



Madelon Hulsebos



Table Representation Learning

UW Seattle
May 24th, 2024

The Impressive Capabilities of Transformers

Transformers leveraged for applications over images, text, code:



Generates images of dogs

You
what to include in speech with PhD graduation committee after defense ceremony?

ChatGPT
After successfully defending your PhD dissertation, it's customary to give a speech to express gratitude and appreciation to your graduation committee and others who have supported you throughout your academic journey. Here's what you might want to include in your speech:

1. **Acknowledgment of the Committee:** Begin by thanking each member of your graduation committee individually, mentioning their contributions to your research and their guidance during the defense process.
2. **Express Gratitude:** Express your sincere gratitude to your advisor, committee members, mentors, colleagues, friends, and family for their support, encouragement, and assistance throughout your PhD journey.
3. **Reflect on the Journey:** Reflect briefly on your academic journey, highlighting key milestones, challenges you've overcome, and significant moments of growth and learning.
4. **Highlight Contributions:** Summarize the main contributions of your research and its potential impact in your field or related areas. This is an opportunity to showcase the significance of your work and its relevance.

Helps writing graduation speech

```
sentiment.ts | ⚡ write.sql.go | 📈 parse.expenses.py | 📩 addresses.rb
```

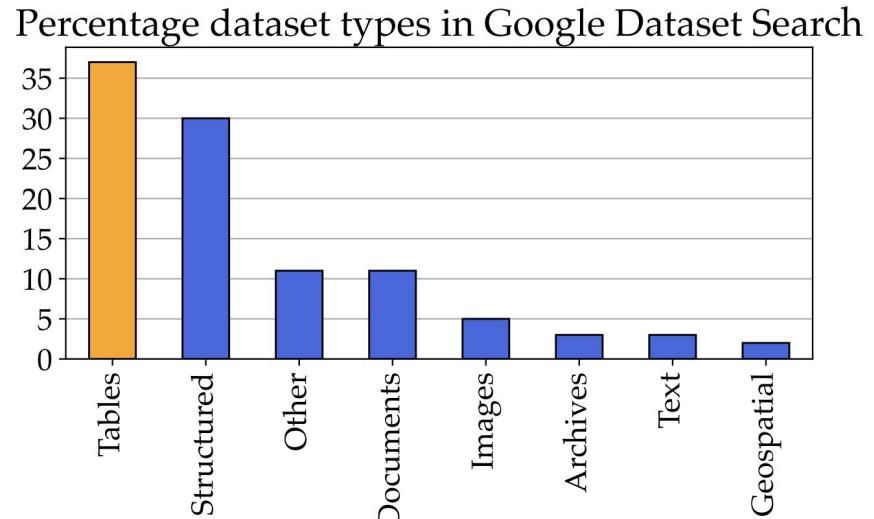
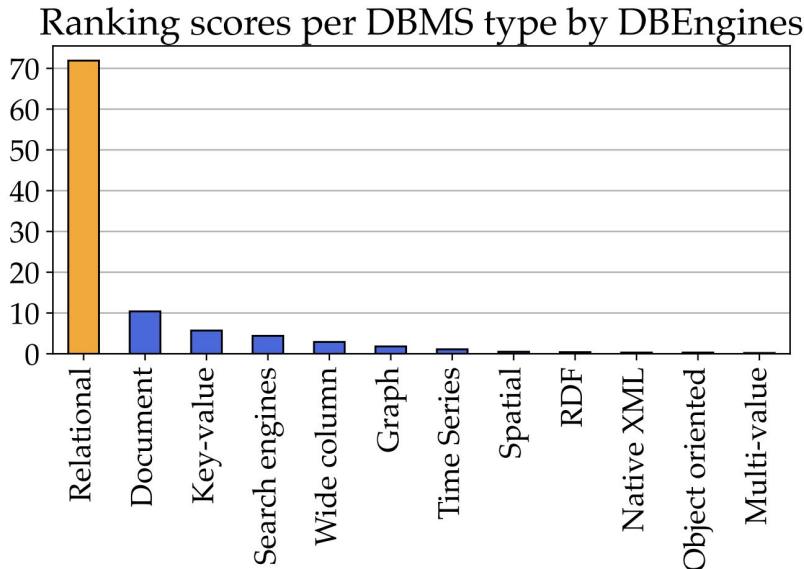
```
1 #!/usr/bin/env ts-node
2
3 import { fetch } from "fetch-h2";
4
5 // Determine whether the sentiment of text is positive
6 // Use a web service
7 async function isPositive(text: string): Promise<boolean> {
8   const response = await fetch(`http://text-processing.com/api/sentiment/`, {
9     method: "POST",
10     body: `text=${text}`,
11     headers: {
12       "Content-Type": "application/x-www-form-urlencoded",
13     },
14   });
15   const json = await response.json();
16   return json.label === "pos";
17 }
```

