# Accommodating LLM Service over Heterogeneous Computational Resources

**Binhang Yuan**

12.04.2024

# Amazing Progress of ML/AI

# The challenge of Today:
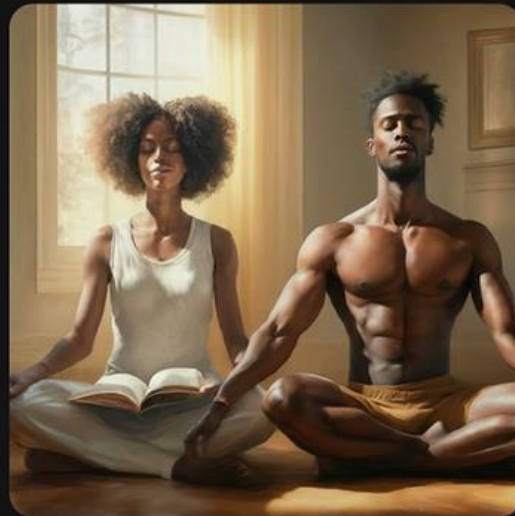
(Million $)

**Building ML Applications at SOTA scale is expensive!**

**Further scaling is facing non-linear bottlenecks.**

# *Optimizing Communications for Distributed and Decentralized LLM Service.*

# Communication Bottlenecks across Infrastructure

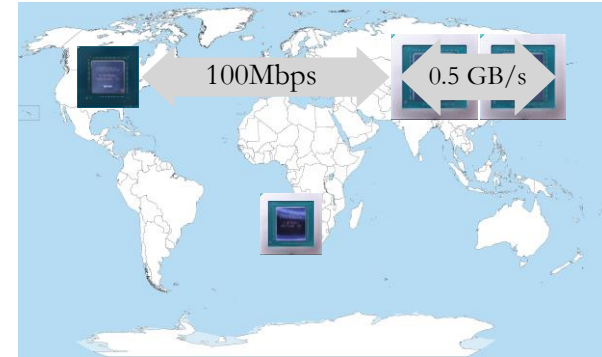communication becomes slower, open up more choices (and some can be cheaper)



| Data Center | (Multi-cloud) Spot Instances | Serverless Environment | Decentralized Network |

**The more we can optimize communications, the more choices we have when building our infrastructure.**

# From Cloud to Decentralized Compute Resource

# _Accommodate LLM training through heterogeneous network._

# Pipeline Parallelism



1. How to schedule the communication to accommodate the decentralized connections?
2. How to compress forward activations and backward gradients?

# Decentralized Training of Foundation Models

- Decentralized training of FM: the network is 100×

  slower, but the pre-training throughput is only

  1.7~3.5× slower!

- Decentralized fine-tuning of FM: **AQ-SGD**

  communication-efficient pipeline training with

  activation compression.



[NeurIPS 2022-(a)]



[NeurIPS 2022-(b)]

# Accommodate Communication in a Decentralized network

A bi-level scheduling algorithm based on an extended balanced graph partition to estimate the communication cost:

- <u>Data parallel communication cost</u>: nodes handling the same stage need to exchange gradients;
- <u>Pipeline parallel communication cost</u>: nodes handling nearby stages for the same micro-batch need to communicate activation in the forward propagation and gradients of the activation in the backward propagation.



(d) perfect matching corresponds to how devices in $C_i$ and devices in $C_j$ communicate in a pipeline.

(a) Communication Topology Graph $\mathbf{G}$ over $N$ devices

(b) Each partition $\mathbf{C_i}$ deals with one stage, running data parallel within each partition

(c) Coarsened graph $\widehat{\mathbf{G}}$ decoding the cost of pipeline parallel

(e) Open-loop-traveling-salesman provides a pipeline structure

# AQ-SGD

$$\min_{x \in \mathbb{R}^d} f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} F(b(a(\xi, x^{(a)}), x^{(b)}))$$

Direct quantization only works to some degree.



