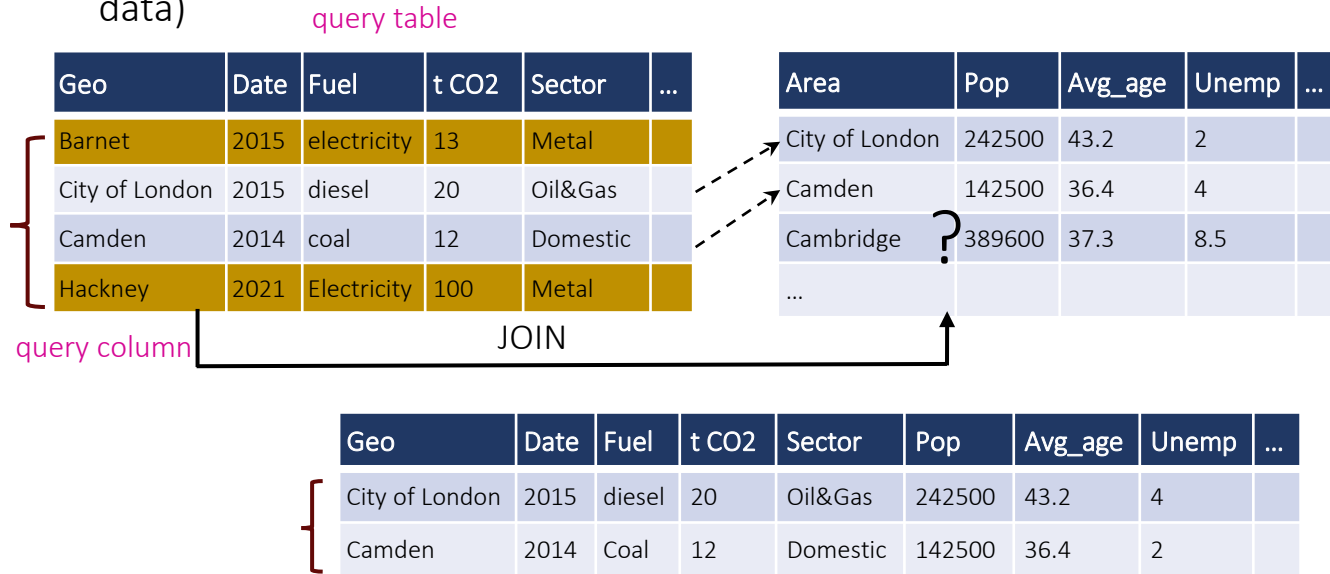
# OVERVIEW: JOIN IN DATABASES

- In databases, we often know which columns to join: join on primary/foreign keys



# JOIN IN DATA LAKES

- Not obvious which table to join on: makes discovery a [search problem](#)
- Joins might not be possible on all query's tuples: smaller result set than query (incomplete data)





# JOINABILITY MEASURE

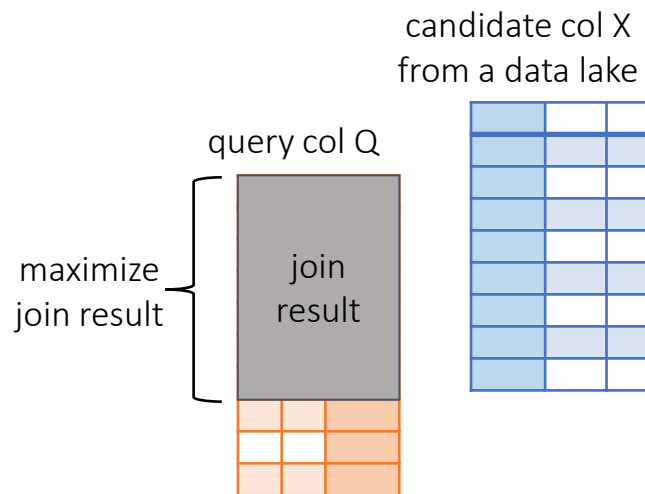
- If columns and query are considered as sets of values, maximize set overlap of query and candidate.

$$\text{Overlap}(Q, X) = |Q \cap X|$$

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

- Another popular set similarity measure is

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





# SEMANTIC JOINABILITY

- Syntactic measures become ineffective for joining dirty and heterogenous data in the wild.

query table

Overlap(Geo,Area) = 1

...	Date	Fuel	t CO2	Sector	Geo	Area	Pop	Avg_age	F.Unemp	Unemp	...
	2015	electricity	130	Domestic	Blaine	LA	8800	43.2	-	-	
	2015	diesel	200	Transport	LA	Big Apple	242500	36.4	62.9	4	
	2014	coal	125	Domestic	NYC	Blain	389600	37.3	66	8.5	
						...					

- Semantic overlap** extends syntactic overlap for effective search despite semantic and syntactic heterogeneity of tuples.

