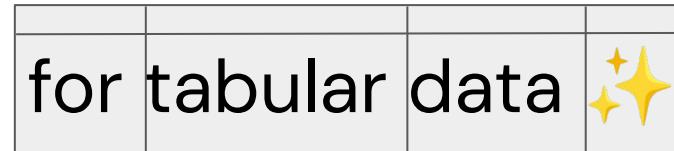


Representation Learning and Generative Models



Madelon Hulsebos (CWI)
TU Berlin
29 January

Agenda of today's lecture

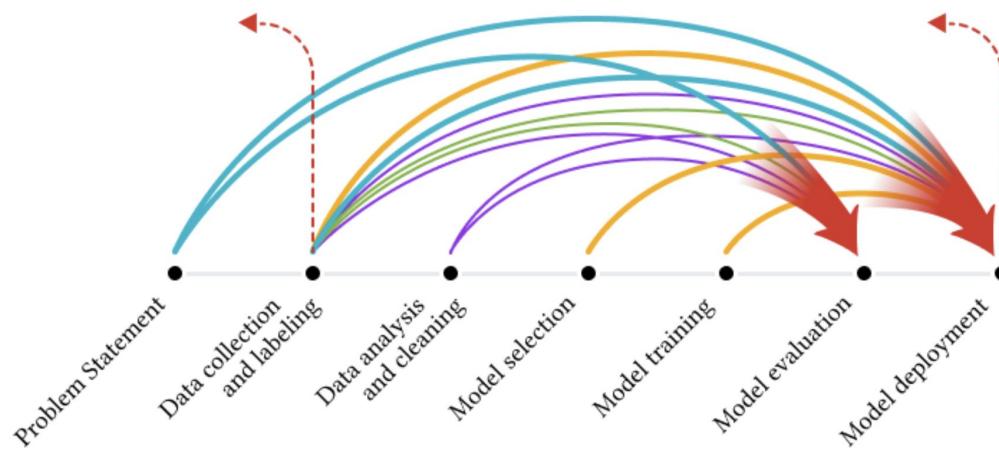
- Why ML for tables?
- Table Representation Learning
 - Background
 - TRL for “data work”
 - TRL for data insights
- Generative models and tabular data
 - Representation learning versus generative models
 - LLMs for (tabular) predictive modeling
 - Agentic data science systems
- Where are we, and where do we go?

Recap

ML pipelines: from raw data to analysis insights of ML model predictions.

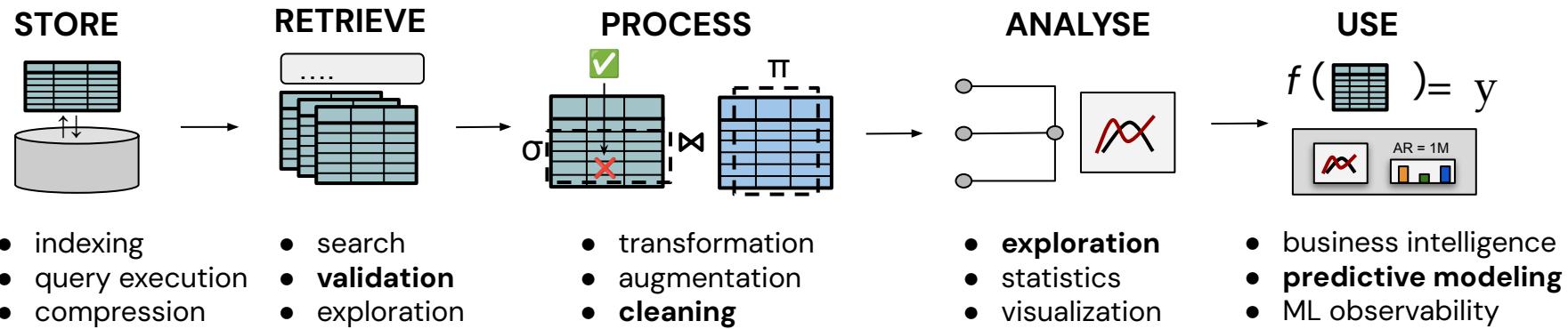
80% data work → what happens here, has huge impact!

20% model work



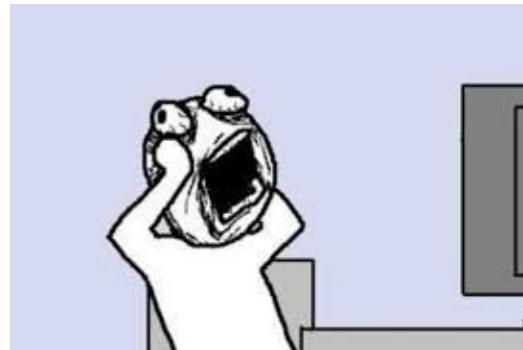
From "Everyone wants to do the model work, not the data work": Data Cascades in High-Stakes AI, Sambasivan et al., SIGCHI, 2021

Breaking down a Data Science pipeline



So much to think about!

So much to go wrong!



Then, I realized: **everyone is doing this....**

What if.... we could use ML to help us do the *data work*, for ML?

→ let's make data work, model work

→ ML for Data Engineering for ML

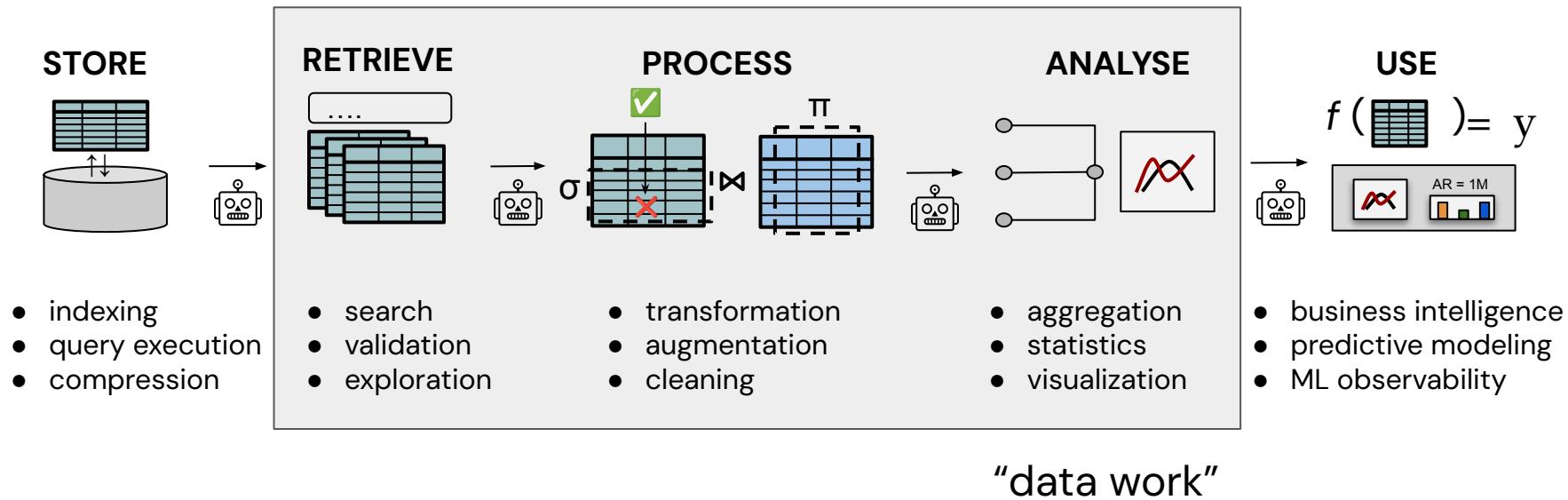


Why would this work?

- “everyone is doing this” = we have data, e.g. CSVs + jupyter notebooks
- If we have data, we can “learn”! ?

Automating data work for ML (predictive modeling)

The ambition...



