# What are we asking from **Tabular Data**?
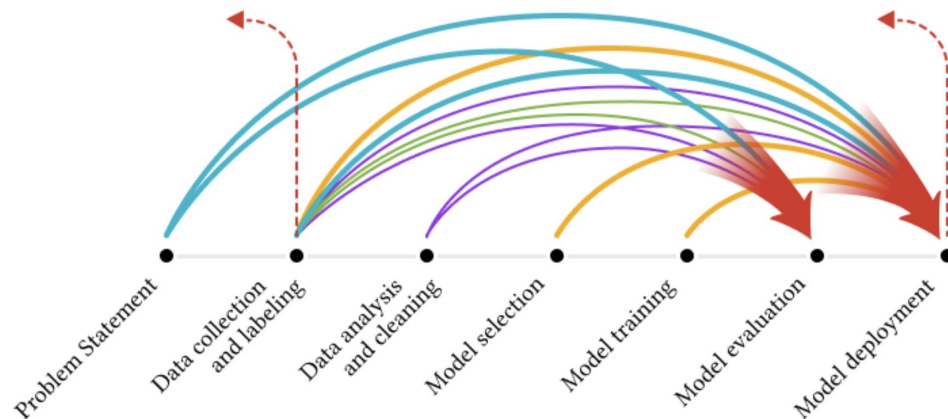
**Madelon Hulsebos**

6 December, Copenhagen

AI for Tabular Data workshop @ EurIPS 2025

As a data scientist,
in the "real world",
I realized 3 things…

# **Realization 1:** most of my work was data work



Data science = "80% data work, 20% model work"

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

# Realization 2: everyone is doing the same....



**over, and over, and over again…**

There must be latent patterns to *learn* (data, code, etc).

# **Realization 3:** tables prevail in the org data landscape

### Popularity of database systems

Bar chart showing popularity values (approximate):
- Relational: 72 (highlighted orange)
- Document: 10
- Key-value: 6
- Search engines: 4
- Wide column: 3
- Graph: 2
- Time Series: 1
- Spatial: <1
- RDF: <1
- Native XML: <1
- Object oriented: <1
- Multi-value: <1

### Distribution Google Dataset Search

Bar chart showing distribution values (approximate):
- Tables: 37 (highlighted orange)
- Semi-Structured: 30
- Other: 11
- Documents: 11
- Images: 5
- Archives: 3
- Text: 3
- Geospatial: 2

For a reason: **tables serve high-value decisions**

# Surprisingly, "tables" ignored as modality in neural AI

**Text**, **Images**, **Code**…

**Tables**?!

# Did a PhD on table semantics, but had a larger vision.

**The Table Representation Learning workshop @ NeurIPS 2022 was born.**

We received what we expected:

- – Tabular QA / text-to-SQL,
- – Synthetic data generation,
- – Data preparation, etc.


And… neural models for predictive tabular ML.

# Anecdotes on neural models for tabular predictive ML

- Received neural predictive ML papers, rejected from main ML confs(?)
- Loud "pro-XGBoost" camp vs. small "pro-neural models" camp.
- Betted on pre-trained neural models for tables in '18: let's empower vision!
- And, 1 paper intro'd a pre-trained cross-table model: TabPFN. It went viral:



Frank's TabPFN tweet

Congrats team TabPFN, TabICL, ConTexTab, etc for pushing through and heading leaderboards!

# Tables weren't really cool in AI

But something happened in a tiny room in New Orleans at TRL @ NeurIPS 2022.

Great vibes, a wildly diverse community, trying to connect the dots.

## Tables were Back.

# Fast-forward to 2025.

**Tabular AI is the "new hot topic" (quote CV researcher)**

**Tabular AI is Europe-led** 🇪🇺 (but let's diversify).

> From an anonymous peer:
>
> ```
> > is there no tab workshop this year at main conf bc you all took to eurips?
>
> > either way wish I could be at the workshop :)
> ```

**And we're only just beginning.**

So, what *are we asking* from tabular data?

