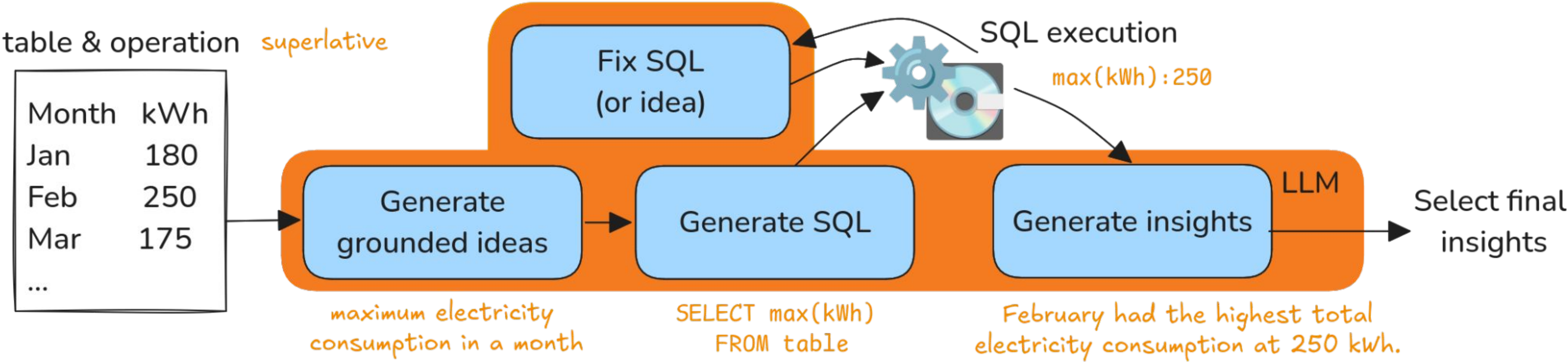# Always use fresh data to evaluate your LLMs! SQL helps analyze tables, but they must be clean.
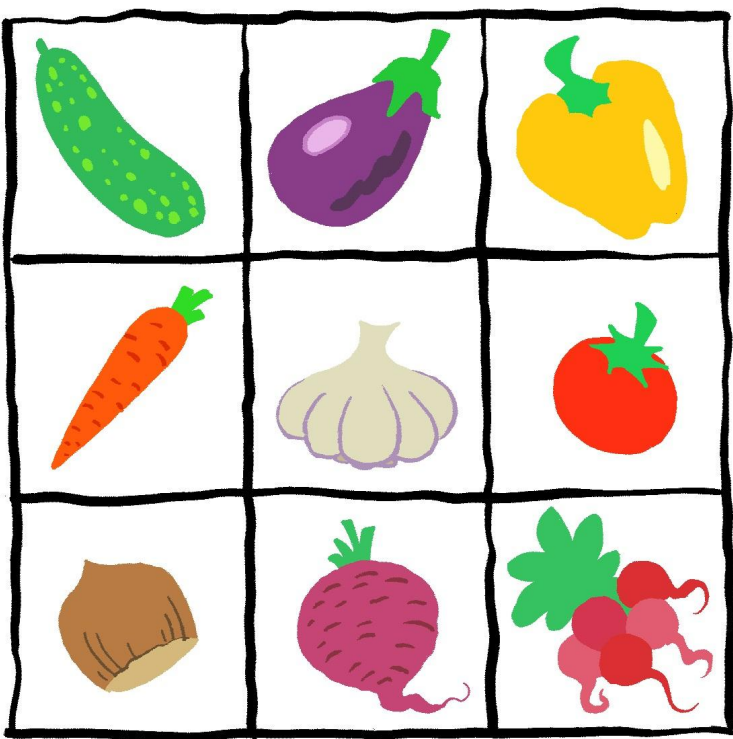
# Generating Data Insights with LLMs by Querying Tables

*K. Onderková, O. Plátek, Z. Kasner, M. Lango, O. Dušek*
*[onderkova, oplatek, kasner, lango, odusek]\*@ufal.mff.cuni.cz*

**Charles University**
**ÚFAL**

table & operation — superlative

| Month | kWh |
|-------|-----|
| Jan | 180 |
| Feb | 250 |
| Mar | 175 |
| ... | |

