# Towards **Table** Representation Learning for End-to-End Data Management and Analysis

Hasso Plattner Institute, Potsdam
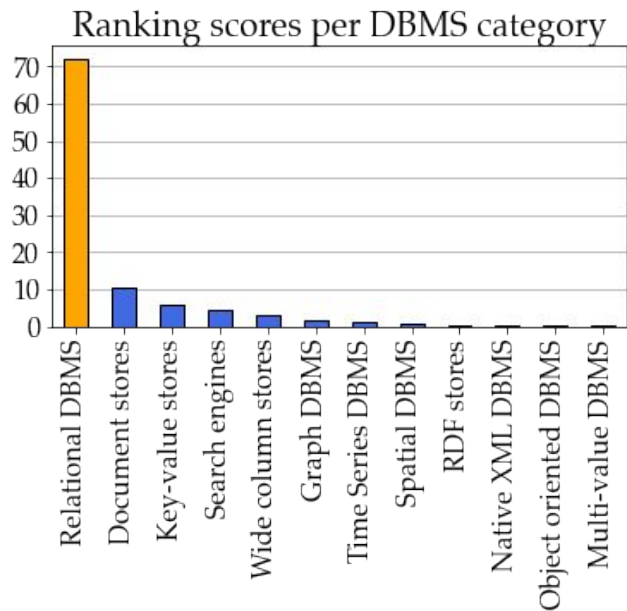06/03/2023

Madelon Hulsebos
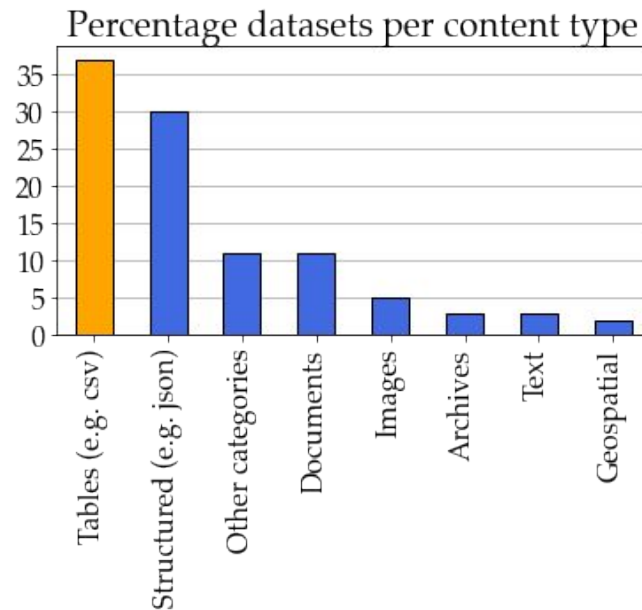
INDE lab

# Tables are **everywhere**

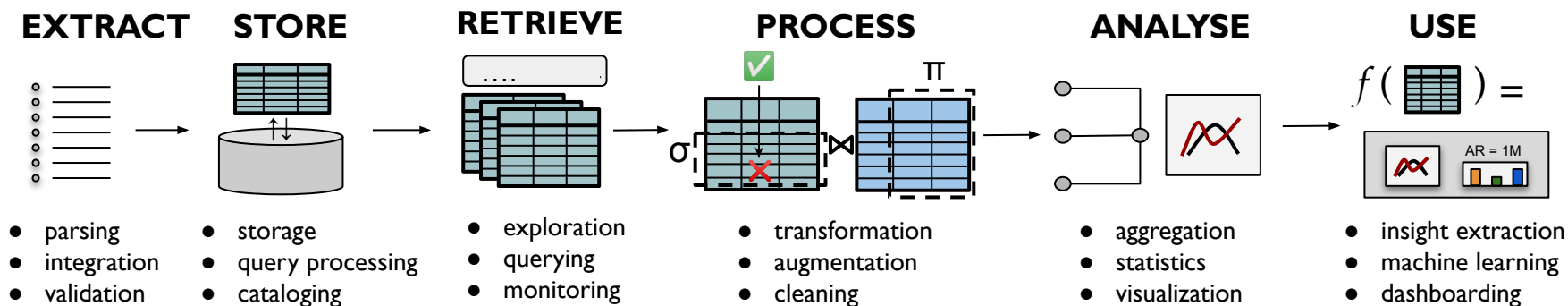Databases, (web) pages, documents, spreadsheets…



From DB-Engines Ranking by Category (Jan '23)



From Google Dataset Search by the Numbers (Benjelloun et al., '20).