# Tabular pipelines are multi-faceted



*Replace "model" with "tool/dashboard" and it's BI.

# LLM-enthusiasts make believe that DS is solved…

# It's not

| Model | SDE (69 tasks) | | PM (28 tasks) | | DS (14 tasks) | |
|---|---|---|---|---|---|---|
| | **Success** | **Score** | **Success** | **Score** | **Success** | **Score** |
| | *Closed model APIs* | | | | *API-based M* | |
| 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* | | | | *Open-weights* | |
| 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 |

(TheAgentCompany: Benchmarking LLM Agents on Consequential Real World Tasks, Xu et al., 2024)

14

It all starts
with retrieving
the right data

# What table to use for my task?

**In 2025**, we get "AGI" but *it still takes weeks to find the right dataset.*



*Task-driven* query

*Result-driven* query suggestions

*Task-driven* dataset relevance

*Data-driven* column suggestions

**Proactive search** makes it easier to find more relevant datasets, with higher success rate.

| Condition | Ease-of-use | Relevance | # Successes |
|---|---|---|---|
| (A) Kaggle | $\mu$=3.08; $\sigma$=0.51 | $\mu$=3.25; $\sigma$=1.05 | 7 *of* 12 |
| (B) Semantic Baseline | $\mu$=3.75; $\sigma$=0.45 | $\mu$=3.25; $\sigma$=0.86 | 6 *of* 12 |
| (C) DATASCOUT | $\mu$=4.75; $\sigma$=0.45 | $\mu$=3.67; $\sigma$=0.78 | 10 *of* 12 |

*Rethinking Dataset Discovery with DataScout*
Lin, R., Chopra, B., Lin, W., Shankar, S., Hulsebos, M., Parameswaran, A., 2025.

**UC Berkeley**  **CWI**

# How to query tables end-to-end? ~RAG over tables

*"What is the highest eligible free rate for K-12 students in the schools in Alameda County?"* → interpret & contextualize query → **retrieve & re-rank table(s)** → table reasoning / query execution → 1.0

How to best retrieve tables in end-to-end tabular QA? Some findings:

- BM25 **less effective than for text**, requires highly descriptive metadata.

- Retrieval over generated summaries helps (but tricky for enterprise-y tables)

- **Row-level retrieval most effective for relational DBs**, but infeasible (many, large tables).

There's a very long road ahead.        Imagine predictive questions over data lakes?

*TARGET: benchmarking Table Retrieval for Generative Tasks*
Ji, X., Parker, G., Parameswaran, A., Hulsebos, M., 2024

*Metadata Matters in Dense Table Retrieval*
Gomm, D., Hulsebos, M., 2025

UC Berkeley    CWI

# Are We Asking the Right Questions?

In end–to–end tabular data analysis, we desire queries that perfectly express the *insight need*: the data, the operation, and output.

Such queries are **platonic**, we won't see them.

But we need to understand **queries for e2d tabular data analysis** to build better systems and eval capabilities. (*join Daniel's talk for more*)

*Are We Asking the Right Questions? On Query Ambiguity in Tabular Data Analysis*
Gomm, D., Wolff, C., Hulsebos, M., 2025

Tables are complex,
they require context

# What is this table about?



Pre crop?
Pre-pre crop?
Pre-pre-pre?

A useful signal for data viz, data prep (integration, missing val mech.), ML, etc.

*Sherlock*: (tiny!) neural model surfaces **table semantics** by column embeddings.

*Sherlock: A Deep Learning Approach to Semantic Data Type Detection,*
Hulsebos et al., 2019.

# Is this data sensitive?

Understudied. Beyond fixed types (pii) sensitivity depends on context. Requires **context beyond the table**.

Input
- Table & metadata
- *Domain-specific docs*

data governance policy documents

**Domain Contextualization: *retrieve-then-detect***

*Retrieve* relevant context

relevant context

*Detect* non-personal sensitivity

Result:
- Better precision.
- Context-grounded (LLM) explanations.