Fix SQL (or idea) → SQL execution — max(kWh):250

Generate grounded ideas → Generate SQL → Generate insights — LLM → Select final insights

maximum electricity consumption in a month

SELECT max(kWh) FROM table

February had the highest total electricity consumption at 250 kWh.
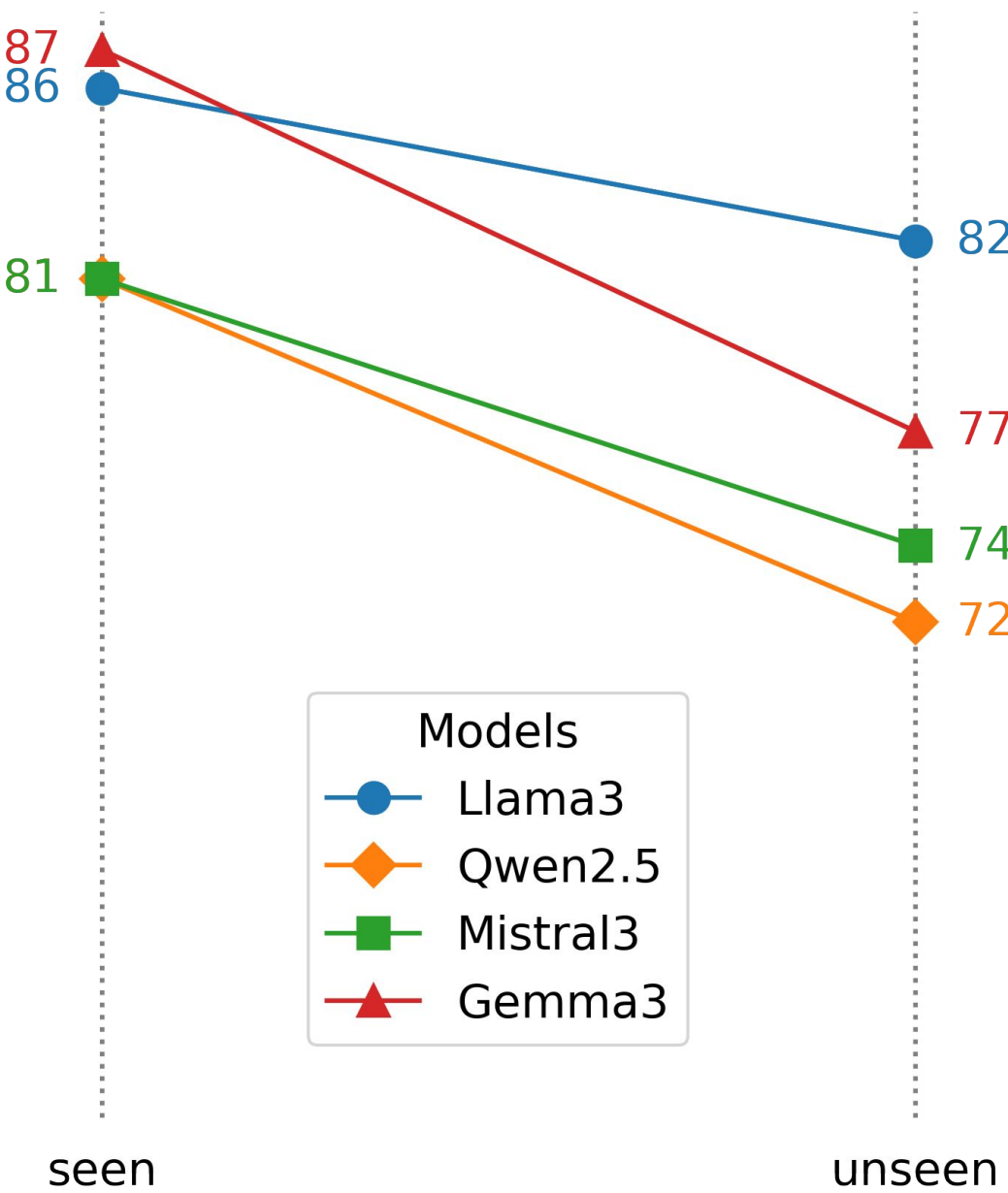
## FreshTab



### Problem
- LLMs **memorized** existing **datasets** (LogicNLG, LoTNLG)
- Datasets contain mostly sports, **no domain** distinction
- They are usually **English only**

### Solution
- **Dynamic dataset construction** based on creation date
  - **Query** Wikidata and **Wikipedia** for list of Wikipages
  - **Scrape**, clean and pick a **table** per Wikipage
  - Filter tables for domain balance

### Results
- LLMs' **performance** visibly **decreases** with unseen data
- FreshTab can **differentiate LLMs** by their **ability to generalize**
- Domain split helps us pinpoint where the models have deficit



TAPEX

| Qwen2.5 Mistral3 Gemma3 Llama3 | sport | politics | culture | other |
|---|---|---|---|---|
| Llama3 | 85 | 81 | 79 | 86 |
| Gemma3 | 84 | 83 | 72 | 68 |
| Mistral3 | 84 | 75 | 75 | 66 |
| Qwen2.5 | 73 | 75 | 75 | 62 |