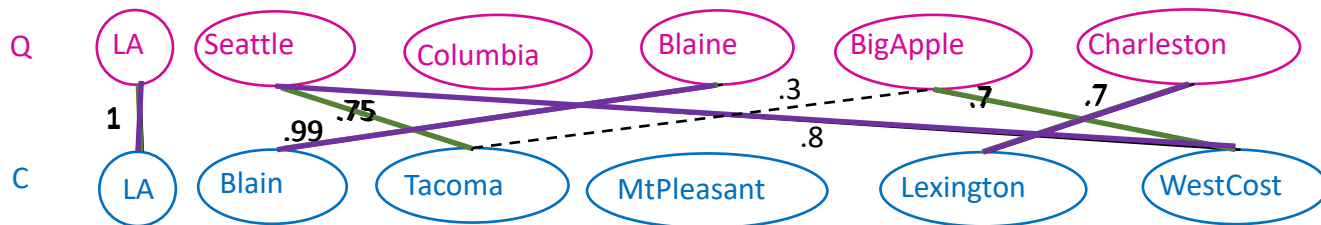
# SEMANTIC OVERLAP MEASURE

...	Q
	LA
	Blaine
	BigApple
	...

...	C
	LA
	Blain
	Tacoma
	...

Q = {LA, Seattle, Columbia, Blaine, BigApple, Charleston}  
 C = {LA, Blain, Appleton, Tacoma, Lexington, WestCoast}

SEMA-JOIN  
 [He, Gunjam+VLDB'15]



Semantic overlap provides a small number of joinable columns. The score of the bipartite graph of Q and C with node user-specified tuple similarity function and threshold.

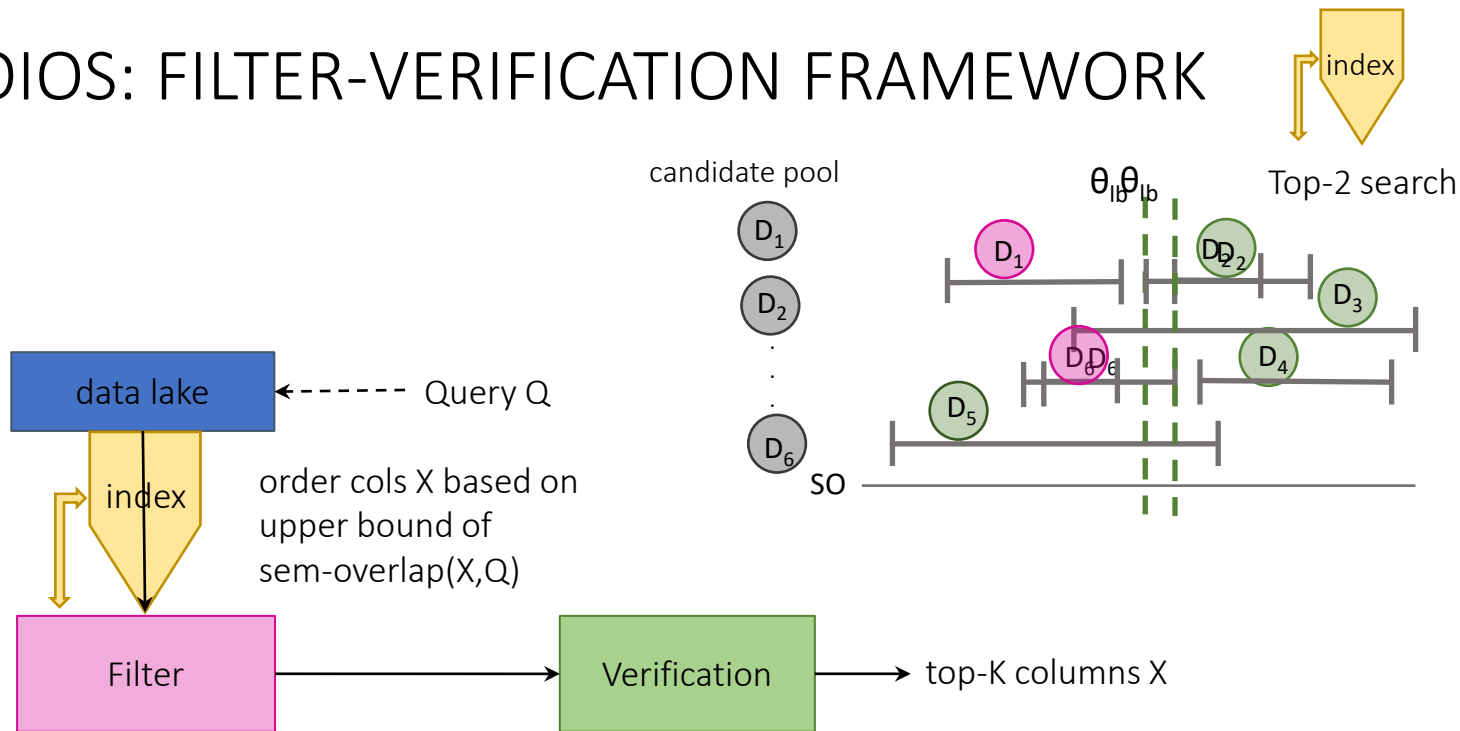
- Semantic similarity: Cosine on embeddings of col. values, etc.
- Join with semantic overlap results in larger join result than syntactic join.
  - Syntactic similarity: Jaccard on n-grams of col. values
- $|Q \cap C| \leq SO(Q, C)$ 
  - Pruned by similarity threshold  $\alpha$
- Discovery requires finding joinable columns and the best way of joining them.