Copilot

Completes code

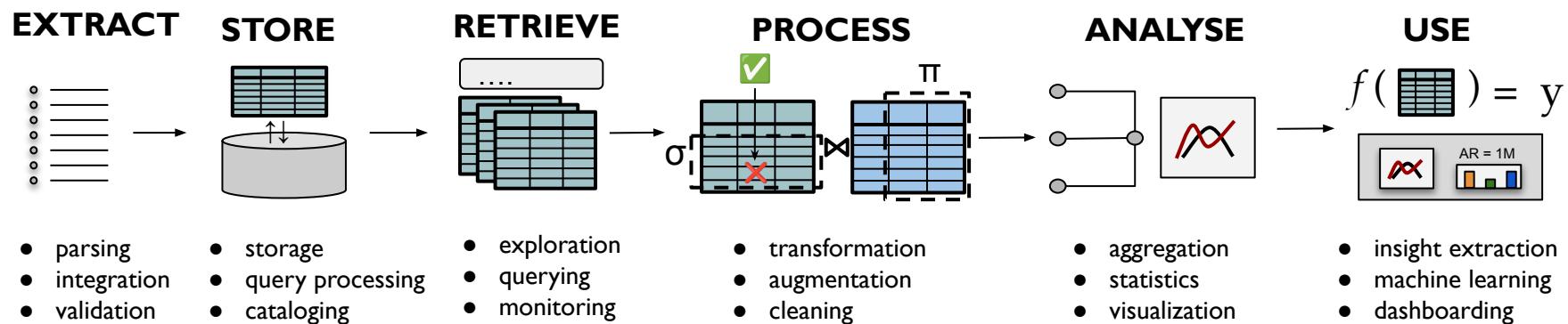
What about **tables**?

Tables Dominate the Data Landscape



Application Potential of Table Representations

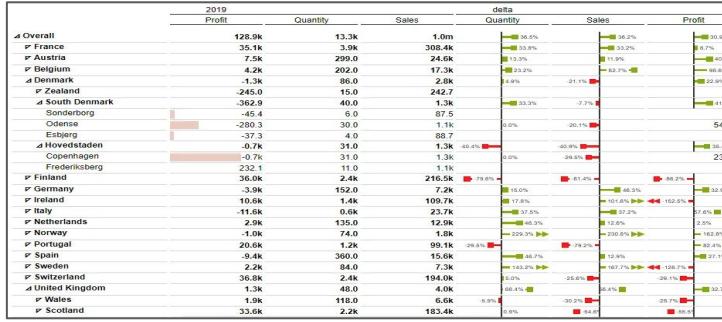
Many tasks in **high-value use-cases** operate over tables, e.g. *data analytics*.



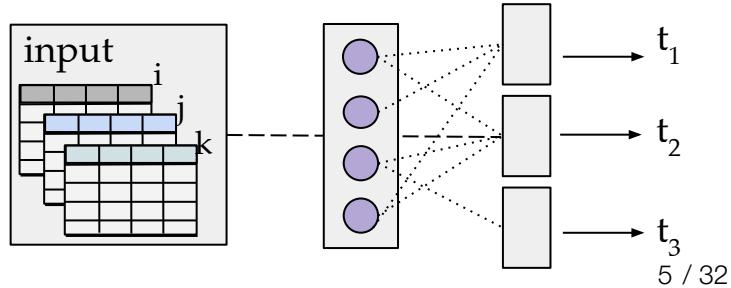
Rich and challenging!

Tables come with diversity in structure, dimensions, content, and semantics...

Nr	ID	seed rate	yield	crop	cultivar	pre crop	pre-pre crop	pre-pre-pre	soil type	precipita	tempera	comment
1	68	91	winter wheat	sugar beets	beans		sandy loam, loe 636	9.6	wb, sg,			
2	68	106	winter wheat	sugar beets	rotation fallow		sandy loam, loe 636	9.6	cultivation			
3	68	97	winter wheat	sugar beets	fallow land (5,5y)		sandy loam, loe 636	9.6	1993-1996			
4	136	95	winter wheat	oats	sugar beets		sandy loam, loe 636	9.6				
5	136	96	winter wheat	potatos	sugar beets		sandy loam, loe 636	9.5	cultivation			
6	136	107	winter wheat	sugar beets	maize		sandy loam, loe 636	9.5	1991-1994			
7	136	107	winter wheat	sugar beets	summer wheat	maize	sandy loam, loe 636	9.5				
8	136	82	winter wheat	oats	sugar beets	sugar beets	sandy loam, loe 636	9.5	organic			
9	136	77	winter wheat	potatos	sugar beets		sandy loam, loe 636	9.5	organic			
10	136	85	winter wheat	sugar beets	maize	maize	sandy loam, loe 636	9.5	organic			
11	136	84	winter wheat	sugar beets	summer wheat	sugar beets	sandy loam, loe 636	9.5	organic			
12	57 371	98	winter wheat	Sperber	sugar beets	winter barley	winter wheat	sandy loam, loe 635	wb, ww			
13	57 365	98	winter wheat	Sperber	potatos	sugar beets	summer barley	sandy loam, loe 635	cultivation, weed			
14	57 365	105	winter wheat	Sperber	sugar beets	maize	maize	sandy loam, loe 635	1987-1992			
15	57 365	97	winter wheat	Sperber	sugar beets	winter wheat	sugar beets	sandy loam, loe 635				
16	39 433	90	winter wheat	Okapi	summer barley		sandy loam, loe 690	8.5	oats, cultivation, weet			
17	39 433	100	winter wheat	Okapi	oats		clay, silt	690	8.5	1982-1986		
18	39 433	97	winter wheat	Okapi	winter wheat		clay, silt	690	8.5			



Goal TRL: map **tables** to some consistent input.
 Learn some representation that helps detect patterns relevant to given task(s).



Outline for today



Images,
videos,
text...

1. Neural Models for Table Understanding
2. Resources for Table Representation Learning
3. Retrieval systems for structured data

Tables

Neural Models for **Table Understanding**

Column type detection: why?

