

# Accommodating LLM Training over Decentralized Computational Resources

Binhang Yuan


06.11.2024

# Amazing Progress of ML/AI




*“space robot studying a book in front of Stanford”*




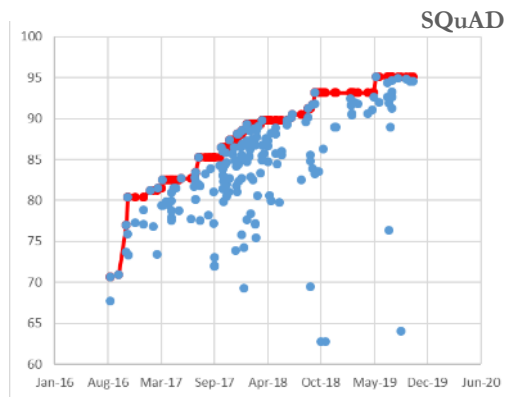
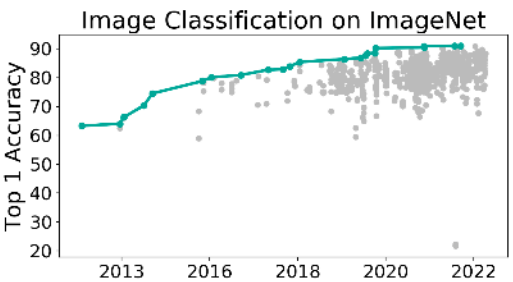


Write a haiku from the perspective of a copywriter who is feeling sad that AI might diminish the value of the written word



Words on a screen,  
Once valued, now just a blur  
Machine takes the pen.





# The challenge of Today:


(Million \$)

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

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

# Bottleneck: Communications & Data Movement

Distributed training at scale is communication-intensive.

  
GPT-3

6.7B Parameters  
 $1.20E+22$   
Floating Point Ops.

---

**32 Machines, 4x A100 GPU each**


Each machine send+recv **4PB** data

100Gbps = **93h** Communication Time

10Gbps = **930h** Communication Time

~**200h** Computation Time



  
GPT-3

175B Parameters  
 $3.14E+23$   
Floating Point Ops.

---

**196 Machines, 8x A100 GPU each**

Each machine send+recv **12PB** data

100Gbps = **279h** Communication Time

10Gbps = **2790h** Communication Time

~**400h** Computation Time



***(Future)** 10 $\times$  further scaling requires fast connections between 10 $\times$  machines. Becoming challenging even for data center.*

***(Today)** Model training today is largely restricted to centralized data centers with fast network connections. Hard to use cheaper alternatives (Non 1st tier clouds, Spot Instances, Volunteer Computes, etc.).*



**NVIDIA DGX SuperPOD:**  
Up to **256** GPUs

*Optimizing Communications for  
Distributed and Decentralized Learning.*



# Communication Bottlenecks across Infrastructure

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



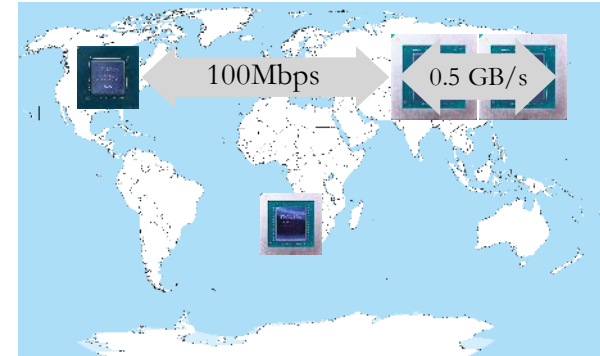
Data Center



(Multi-cloud) Spot Instances



Serverless Environment



Decentralized Network

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

$$\min_x \mathbb{E}_\xi f(\xi, x)$$

$$\min_x \mathbb{E}_\xi f(\xi, x)$$

### Data

- (ImageNet) 1.3M Images (est. 160+ GB)
- (GPT-3) 300 Billion Tokens (est. 2+ TB)

### Model

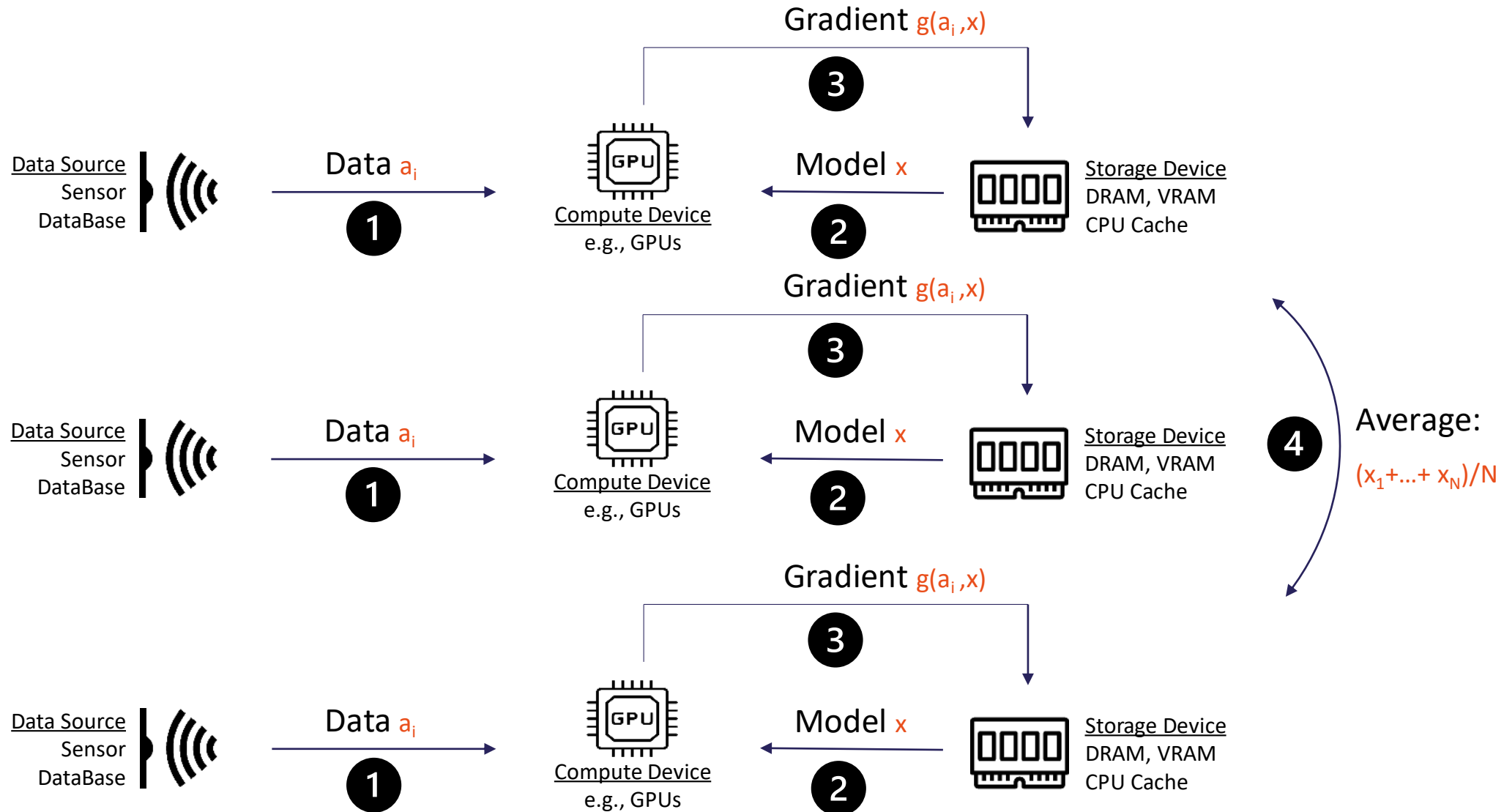
- (GPT-2) 1.3 Billion Parameters (2.6 GB fp16)
- (GPT-3) 175 Billion Parameters (350GB fp16)