$$\text{score}(M) = 4.14$$

# TOP-K SEMANTIC OVERLAP SEARCH

- Semantic overlap can be expensive
  - Bipartite graph matching:  $O(n^3)$ ,  $n$  is col. size [Kuhn'1995]
  - Bipartite graph construction:  $O(n^2)$
- **Problem.** Given a column  $Q$  and parameter  $K$ , find the **top- $K$**  columns based on the **semantic overlap** measure.
- Search complexity:  $O(mn^3)$ ,  $n$  is the size of cols. and  $m$  is the number of sets
- Linear scan over all datasets and computing graph matching is infeasible in practice for data lakes of tens of thousands of datasets.
- **Solution.** **KOIOS** is an exact and efficient top- $K$  join search algorithm with semantic overlap.

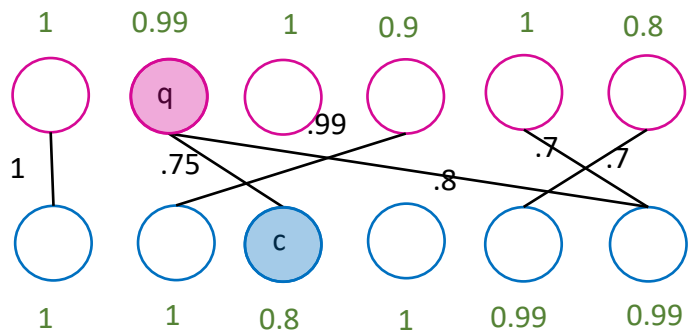
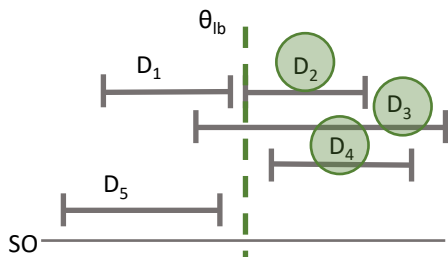
# KOIOS: FILTER-VERIFICATION FRAMEWORK



- Upper- and lower-bound filters
- A partitioning scheme for efficient filtering
- Early termination of bipartite graph matching



# VERIFICATION: EARLY TERMINATION OF MATCHING



- Hungarian algorithm assigns a valid labeling function  $label: nodes \rightarrow R$ , s.t. for two nodes  $q$  and  $c$ ,  $label(q) + label(c) \geq edge\_weight(q, c)$
- The algorithm improves on the labeling function iteratively. At each iteration:

$$\text{bipartite matching score}(Q, C) \leq \sum \text{node labels}$$

$$\text{upper bound}(Q, C) = \sum \text{node labels}$$

- Terminate matching prematurely.

# EVALUATION

datasets statistics

Dataset	#Sets	Max Card.	Avg. Card.	#Unique Elements
DBLP	4,246	514	178.7	25,159
OpenData	15,636	31,901	86.4	179,830
Twitter	27,204	151	22.6	72,910
WDC	1,014,369	10,240	30.6	328,357