Essential understanding of a table comes through its columns.

name	salary	country
name	salary	cntr

Looks easy, but....

- Undescriptive header?
- Messy and heterogeneous values?
- Unknown types?

Semantic column types dictate operations to perform on them:

name	salary	cntr



naam	status	land

name
Xi
carl
sara

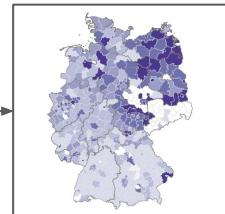
name
Xi
Carl
Sara

Join tables on “name” and “country” columns

Capitalize “name” columns

name	salary	cntr

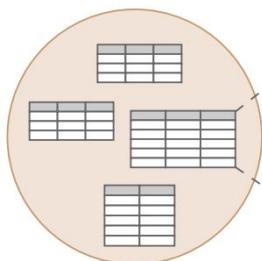
Plot “country” data



Sherlock: Column Type Detection with DL

Prior: string matching (col name/values) w/ regex or dict: **robust?** **scale?** **accuracy?**

1. Source Corpus



2. Sampled Dataset and Features

Column Header		Date
Country	Name	Exact Matching
Nigeria	M Buhari	08-18-2018
Germany	A Merkel	06-14-2015
...
Canada	J Trudeau	07-20-2017

Column Values Feature Extraction

Feature Categories

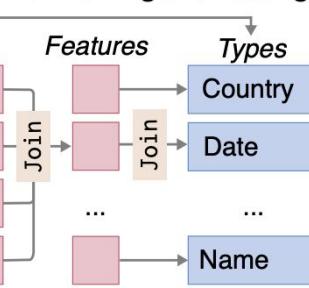
Character Distributions

Word Embeddings

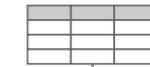
Paragraph Vectors

Global Statistics

3. Training and Testing Set



4. Semantic Type Detection

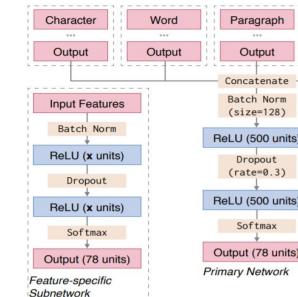


Unseen Data

Model

Prediction
Location: 0.9
Name: 0.7
Year: 0.8

Predicted Types and Confidences



How well does **Sherlock** detect types?

78 semantic types (name, address, etc).

Method	F ₁ Score	Runtime (s)	Size (Mb)
<i>Machine Learning</i>			
Sherlock	0.89	0.42 (± 0.01)	6.2
Decision tree	0.76	0.26 (± 0.01)	59.1
Random forest	0.84	0.26 (± 0.01)	760.4
<i>Matching-based</i>			
Dictionary	0.16	0.01 (± 0.03)	0.5
Regular expression	0.04	0.01 (± 0.03)	0.01
<i>Crowdsourced Annotations</i>			
Consensus	0.32 (± 0.02)	33.74 (± 0.86)	-

Examples of misclassifications.

Examples	True type	Predicted type
<i>Low Precision</i>		
81, 13, 3, 1	Rank	Sales
316, 481, 426, 1, 223	Plays	Sales
\$, \$\$, \$\$\$, \$\$\$\$, \$\$\$\$\$	Symbol	Sales
<i>Low Recall</i>		
#1, #2, #3, #4, #5, #6	Ranking	Rank
3, 6, 21, 34, 29, 36, 54	Ranking	Plays
1st, 2nd, 3rd, 4th, 5th	Ranking	Position

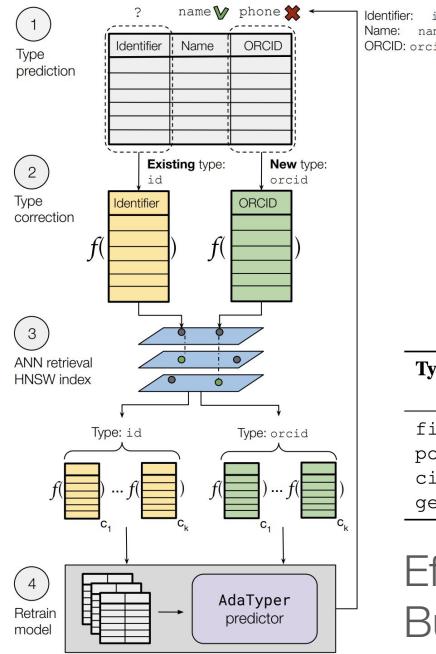
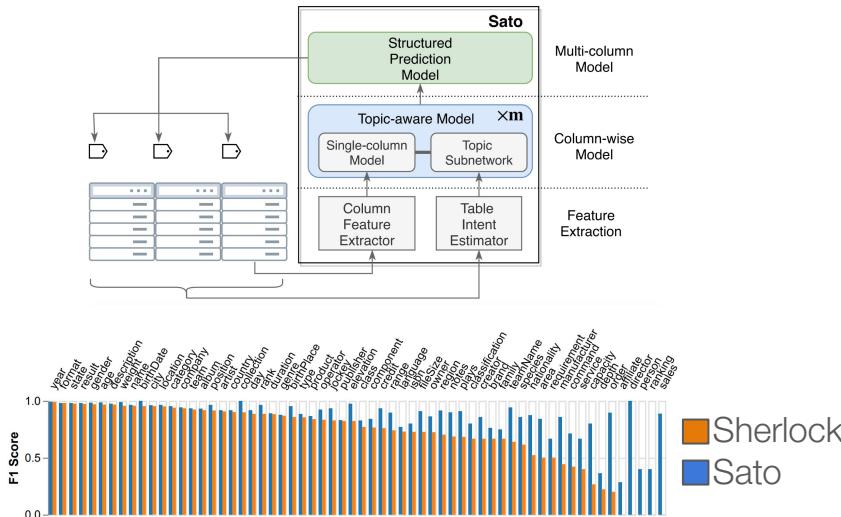
Challenges:

- Numeric data
- Non-mutually exclusive types

Deployed (often in healthcare), benchmarked against, and extended...