Converge better

Quantization Error of Diff smaller

Model change smaller

Activation Diff Smaller

Activation change smaller

- **(A1: Lipschitz assumptions)** We assume that $\nabla f$, $\nabla(f \circ b)$ and $a$ are $L_f$, $L_{f \circ b}$-, and $\ell_a$-Lipschitz, respectively, recalling that a function $g$ is $L_g$-Lipschitz if

$$\|g(x) - g(y)\| \le L_g \|x - y\|, \quad \forall x, \forall y.$$

  Furthermore, we assume that $a$ and $f \circ b$ have gradients bounded by $C_a$ and $C_{f \circ b}$, respectively, i.e. $\|\nabla a(x)\| \le C_a$, and $\|\nabla(f \circ b)(x)\| \le C_{f \circ b}$.

- **(A2: SGD assumptions)** We assume that the stochastic gradient $g_\xi$ is unbiased, i.e. $\mathbb{E}_\xi[g_\xi(x)] = \nabla f(x)$, for all $x$, and with bounded variance, i.e. $\mathbb{E}_\xi \|g_\xi(x) - \nabla f(x)\|^2 \le \sigma^2$, for all $x$.

**Theorem 3.1.** *Suppose that Assumptions A1, A2 hold, and consider an unbiased quantization function $Q(x)$ which satisfies that there exists $c_Q < \sqrt{1/2}$ such that $\mathbb{E}\|x - Q(x)\| \le c_Q \|x\|$, for all $x$.[1] Let $\gamma = \frac{1}{3(C + 3L_f)\sqrt{T}}$ be the learning rate, where*

$$C = \frac{4 c_Q \ell_a (1 + C_a) L_{f \circ b} N}{\sqrt{1 - 2 c_Q^2}}.$$

*Then after performing $T$ updates one has*

$$\frac{1}{T} \sum_{t \in [T]} \mathbb{E}\|\nabla f(x_t)\|^2 \lesssim \frac{(C + L_f)(f(x_1) - f^*)}{\sqrt{T}} + \frac{\sigma^2 + (c_Q C_a C_{f \circ b})^2}{\sqrt{T}}. \quad (3.1)$$

11

# ***<u>LLM service is NOT all about training.</u>***

*<u>"90% of the machine learning demand in the cloud is for inference."</u>*

-- AWS Report

# FlexGen

## *High-Throughput Generative Inference of Large Language Models with a Single GPU*

- *OPT-175B Scale Inference on a single GPU*:
  - 6.5K stars on Github;
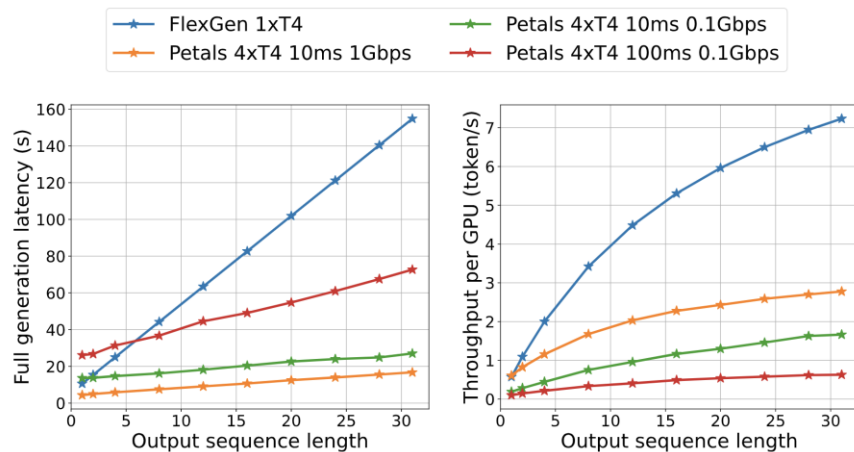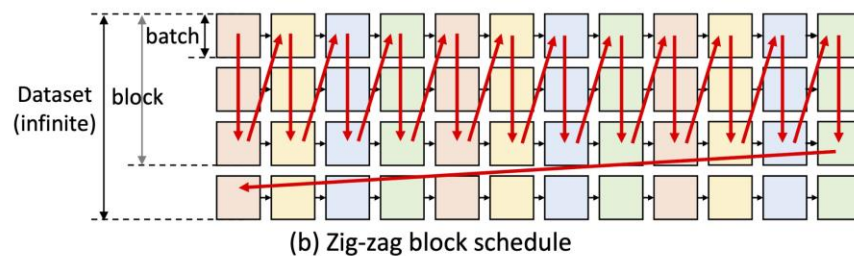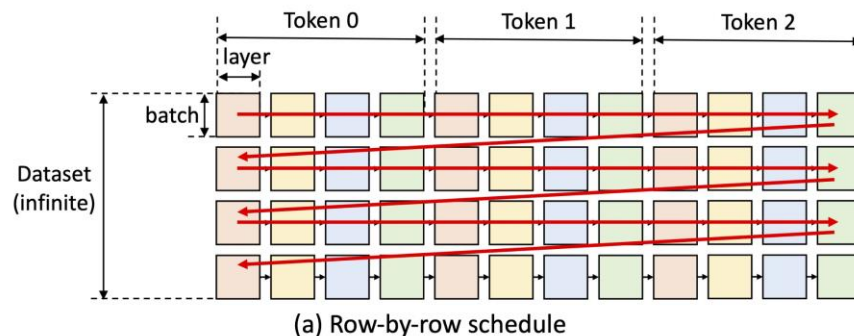  - Top discussion on Hacker News;
  - High throughput scenario: *1 token/s.*