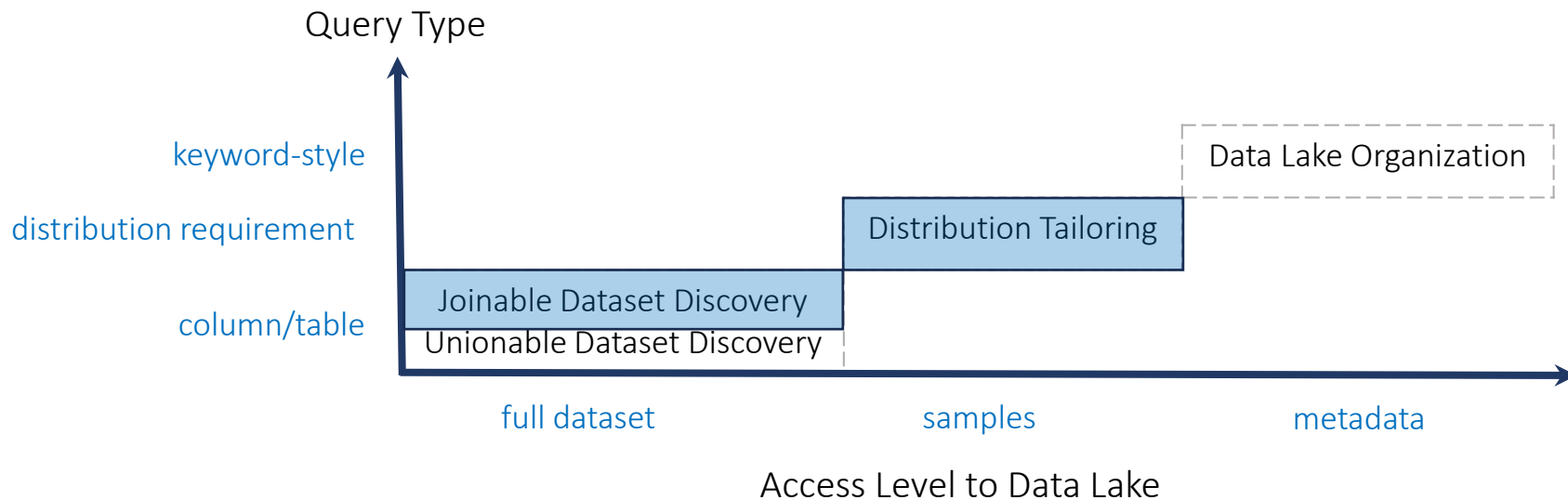
comparison to SOTA

Dataset	KOIOS Response Time (s)	SOTA Response Time (s)	KOIOS Mem (MB)	SOTA Mem (MB)
DBLP	0.83	211	0.83	11
OpenData	18.6	101	18.6	102.5
Twitter	0.7	518	0.7	10
WDC	147	1062	147	885

- KOIOS achieves at least 5X speed up over the SOTA on massive data lakes.
- Even better speedup for medium and large queries compared to the SOTA.

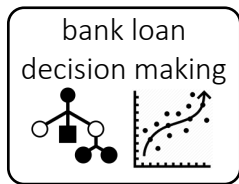


# DATASET DISCOVERY



# DISCOVERY + DISTRIBUTION REQUIREMENT

- Scenarios with distribution requirements
  - Representation in test data
  - Seeking data with a sufficient representation to avoid overfitting
  - Selection bias leads to flawed and unreliable outcomes.
- Nearest neighbor search index on histograms of groups in datasets [Mao+, AAAI'17]
  - Non-existent results
  - Need to know the group of each tuple apriori at index time.
- Combine multiple datasets to get the distribution we want!

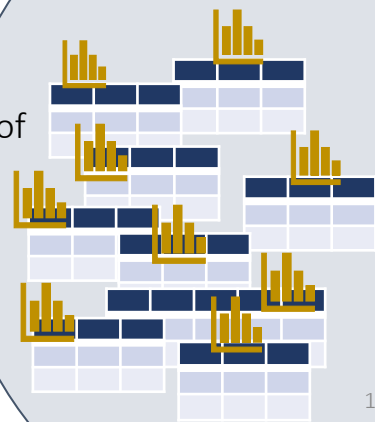


test

At least 30% African American  
female employees from Texas w/  
salary  $\geq$  55K and at least 30% white  
male employees w/ salary  $\geq$  55K