Extending Sherlock: **Contextual** and **Adaptive** Models

SATO: condition Sherlock-predicted column type on preds of neighbor columns



AdaTyper: adapt base type detection model by generating labeled samples from example columns

Type	regex	dictionary	AdaTyper		
			i=0	i=5	$\Delta F1$
first name	0	0.508	-	0.580	+0.580
postal code	0.108	0	-	0.068	+0.068
city	0.195	0	0.296	0.387	+0.091
gender	0	0.098	0.340	0.341	+0.001

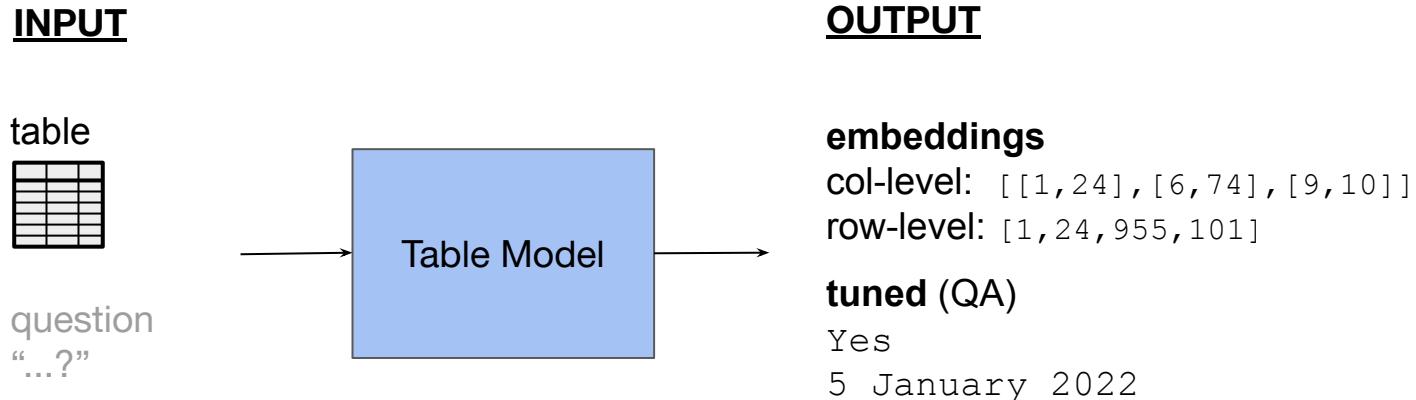
Effective w/ only few examples
But numeric/pattern? Regex..

Sato: Contextual Semantic Type Detection in Tables. Zhang, Suhara, Li, Hulsebos et al, VLDB, 2020.

AdaTyper: Adaptive Semantic Column Type Detection. Hulsebos et al, arxiv, 2023.

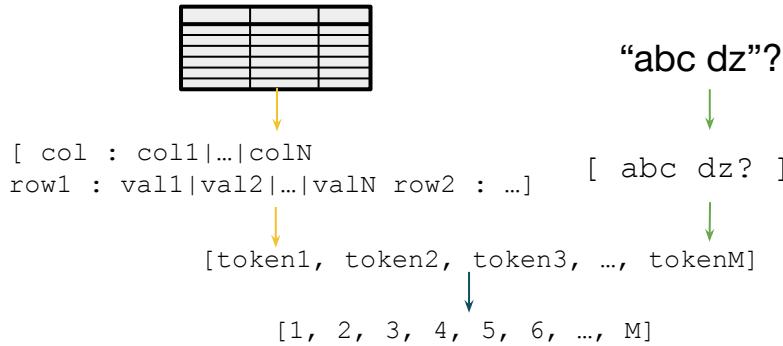
Then came **Transformers for Tables**

Transformers for Tables



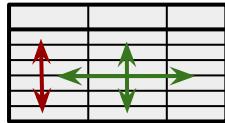
Transformers for Tables

Input



Model

Structured attention:



Vertical
Matrix

Output

Embeddings:

token-level agg to cell-, col-, row-level.

Serialize (e.g. row-wise or col-wise)

Merge tokens table & context (e.g queries)

Map tokens to “token IDs”

Pre-training tasks:

Default: recover column names or cell values.

Efficient: synthesized SQL execution.

Fine-tuning:

Predicting cells+operators, SQL, etc

Transformers Not Always SOTA...

Problem: pretrained table models poor OOD performance (e.g. col type prediction).

