# ChatGPT has been globally exposed to at least 4.7M samples from 263 benchmarks

# 🥬 Leak, Cheat, Repeat: Data Contamination and Evaluation Malpractices in Closed-Source LLMs

Simone Balloccu   Patrícia Schmidtová   Mateusz Lango   Ondřej Dušek

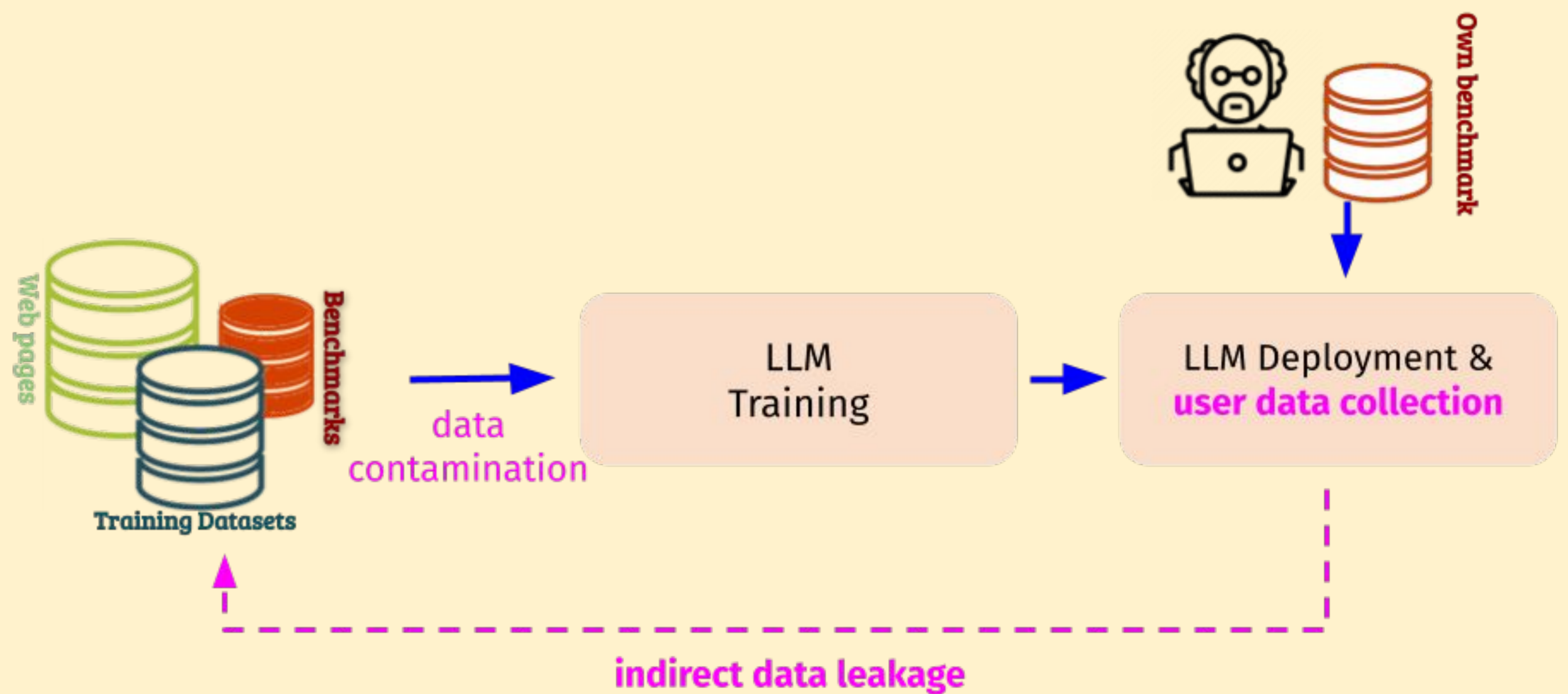{balloccu, schmidtova, *lango, odusek*}@ufal.mff.cuni.cz