Why “ML for Tables”?

From language and vision models → table models



✓

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 of 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:

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

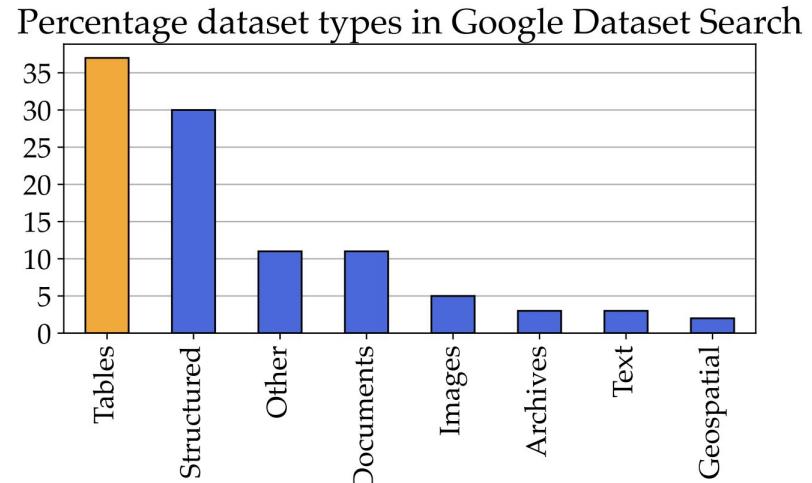
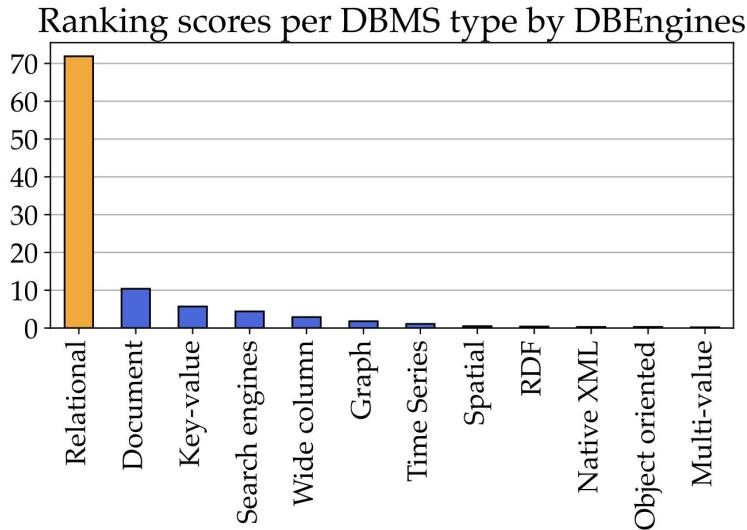
✓

Nr	ID	seed rate	yield	crop	cultivar	pre crop	pre-pre crop	pre-pre-pre
1	68	91	winter wheat	sugar beets	beans			
2	68	100	winter wheat	sugar beets	rotation fallow			
3	68	97	winter wheat	sugar beets	fallow land (5,5y)			
4	136	95	winter wheat	oats	sugar beets			
5	136	96	winter wheat	potatos	sugar beets			
6	136	107	winter wheat	sugar beets	maize			
7	136	107	winter wheat	sugar beets	summer wheat	maize		
8	136	82	winter wheat	oats	sugar beets	sugar beets		
9	136	77	winter wheat	potatos	sugar beets			
10	136	85	winter wheat	sugar beets	maize	maize		
11	136	84	winter wheat	sugar beets	summer wheat	sugar beets		
12	57 371	98	winter wheat	Sperber	sugar beets	winter barley	winter wheat	
13	57 365	98	winter wheat	Sperber	potatos	sugar beets	summer barley	
14	57 365	105	winter wheat	Sperber	sugar beets	maize	maize	
15	57 365	97	winter wheat	Sperber	sugar beets	winter wheat	sugar beets	
16	39 433	90	winter wheat	Okapi	summer barley			
17	39 433	100	winter wheat	Okapi	oats			
18	39 433	97	winter wheat	Okapi	winter wheat			

?

Tables are Everywhere

Data modalities in the real-world data landscape



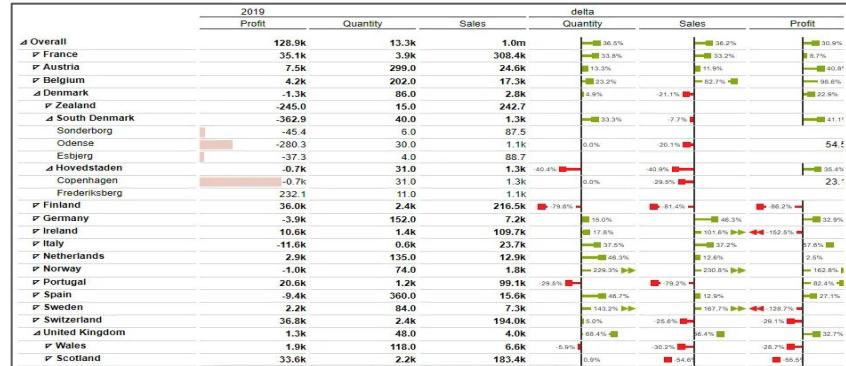
For a reason: tables serve high-value applications, e.g. data analysis & predictive modeling

A Challenge of Heterogeneity...

Tables store lots of *structured, fresh, domain* data!

Tables come in all shapes, semantics and sizes...

crop rotation : Tabelle												
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	100	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		
6	136	107	winter wheat		sugar beets	maize			sandy loam, loe 636	9.5	1991-1994	
7	136	107	winter wheat		sugar beetsn	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	9.5	wb, ww	
13	57 365	98	winter wheat	Sperber	potatos	sugar beets	summer barley		sandy loam, loe 635	9.5	cultivation, weed	
14	57 365	105	winter wheat	Sperber	sugar beets	maize	maize		sandy loam, loe 635	9.5	1987-1992	
15	57 365	97	winter wheat	Sperber	sugar beets	winter wheat	sugar beets		sandy loam, loe 635	9.5		
16	39 433	90	winter wheat	Okapi	summer barley				sandy loam, loe 690	8.5	oats, cultivation, weed	
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	

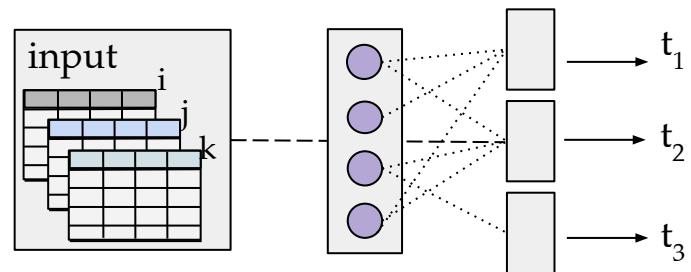


Challenging... how to deal with variation?

Representation Learning for tabular data

Table Representation Learning

Map each **table** to some consistent input.
Learn some **representation** that helps
detect patterns relevant to given task(s).



TRL for data work

Table Semantics Are All You Need

A table's understanding comes through its columns.