### Compute

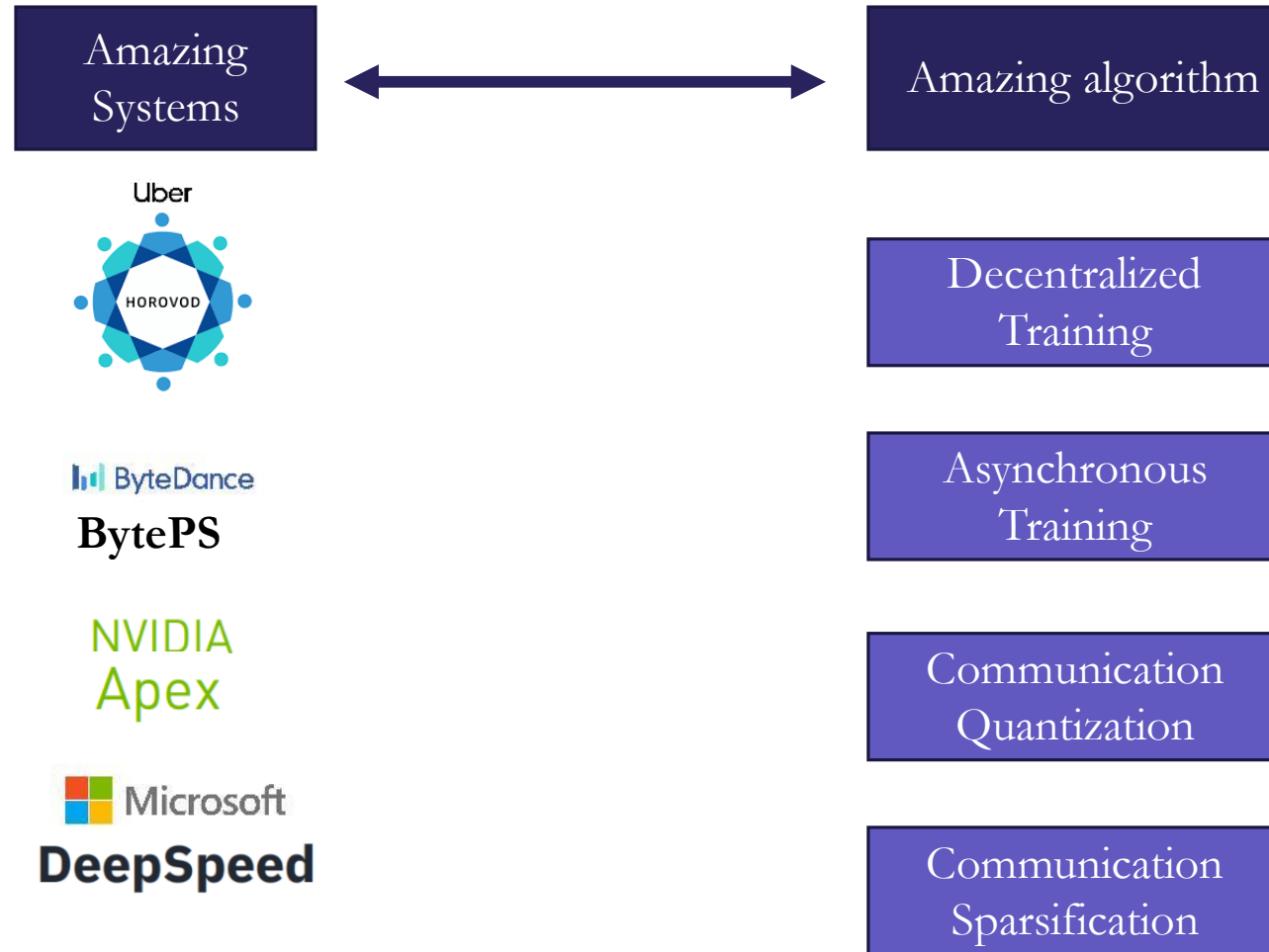
- (GPT-2) est. 2.5 GFLOPS/token
- (GPT-3) est. 0.4 TFLOPS/token



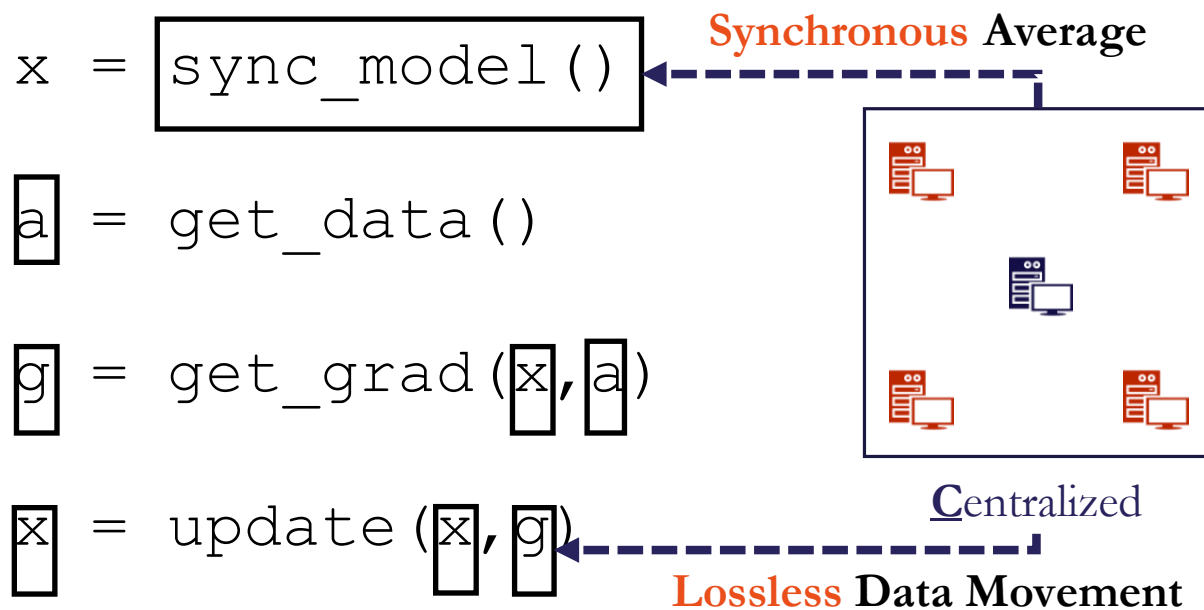
# Data Parallel SGD



# System Optimizations and Relaxed Algorithms



# Baseline: Centralized, Synchronous, Lossless, SGD



## Idea

- Distribute batch gradient calculation to multiple workers;
- Synchronize workers with a central server (or AllReduce).