Models: Llama3, Qwen2.5, Mistral3, Gemma3

### Observations
- Llama 3.3 70b seems to generalize best to new data
  - most balanced performance across domains
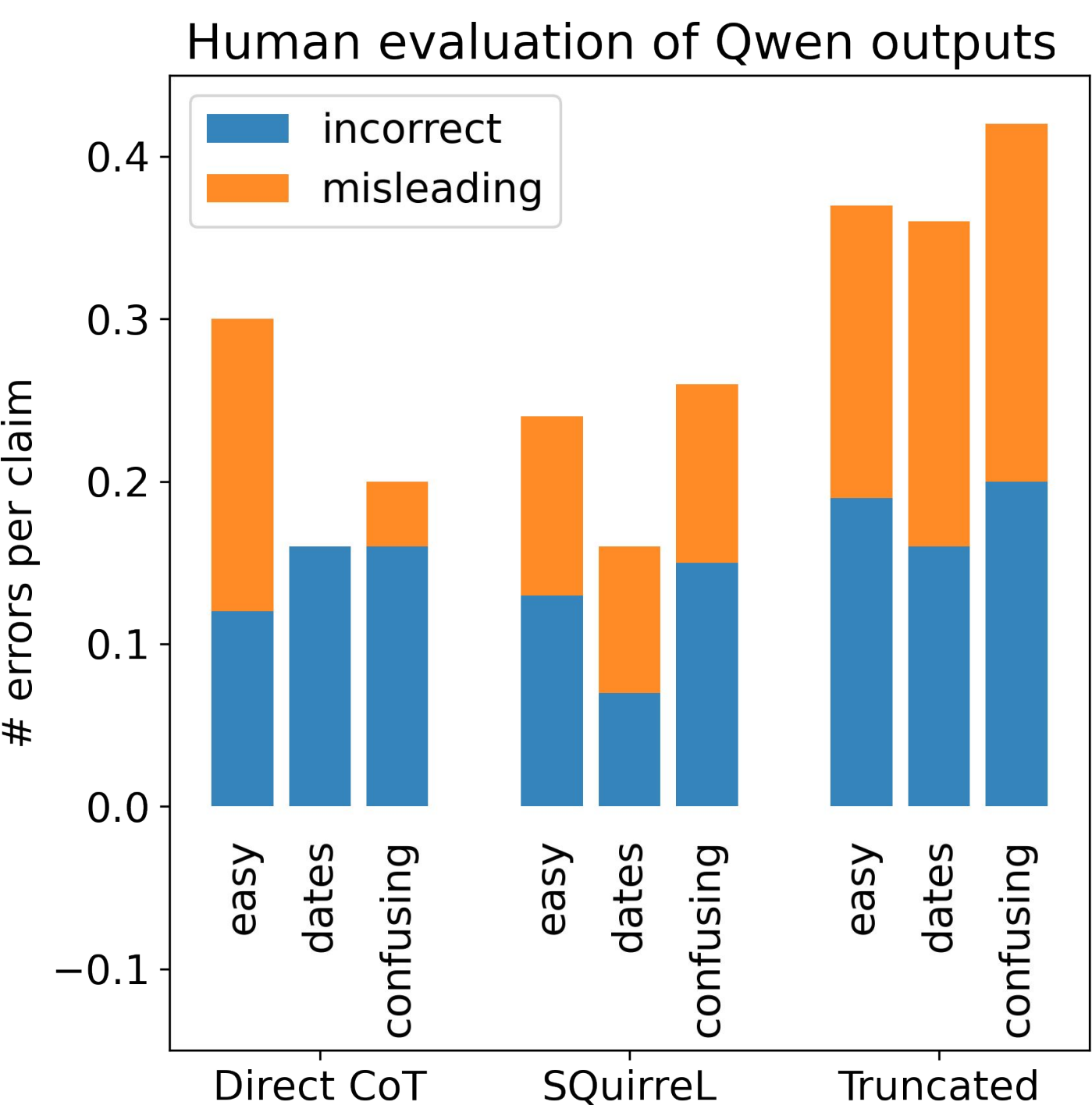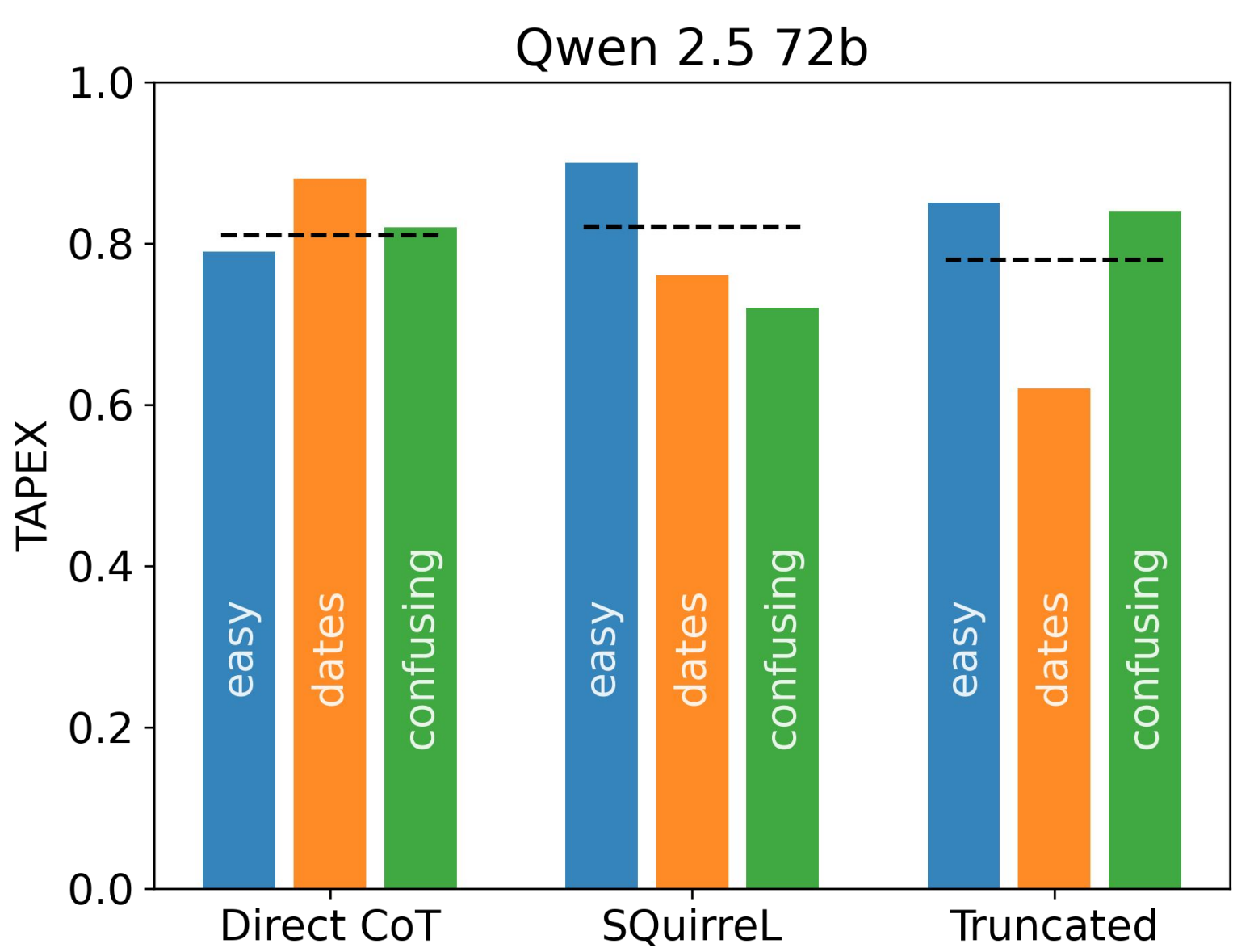- Sport domain is easiest, as expected from lot of training data

## SQuirreL



From real a human!

### Problem
- Symbolic logical forms are **hard to scale**
- LLMs are **hallucinating** and weak at logic

### Solution
- SQuirreL: **Ground** LLMs **in data** retrieved with SQL
  - as a bonus, it is nicely **checkable**
  - can work with truncated tables **on bigger data**

### Results
- Factuality does not improve. **Why?**
- Some **tables** are **not well defined** and are confusing
  - We can detect those tables
- SQL needs well-defined tables to work



Qwen 2.5 72b

Human evaluation of Qwen outputs (incorrect / misleading)

### Observations

| unnamed | entering | advancing |
|---------|----------|-----------|
| First round (14 teams) | 14 Second League teams | - |
| Last 16 (16 teams) | 5 Premier League teams | 11 winners |
| Semi-finals (4 teams) | - | 8 winners |

| event | men |
|-------|-----|
| 5000 metres | 13:01.00 |
| High jump | 2.33 |
| Decathlon | 8550 |

| candidate | votes |
|-----------|-------|
| dog | 60 |
| cat | 70 |
| Total | 130 |