ML task:

- given table T ,
- predict semantic column types C ,
- with each c in C from preset ontology.

name	salary	country
name	sal	cntr

Semantic column types dictate operations *sensible* to perform on them:

name	salary	cntr



naam	status	land

Inform semantic join on tables

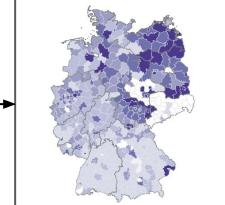
name
Xi
Carl
sara

name
Xi
Carl
Sara

Capitalize "name" columns

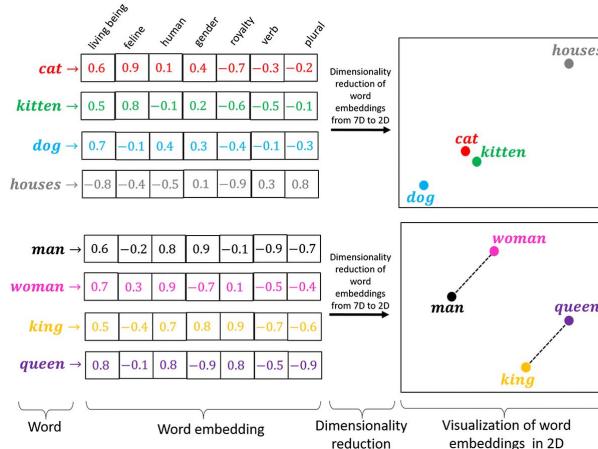
name	salary	cntr

Plot "country" data

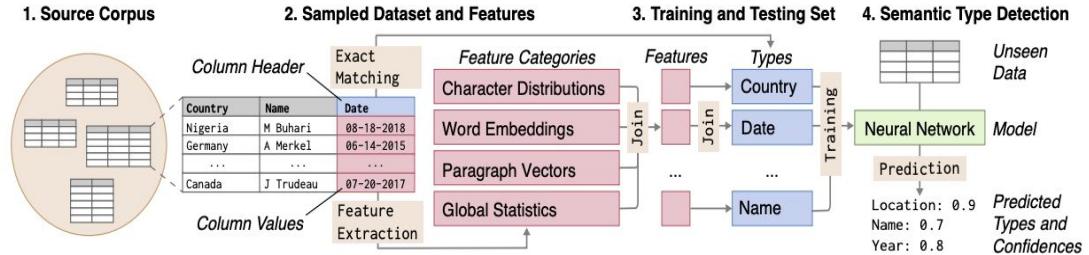


Semantic column types

Word embeddings: represent “words” in numeric vector space reflecting semantics



Starting point: treat a column as set of strings



Fast Forward to Transformers

- Bottleneck of existing Deep Learning models:
 - Don't take in much context (1 input column -> 1 output label)
 - Not very scalable
- Transformer architecture: **attention** mechanism enables “contextual learning” in parallel!

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*

Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukasz.kaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Start of LMs.. read!

Transformers for Tables

High level pipeline

More context, more efficient → lower level training!

Example task: question answering over tables

INPUT

table



question

“...?”



OUTPUT

embeddings

`[[1,24],[6,74],[9,10]]
[1,24,955,101]`

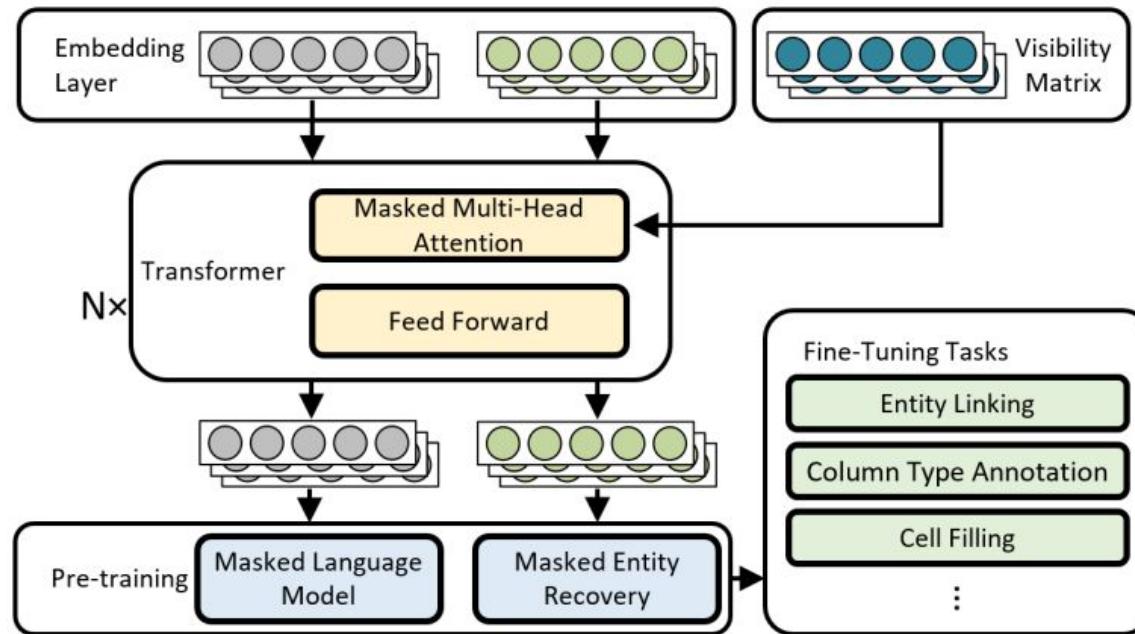
label (tuned model)

Yes

5 January 2022

“Table Model”: TURL

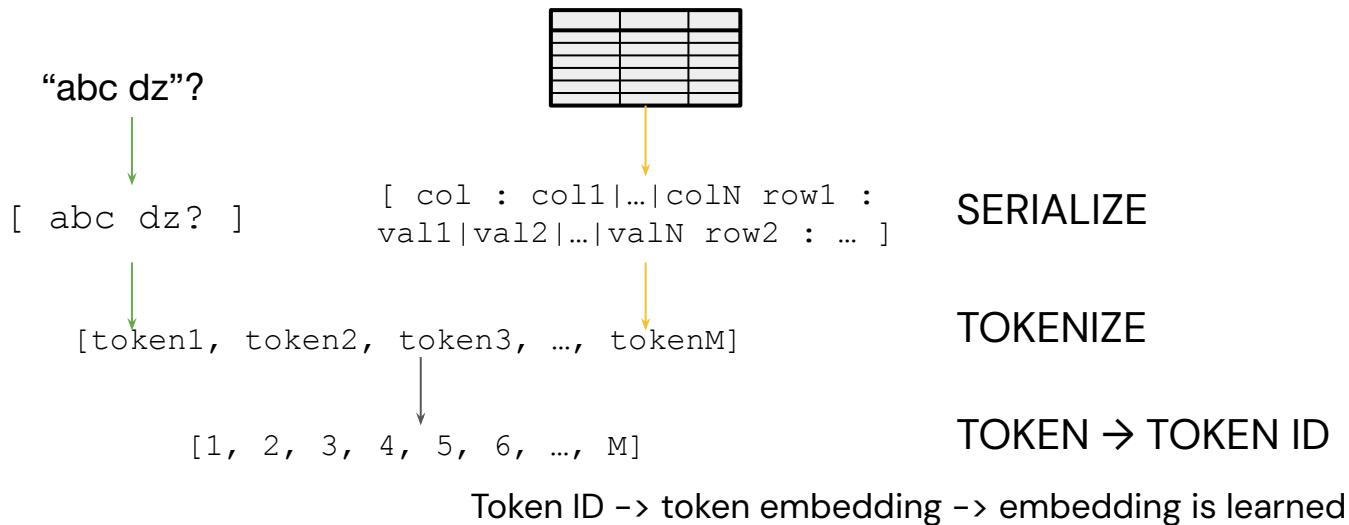
Low level inputs
and architecture



Transformers for Tables

Input: from Table to Tokens

- Sample, serialize, and tokenize(+pad) table.
- Table can be aligned w/ metadata or other input (if any).
- Many variations for serialization (e.g. row-wise, w/ SEP tokens etc).



Transformers for Tables

Architecture: Learning Adjusted for Table Structure

- Attention learns across all tokens (context) in input text.
- But Mrinal has not much to do with Goopy → structural attention:
→ Structural attention: vertical = across column or matrix = across row/col (TURL).