Imagine in healthcare: contextualizing **EHR tables**, in **patient–doctor transcripts**, **medicine docs**, etc?

*Towards Contextual Sensitive Data, Telkamp and Hulsebos, 2025.*

21

But eventually,
we want that insight

# What is the answer to my question?

**Typical questions that drive decisions:**

– Analytical questions → *e.g.* what has happened?
– Predictive queries → *e.g.* what might happen?

Some stats from Gael's keynote at TRL @NeurIPS '24 resonated.

**Pypi #downloads last month** (updated)**:**

– **Scikit-learn: 189,875,197**
– **Pandas: 5O2,255,215**

# Analytical questions: we can use SQL

Why enter the *text-to-SQL* game?

- Saturated (but not solved), everyone just throwing LLMs at it.
- SQL made for humans not machines: it's too flexible, different abstraction?

Then, Cornelius (PhD) convinced me that we need **SQaLe to specialize**.

- Large text-to-SQL dataset
- 500K+ (question, sql, schema)
- Grounded in *real* schemas



*SQaLe: A large text-to-SQL corpus grounded in real schemas, Wolff, Gomm, and Hulsebos, 2025.*

# Specialization over general-purpose

Initial Qwen model *fine-tuned* for text-to-SQL on **SQaLe**:



0.5 **billion** params vs 1.8 **trillion** params

*Find Cornelius in the poster session for more.

# Executability is not accuracy, sure.

*"For each customer in the gold or platinum loyalty tier who has placed at least two delivered orders in the last six months, return their full name, email, city and state, the total number of delivered orders, the total amount they paid based only on captured payments, their most frequently ordered product category in that period, and whether any of the products they ordered currently have low inventory in any warehouse (defined as quantity on hand minus quantity reserved less than 10), sorted by total captured payment amount from highest to lowest."*

Realistic phrasing? Maybe not: inherited from source.
But look at that complex SQL! Promising start.

```sql
WITH gold_platinum_customers AS
  (SELECT c.id,
          c.first_name,
          c.last_name,
          c.email,
          c.city,
          c.state
   FROM customers c
   JOIN orders o ON c.id = o.customer_id
   WHERE c.loyalty_tier IN ('gold',
                            'platinum')
     AND o.status = 'delivered'
     AND o.delivered_at >= date('now', '-6 months')
   GROUP BY c.id
   HAVING COUNT(*) >= 2),
     order_details AS
  (SELECT gpc.id,
          gpc.first_name,
          gpc.last_name,
          gpc.email,
          gpc.city,
          gpc.state,
          o.order_number,
          SUM(p.amount) AS total_paid,
          p.status AS payment_status
   FROM gold_platinum_customers gpc
   JOIN orders o ON gpc.id = o.customer_id
   LEFT JOIN payments p ON o.id = p.order_id
   AND p.status = 'captured'
   WHERE o.delivered_at >= date('now', '-6 months')
     AND o.status IN ('delivered',
                      'completed')
   GROUP BY gpc.id,
            o.id),
     product_category AS
  (SELECT od.id,
          od.first_name,
          od.last_name,
          od.email,
          od.city,
          od.state,
          p.category,
          ROW_NUMBER() OVER (PARTITION BY od.id
                             ORDER BY COUNT(*) DESC) AS rn
   FROM order_details od
   JOIN order_items oi ON od.order_number = CAST(oi.order_id AS TEXT)
   JOIN products p ON oi.product_id = p.id
   GROUP BY od.id)
SELECT gpc.first_name,
       gpc.last_name,
       gpc.email,
       gpc.city,
       gpc.state,
       COUNT(od.order_number) AS total_orders,
       SUM(od.total_paid) AS total_captured_payment_amount
FROM gold_platinum_customers gpc
JOIN order_details od ON gpc.id = od.id
AND od.payment_status = 'captured'
LEFT JOIN product_category pc ON od.id = pc.id
AND pc.rn = 1
WHERE pc.category IS NOT NULL
GROUP BY gpc.id
ORDER BY total_captured_payment_amount DESC;
```
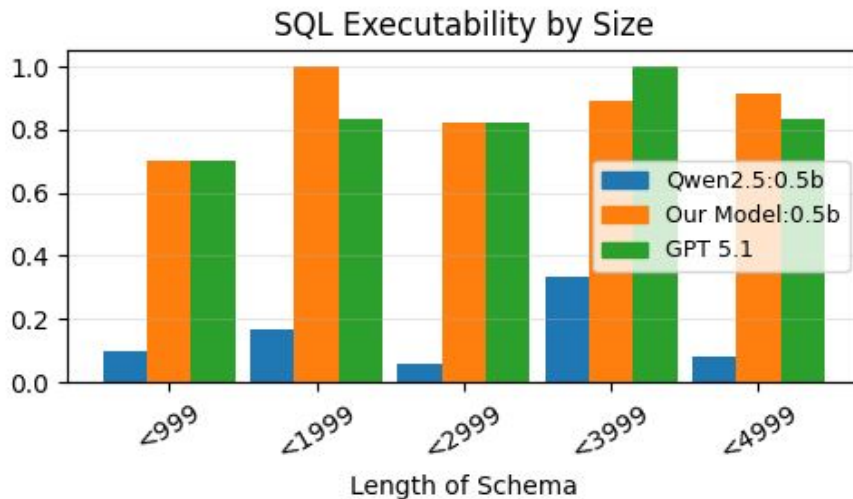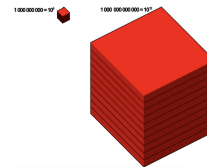
26