## Mathematical Formulation

$$x_{t+1} = x_t - \gamma \sum_{i=1..n} g_i(x_t; a_i)$$

## Convergence

$$O(1/\sqrt{nT})$$

Goal 1: Keep This Similar

## System Profile

Goal 2: Make this Faster



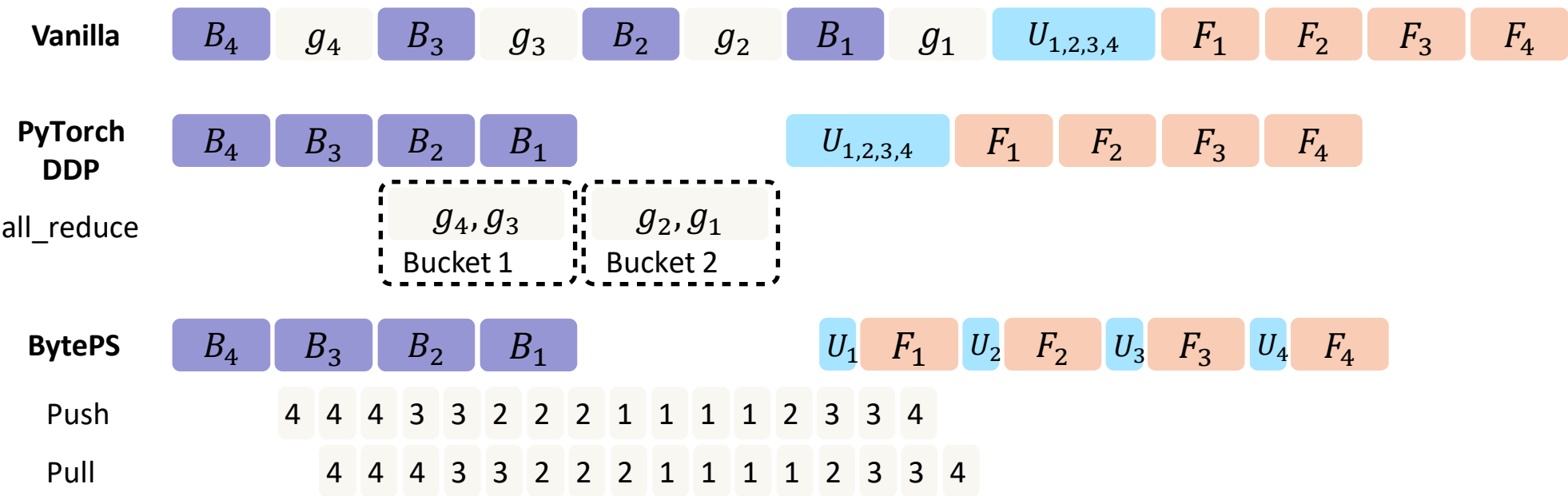
■ Computation  
■ Communication

# System Optimizations

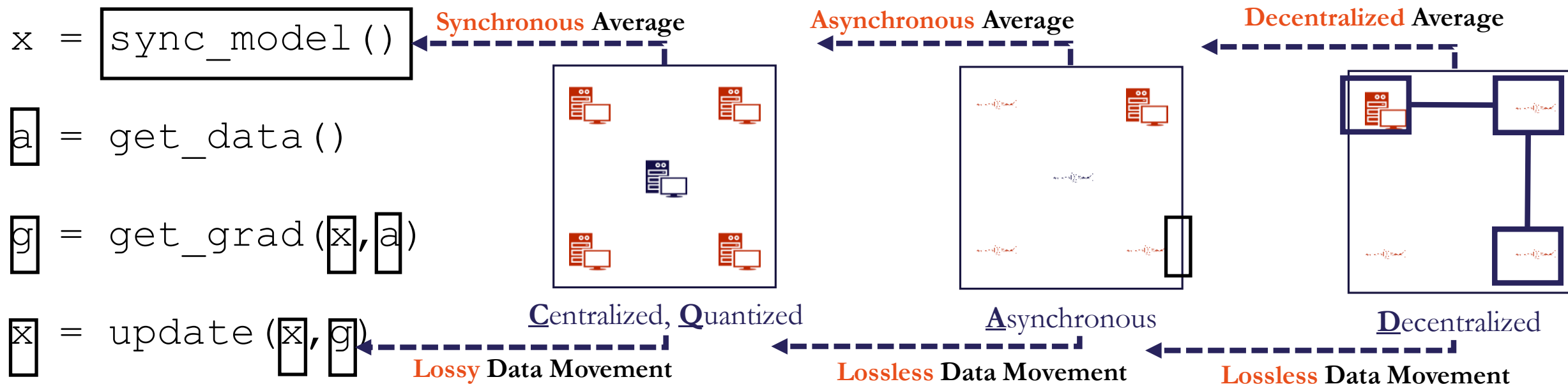


Existing Systems:

Optimize the standard DP-SGD computation:



# Relaxed Algorithms



Mathematical Formulation

$$x_{t+1} = x_t - \gamma \mathcal{Q} \left( \sum_{i=1..n} \mathcal{Q}(g_i(x_t, a_i)) \right)$$

$$x_{t+1} = x_t - \gamma g(x_{t-\tau_t}; a_i)$$

staleness caused by async

$$x_{t+1,i} = \frac{x_{t,i-1} + x_{t,i} + x_{t,i+1}}{3} - \gamma g(x_{t,i}; a_i)$$

Convergence

$$O(1/\sqrt{nT} + \epsilon/\sqrt{T})$$

Quantization error:  $\epsilon$

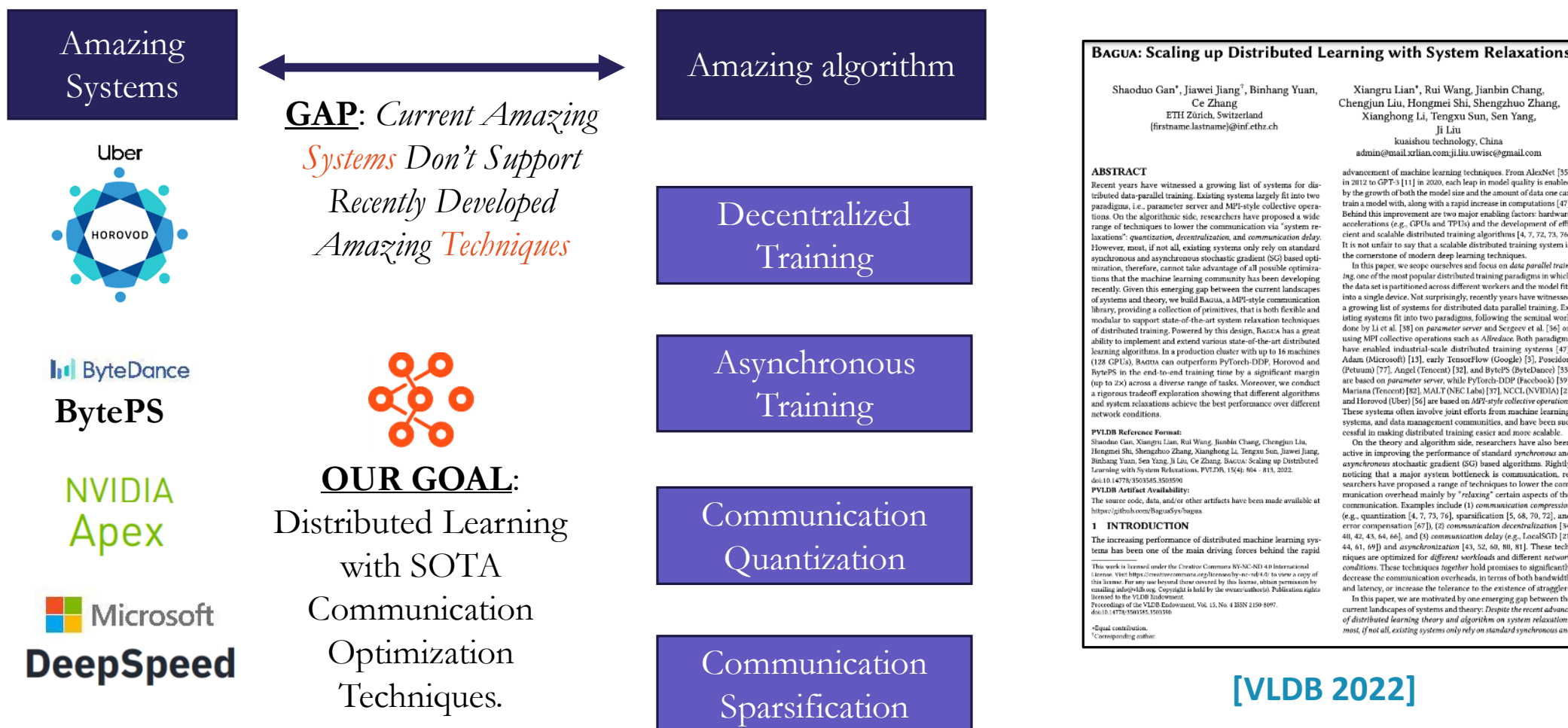
$$O(1/\sqrt{nT} + \tau/T)$$


$$O(1/\sqrt{nT} + \rho/T^{1.5})$$

$\rho$ : network topology constant

# Attempt 1

## Automatic System Optimization for Relaxed Algorithms





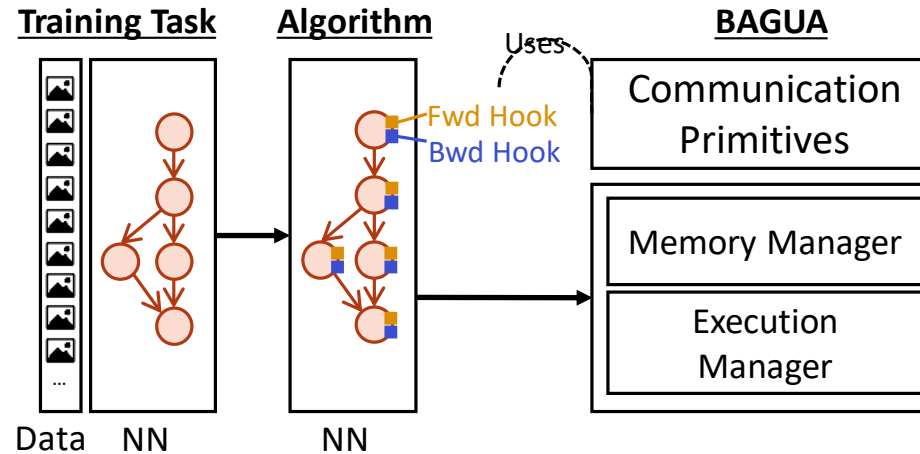
It is not easy to translate *algorithmic flexibility* into *system performance gain*.