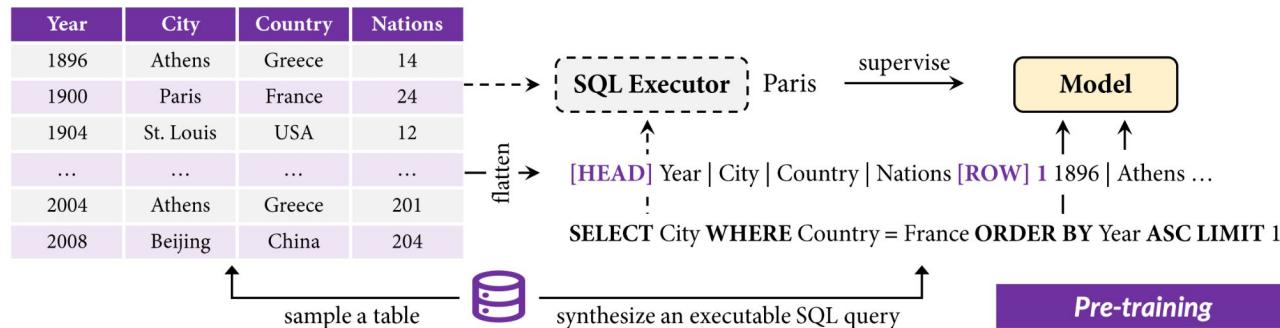
Transformers for Tables

Pre-training: Table-specific Tasks

Pretraining because the goal is to obtain “generic model” that can be “fine-tuned” for various tasks (using “embeddings” of inputs)

- Pre-training tasks:

- Typical: recovering (predicting) column names or cell values.
- Efficient: (synthesized) SQL execution (\downarrow TaPEX [1]).

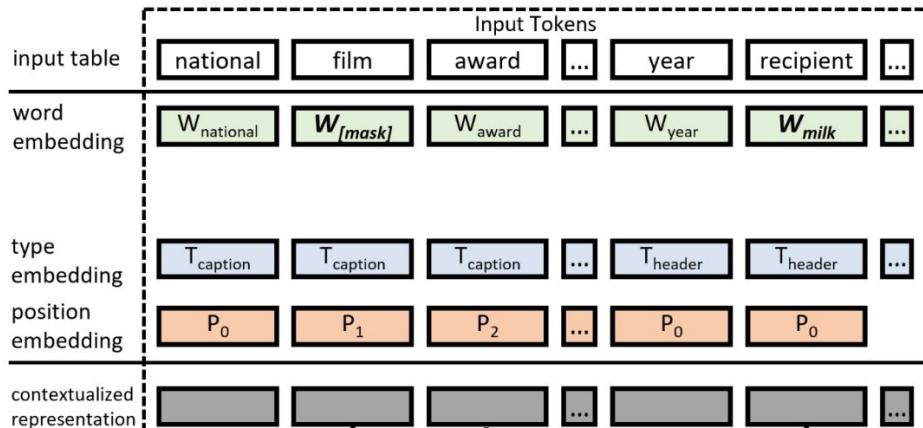


*Liu et al, TAPEX: Table pre-training via learning a neural SQL executor, ICLR, 2022.

Transformers for Tables

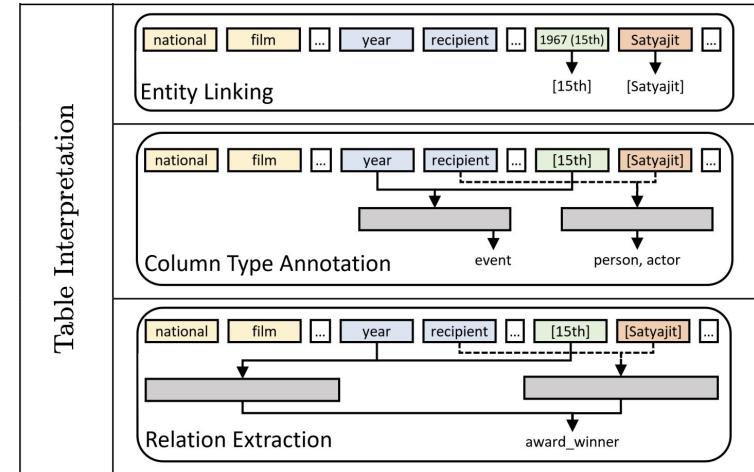
Output: Embeddings or Predictions

Embeddings



Typically aggregated from token-level embeddings
to cell/row/col level

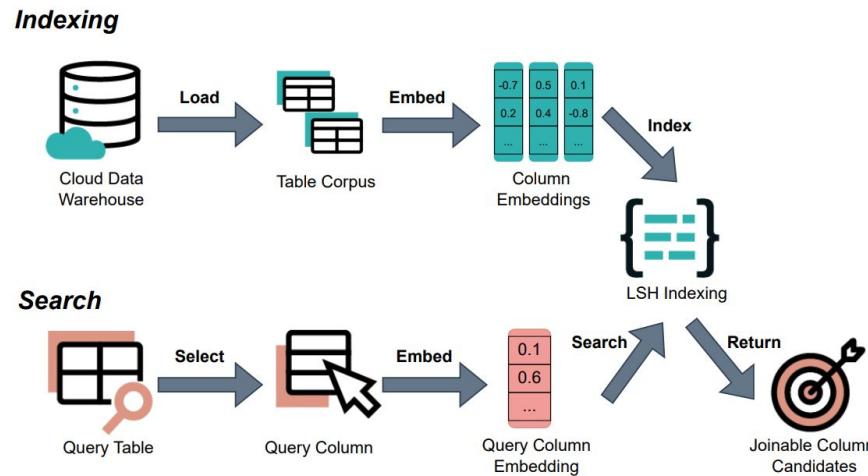
Predictions



These are “fine-tuned” from embeddings
to explicit labels

Representation Learning for Join Discovery

Task: given input table, find joinable tables for given column.



But embeddings used for retrieval, correlation prediction, data validation, etc.

Cong et al., WarpGate: A Semantic Join Discovery System for Cloud Data Warehouses, CIDR, 2023

TRL for data insights

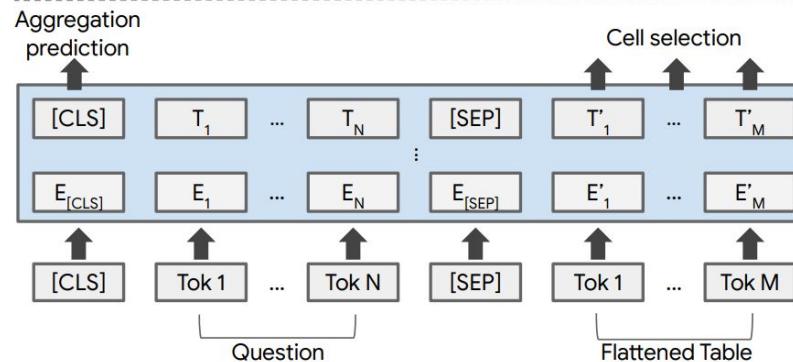
TRL for Question Answering over Tables

TaPas: TRL model for QA predicting operator + cell span

op	$P_a(op)$	$\text{compute}(op, P_s, T)$
NONE	0	-
COUNT	0.1	.9 + .9 + .2 = 2
SUM	0.8	.9×37 + .9×31 + .2×15 = 64.2
AVG	0.1	64.2 ÷ 2 = 32.1

$$S_{\text{pred}} = .1 \times 2 + .8 \times 64.2 + .1 \times 32.1 = 54.8$$

Rank	...	Days	P_s
1	...	37	0.9
2	...	31	0.9
3	...	17	0
4	...	15	0.2
...	0



Herzig et al, TAPAS: Weakly Supervised Table Parsing via Pre-training, ACL, 2020.