# Tables are **driving many analysis pipelines**

End-to-end pipelines involve tons of applications.

**EXTRACT**
- parsing
- integration
- validation

**STORE**
- storage
- query processing
- cataloging

**RETRIEVE**
- exploration
- querying
- monitoring

**PROCESS**
- transformation
- augmentation
- cleaning

**ANALYSE**
- aggregation
- statistics
- visualization

**USE**
- insight extraction
- machine learning
- dashboarding

As w/ images and text: can we learn table representations to fuel these pipelines?
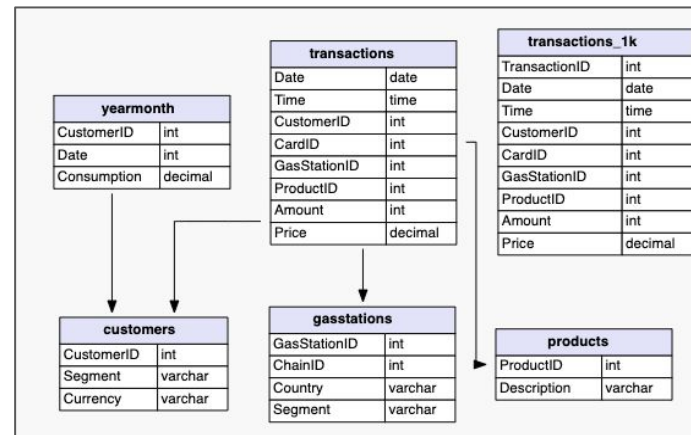
# Tables are **rich and challenging**

Content:   measurements, messy, heterogeneous dtypes.

Structure: columns, rows, cells, headers, hierarchical.

Context:   relations, constraints, metadata.

Usage:      analyses, ml models, visualizations.



CTU Prague Relational Learning Repository



First table when searching "crop data"



Published by a Tableau user

# Today: **learning over tables**

How to **represent** a table?

How to **understand** them?

How to **adapt** table models?

How to find good **data**?

And beyond…



Images, videos, text…

Tables

# 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?

As in other type systems, semantic column types dictate operations to perform on them.

| name | salary | cntr |
|------|--------|------|
|      |        |      |
|      |        |      |
|      |        |      |
|      |        |      |

⋈

| naam | status | land |
|------|--------|------|
|      |        |      |
|      |        |      |
|      |        |      |
|      |        |      |

Join tables on "name" and "country" columns

| **name** |
|----------|
| Xi Yu |
| carl bert |
| Sara zi |

→

| **name** |
|----------|
| Xi Yu |
| Carl Bert |
| Sara Zi |

Capitalize "name" columns

| name | salary | cntr |
|------|--------|------|
|      |        |      |
|      |        |      |
|      |        |      |
|      |        |      |

→ 

Plot "country" data

# Column type detection: **how?**

**Matching** header or values by ① matching column values, ② aggregating to types.

In commercial systems (e.g. Tableau):

- Preset regular expressions.
- Preset type:values dictionary.

SOTA:

- Ontology-based [1].
- Extracted rules from GitHub [2].

[1] Recovering semantics of tables on the web. Petros et al, 2011
[2] Synthesizing type-detection logic for rich semantic data types using open-source code, Cong Yan and Yeye He, 2018.

What if we remove column names?



Detected Types With Column Headers

| Country/Region | String | Latitude | Longitude | Country/Region | String |
|---|---|---|---|---|---|
| country-capitals.csv Country Name | country-capitals.csv Capital Name | country-capit... Latitude | country-capital... Longitude | country-capitals.csv Country Code | country-capitals.csv Continent Name |
| Aruba | Oranjestad | 12.517 | -70.033 | AW | North America |
| Australia | Canberra | -35.267 | 149.133 | AU | Australia |
| Austria | Vienna | 48.200 | 16.367 | AT | Europe |

Detected Types Without Column Headers

| String | String | Decimal | Decimal | String | String |
|---|---|---|---|---|---|
| country-capitals-edite... F1 | country-capitals-edi... F2 | country-capit... F3 | country-capital... F4 | country-capitals-edite... F5 | country-capitals-edited... F6 |

Remove Headers

# Column type detection: **Sherlock**

**Scale**, **robustness**, **accuracy**?



Published at KDD 2019

# Can **Sherlock** detect types?

Evaluated on >600K columns from Web tables.