on histograms of  
target groups

Index

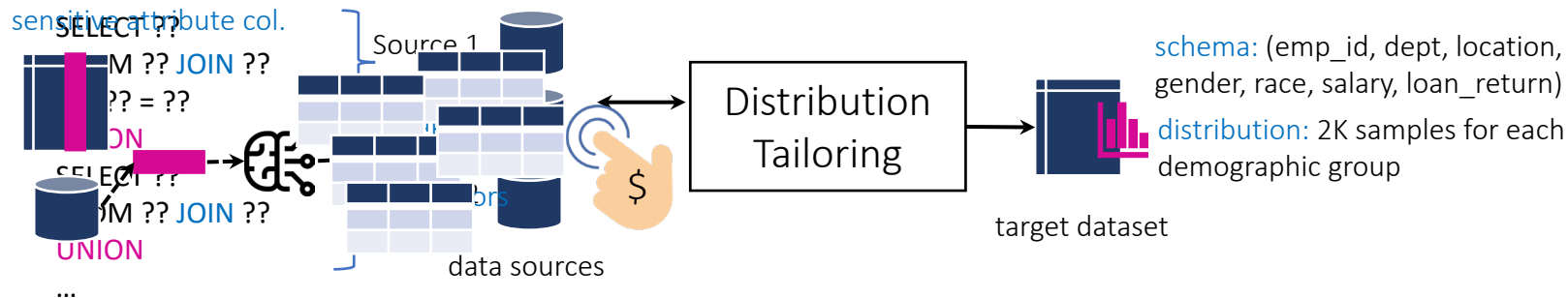


# DISTRIBUTION-AWARE DISCOVERY

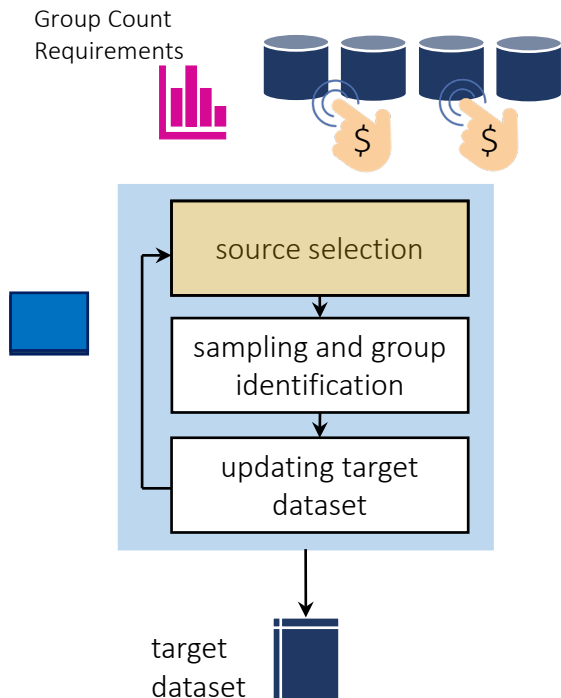
How to construct a dataset that satisfies group distribution requirements from multiple sources in a cost-effective manner?

# QUERY/DATA/COST MODELS

- **Target dataset:** schema + count requirements specified over groups
  - Schema: description of columns
- **Data sources**
  - Data lake tables with the same schema as target schema
  - Project-join views over a database/data lake: join is expensive to execute, resort to tuple sampling from joins [Zhao+, SIGMOD'18; Li+, SIGMOD'16; Haas+, SIGMOD'99]
  - Other sources: crowd-sources, data providers – data market setting [AN, VLDB'22]: monetary cost for purchasing data



# DATA DISTRIBUTION TAILORING (DT)



- **Problem.** Given **sources** with their **costs**, and **minimum count requirements** on the **groups**, select samples from sources s.t. the union of samples fulfills the count requirements, while the **expected total query cost is minimized**.
- **Solution.** Iterative sampling: find a sequence of sources to sample, until distribution requirement is satisfied.
- The cost optimality depends on how much DT knows about group distributions in sources.

# DT STRATEGIES FOR (UN)KNOWN GROUP DISTRIBUTIONS

- Known distributions
  - Dynamic programming solution
    - Pseudo-polynomial time and space complexity
    - Not practical for large number of groups and count requirements
  - Optimal strategy for binary groups and sources with equal costs
  - Practical strategy for m-ary groups and sources with arbitrary costs
- Unknown distributions
  - Budget allocation strategy based on multi-armed bandit

Tailoring Data Source Distributions for Fairness-aware Data Integration. F. Nargesian, A. Asudeh, H. V. Jagadish, VLDB, 2021.

Data Distribution Tailoring Revisited: Cost-Efficient Integration of Representative Data. J. Chang, B. Cui, F. Nargesian, A. Asudeh, H. V. Jagadish, VLDBJ, 2024.



UR PhD  
student



UR Undergrad  
student

# UNKNOWN DISTRIBUTIONS

No information about group distribution

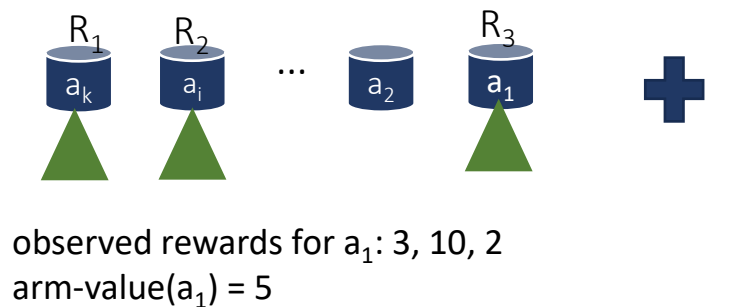


- Learning group distribution and source goodness as we sample.
- **Solution.** Applying Multi-armed Bandit (MAB)

# OVERVIEW: MULTI-ARMED BANDIT

- Given  $k$  arms a time horizon (budget)  $T$ , at each timestep  $t=1,\dots,T$ , we choose an arm, and receive a **real-valued reward**  $R_t$ .

- Reward** depends on the arm and is iid.
- Distributions of rewards are unknown.
- Nevertheless, we must **maximize our total reward**.
- As selecting arms, form estimates for an arm's value: e.g., the **average of sample rewards** from the arm.

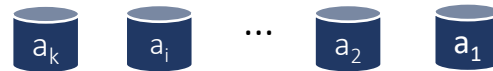


[Sutton&Barto, 1998]



# OVERVIEW: MULTI-ARMED BANDIT

- Maximize the total reward.
- Both try arms to learn their values (**explore**) and prefer those that appear best (**exploit**).
  - Never stop exploring; maybe explore less with time; or not!
- **Regret** is the opportunity loss for one step: difference of obtained reward and optimal reward
- **Goal**. Minimize total regret  $\sim$  maximize cumulative reward



$$R_1 + R_2 + R_3$$

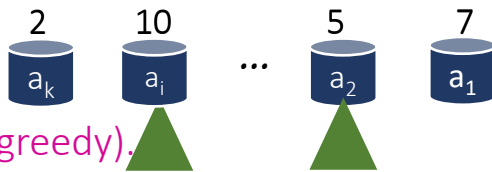
$$\text{Regret}(T) = \text{OPT reward @T} - \text{Planner reward @T}$$

observed rewards so far for  $t=1,..,4$ : 3, 10, 5, 10  
should have pulled the arm with reward 10 all along  
**regret**:  $(10-3)+(10-10)+(10-5)+(10-10) = 12$