# Bagua: System Design & Implementation



- A modular design to accommodate the diversity of different algorithms and communication patterns.
- An optimization framework that applies automatically to an algorithm implemented in BAGUA.

End user: simply wrap up your training script with BAGUA. Specify the algorithm you want to use



## MPI-Style

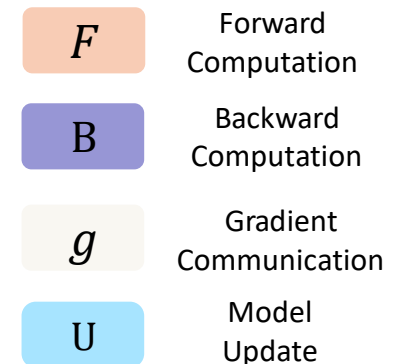
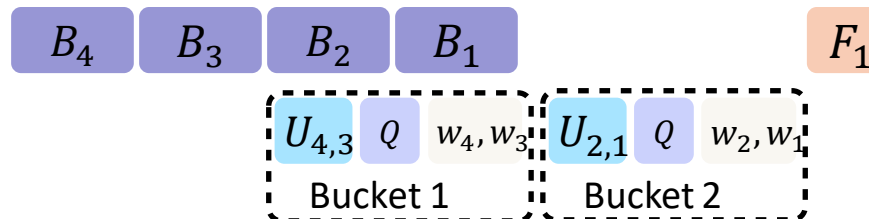
FCS: Full Prec., Centralized, Sync  
FDS: Full Prec., Decentralized, Sync  
LCS: Low Prec., Centralized, Sync  
LDS: Low Prec., Decentralized, Sync  
...

Optimizer: automatically optimize communication and computations

## E.g., Decentralized, Low Precision Alg.



**Automatic**



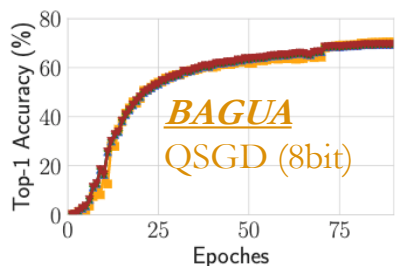
```
1 import torch
2 from bagua import bagua_init, DefaultAlgo
3
4 def main():
5     args = parse_args()
6
7     # define model and optimizer
8     model = MyNet().to(args.device)
9     optimizer = torch.optim.SGD(model.parameters(), lr=args.lr)
10    # transform to BAGUA wrapper
11    model, optimizer = bagua_init(model, optimizer, DefaultAlgo,
12                                  is_intra)
13
14    # train the model over the dataset
15    for epoch in range(args.epochs):
16        for b_idx, (inputs, targets) in enumerate(train_loader):
17            outputs = model(inputs)
18            loss = torch.nn.CrossEntropyLoss(outputs, targets)
19            optimizer.zero_grad()
20            loss.backward()
21            optimizer.step()
```



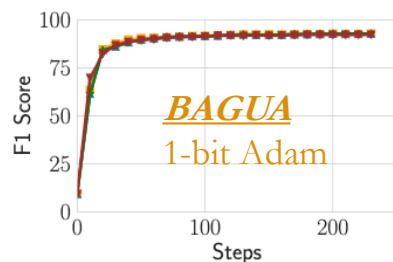
# Bagua Results



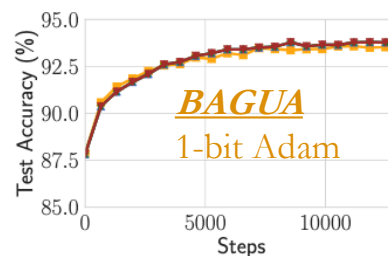
**Setup:** 16 machines, each 8 V100 GPUs. Connected via {10Gbps, 25Gbps, 100Gbps} networks.



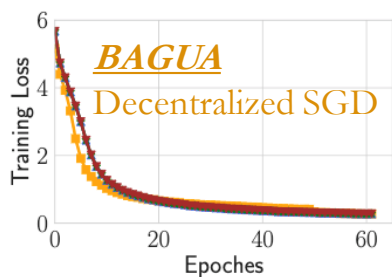
(a) VGG16



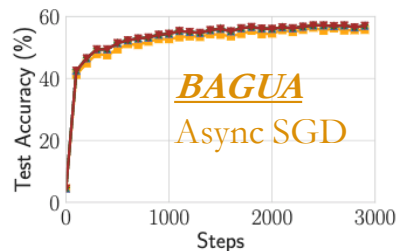
(b) BERT-LARGE Finetune



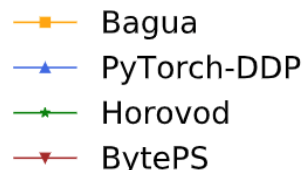
(c) BERT-BASE Finetune



(d) Transformer



(e) LSTM+AlexNet



Network Conditions	VGG16	BERT-LARGE	BERT-BASE	Transformer	LSTM+AlexNet
100 Gbps	1.1×	1.05×	1.27×	1.2×	1.34×
25 Gbps	1.1×	1.05×	1.27×	1.2×	1.34×
10 Gbps	1.94×	1.95×	1.27×	1.2×	1.34×

**Significant speed-up over {Torch-DDP, Horovod 32bits, Horovod 16bits, BytePS}**

*Supporting a diverse set of algorithms can provide significant improvements over existing systems.*

**Same Convergence with Relaxed Algorithms**

# From Cloud to Decentralized Compute Resource

Instance Size	vCPUs	Instance Memory (GiB)	GPU – A100	GPU memory	Network Bandwidth (Gbps)	GPUDirect RDMA	Storage (GB)	Bandwidth (Gbps)	Price/hr
p4d.24xlarge	96	1152	8	320 GB HBM2	400 ENA and EFA	Yes	60 NV		
p4de.24xlarge (preview)	96	1152	8	640 GB HBM2e	400 ENA and EFA	Yes	60 NV		



This is \$4.09/hour for an A100 GPU.



Interruptible ☒ On-Demand #GPUs: ANY 0X 1X 2X 4X 8X 8X+

m.7424	datacenter-40660	Netherlands, NL	Motherboard	↑628 Mbps ↓602 Mbps	0 ports	verified	\$0.500/hr
	1x A100 SXM4		PCIe 4.0, 16x 20.8 GB/s				
	19.5 TFLOPS	80 GB	AMD EPYC 7542 ...				
	Max CUDA: 11.7	1401.7 GB/s	24.0/24 cpu 129/129 GB	Storage 583 MB/s	270.0 GB		

Type #5995170

m.7207

host-33081

Not Specified PCIe 4.0, 16x 19.8 GB/s | ↑11 Mbps ↓321 Mbps | 250 p |  |  ||  | 1x A100 SXM4 |  |  |  |  |  |  |
|  | 19.5 TFLOPS | 39 GB | AMD EPYC 7763 ... |  |  |  |  |
|  | Max CUDA: 11.8 | 1140.6 GB/s | 64.0/256 cpu 121/483 GB | nvme 1008 MB/s | 813.1 GB |  |  |

Type #5105289

m.5308

host-33081

Texas, US 08XP3P | ↑11 Mbps ↓317 Mbps | 4 ports |  |  ||  | 1x A100 SXM4 |  | PCIe 4.0, 16x 22.5 GB/s |  |  |  |  |
|  | 19.5 TFLOPS | 40 GB | AMD EPYC 7513 ... |  |  |  |  |
|  | Max CUDA: 11.7 | 1130.8 GB/s | 32.0/128 cpu 64/258 GB | 1218 MB/s | 238.4 GB |  |  |