Works Well... for Basic Cases

Example with TaPas:

Actor	Age	Number of movies
Brad Pitt	56	87
Leonardo Di Caprio	45	53
George Clooney	59	69

How old is Leonardo Di Caprio? AVERAGE 45

What is Leonardo Di Caprio his age? AVERAGE 45

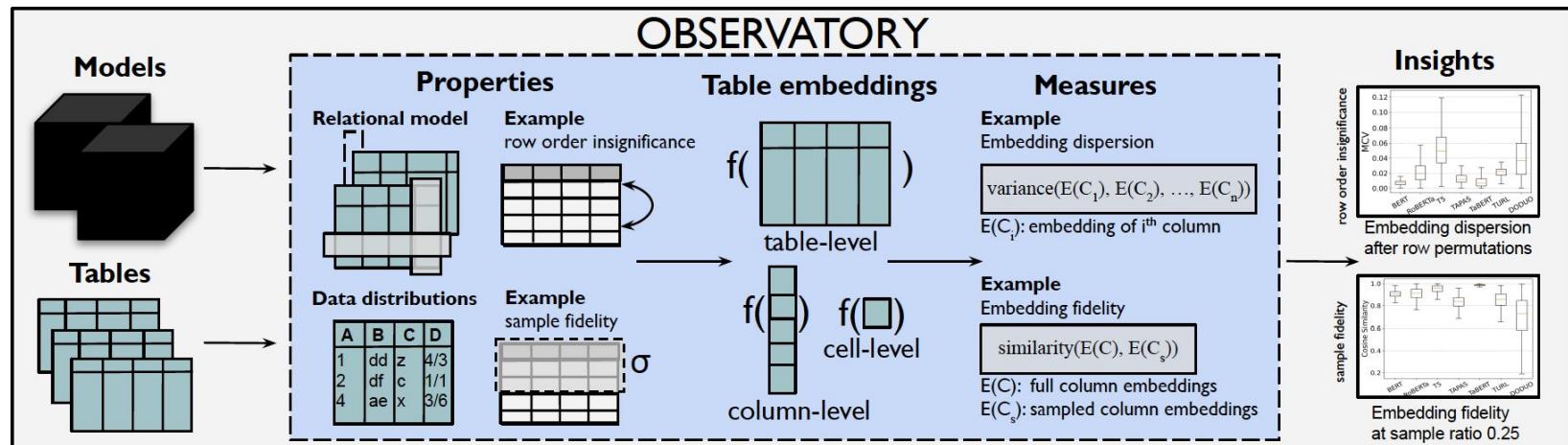
What is the sum of the number of movies? SUM $87, 53, 69 = 209$

How many movies are there in total? COUNT 87, 53, 69 = 3

Do Table Embeddings capture Relational Properties?

Tables ≠ natural language ?

Studying neural table embeddings through **Codd's relational model**.



Observatory: Characterizing Relational Table Embeddings. Cong, Hulsebos, Sun, Groth, Jagadish, VLDB, 2024.

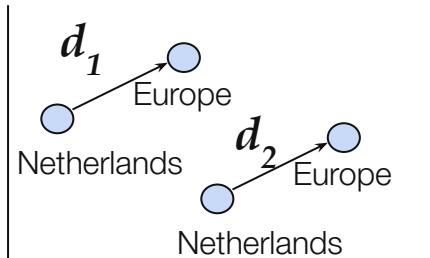
Example Property: Functional Dependencies

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

We argue that:

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

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



Current Architectures Fall Short...

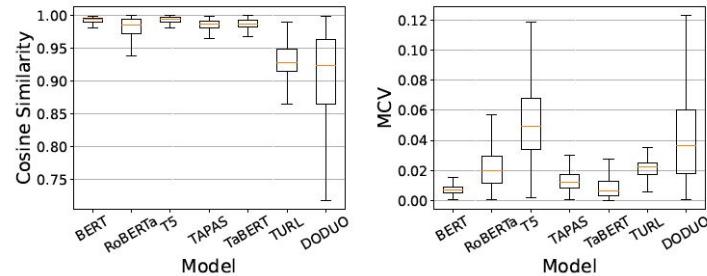
Turns out, most models do not preserve FDs!

We also consider simpler 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 downstream tasks: **row shuffling affects 34% semantic column types!**

Generative Models for tabular data

Representation Learning vs Generative Models

From **predicting** labels, e.g.:

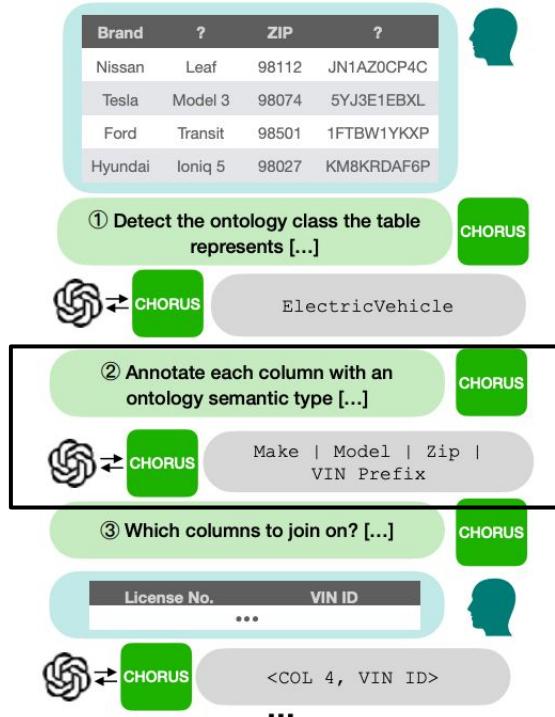
- semantic types / relations between columns,
- cell span + aggregations (QA),
- *True or False* (fact verification).

To **generating** answers...!

- Input: table, context, query (question / task / anything)
- Output: anything (e.g. code, or explicit answers)
- *Underlying* mechanism: **next-token prediction**

“Underlying” because “underlying” LM might be “tuned” to predict discrete labels.

Can LLMs help with data discovery?



Consider this example. Input:

```

Name, Famous Book, Rk, Year

Fyodor Dostoevsky, Crime and Punishment, 22.5, 1866

Mark Twain, Adventures of Huckleberry Finn, 53, 1884

Albert Camus, The Stranger, -23, 1942

```

Output:

`dbo:author, dbo:title, Unknown, dbo:releaseDate`.

For the following CSV sample, suggest a DBpedia.org Property for each column from the `dbo:` namespace.
``` [...] ```

(b) Column-type annotation

Typical format:

- Instructions
- Example (input, output)

# Transformers Not Always SOTA

Problem: pretrained TRL models poor OOD performance on *col type prediction*.

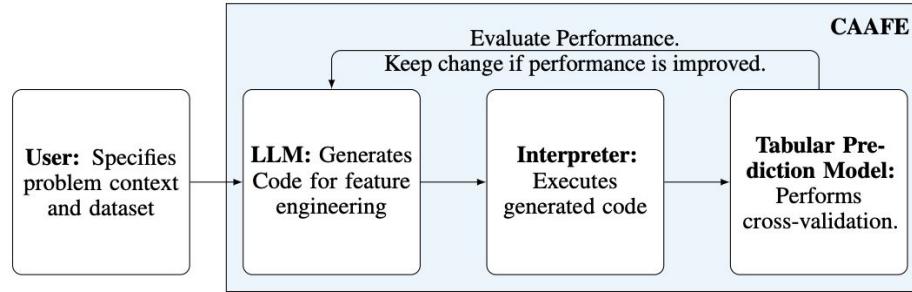
Just use generative model? Sherlock outperforms LLMs.

	<b><math>F_1</math>-score</b>	<b>Precision</b>	<b>Recall</b>
DoDuo-VizNet*	0.900	90.3%	89.9%
Sherlock*	<b>0.930</b>	<b>92.2%</b>	<b>93.1%</b>
TaBERT	0.380	38.9%	38.3%
DoDuo-Wiki	0.815	82.6%	81.4%
<b>CHORUS</b>	<b>0.865</b>	<b>90.1%</b>	<b>86.7%</b>

LLM analyses show tables best formatted w HTML tags, but many challenges.  
Messy data? Large tables? Full DBs? Non-descriptive headers? Numeric data?

# Beyond SQL: feature engineering!

General approach:



Example, binning:

```
Feature: AgeGroup (categorizes passengers into age groups)
Usefulness: Different age groups might have different likelihoods
of being transported.
Input samples: 'Age': [30.0, 0.0, 37.0]
bins = [0, 12, 18, 35, 60, 100]
labels = ['Child', 'Teen', 'YoungAdult', 'Adult', 'Senior']
df['AgeGroup'] = pd.cut(df['Age'], bins=bins, labels=labels)
df['AgeGroup'] = df['AgeGroup'].astype('category')
```

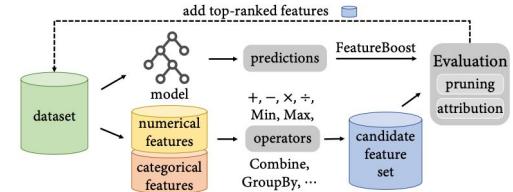
# Automated features helpful, but minimal gain

No feature engineering at all...



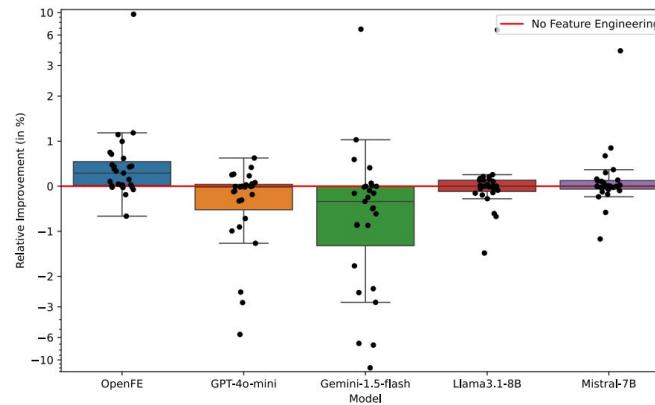
		No FE	Baselines				CAAFE GPT-4
			DFS	AutoFeat	FETCH	OpenFE	
Log. Reg.	Mean	0.749	0.764	0.754	0.76	0.757	0.763
	Mean Rank	27.4	23.6	26.2	25.2	25	24.8
Random Forest	Mean	0.782	0.783	0.783	0.785	0.785	0.79
	Mean Rank	23.4	22.1	21.8	23.5	22.3	23.1
ASKL2	Mean	0.807	0.801	0.808	0.807	0.806	0.815
	Mean Rank	12.2	12.9	12.6	13.4	13.5	10.9
Autogluon	Mean	0.796	0.799	0.797	0.787	0.798	0.803
	Mean Rank	17.6	15.4	16.4	17.6	16.6	15.8
TabPFN	Mean	0.798	0.791	0.796	0.796	0.798	0.806
	Mean Rank	13.9	15	14.8	16.5	13.9	12.9

Generating feature pool,  
prune feature candidates



# What's going on?

Language models engineer too many simple features...



Better with domain expertise? Or much better training (data, tricks)....

# Make LMs significantly better for tabular tasks?

Train LLMs on tabular tasks, **at scale**:

- GPT-based model trained on 86B tokens
- >593.8K table+language samples for training encoder
- >2.36M query+table+output tuples for fine-tuning

Scale of evaluation:

- **23 benchmarking metrics**
- TableGPT2 **7B** model: +35.20% improvement
- TableGPT2 **72B** model: +49.32% improvement



## TableGPT2: A Large Multimodal Model with Tabular Data Integration

Aofeng Su, Awonen Wang, Chao Ye, Chen Zhou, Ga Zhang, Gang Chen, Guangcheng Zhu, Haobo Wang, Haokai Xu, Hao Chen, Haoze Li, Haoxuan Lan, Jiaming Tian, Jing Yuan, Junbo Zhao, Junlin Zhou, Kaizhe Shou, Liangyu Zha, Lin Long, Liyao Li, Pengzu Wu, Qi Zhang, Qingyi Huang, Saisai Yang, Tao Zhang, Wentao Ye, Wufang Zhu, Xiaomeng Hu, Xijun Gu, Xinjie Sun, Xiang Li, Yuhang Yang, Zhiqing Xiao

Authors are ordered alphabetically by the first name.

Zhejiang University Institute of Computing Innovation, Zhejiang University

### Abstract

The emergence of models like GPTs, Claude, LLaMA, and Qwen has reshaped AI applications, presenting vast new opportunities across industries. Yet, the integration of tabular data remains notably underdeveloped, despite its foundational role in numerous real-world domains.

This gap is critical for three main reasons. First, database or data warehouse data integration is essential for advanced applications; second, the vast and largely untapped resource of tabular data offers immense potential for analysis; and third, the business intelligence domain specifically demands adaptable, precise solutions that many current LLMs may struggle to provide.

In response, we introduce TableGPT2, a model rigorously pre-trained and fine-tuned with over **593.8K** tables and **2.36M** high quality query-table-output tuples, a

# Impression of scale

Benchmark	Metric	GPT-4o	TableLLM (Qwen2)	TableLLM (CodeQwen)	TableLLM (LLaMA3)	TableLLM (LLaMA3.1)	TableLLM (DeepSeek)	TableLLM-13B	DeepSeek-lite	Yi-Coder	Qwen2.5-Coder	Qwen2.5-Instruct	TableGPT2-7B	TableGPT2-72B
<i>Table Understanding</i>														
Col Type Annot.	F1	31.75	10.10	5.71	1.47	1.59	6.04	12.70	20.58	5.38	32.59	22.19	<b>85.88</b>	85.67
Relation Extract.	F1	52.95	1.60	3.79	2.39	2.00	3.34	18.16	8.67	2.25	31.00	15.92	<b>83.35</b>	79.50
Entity Linking	Acc	90.80	47.10	39.70	0.20	0.60	15.50	66.25	70.15	41.75	71.70	82.25	92.00	<b>93.30</b>