Token 0    Token 1    Token 2

layer

batch

Dataset (infinite)

(a) Row-by-row schedule

batch

Dataset | block (infinite)

(b) Zig-zag block schedule

FlexGen 1xT4    Petals 4xT4 10ms 0.1Gbps
Petals 4xT4 10ms 1Gbps    Petals 4xT4 100ms 0.1Gbps

Full generation latency (s) — Output sequence length

Throughput per GPU (token/s) — Output sequence length

FlexGen: High-Throughput Generative Inference of Large Language Models with a Single GPU

Ying Sheng [1]  Lianmin Zheng [2]  Binhang Yuan [3]  Zhuohan Li [2]  Max Ryabinin [4 5]
Beidi Chen [6 7]  Percy Liang [1]  Christopher Ré [1]  Ion Stoica [2]  Ce Zhang [3]

**Abstract**

The high computational and memory requirements of large language model (LLM) inference make it feasible only with multiple high-end accelerators. Motivated by the emerging demand for latency-insensitive tasks with batched processing, this paper initiates the study of high-throughput LLM inference using limited resources, such as a single commodity GPU. We present FlexGen, a high-throughput generation engine for running LLMs with limited GPU memory. FlexGen can be flexibly configured under various hardware resource constraints by aggregating memory and computation from the GPU, CPU, and disk. By solving a linear programming problem, it searches for efficient patterns to store and access tensors. FlexGen further compresses the weights and the attention cache to 4 bits with negligible accuracy loss. These techniques enable FlexGen to have a larger space of batch size choices and thus significantly increase maximum throughput. As a result, when running OPT-175B on a single 16GB GPU, FlexGen achieves significantly higher throughput compared to state-of-the-art offloading systems, reaching a generation throughput of 1 token/s for the first time with an effective batch size of 144. On the HELM benchmark, FlexGen can benchmark a 30B model with a 16GB GPU on 7 representative sub-scenarios in 21 hours. The code is available at https://github.com/FMInference/FlexGen.

Figure 1. The total latency for a block and throughput trade-offs of three offloading-based systems for OPT-175B (left) and OPT-30B (right) on a single NVIDIA T4 (16 GB) GPU with 208 GB CPU DRAM and 1.5TB SSD. FlexGen achieves a new Pareto-optimal frontier with 100× higher maximum throughput for OPT-175B. Other systems cannot further increase throughput due to out-of-memory issues. "(c)" denotes compression.