44.4 DLPerf

48.9 DLP\$/hr

Reliability 99.69%

MAKE BID

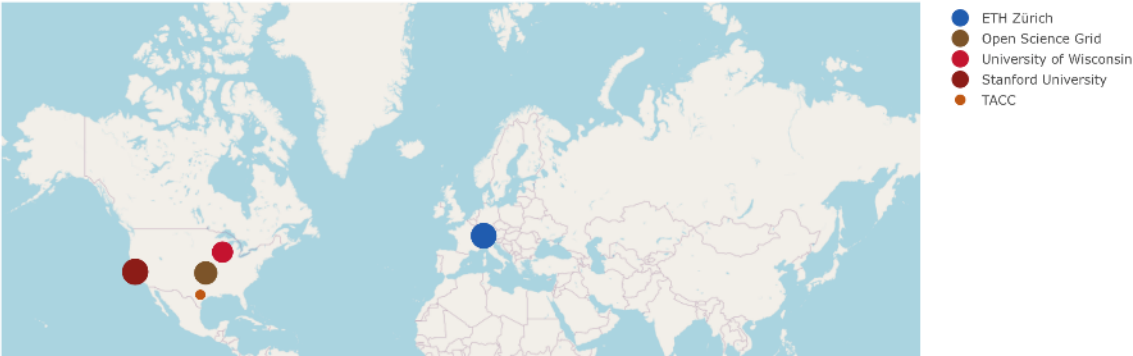
This is what you can get from a decentralized GPU pool!

Available TFlops  
71,509 TFlops

Available GPUs  
# 536

Total TFlops  
124,428 TFlops

Status Global View



# Attempt 2

These algorithmic  
building blocks need  
to *be put together!*

## CocktailSGD: Fine-tuning Foundation Models over 500Mbps Networks

Jue Wang<sup>\*1</sup> Yucheng Lu<sup>\*2</sup> Binhang Yuan<sup>1</sup> Beidi Chen<sup>3</sup> Percy Liang<sup>4</sup> Christopher De Sa<sup>2</sup> Christopher Re<sup>4</sup>  
Ce Zhang<sup>1</sup>

### Abstract

Distributed training of foundation models, especially large language models (LLMs), is communication-intensive and so has heavily relied on centralized data centers with fast interconnects. *Can we train on slow networks and unlock the potential of decentralized infrastructure for foundation models?* In this paper, we propose COCKTAILSGD, a novel communication-efficient training framework that combines three distinct compression techniques—random sparsification, top-K sparsification, and quantization—to achieve much greater compression than each individual technique alone. We justify the benefit of such a hybrid approach through a theoretical analysis of convergence. Empirically, we show that COCKTAILSGD achieves up to 117× compression in fine-tuning LLMs up to 20 billion parameters without hurting convergence. On a 500Mbps network, COCKTAILSGD only incurs ~ 1.2× slowdown compared with data center networks.

merely 10 billion tokens would require 6 petaflops-days: 8 A100 GPUs running at 50% capacity for 5 days!

When training foundation models in a distributed way, *communication* is the key bottleneck in scaling. As an example, fine-tuning GPT-J-6B over 10 billion tokens with a batch size of 262K tokens over 4 machines (each with 2 A100 GPUs) would require 915.5 TB data being communicated during the whole training process! The computation time for such a workload is 114 hours, which means that we need to have at least 20 Gbps connections between these machines to bring the communication overhead to the same scale as the computation time. Not surprisingly, today's infrastructure for training and fine-tuning foundation models are largely *centralized*, with GPUs connected via fast 100Gbps–400Gbps connections (Microsoft, 2020).

Such a heavy reliance on centralized networks increases the cost of infrastructure, and makes it incredibly hard to take advantage of cheaper alternatives, including tier 2 to tier 4 clouds, spot instances and volunteer compute. For example, while volunteering compute projects such as Folding@Home can harvest significant amount of computes for embarrassingly parallelizable workloads (e.g., 2.43exaflops in April 2020 (Larson et al., 2009)), it is challenging to harvest these cycles for foundation model training due to the communication bottleneck. Recently, there has been an exciting collection of work focusing on the decentralized training of neural networks, including those that are algorithmic (Lian et al., 2017; Ryabinin & Gusev, 2020; Diskin et al., 2021; Ryabinin et al., 2021; Yuan et al., 2022; Jue et al.) as well as system efforts such as Training Transformer Together (Borzunov et al., 2022b), and PETALS (Borzunov et al., 2022a). However, despite of these recent efforts, communication is still a significant bottleneck, and one can only compress the communication by at most 10-30× in these recent efforts without hurting convergence. To fully close the gap between centralized infrastructure (100Gbps) and decentralized infrastructure (100Mbps-1Gbps), we need to decrease the communication overhead by at least 100×!

Luckily, there have also been rapid development of communication-efficient optimization algorithms and these efforts provide the foundational building blocks of this paper. Researchers have proposed a wide range of

### 1. Introduction

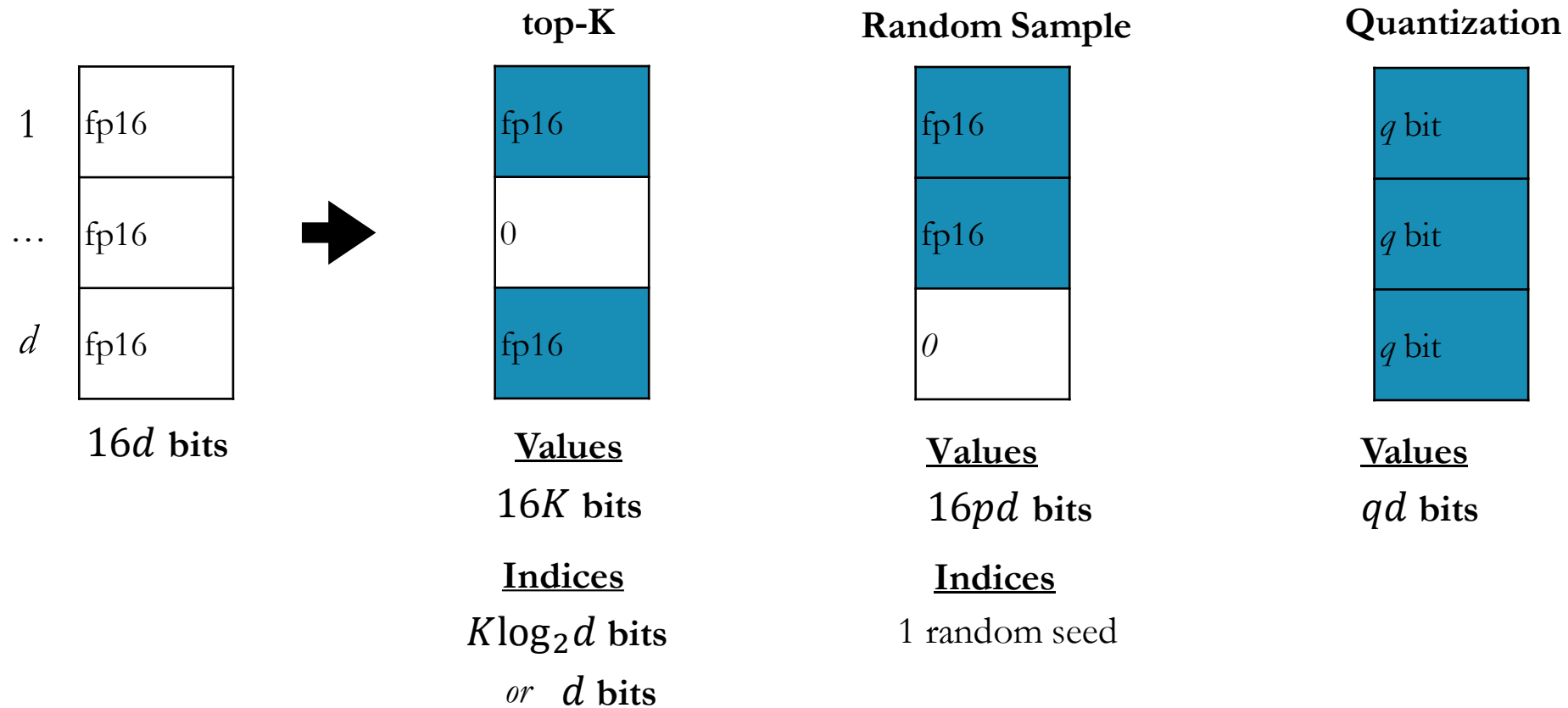
In recent years, foundation models (Bommasani et al., 2021), including large language models (Brown et al., 2020; Chowdhery et al., 2022; Bommasani et al., 2021; Zhang et al., 2022; Liang et al., 2022; Scao et al., 2022), have enabled rapid advancement for various machine learning tasks, especially in natural language processing (Brants et al., 2007; Austin et al., 2021). Such a significant improvement on quality has been fueled by an ever-increasing amount of data and computes that are required in training these models (Kaplan et al., 2020). Today, training even modest scale models requires a significant amount of compute: For example, fine-tuning GPT-J-6B (6 billion parameters) over