CHARLES UNIVERSITY   ÚFAL
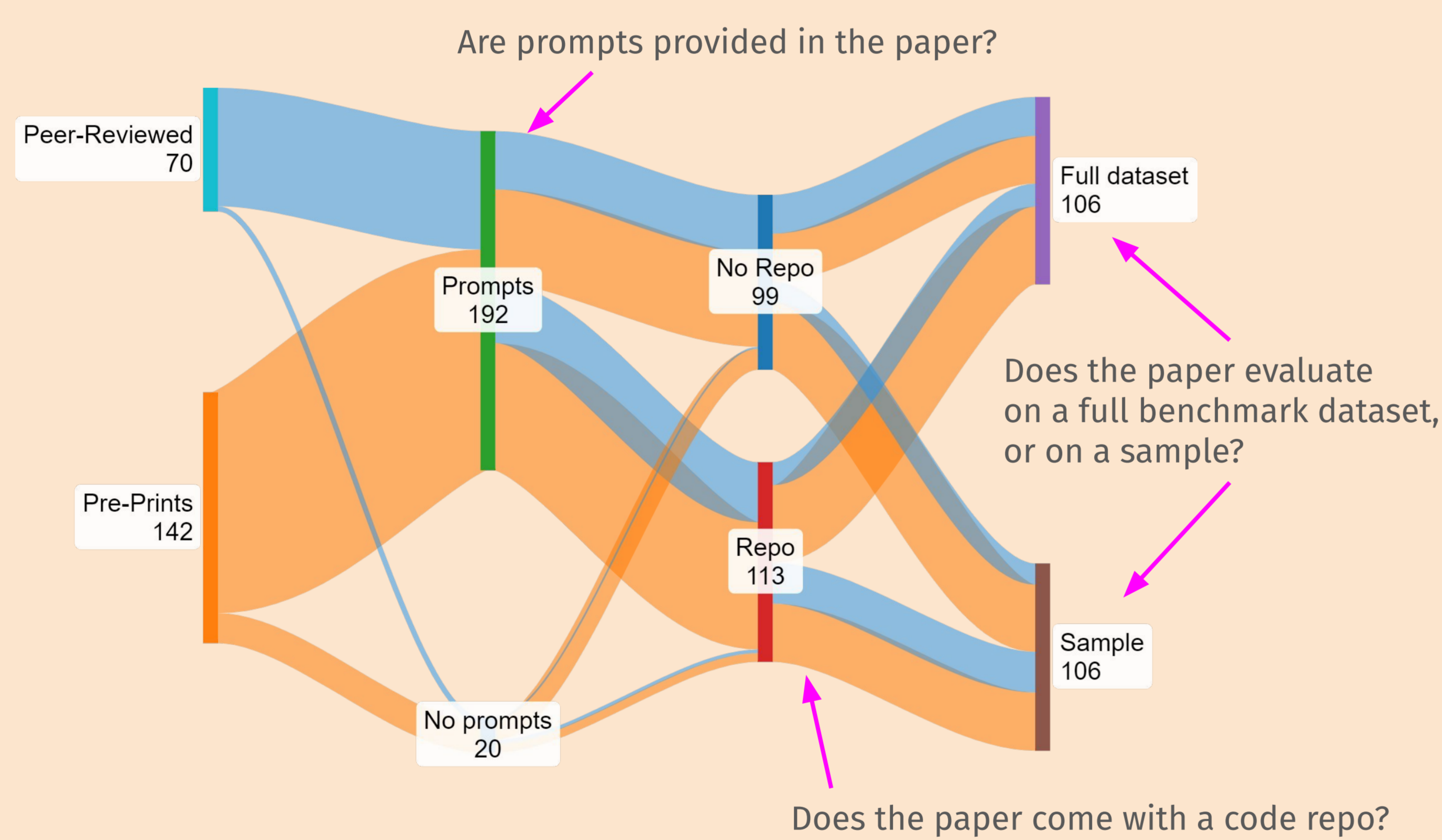
## Data leakage in closed-source LLMs

**Motivation**

- some popular LLMs are **closed-source**
- lack of information on training data raises questions about the **credibility of LLM performance evaluation**
- several attempts to address this issue overlook the problem of **indirect data leaks**

**Methodology & Findings**

- We analyse **255 papers** evaluating OpenAI's GPT-3.5 and GPT-4 on a variety of tasks
- We conclude that **~42% of the relevant reviewed papers leaked data** to GPT-3.5 and GPT-4, for a total of **~4.7M samples** across 263 benchmarks