**1. Introduction**

In recent years, large language models (LLMs) have demonstrated strong performance across a wide range of tasks (Brown et al., 2020; Bommasani et al., 2021; Zhang et al., 2022; Chowdhery et al., 2022). Along with these unprecedented capabilities, generative LLM inference comes with unique challenges. These models can have billions, if not trillions of parameters (Chowdhery et al., 2022; Fedus et al., 2022), which leads to extremely high computational and memory requirements to run. For example, GPT-175B requires 325GB of GPU memory simply to load its model weights. Fitting this model onto GPUs would require at least five A100 (80GB) GPUs and complex parallelism strategies (Pope et al., 2022; Aminabadi et al., 2022). Thus, lowering LLM inference resource requirements has recently attracted intense interest.
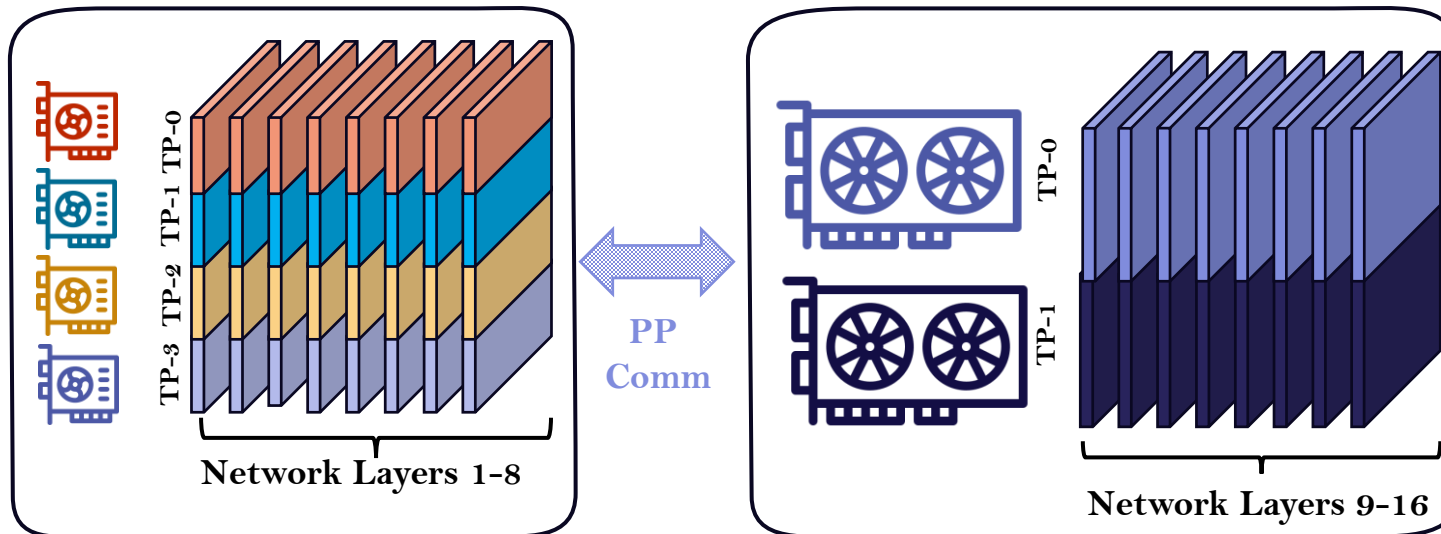
In this paper, we focus on a setting that we call *throughput-oriented generative inference*. In addition to interactive use cases such as chatbots, LLMs are also applied to many "back-of-house" tasks such as benchmarking (Liang et al., 2022), information extraction (Narayan et al., 2018), data wrangling (Narayan et al., 2022), and form processing (Chen et al., 2021). One key characteristic of these tasks is that they often require running LLM inference in batches over a large number of tokens (e.g., all the documents in a company's

**[ICML 2023 Oral]**

# HexGen

## *Generative Inference of Foundation Model over Heterogeneous Environment*

- An implementation that accommodates tensor model parallelism and pipeline parallelism.

- A scheduling algorithm that optimizes pipeline partitions and parallel strategies over heterogeneous GPUs.



Network Layers 1-8

PP Comm

Network Layers 9-16



[Preprint: arxiv. 2311.11514]

# Summary

- ***Communication*** is a key bottleneck of distributed learning, both for centralized data center network and decentralized environments.

- We can develop ***Algorithms*** to alleviate communication bottlenecks:

  - *LLM Training*: *system scheduling and algorithm relaxation.*

  - *LLM Inference*: *latency and throughput orientated scenarios.*

Personal page:
https://binhangyuan.github.io/site/

*Thank you!*