Row Pop.	MAP	53.40	2.20	5.14	1.93	6.23	3.13	14.25	1.20	1.00	13.23	12.30	<b>59.97</b>	55.83
<i>Question Answering</i>														
HiTab	Exec Acc	48.40	11.74	0.00	0.00	0.00	39.08	6.30	0.76	0.00	1.70	10.73	70.27	<b>75.57</b>
FetaQA	BLEU	21.70	12.24	8.69	2.42	3.10	7.94	10.83	15.08	11.17	13.00	16.91	28.97	<b>32.25</b>
HybridQA	Acc	58.60	27.12	20.14	27.35	27.61	19.53	51.88	42.58	29.83	51.10	51.13	53.17	<b>56.41</b>
WikiSQL	Acc	47.60	46.50	37.20	39.26	39.00	36.14	41.10	38.30	25.34	46.90	47.42	53.74	<b>57.32</b>
WikiTQ	Acc	68.40	64.16	36.05	34.95	38.84	36.05	66.30	47.65	43.37	<b>74.50</b>	68.55	61.42	71.45
<i>Fact Verification</i>														
TabFact	Acc	74.40	72.00	53.20	40.06	27.13	60.76	68.95	62.27	79.6	77.26	84.60	77.80	<b>85.43</b>
FEVEROUS	Acc	71.60	20.10	46.90	51.50	42.30	18.39	21.45	7.80	38.10	60.70	63.30	<b>78.05</b>	76.80
<i>Table to Text</i>														
ToTTo	BLEU	12.21	6.95	3.10	5.50	6.23	3.81	5.36	8.76	2.64	10.50	11.91	14.10	<b>22.69</b>
<i>Natural Language to SQL</i>														
BIRD(dev)	Exec Acc	-	9.13	7.37	1.83	2.48	0.39	0.72	25.10	24.19	27.18	18.97	31.42	<b>38.40</b>
BIRD(dev-knowledge)	Exec Acc	-	15.45	18.19	3.39	3.72	0.39	1.83	36.51	39.96	42.96	31.42	49.28	<b>60.76</b>
Spider(dev)	Exec Acc	-	42.26	32.88	12.86	18.96	2.71	4.26	66.44	58.12	70.99	61.70	76.31	<b>79.40</b>
Spider(test)	Exec Acc	-	40.29	34.93	12.02	16.35	7.33	2.93	66.65	56.87	69.73	60.18	74.38	<b>78.48</b>
<i>Holistic Table Evaluation</i>														
TableBench	DP	-	26.62	26.44	26.71	26.73	26.15	3.88	29.60	21.94	28.67	25.18	32.03	<b>38.90</b>
	TCoT	-	37.08	31.33	29.79	30.01	28.65	3.85	30.93	22.8	36.25	29.77	42.34	<b>50.06</b>
	SCoT	-	14.11	17.78	9.60	12.38	22.39	2.88	22.61	8.43	25.95	24.35	25.01	<b>30.47</b>
	PoT@1	-	21.05	26.39	31.96	25.80	28.39	2.94	10.90	11.36	16.15	22.58	<b>33.52</b>	28.98

# LLMs for predictive modeling

# Problem setting

$$f(\mathbf{X}) \rightarrow \mathbf{y}$$

Where  $\mathbf{X}$  are *features*, and  $\mathbf{y}$  is the *target* to predict.

Quite similar to missing value imputation, where  $\mathbf{y}_{\text{test}}$  are missing values?

# General approach

Given **generative model M** and a **table T**:

- Serialize rows in **T** into “sentences”
- Template the prediction “task”
- Fine-tune **M** on train set (where to-be-predicted labels are provided)
- Evaluate on test set (labels to-be-predicted)

→ Still “generation”, so prone to errors in hallucination, formatting, etc.

# TableLLM: few-shot LLMs for predictive modeling

## 1. Tabular data with $k$ labeled rows

age	education	gain	income
39	Bachelor	2174	$\leq 50K$
36	HS-grad	0	$> 50K$
64	12th	0	$\leq 50K$
29	Doctorate	1086	$> 50K$
42	Master	594	

## 2. Serialize feature names and values into natural-language string with different methods

### Manual Template

The age is 42. The education is Master. The gain is 594.

### Table-To-Text

The person is 42 years old. She has a Master. The gain is 594 dollars.

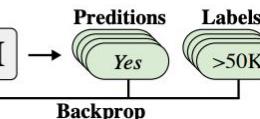
### LLM

The person is 42 years old and has a Master's degree. She gained \$594.

## 3. Add task-specific prompt

Does this person earn more than 50000 dollars? Yes or no? Answer:

## 4a. Fine-tune LLM using labeled examples

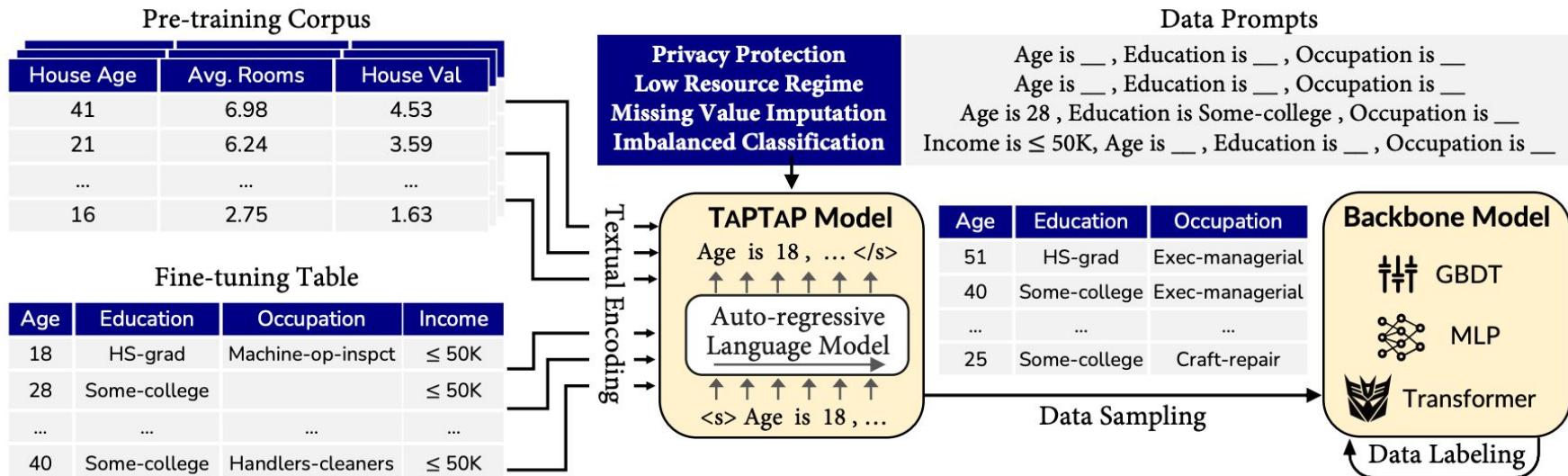


## 4b. Use LLM for prediction on unlabeled examples

The age is 42. The education is Master. The gain is 594.  
Does this person earn more than 50000 dollars? Yes or no?  
Answer:



# TapTap: Language Models for Predictive Models



Note: data augmentation with TapTap, improving robustness to invariant row/col order!

# Agentic Data Science

# Agentic Systems for Data Science

“Agentic”: the LLM-system has some “agency”, i.e. it plans what to do.