78 semantic types (`name`, `address`, etc).

| Method | $F_1$ Score | Runtime (s) | Size (Mb) |
|---|---|---|---|
| *Machine Learning* | | | |
| 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 |
| *Matching-based* | | | |
| Dictionary | 0.16 | 0.01 (±0.03) | 0.5 |
| Regular expression | 0.04 | 0.01 (±0.03) | 0.01 |
| *Crowdsourced Annotations* | | | |
| Consensus | 0.32 (±0.02) | 33.74 (±0.86) | – |

Current usage:

- Adopted in industry: health tech and fashion (e.g. data integration).
- People contributed bugfixes, speedups.
- Was extended to SATO (w context).
- Research benchmarks (competitive!).

Paper, model, data and code: https://sherlock.media.mit.edu

# In the wake of **Sherlock**

Pre-trained models for table understanding: large-scale training without ground-truth labels.

Industry feedback Sherlock: nice but **data mismatch**, cannot add **custom types.**

① How to transfer to new data domains?

② How to detect new types?

# What **data** do we need?

① How to transfer to new data domains?  → Why asked?

Tables needed:

- Large to facilitate learning  → WebTables [3] ✅
- Table semantics (e.g. col types) → WebTables ✅
- DB-like table content and structure (semantics, dtypes, size)  ❓
- Coverage to generalize across domains  ❓

| President | Party | Term as President | Vice-President |
|---|---|---|---|
| 1. George Washington (1732-1799) | None, Federalist | 1789-1797 | John Adams |
| 2. John Adams (1735-1826) | Federalist | 1797-1801 | Thomas Jefferson |
| 3. Thomas Jefferson (1743-1826) | Democratic-Republican | 1801-1809 | Aaron Burr, George Clinton |
| 4. James Madison (1751-1836) | Democratic-Republican | 1809-1817 | George Clinton, Elbridge Gerry |
| 5. James Monroe (1758-1831) | Democratic-Republican | 1817-1825 | Daniel Tompkins |
| 6. John Quincy Adams (1767-1848) | Democratic-Republican | 1825-1829 | John Calhoun |
| 7. Andrew Jackson (1767-1845) | Democrat | 1829-1837 | John Calhoun, Martin van Buren |
| 8. Martin van Buren (1782-1862) | Democrat | 1837-1841 | Richard Johnson |
| 9. William H. Harrison (1773-1841) | Whig | 1841 | John Tyler |
| 10. John Tyler (1790-1862) | Whig | 1841-1845 | |
| 11. James K. Polk (1795-1849) | Democrat | 1845-1849 | George Dallas |
| 12. Zachary Taylor (1784-1850) | Whig | 1849-1850 | Millard Fillmore |
| 13. Millard Fillmore (1800-1874) | Whig | 1850-1853 | |
| 14. Franklin Pierce (1804-1869) | Democrat | 1853-1857 | William King |
| 15. James Buchanan (1791-1868) | Democrat | 1857-1861 | John Breckinridge |

Table from a Web page.

[3] WebTables: exploring the power of tables on the web, Cafarella et al., 2008

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 | | potatoes | 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 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 | | potatoes | 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 | | sandy loam, loe | 635 | | cultivation, weed |
| 13 | 57 | 365 | 98 | winter wheat | Sperber | potatoes | sugar beets | summer barle | 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, 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 | |

Table with crop data, first result "example database table".

# Can we use GitHub CSV files?



Result from GitHub code search when querying for CSV files containing "id".