<sup>\*</sup>Equal contribution <sup>1</sup>ETH Zürich, Switzerland <sup>2</sup>Cornell University, USA <sup>3</sup>Carnegie Mellon University, USA <sup>4</sup>Stanford University, USA. Correspondence to: Jue Wang <juewang@inf.ethz.ch>.

Proceedings of the 40<sup>th</sup> International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

[ICML 2023]

# Three Methods of Compression



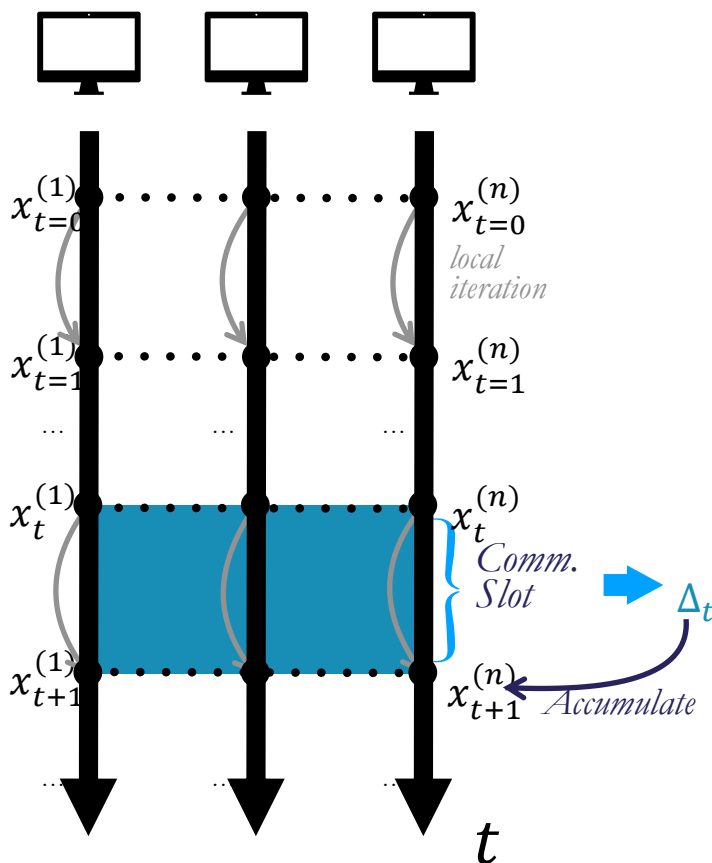
Expensive to compute  
and to encode Indices

Might not keep top  
values as in Top-K

Only provide up to 16x  
compression; hard to go aggressive

It is very hard to reach *100X*  
*compression* ratio with a single method.

# CocktailSGD: Mixture of Compression Methods



*Idea: A Mixture of communication compression techniques.*

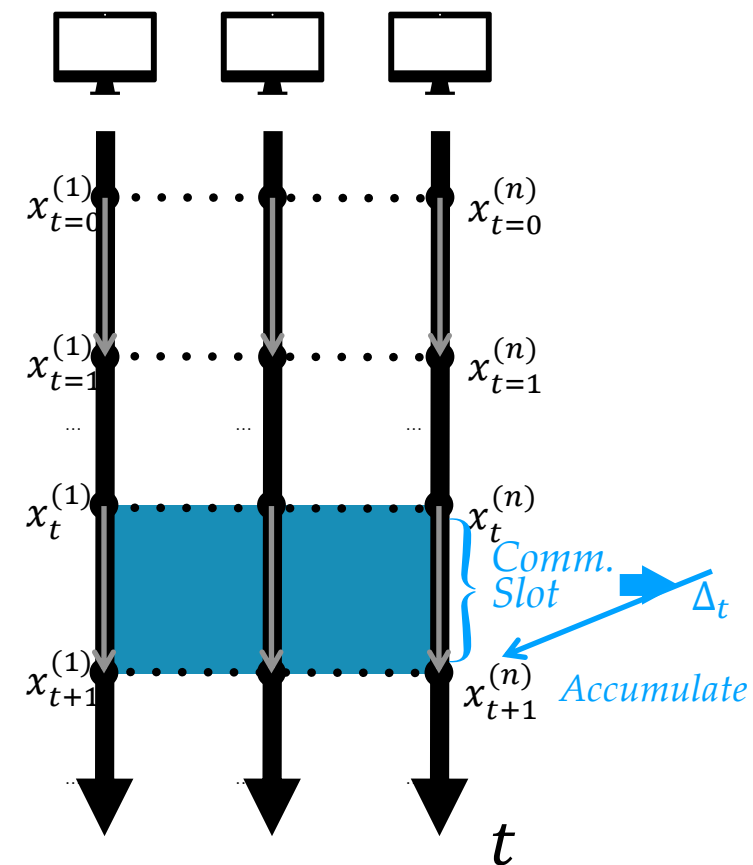
Looking at  $\Delta_t$ :

- It has 1-step staleness ——— // asynchrony
- At  $t$ , randomly pick  $p\%$  parameters to communicate ——— // local training: compress  $\sim \frac{1}{p\%} \times$
- For selected parameters, let  $\delta_t^{(i)}$  be local model updates since last communication:
  - $\tilde{\delta}_t^{(i)} = \text{top-}K\%(\delta_t^{(i)})$  ——— // topK: compress  $\sim \frac{1}{K\%} \times$
  - $\hat{\delta}_t^{(i)} = \text{Quantize}(\tilde{\delta}_t^{(i)}, q \text{ bits})$  ——— // Quantization: compress  $\sim \frac{16}{b} \times$
- Communicate:  $\Delta_t = \sum_i \hat{\delta}_t^{(i)}$

As long as **Communication** fully fills the **Comm. Slot**, no slow down caused by communication.