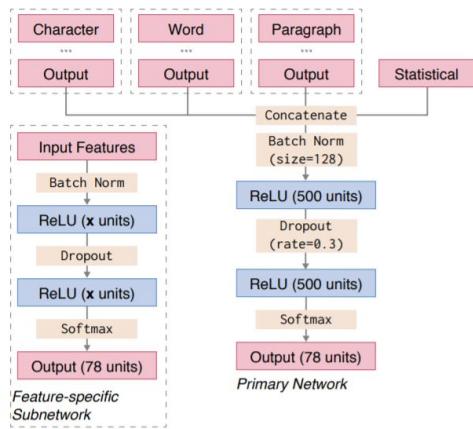
Just use GPT for col type prediction? [1]

	F_1 -score	Precision	Recall	
DoDuo-VizNet*	0.900	90.3%	89.9%	
Sherlock*	0.930	92.2%	93.1%	← Specific DL model
TaBERT	0.380	38.9%	38.3%	
DoDuo-Wiki	0.815	82.6%	81.4%	
CHORUS	0.865	90.1%	86.7%	← GPT-based

[1] CHORUS: Foundation Models for Unified Data Discovery and Exploration. Kayali, et al. VLDB, 2024.

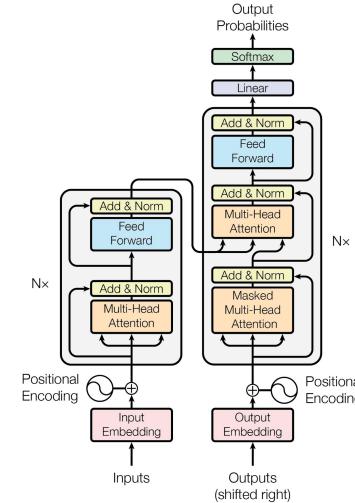
LLMs for tables: Overkill or Unutilized Potential?

Sherlock col type pred model:



hundreds of params
VS
billions/trillions of params

Transformer architecture:



But LLMs are promising... if they'd work well for tables:

How to handle messy data, large tables, full DBs, vague headers, numeric data?

Resources for **TRL**

What Data Do We Need?

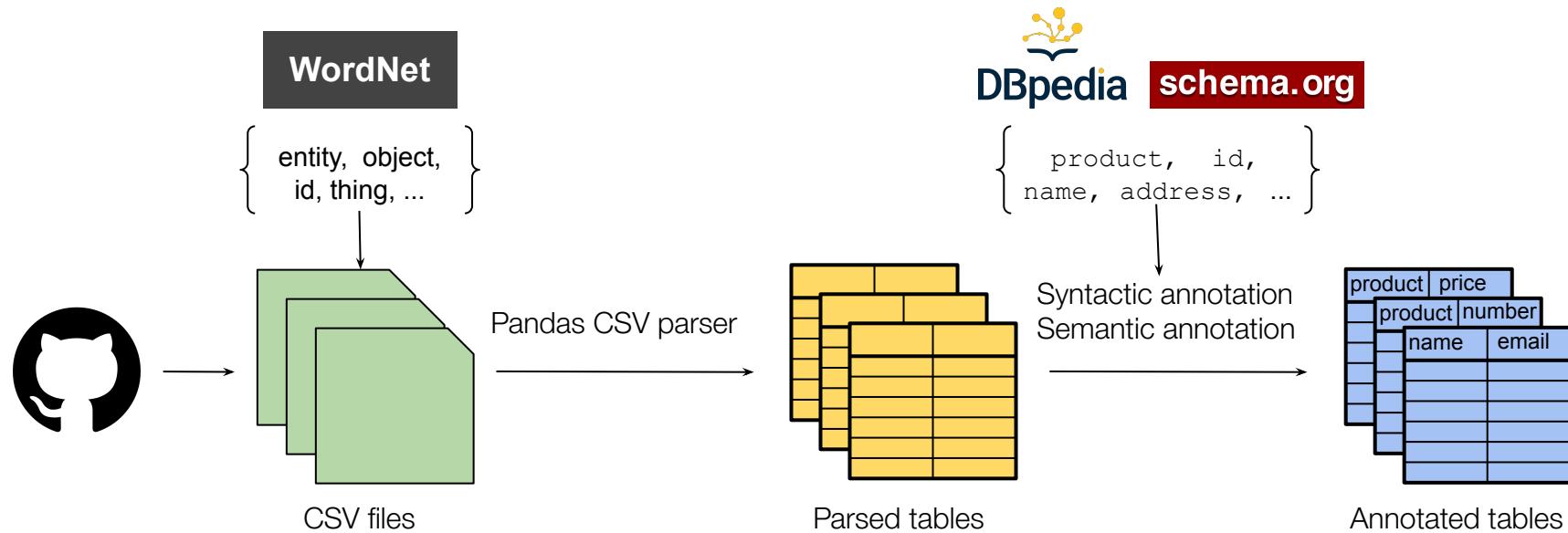
- Web/WikiTables → Web applications. Web tables * DB tables...
- Data tasks on offline tables? GitHub as a data source?