[Sutton&Barto, 1998]

# OVERVIEW: EXPLORATION, EXPLOITATION, AND REGRET

estimated action-values:



- **Explore** at time  $t$ : select a **random** arm.
- **Exploit** at time  $t$ : select the arm with the **best** estimated value so far (**greedy**).
- $\epsilon$ -greedy. At each step explore with **prob.  $\epsilon$  (exploration rate)** and exploit with **prob.  $(1-\epsilon)$** 
  - Linear regret
    - Assume perfect estimates. We potentially pull imperfect arm with  $\epsilon$  prob., resulting in expected  $\epsilon.T$  regret.
- **Decaying exploration rate**. More exploration at the beginning.
  - E.g., if new to a city, extensively explore restaurants at the beginning, explore less later.
  - Regret  $O(T^{2/3} \log T^{1/3})$
  - Can be brought down by Upper Confidence Bound (UCB) to  $O(T^{1/2} \log T^{1/2})$

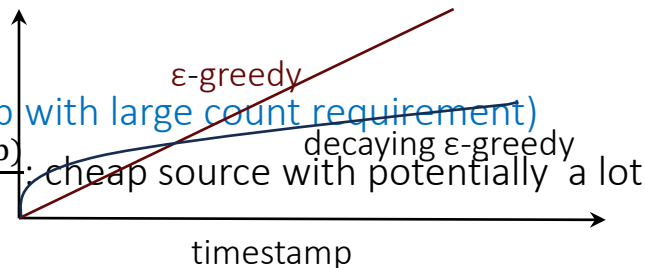
# MAB FOR UNKNOWN DISTRIBUTION DT



- Each source is an arm.
- Apply  $\epsilon$ -greedy.
- With each sample from a source, we learn about the distribution of groups in that source.

- Greedy strategy at time  $t$ .

- First, select a hard-to-find group (rare group, total group with large count requirement)
  - Greedy action.  $\arg\max_{source} \frac{\text{ratio(hard to find group)}}{\text{source cost}}$ : cheap source with potentially a lot of samples of hard-to-find-group.



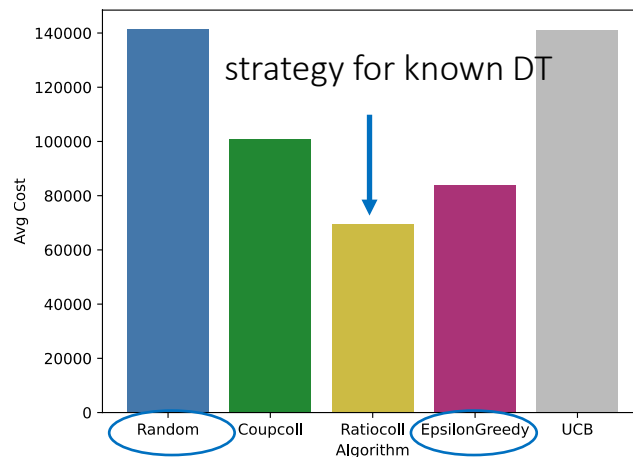
- Regret is a proxy for cost overhead.
- Regret. DT with  $\epsilon$ -greedy strategy with exploration rate  $\sqrt[3]{\ln t / t}$  at time  $t$  has regret of  $O\left(t^{\frac{2}{3}} \log t^{\frac{1}{3}}\right)$  -- for sources with equal costs and strategy *ratiocoll*.

# PRACTICAL STRATEGY

- Explore-then-Exploit
  - Crudely approximate budget  $T \approx \sum group\_count\_requirements$
  - Randomly sample sources for  $\alpha T$  iterations
    - Performing all explorations at the beginning allows sampling to be done in batches and in parallel!
  - Be greedy in the rest of time steps.

# EVALUATION: KNOWN AND UNKNOWN DISTRIBUTIONS

Known/Unknown DT on Flights dataset



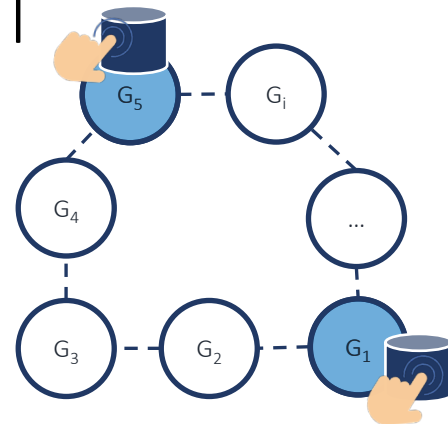
- $\epsilon$ -greedy outperforms random sampling and is in competition with strategy for. Known distribution.