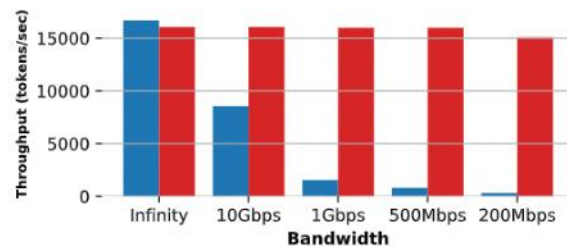
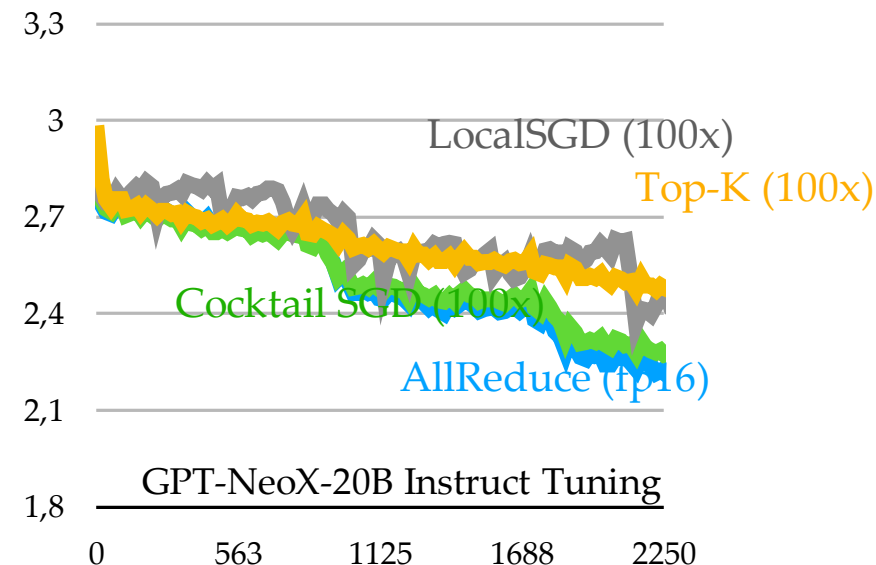
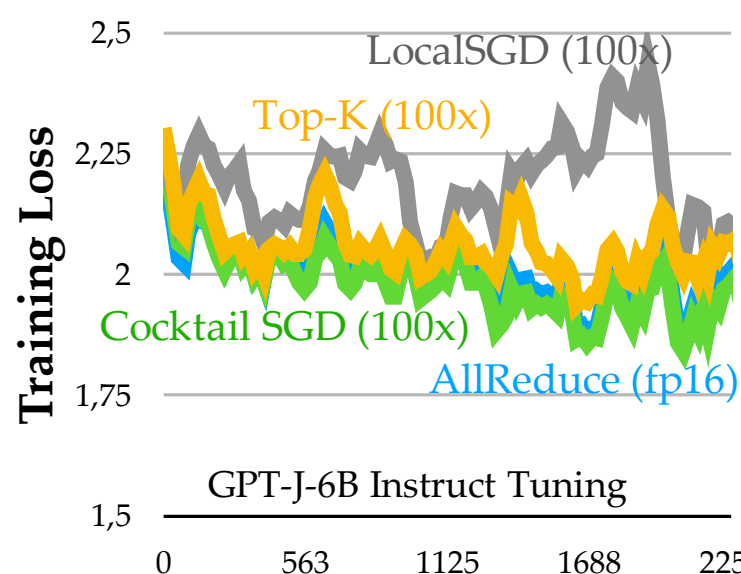
# “Cocktail SGD”: Data Parallel over 1Gbps

AI Infra

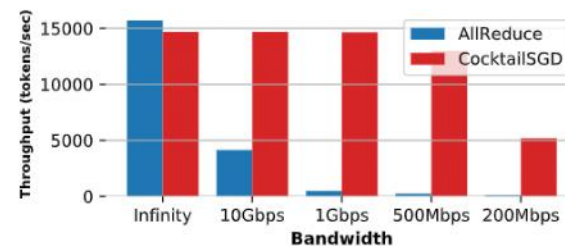


As long as **Communication** fully fills the **Comm. Slot**, no slow down caused by communication.

Different communication compression techniques complement each other and compose well!



(b) GPT-J-6B



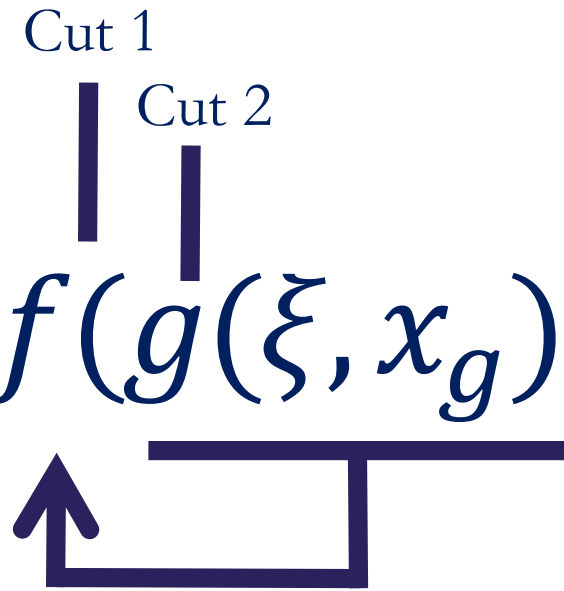
(c) GPT-NeoX-20B

Data parallel  
over ~500 Mbps  
network!



Large language model training goes  
*beyond data parallelism.*



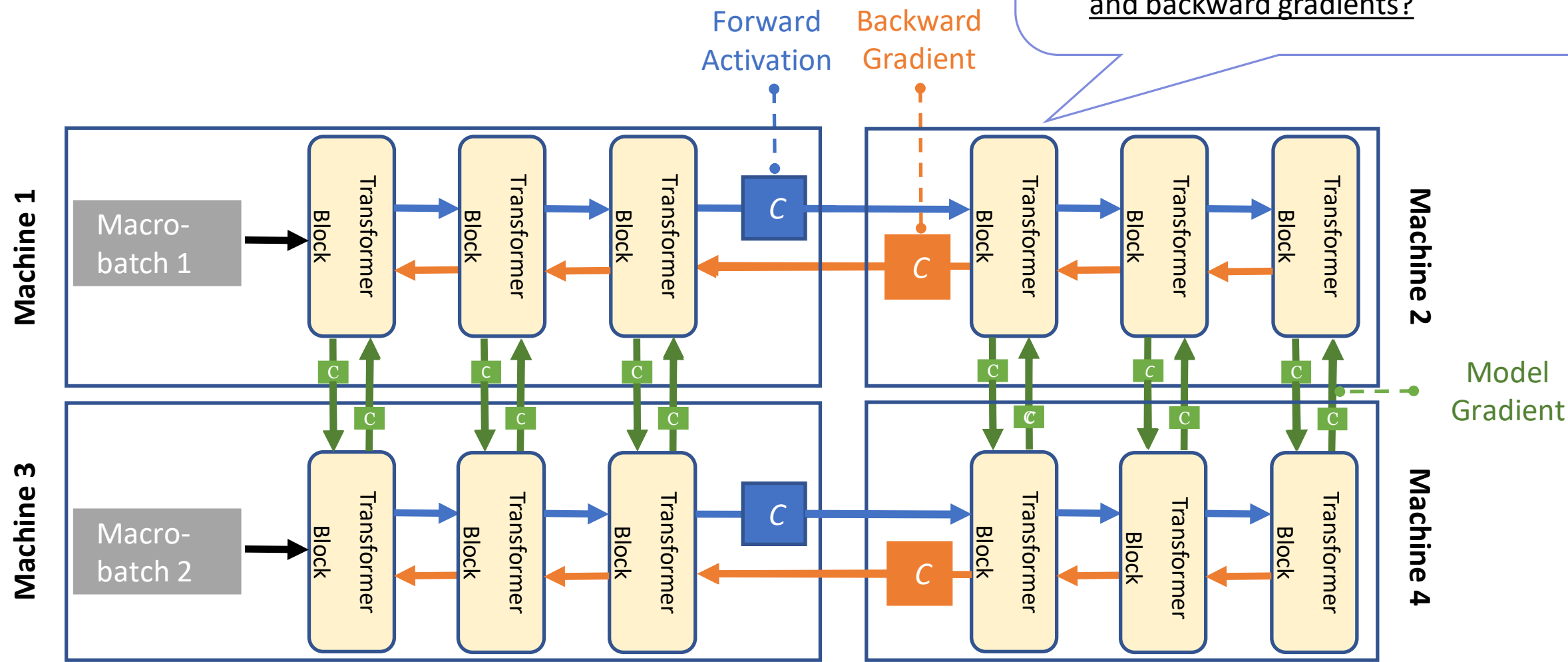
$$\min_x \mathbb{E}_\xi f(\xi, x) \quad \Rightarrow \quad \min_{x_f, x_g} \mathbb{E}_\xi f(g(\xi, x_g), x_f)$$