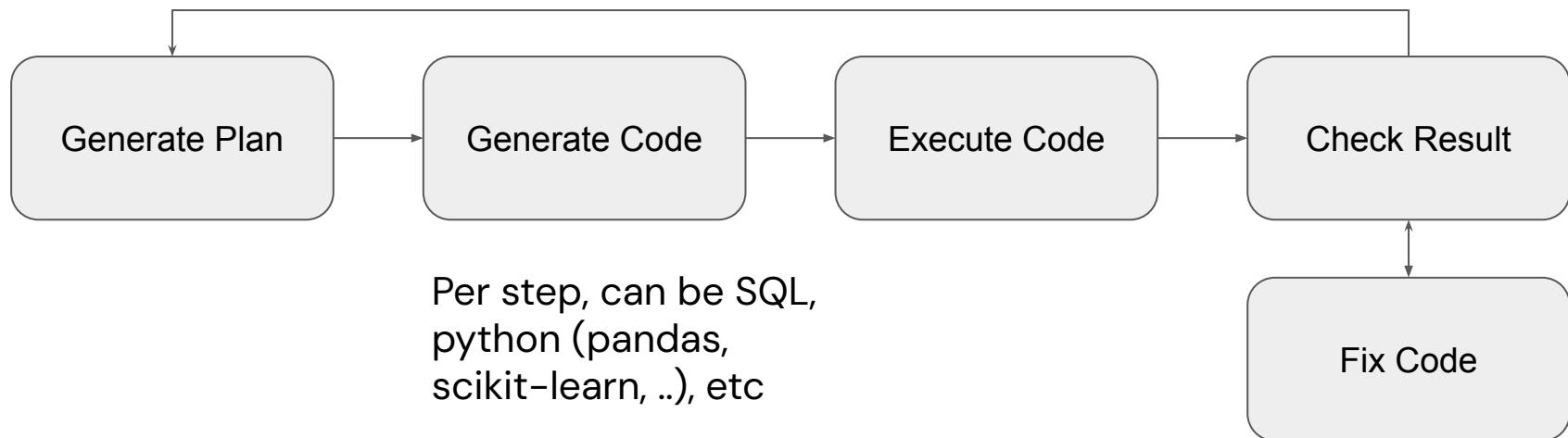
The screenshot shows a user interface for a "Data Science Agent". At the top, there's a purple header bar with a flask icon and the text "Data Science Agent" followed by a button labeled "Experiment". Below the header, there are two buttons: "README.md" and "Playground", with "Playground" being underlined, indicating it's the active tab. The main content area displays a file named "2018\_Central\_Park\_Squirrel\_Census\_-\_Squirrel\_Data\_20240501.csv". Below the file name is a descriptive text: "Analyze the proportion of adult and juvenile animals in the census data. Are there any spatial patterns in age distribution?". A call-to-action button labeled "2018\_Central\_Park\_Squirrel\_Census\_-\_Squirrel\_Data\_20240501.csv" is centered below the file. At the bottom, there's a "Plan" button with a dropdown arrow next to it.

**Pipeline of 8 steps, automated!** (8, but didn't even train/eval an ML model)

# One step, e.g. “data cleaning”

Data cleaning:

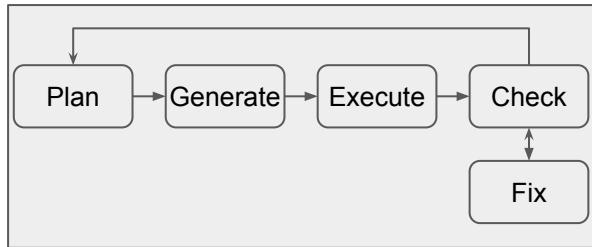
- Remove invalid values
  - Remove outliers
  - Impute missing values
  - ...
- Nice, an LLM can reason about what “invalid” would mean, examples?



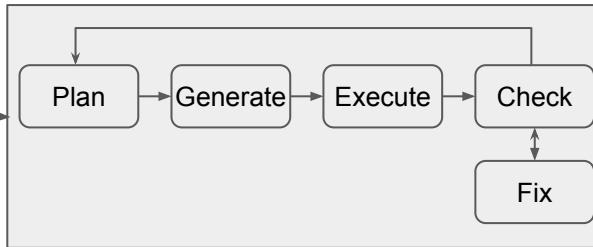
# What can possibly go wrong?

Analyze the proportion of adult and juvenile animals in the census data.  
Are there any spatial patterns in age distribution?

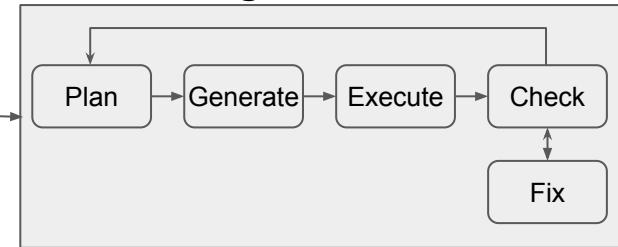
Validate data



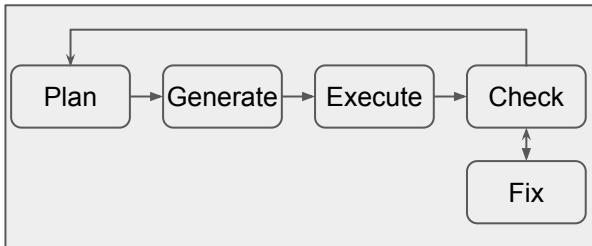
Clean data



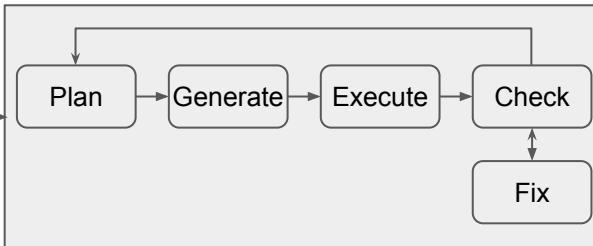
Integrate data



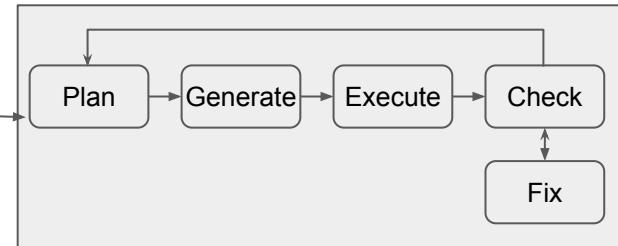
Summarize data



Aggregate data

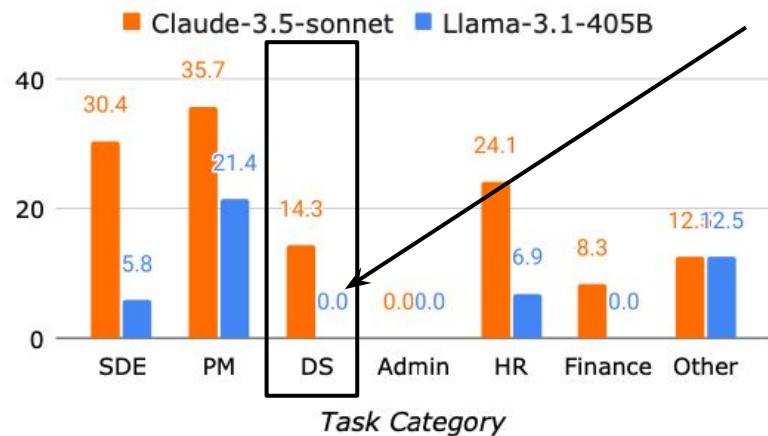


Visualize data



# How do agentic DS systems perform?

Realistic data science (DS) task



(b) Success rate across task categories

Model	SDE (69 tasks)		PM (28 tasks)		DS (14 tasks)	
	Success	Score	Success	Score	Success	Score
Closed model APIs						
Claude-3.5-Sonnet	30.43	38.02	35.71	51.31	14.29	21.70
Gemini-2.0-Flash	13.04	18.99	17.86	31.71	0.00	6.49
GPT-4o	13.04	19.18	17.86	32.27	0.00	4.70
Gemini-1.5-Pro	4.35	5.64	3.57	13.19	0.00	4.82
Amazon-Nova-Pro-v1	2.90	6.07	3.57	12.54	0.00	3.27
Open-weight models						
Llama-3.1-405b	5.80	11.33	21.43	35.62	0.00	5.42
Llama-3.3-70b	11.59	16.49	7.14	19.83	0.00	4.70
Qwen-2.5-72b	7.25	11.99	10.71	22.90	0.00	5.42
Llama-3.1-70b	1.45	4.77	3.57	15.16	0.00	5.42
Qwen-2-72b	2.90	3.68	0.00	7.44	0.00	4.70

# Some suggestions...

- Errors are costly!
  - Need for human interaction – “how” is an open question (reviewing code 🤔?)
  - Better interpretation (refinement) of input query.
- Generalizability is key
  - Robustness to variation (data, workflow needs)
  - Need to acknowledge limitations (current demo mode: “can do, will do!”)

But... promising!

# Key take-aways

- Potential of TRL & generative models for tables for data work!
- LLMs can do predictive modeling, reasonably
- We're moving towards agentic Data Science systems

## More attention needed to:

- We need representative and large-scale datasets (hard to get!)
- We need specialized “tricks” (e.g. architecture, pretraining, tokenization, etc)
- We need better domain context and ways to fetch human guidance

Got interested?

Join us for a workshop on this topic on 27 February in Amsterdam:

## ***ELLIS workshop on Representation Learning and Generative Models for Structured Data***



e l l i s



UNIT  
AMSTERDAM

<https://sites.google.com/view/rl-and-gm-for-sd/home>

Or check: <https://www.madelonhulsebos.com/upcoming/>

# Thank you!



Madelon Hulsebos  
TRL Lab @ CWI

<https://www.madelonhulsebos.com>  
@madelonhulsebos on Bluesky