# EPIC DATA lab

# Prompt Compression for retrieval-augmented generation

parth asawa, Shreya shankar, Aditya parameswaran

**Exploring applications of compression techniques to decrease cost and latency in retrieval-augmented generation.**

## Background

- **Problem**
  - LLMs are extremely powerful at generative tasks, though are often expensive and high-latency.
  - LLM API inference time scales quadratically and costs scale linearly with input length.

- **Goal**
  - Explore the effectiveness of various types of compression techniques from both the databases and information theory world on reducing input size while maintaining accuracy.
  - Propose methods for compression over both structured and unstructured data corpuses.

## Prior Prompt Compression Work

*LLMLingua: Compressing Prompts for Accelerated Inference of Large Language Model,* Jiang, et al.



Figure 1: Framework of the proposed approach *LLMLingua*.

## Structured Data Corpuses

- **Case Study:**
  - Compression with 3 strategies:

*Column Filtering*

*NL2SQL*



NL question

*Columnar Compression*



## Use Cases & Workloads

- We explore use cases of compression in web tables, tabular datasets, unstructured text, and more.

Compositional Semantic Parsing on Semi-Structured Tables



*Panupong, et al.*          *Islam, et al. (FinanceBench)*
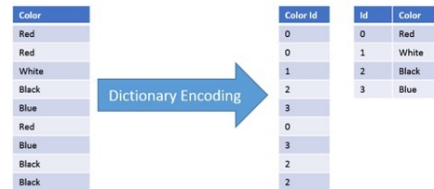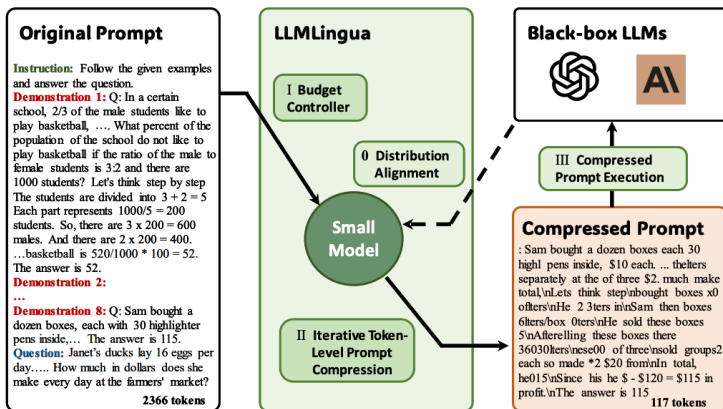
## Unstructured Data Corpuses

- General-purpose large language models don't model the importance of words/tokens in specific domains well.

- **Domain-specific compression:** use fine-tuned models.
  - BloombergGPT: A Large Language Model for Finance
  - Med-PaLM: A Medical Large Language Model



## Other Considerations & Future Work

- Identifying best compression strategies automatically over long mixed-structure documents.

- Revisiting traditional compression techniques like TF-IDF based stop-word removal.

- What are the most important data workloads or problems for which people would find this work impactful?