What *should we ask* from tabular data, next?

# Generalization versus specialization?

Deriving insights from tabular data to inform decisions is an inherently multi-faceted problem, **perfect for smaller specialized models**, but still looking for the right level of abstraction, which may vary.

Eventually tabular systems are a mix and match:

- Databases and other tools for analytical queries, data prep, etc.
- Some ML models for statistical reasoning
- Some human, document, and LLM contextualization
- Some agentic capabilities
- Some fluid interfaces that connect things together

# Still, many open questions

What is a table? What is a relation? How to deal with different representations?

What data do we encounter in practice?

What queries are asked? What can and should be asked? → check Daniel's talk

How to contextualize tabular data w/ domain knowledge from docs / human input?

How do we want to enable interaction with tabular data?

What is the ideal scope for foundation models? When to generalize vs specialize?

How do we go full-cycle from data collection to decision? What can we "agentize"?
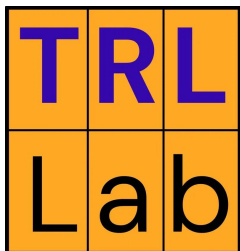
# Where does the community stand?

Excitement:

- **Tabular AI** has gained momentum – it's recognized widely!
- Academic / VC funding is flowing.
- There's much more to **explore and exploit**.

❤ **Tabular AI community**!

You, too? We're looking for committed members to foster the community.
*Reach out if you want to contribute, stay tuned for more initiatives :)*.

# Reach out



**TRL Lab**

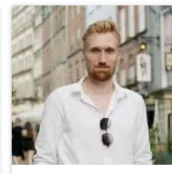| | |
|---|---|
| **Madelon Hulsebos** | Faculty |
| **Xue (Effy) Li** | Postdoctoral researcher |
| **Daniel Gomm** | PhD researcher |
| **Cornelius Wolff** | PhD researcher |
| **Jan Henrik Bertrand** | Research student (ELLIS MSc Honours Program) |
| **Wojciech Kosiuk** | Research student (ELLIS MSc Honours Program) |

`madelon@cwi.nl`

`https://trl-lab.github.io/`