The screenshot shows a GitHub search results page with the following details:

- Search Query:** extension:"csv" "id"
- Results Count:** 15,768,996 code results
- Repository:** Kreef123/Sendy-Logistics-Challenge
- File:** data/Riders.csv
- Preview of CSV Data:**

1	Rider	Id	No_of_Orders	Age	Average_Rating	No_of_Ratings
2	Rider	Id	396,2946,2298,14,1159			
3	Rider	Id	479,360,951,13.5,176			
4	Rider	Id	648,1746,821,14.3,466			
5	Rider	Id	753,314,980,12.5,75			
6	Rider	Id	335,536,1113,13.7,156			
- CSV Link:** CSV
- Last Indexed:** 27 Mar 2021

GitTables



Properties and Use-cases of GitTables

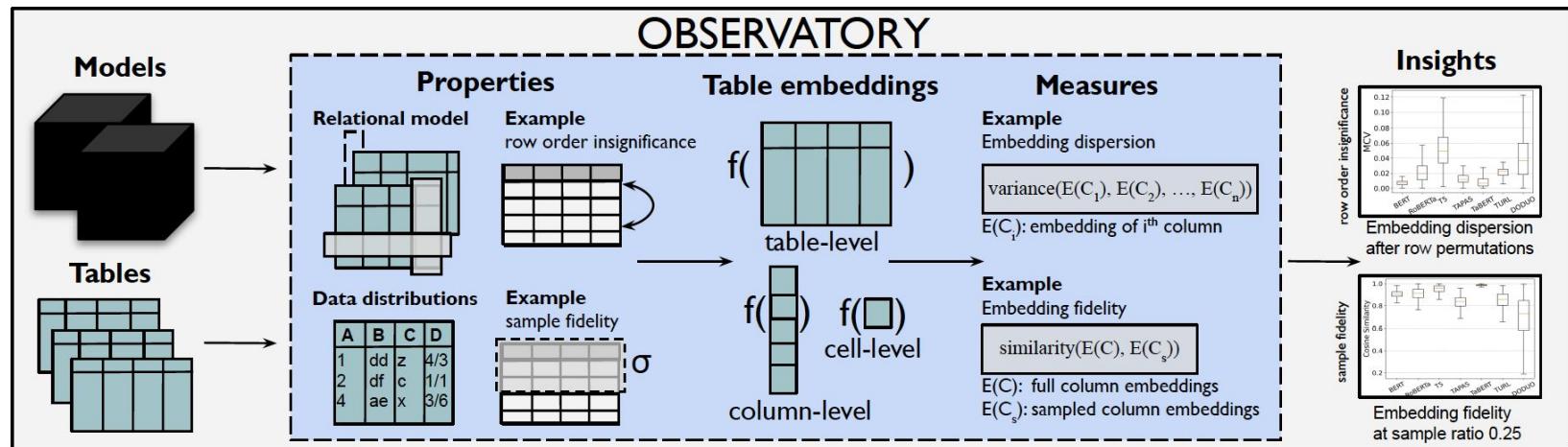
- **>1M tables** and **800K CSV** files.
- More representative: wider+taller, and IDs most common attribute.
- Usage shown for semantic **column type detection** and **schema completion**:

Header prefix	Suggested completion
payment_id, customer_id	review_id , product_id , product_parent , product_title , ...
id, company	ReceivablePaymentHeader , ReceivablePayment , Status , Customer , BankEntity , ...
id, name, location	phone , email , uid , active , ad_organization_id , ...

Also used for join discovery, CSV parsing, KG enhancement, etc.

Do Current Models Capture Relational Properties?

Neural table embeddings through the lens of Codd's relational model.



Example Property: Functional Dependencies

Given table T with FD: $X=\text{country} \rightarrow Y=\text{continent}$

We argue that:

- FD relation interpretable as *translation* between embeddings $E(\pi X(s))$ and $E(\pi Y(s))$
- Model f preserves FD if $d(E(\pi X(s)), E(\pi Y(s))) = d(E(\pi X(t)), E(\pi Y(t)))$ where d preserves norm+direction (L1/L2-norm).

ID	name	country	continent
1	Kathryn	Netherlands	Europe
2	Oscar	Netherlands	Europe
3	Lee	Canada	North America
4	Roxanne	USA	North America
5	Fern	Netherlands	Europe
6	Raphael	USA	North America
7	Rob	USA	North America
8	Ismail	Canada	North America

Measure the average group-wise variance over all n “FD-groups”:

$$\overline{s^2} = \frac{1}{n} \sum_{j=1}^n \frac{\sum_{i=1}^{m_{G_j}} \|d_{ji} - \bar{d}_j\|_2^2}{m_{G_j} - 1}$$

$\overline{s^2}$ approaches 0 if the translation between group-wise FD value pairs in X (country) and Y (continent) is approx. equal. At least $\overline{s^2}$ is smaller than in non-FD value pairs.



Current Architectures Often Fall Short...

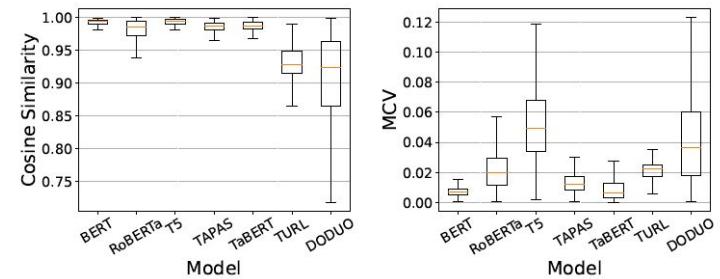
Turns out, most models do not preserve FDs!

Also more fundamental properties:

A *relation* then consists of a set of tuples, each tuple having the same set of attributes. If the domains are all simple, such a relation has a tabular representation with the following properties.

- (1) There is no duplication of rows (tuples).
- (2) Row order is insignificant.
- (3) Column (attribute) order is insignificant.
- (4) All table entries are atomic values.