# The birth of **GitTables**



WordNet

{ entity, object, id, thing, ... }

querying CSV files

**GitHub repositories**

**CSV files**

CSV parsing
table filtering
content curation

**Curated tables**

DBpedia  schema.org

{ product,id, name,address,... }

syntactic annotation
semantic annotation

| product | price |
| --- | --- |

| product | number |
| --- | --- |

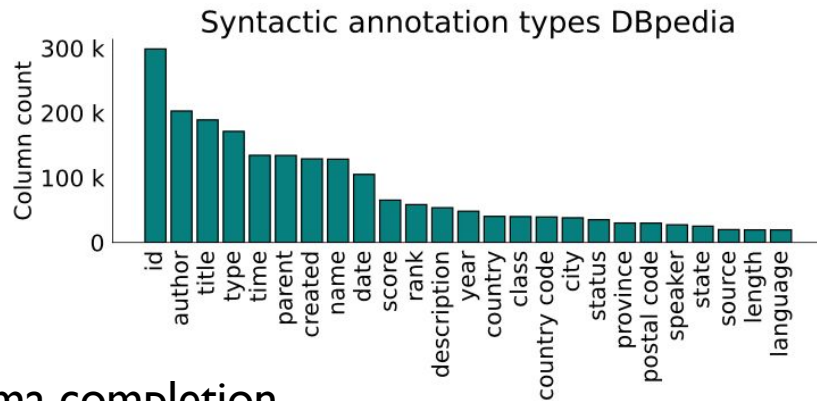| name | email |
| --- | --- |

**Annotated tables**

INDE lab

Published at SIGMOD 2023

# What can we do with **GitTables**

We publish >1M tables, also underlying >800K CSV files.

| Data type | GitTables | WDC WebTables |
|-----------|-----------|---------------|
| Numeric   | 57.9%     | 51.4%         |
| String    | 41.6%     | 47.4%         |
| Other     | 0.5%      | 1.2%          |



Syntactic annotation types DBpedia

We show:      ML for type detection and schema completion.

Other use:      join discovery, schema matching, benchmarking.

General Table Representations? E.g. parsing, compression, error repair?

Paper, data and code: https://gittables.github.io

# Adaptive type detection: **AdaTyper** [WIP]

② How to detect new types?

**Current**: by user-provided value dict or regular expression.
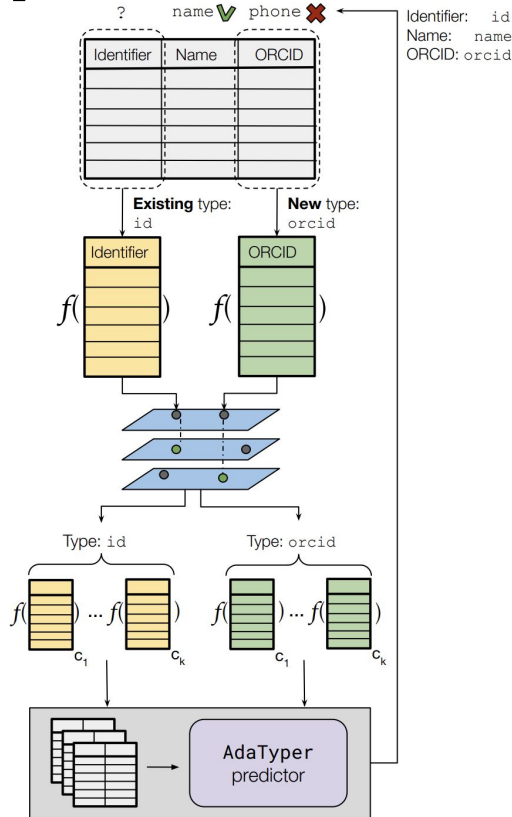
**Interactive adaptation by example**:

1. Predict initial column type.
2. User corrects with (new) type.
3. Embed example column.
4. Retrieve similar col embeddings from HNSW index [4].
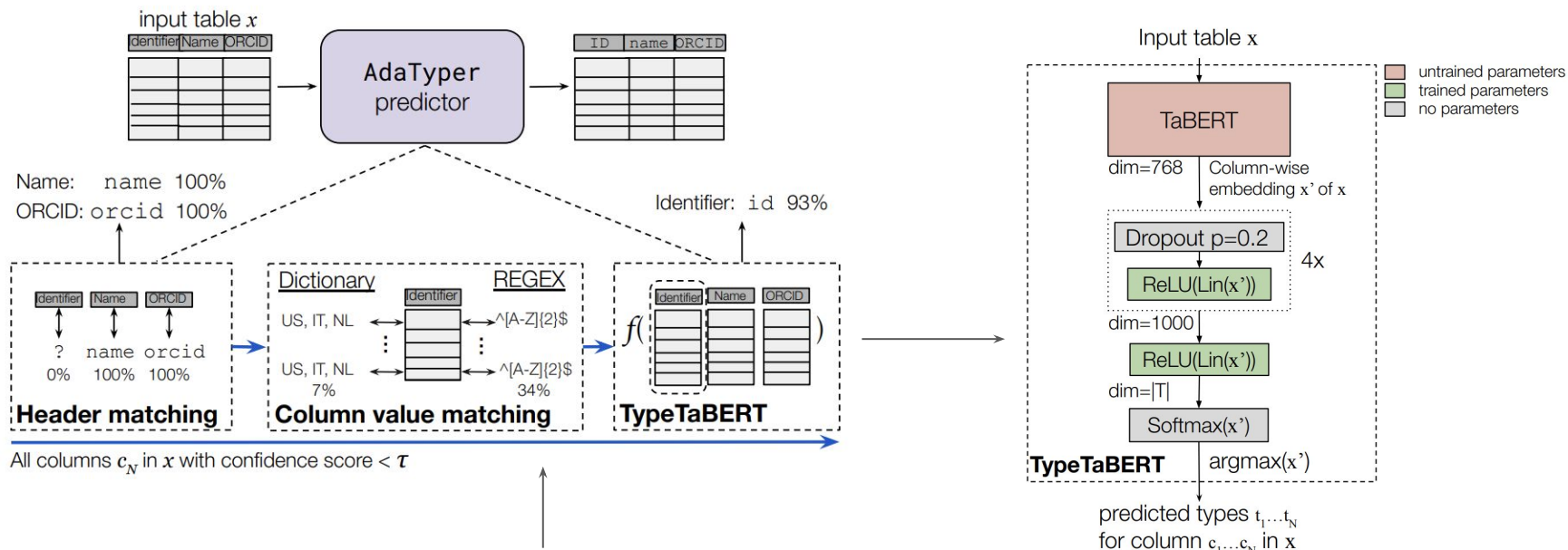5. Retrain type prediction model.



[4] Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs, Malkov and Yashunin, 2018

# **AdaTyper** predictor

**Hybrid** type detection pipeline enabling **different adaptation methods**.



So, we can still adapt through regular expressions….

# How well does **AdaTyper** adapt?

Measuring **performance after *x* examples** of new type

- Human annotated tables from Prague Relational Learning Repo (not used for training!).
- High precision.
- WIP: low recall, increase -> drop: issues w example diversity and label errors?
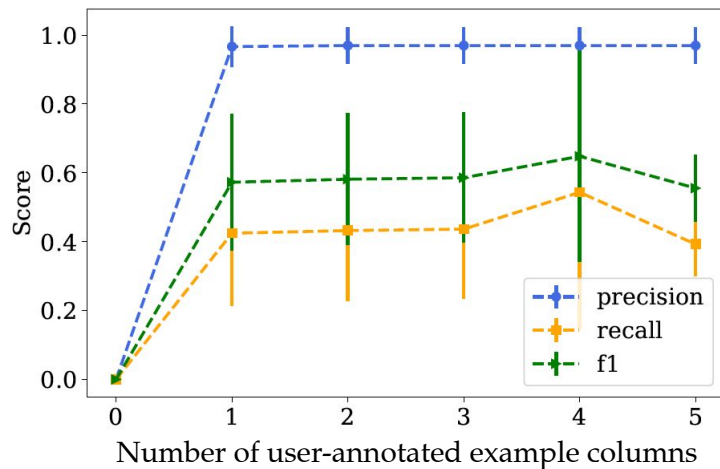
# Table on the **Horizon**

1. What *do* table representations capture? Now blindly adopting models for any task.

2. What *can* table representations capture within E2E pipeline?

    - Left and right, from storage & query optimization to analysis recommendation
    - Developing new neural architectures aligned with data management tasks
    - Contextualizing tables w.r.t. downstream usage

3. Table-specific deployment challenges.

# **Interested**?

New research area with many challenging problems and impactful applications!

Exciting community spanning different communities (e.g. NLP, DB, ML). Take part:

1.  **Join**: Dedicated TRL Slack space → reach out [m.hulsebos@uva.nl](mailto:m.hulsebos@uva.nl)!

2.  **Learn**: SIGMOD '23 Tutorial "Models and Practice of Neural Table Representations".

3.  **Contribute**: hopefully 2nd [Table Representation Learning workshop](#) at NeurIPS '23.

Ideas for TRL applications, challenges, questions → [m.hulsebos@uva.nl](mailto:m.hulsebos@uva.nl)?