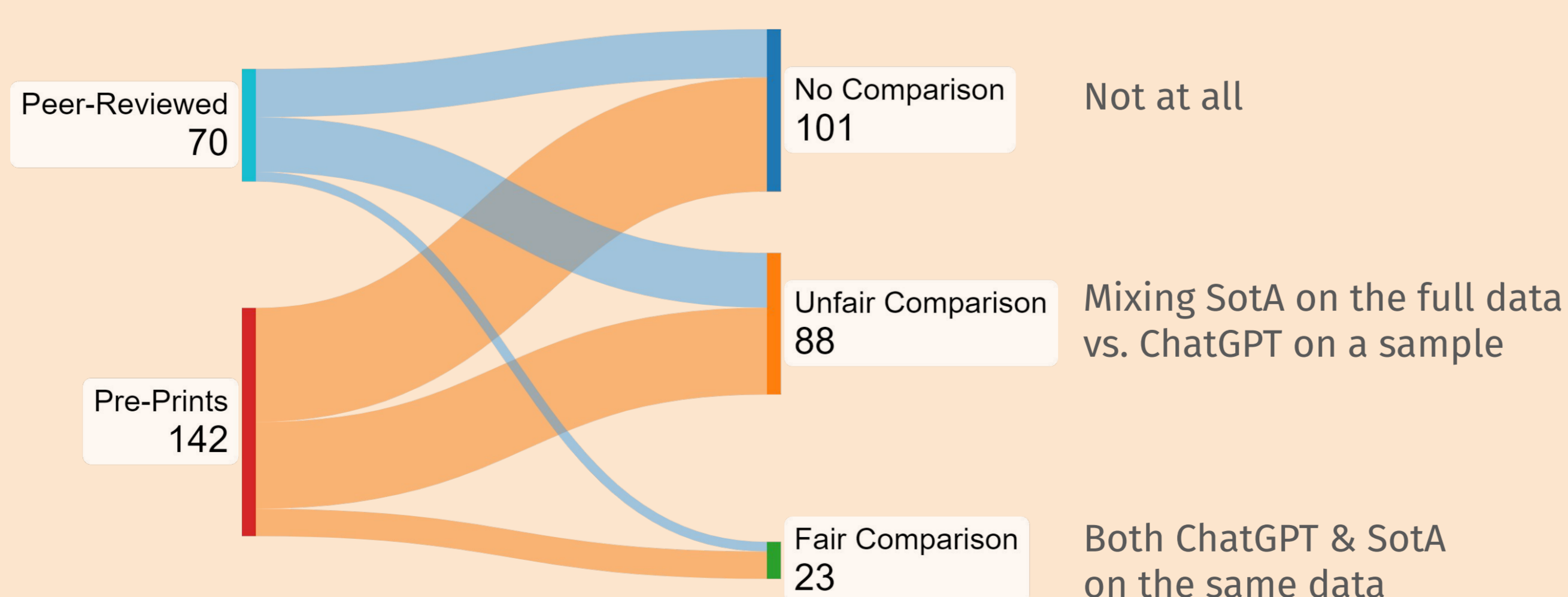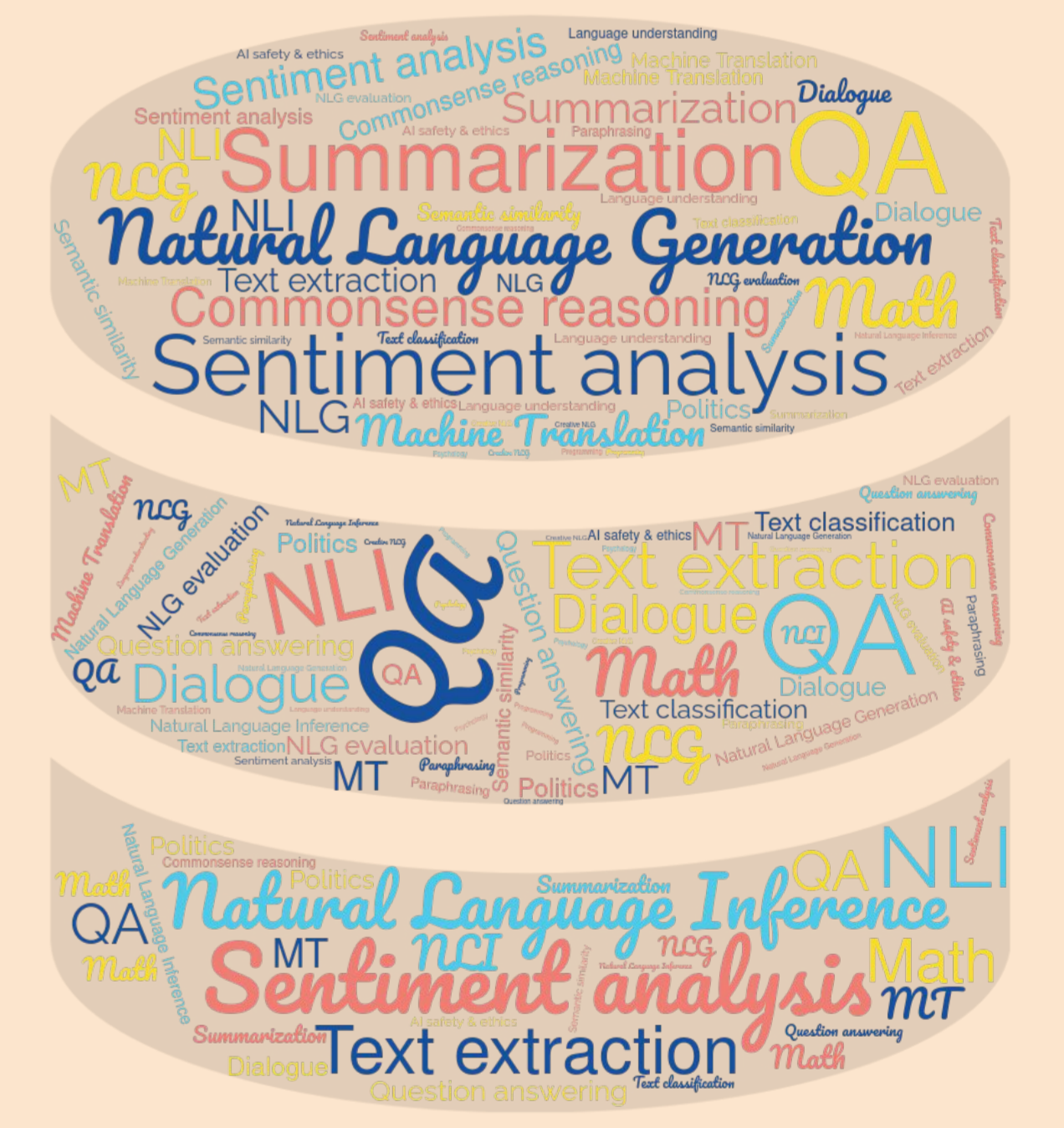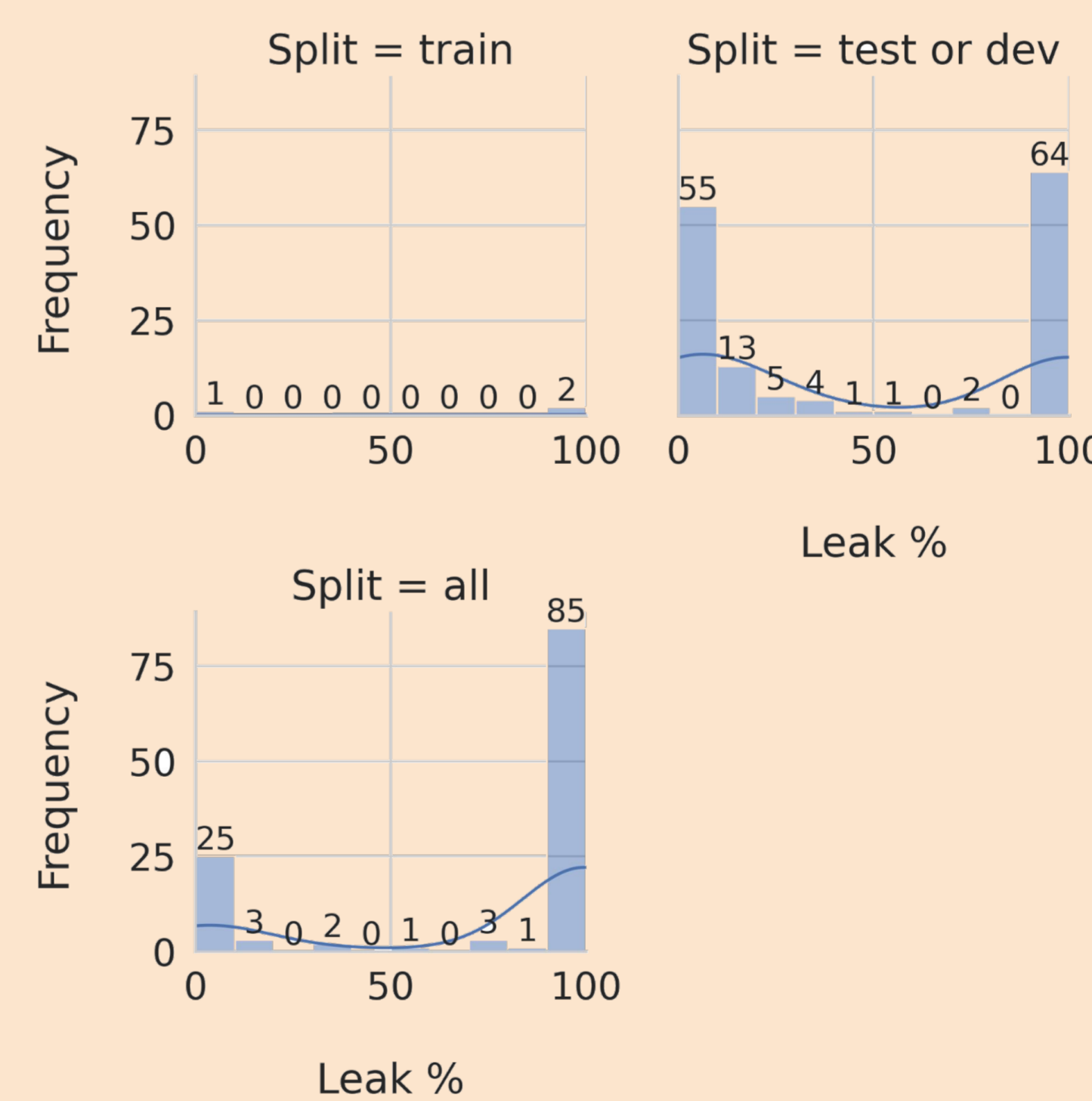Web pages | Benchmarks | Training Datasets → data contamination → LLM Training → LLM Deployment & **user data collection** | Own benchmark | **indirect data leakage**