Measure by avg cosine similarity of col embeddings across row permutations.



row order robustness

Impact on downstream tasks! **Row shuffling affects 34% col type predictions.**

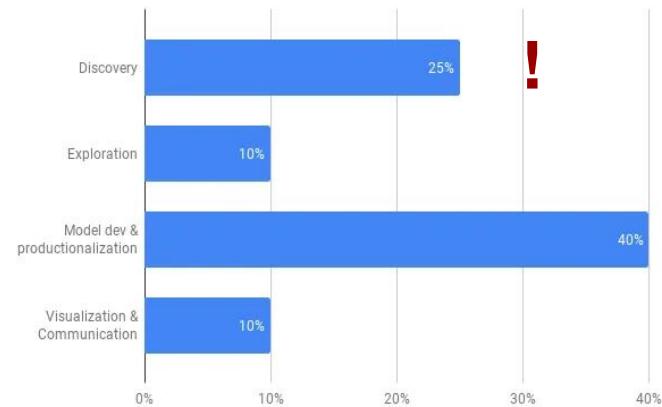
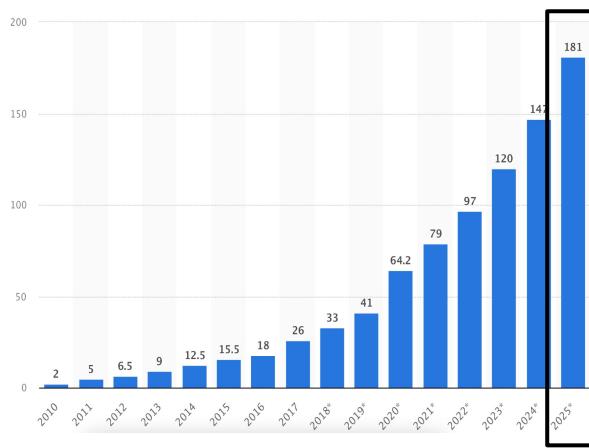
[2] Extending the Database Relational Model to Capture More Meaning, Codd, 1979.

Embeddings -> retrieval systems for structured data

Use-case: dataset search for analytics/ML

Immense growth of data → desire for insights

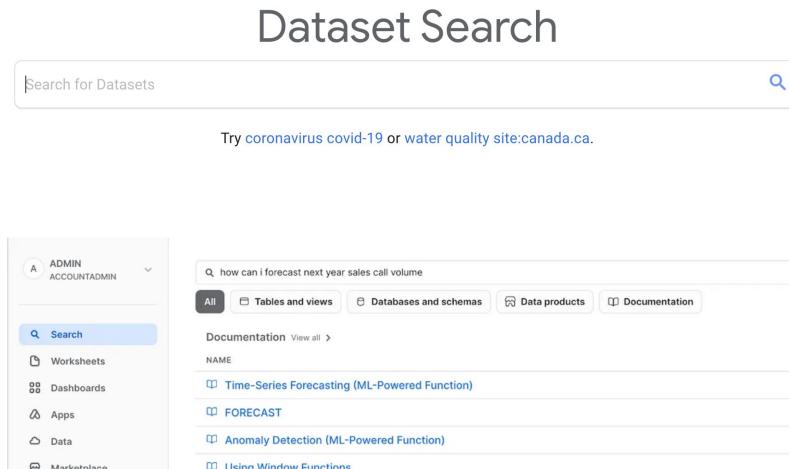
Finding the right dataset = still time-consuming



Research versus Practice

Systems in **industry** vs focus in **research**

Dataset Search



Search for Datasets

Try [coronavirus covid-19](#) or [water quality site:canada.ca](#).

ADMIN ACCOUNTADMIN

Search Worksheets Dashboards Apps Data Marketplace

Documentation View all >

NAME

- Time-Series Forecasting (ML-Powered Function)
- FORECAST
- Anomaly Detection (ML-Powered Function)
- Using Window Functions



“Basic” dataset search (e.g. keyword search)

Method	Task	Rep. Learning	ANN Index
Octopus [18]	KS	✗	✗
G.D.S. [2]	KS	✗	✗
Aurum [13]	KS	✗	LSH
LSH-Ensemble [3]	Join	✗	LSH
Juneau [4]	Join	✗	✗
JOSIE [5]	Join	✗	✗
MATE [6]	Join	✗	XASH
DeepJoin [7]	Join	✓	HNSW
D ³ L [14]	Union, Join	✗	LSH
Starmie [8]	Union, Join	✓	LSH, HNSW
TUS [9]	Union	✗	LSH
SANTOS [10]	Union	✗	✗
TURL [12]	TU	✓	✗
Sherlock [11]	TU	✓	✗
SATO [19]	TU	✓	✗

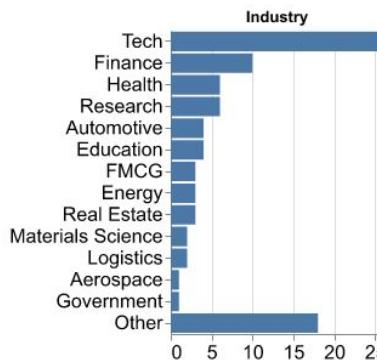
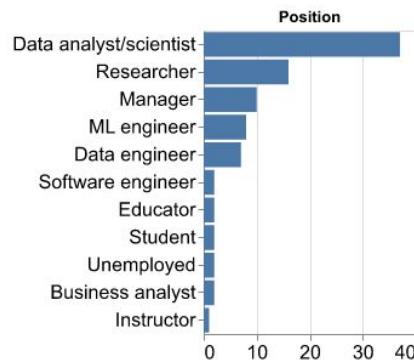
Majority research on data augmentation

Basic Dataset Search is not a Solved Problem!

