

Cut Costs, Not Accuracy: LLM-Powered Data Processing with Guarantees

Sepanta Zeighami, Shreya Shankar, Aditya Parameswaran

LLMs for Data Processing

Example: game reviews dataset Elden Ring is a monumental achievement in open-world ...

Stardew Valley is the ultimate cozy game, much better than Animal Crossing ...



Other examples

- Extract name of police officers with misconduct from court cases
- Find legal contracts that rely on a particular law

Prompt for processing: does the review mention another game more positively?



Low-Cost LLM-Powered Data Processing



User wants outputs from accurate top-of-theline (e.g., gpt-4o) LLM but it's too expensive

We can use cheap potentially inaccurate LLM (e.g., gpt-4o-mini) as long as we provide guarantee most outputs are the same as the top-of-the-line LLM

Low-Cost Pipelines Through Model Cascade



Mode Cascade Framework

<u>Text Dataset</u>



Overview of Existing Approaches



1: Kang et al VLDB 2020

PRISM





PRSIM determines a cascade threshold to save cost while providing theoretical guarantees on meeting the accuracy target

Takes *data and task characteristics* into account to provide high-utility

PRISM Overview



- Iteratively samples records
- Considers label distribution and the accuracy target





- Uses recent statistical tools
- Improves accuracy when variance is low



- Analysis of probability of not meeting the target
- Takes data characteristics into account















Further Details



PRISM also provides guarantees if the user is interested in precision or recall instead of accuracy



Estimation through Hypothesis Testing









Results

- Accuracy target: 90%
- Probability of failure: 10%



Results Across Accuracy Targets



Summary

- PRISM helps reduce costs while processing text data with LLMs
 - Uses cheaper models as much as possible for the specific dataset and query
 - On average 86% lower cost, 118% higher recall and 19% higher precision over state-of-the-art
- Provides theoretical guarantees to avoid performance regression when using small models



Thanks! Q&A