# PLUTUS: UNDERSTANDING DT FOR ML

- DT is available in Apache SystemDS (open-source SystemML from IBM: declarative ML system on Spark)



# KNOWN DISTRIBUTIONS: BINARY GROUPS AND EQUI-COST



- **Given.** The ratio of each group in each source. Same cost for all sources.
- **Algorithm.** At each iteration until all requirements are satisfied:
  - From unsatisfied groups, pick a group to prioritize.
  - For that group pick the cheapest sources for acquiring samples of that group.
  - Sample from that source; update target dataset; update remaining group requirements.
- **Cheapest source for a group:** source that has the highest ratio of that group
- **Group to prioritize/sample for:** group that is minority in its best source

# EQUI-COST BINARY DT: SAMPLE FOR MINORITY GROUP

cost=1



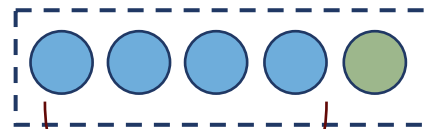
20% of  $G_1$  and 80% of  $G_2$

cost=1



5% of  $G_1$  and 95% of  $G_2$

sample from  $S_1$  until  
got one tuple of  $G_1$



No need to separately  
sample for  $G_2$ , req.  
satisfied

sample from  $S_2$  until  
got one tuple of  $G_2$



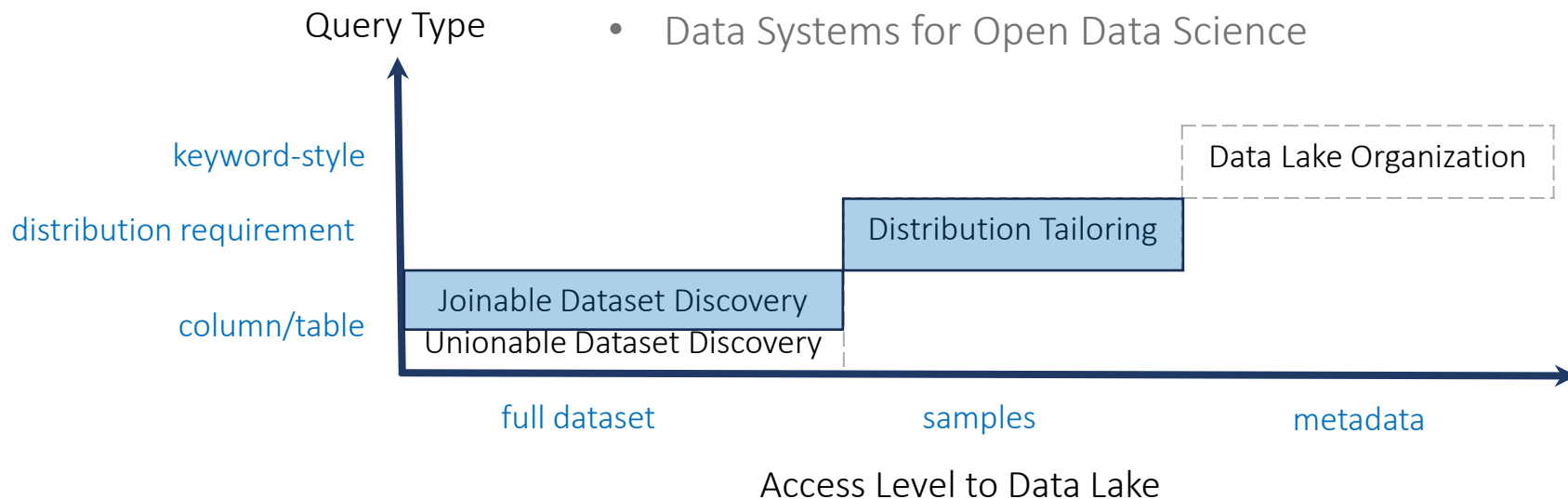
Need to sample  
another source to get  $G_1$

- **Requirement.** Collect at least one tuple of each group.
- Expected cost of getting one sample of  $G_1$  from  $S_1$  is  $100/20=5$  and from  $S_2$  is  $100/5=20$ .
  - Best source for getting  $G_1$  is  $S_1$ . Similarly, best source for  $G_2$  is  $S_2$ .
- $G_1$  is minority in its best source ( $S_1$  has 20% of  $G_1$ ,  $S_2$  only has 5% of  $G_1$ ) → Pick  $G_1$ 
  - Piggybacking: as we are sampling for  $G_1$ , we can fulfill the count requirements of  $G_2$  with no cost.
- Proof by contradiction.