We asked ~~ourselves~~: **why is dataset search *still* so hard in practice?**

89 data practitioners!!

recruited through social media & mailing lists:



We asked:

- What and how they search?
- What challenges they face?
- How they *want* to search?

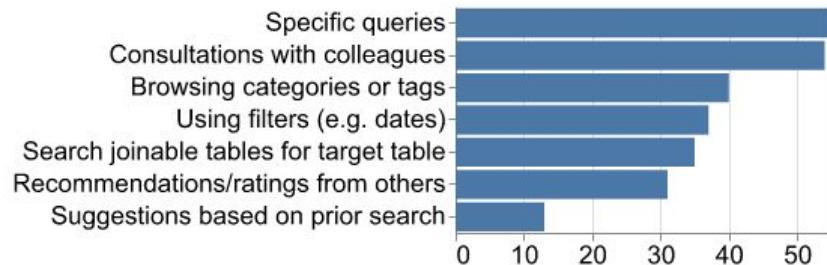


"It Took Longer than I was Expecting:" Why is Dataset Search Still so Hard?. Hulsebos, et al, HILDA@SIGMOD, 2024.

Practitioner's perspective: what and how they search

79% searches for **initial dataset**, 52% for **data enrichment**.

How do you search?

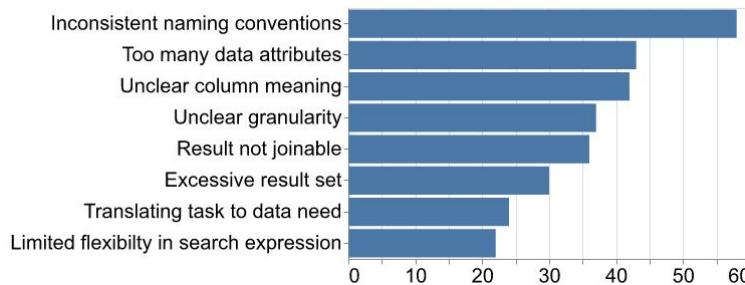


*"Identify the **problem and the data for the problem**, ... then specific keyword or tag search. Also, **identify people** who have worked on similar problems..."*

*"Having **so many tables**, I ask more experienced colleagues **which ones are most inherent to the analysis** I need to do. I then navigate through the categories and tags to look for others."*

Practitioner's perspective: key challenges

Key challenges with existing systems?



*"The biggest challenge I've noticed is **messy variable naming** - it takes me a long time to unpack what each variable means...."*

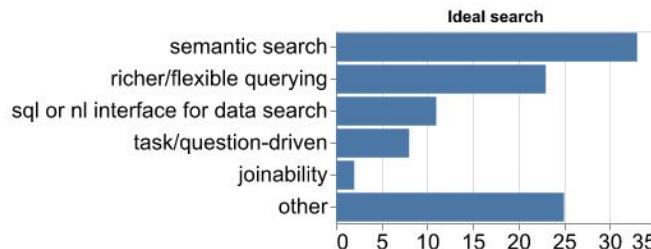
*"**Categorical level of detailing** is required, which isn't possible now."*

*"There are **too many table results** after the initial search...."*

*"Not many features to search/query keywords, a lot of times **changing query still renders same data results...**"*

Practitioner's perspective: ideal search systems

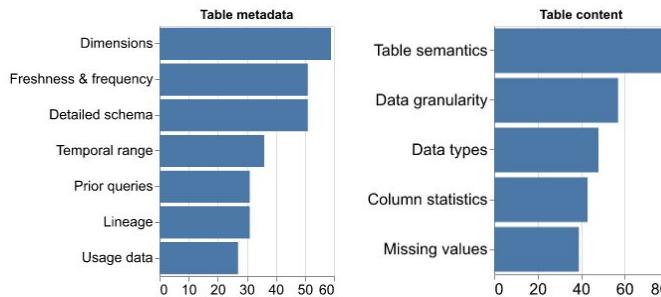
What should search systems facilitate?



"**Topic model search results**, based on sentence similarity with the dataset description."

"Ideally I would have something **across all of the various data sources and tables** and be able to use SQL (**or a trustable NLP solution**) and pull all relevant data and metadata."

What properties to search over?



"Show me **product usage** datasets where the main fact table is **event-level usage** data with **hundreds of millions** of records and there are dimension **tables for user and account**."

"Dataset to <**solve issue of ...**> with columns <1,2,3,...> on <**granularity desired**>"

Desiderata for Dataset Search

Task-driven: explicit **data needs often unknown** requiring back-and-forths w/ experts

Hybrid: search spans **multiple “views”** of a table; raw metadata + embeddings

Iterative: data search queries **don't fit a search bar**; complex process

Comprehensible and diverse results: result sets **hard to digest and navigate**

TBC...

Key takeaways...

- **Tables prevalent** in the data landscape, especially enterprises (eg for analytics).
- Capabilities of **transformer should extend beyond** images & text -> tables & DBs.
- Pre-train & tune table models on **representative data** (eg GitTables).
- **Tables ≠ natural language**: specific challenges and properties (Observatory).
- Important applications of table embeddings, e.g. **retrieval systems**.

Interested?

- Reach out: madelon@berkeley.edu
- TRL papers/resources: madelonhulsebos.com/trl
- Table Representation Learning workshop @ NeurIPS 2024??

Thanks!