



What Format Is Your LLM Dataset?

Large language model (LLM) training datasets are web-scale.

- PiB-level raw data; TiB-level cleaned data;
- Massive data preprocessing.
- Needs to reside on costly high performance storage.

Many training frameworks and datasets are using **JSONL**, but we know that it is not an ideal choice for managing such amount of data.

Data Preparation in Parquet; Why Not Training?



Figure 1. Parquet logical and physical layouts.

A lot of data cleaning & preparation tools use Parquet, which is designed for efficiently processing web-scale datasets.

With Parquet, data preparation reuses existing data infrastructure for analytical workload.

- Designed for scan-based operations.
- Optimized for block storage.
- High compression rate to reduce storage cost.

However, Parquet is inefficient with random access, a key procedure for data shuffling to ensure model accuracy.

Costly Pit-stops between Preparation and Training



The current practice is to convert & pre-shuffle columnar data into other formats like **JSONL**.

- 3× more storage capacity needed on costly performant storage.
- Break single truth of data.
- Additional human efforts to keep up with ever evolving datasets with new data and preprocessing refinement.

Tianle Zhong¹ ¹University of Virginia

Jiechen Zhao²

²University of Toronto



Figure 2. Mmap file I/O

- Distributed memory like Ray and Spark High memory footprint and complexity.
- Memory contention with other key functionalities like model checkpointing. Disk-backed memory mapping (HuggingFace Datasets). Unsatisfactory throughput due to thrashing (page faults).
- Streaming I/O-based local shuffle.
- Obvious model accuracy loss due to limited shuffle quality.

Our Goals

We aim to stay with one consistent columnar format for training data pipeline, and achieves

- Controlled DRAM footprint (against distributed memory). • provide memory resilience for other critical functionalities like tensor offloading and model checkpointing.
- Sufficient throughput (against MMap I/O) to avoid data starvation for GPU utilization.
- High shuffle quality (against local shuffle) to ensure model accuracy.

Chunk Row

Figure 5. A lot of I/O bandwidth is wasted due to the granularity gap.

Figure 4. Disk I/O remains as bottleneck

- Marginal benefits of caching unless the dataset fully cached, considering that access to the same row is exactly once per epoch.
- Significant granularity gap: Fine-grained shuffling vs. Coarse-grained columnar chunk I/O
- Consume rows but have to read chunks, resulting in extremely low effective bandwidth. In other words, a lot of data is fed into CPU but never used by GPU.

Youmu Overview

Figure 6. Youmu system overview.

Youmu: Efficient Columnar Data Pipeline for LLM Training

Qiang Su³ Geoffrey Fox ¹

³The Chinese University of Hong Kong

Figure 3. Streaming shuffle.

Design highlights:

- No cache, only buffer. • From Observation 1.
- Fine-grained page-level I/O. • From Observation 2.
- Practical compatibility.
- Directly works with widely adopted Parquet.

- Friendly I/O size for SSD.
- Many pages to shuffle.
- Improved I/O goodput.

- Python APIs to be directly integrated with PyTorch Dataset interface.
- Support full shuffle by extracting a random row from retrieved page.

Figure 10. Perplexity against training steps

Figure 11. Memory footprint per node against dataset sizes on 16 nodes.

Youmu Design

File	• 1		Pag	geOffset_File	[0] x+y+z
File 0		Col_ck 0	Col_ck 1	Col_ck 2	
	Pageldx 0	Page Addr	Page Addr	Page Addr	
	Pageldx 1	Page Addr	Page Addr	Page Addr	• • • •
		•			
PageOffset_ColCk x			x+y	x+y+z	

Figure 8. A lot of I/O bandwidth is wasted due to the granularity gap

Aggressive Buffer Shuffle

- Based on the fact: page size < buffer size.</p>
- Shuffle rows in buffer upon the arrival of each new page.
- Better shuffle quality than standard buffer shuffle at initialization.

Global Page Index

- A 3-dim index across file, column chunk and page.
- Lightweight initialization since only metadata needed.
- Layered binary search-based index translation.

Implementation

• Rust runtime: Based on official Rust implementation of Apache Parquet and Arrow.

Zero-copy conversion to various dataframe and tensor formats enabled by Arrow memory model ecosystem.

Evaluation

Key Take-away

- A small page (10KB) achieves competitive quality to row shuffling.
- A big page (1MB) still outperforms streaming shuffle.

TL;DR

Youmu can achieve sufficiently low latency to saturate GPU while keeping very low memory overheads