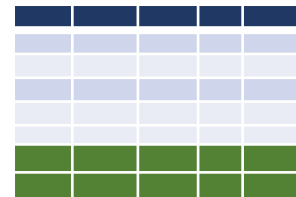
# OUTLOOK

- Open Data for AI
- Open Data for Users
- Data Systems for Open Data Science



# OPEN DATA AND AI

- Sample discovery
  - Discovering novel samples by unionable dataset discovery [NZPM, VLDB'18]
- Feature discovery
  - Pushing down feature selection measures into join discovery
  - More interesting (and rare) relationships: causal dataset discovery
- LLMs and dataset discovery
  - Dataset understanding and training data generation
- Dataset discovery for LLMs
  - Semantic/query-based nearest neighbor search indexes for vector DBs

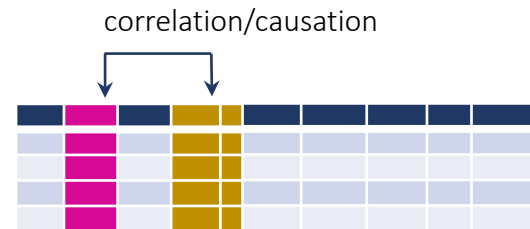


# CAUSAL DATASET DISCOVERY

- Consider tables  $Q$  and  $T$ . Candidate columns  $X \in Q$  and  $Y \in T$  have a **correlation link** over  $Q \bowtie_{K=K'} T$ , if their correlation after the join is higher than a threshold  $\theta$ .

- A correlation link can potentially be

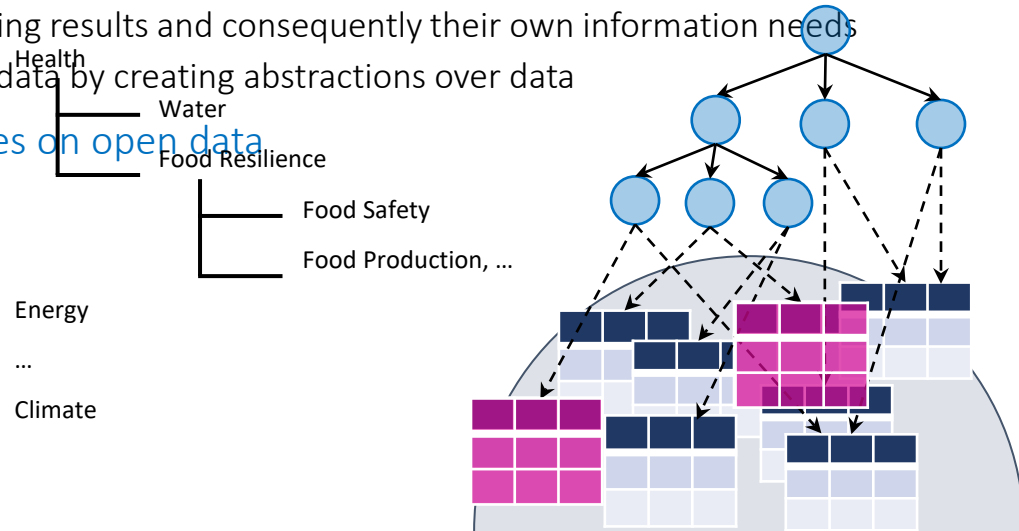
- $X$  causes  $Y$  over join  $Q \bowtie_{K=K'} T$
- $Y$  causes  $X$  over join  $Q \bowtie_{K=K'} T$
- No causal relation exists between  $X$  and  $Y$



- Problem.** Given a query table  $Q$  and a data lake of tables  $L$ , find all tables  $C \in L$  such that  $Q$  and  $C$  have at least one causal link over their join.
- Challenges.** Sparsity of causal relation and limited training data for fine-tuning
- Results.**
  - Role prompting + CoT + Socratic prompting + fine tuning performs the best.
  - LLMs have potential to learn and generalize causality.

# OPEN DATA AND USERS

- Metadata, search, query writing
  - Metadata standardization and taxonomy induction in domains such as social sciences
- User interfaces for search
  - Constructing abstract structures for navigation and exploration [NZGPM, SIGMOD'20&TKDE'24]
  - Assisting users with debugging results and consequently their own information needs
  - Dealing with heterogenous data by creating abstractions over data
- Discovery with private queries on open data





# ACKNOWLEDGEMENT

- Abolfazl Asudeh (UIC)
- Nikolaus Augsten (U. Salzburg)
- Matthias Boehm (TU Berlin)
- H. V. Jagadish (U. Michigan)
- Renée J. Miller (NEU)
- Felix Naumann (HPI)
- Ken PU (UOIT)
- Divesh Srivastava (AT&T)
- Eric Zhu (MSR)
- Leon Bornemann (HPI, PhD student)
- Jiwon Chang (UR, PhD student)
- Tianji Cong (U. Michigan, PhD student)
- Yurong Liu (UR, undergrad)
- Pranay Mundra (UR, MSc.)
- Aidan Sciortino (UR, undergrad)
- Zhengbin Tao (UR, PhD student)
- Draco Xu (UR, undergrad)
- Mengqi Zhang (UR, PhD student)
- Jianhao Zhang (UR, MSc.)

Code and data: [github.com/DataIntelligenceCrew](https://github.com/DataIntelligenceCrew)

THANKS.