## Evaluation reproducibility

Are prompts provided in the paper?

Peer-Reviewed 70 | Pre-Prints 142 | Prompts 192 | No prompts 20 | No Repo 99 | Repo 113 | Full dataset 106 | Sample 106

Does the paper evaluate on a full benchmark dataset, or on a sample?

Does the paper come with a code repo?

## Evaluation fairness

Does the paper compare with previous SotA?

Peer-Reviewed 70 | Pre-Prints 142 | No Comparison 101 — Not at all | Unfair Comparison 88 — Mixing SotA on the full data vs. ChatGPT on a sample | Fair Comparison 23 — Both ChatGPT & SotA on the same data

## Indirectly leaked datasets

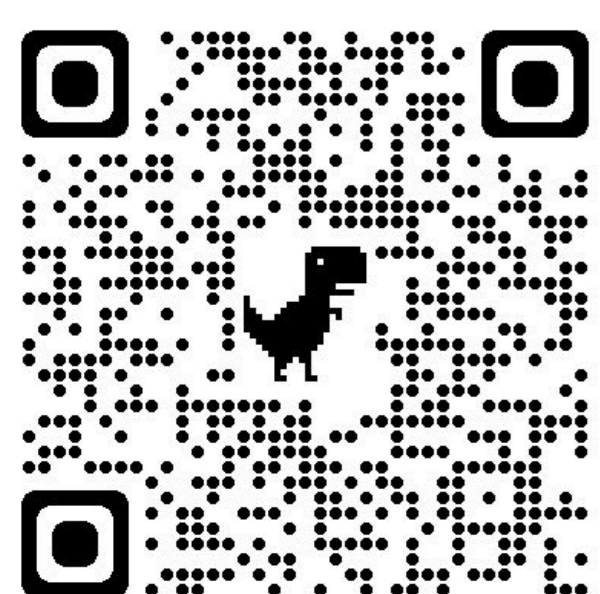Split = train | Split = test or dev | Split = all

Number of times (y) we observed a specific percentage of data leaked (x) from the given split

NLP tasks: size/frequency corresponds to the number of leaked datasets.

## Suggested practices in closed-source LLM evaluation

- **Access models in a non-leaky way:** use the API or opt-out of data collection (check vendor policy)
- **Interpret performance with caution:** incredible performance may be explained by data leakage
- **When possible, avoid using closed-source models**
- **Use a fair comparison:** sufficiently large samples, same for all methods, re-run experiments – do not copy & paste results
- **Make evaluation reproducible:** release code, prompts, model versions, sample selection (if subsampling)
- **Report indirect data leakage**, e.g. on our webpage!