*Forward Activation*

- (GPT-3) 24MB / 1000tokens

# 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)]

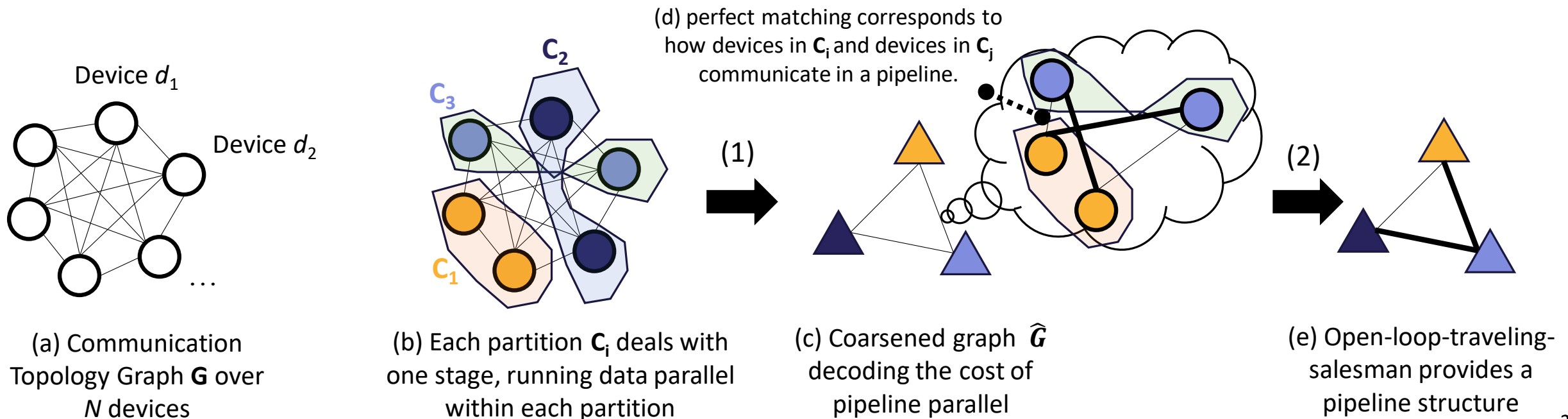


[NeurIPS 2022-(b)]

# Accommodate Communication in a Decentralized network

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

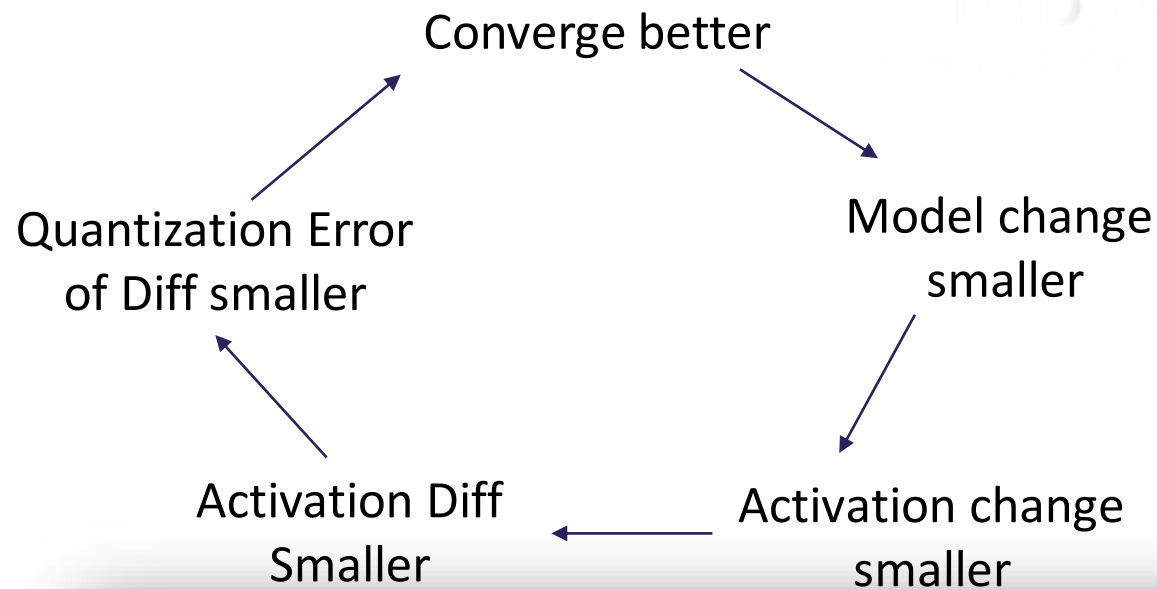
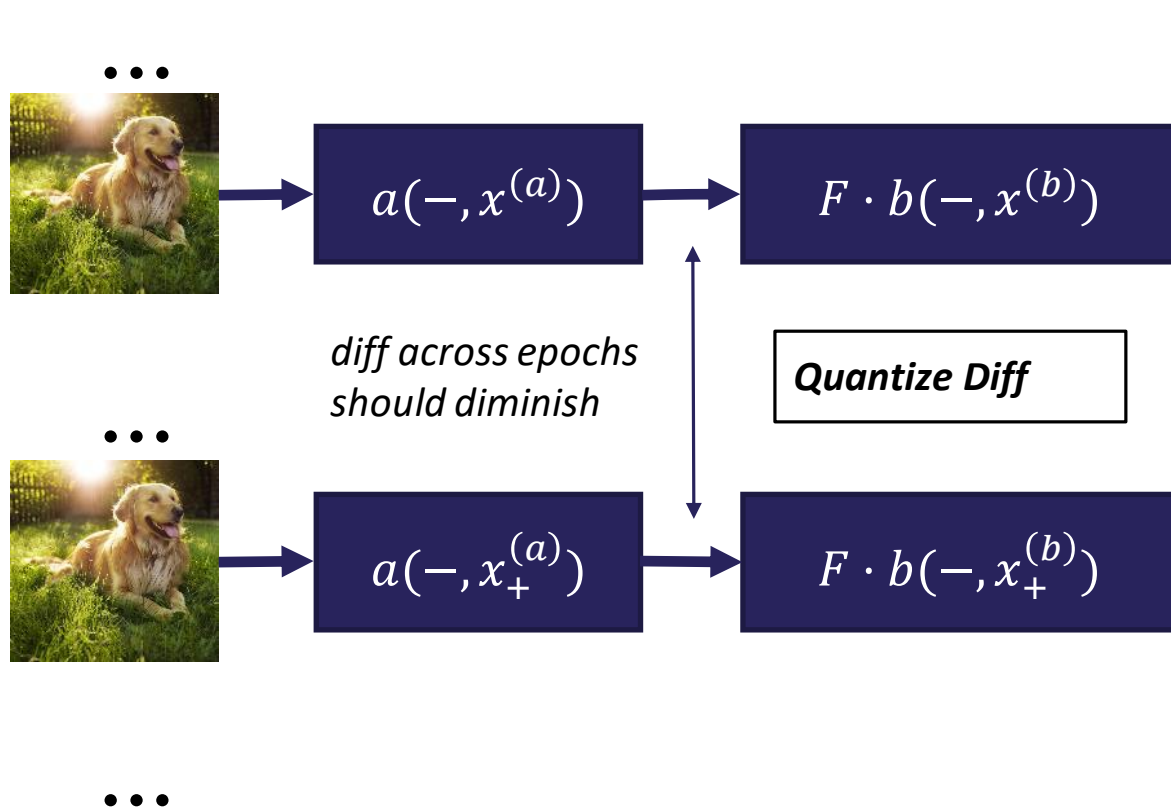
- Data parallel communication cost: nodes handling the same stage need to exchange gradients;
- Pipeline parallel communication cost: nodes handling nearby stages for the same micro-batch need to communicate activation in the forward propagation and gradients of the activation in the backward propagation.



# 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.



- **(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

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

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

- **(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.  $\mathbb{E}_\xi\|g_\xi(x) - \nabla f(x)\|^2 \leq \sigma^2$ , 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  $\mathbb{E}\|x - Q(x)\| \leq c_Q \|x\|$ , for all  $x$ .<sup>1</sup> Let  $\gamma = \frac{1}{3(C+3L_f)\sqrt{T}}$  be the learning rate, where

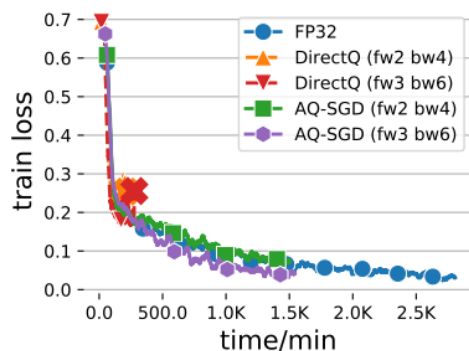
$$C = \frac{4c_Q \ell_a (1 + C_a) L_{f \circ b} N}{\sqrt{1 - 2c_Q^2}}.$$

Then after performing  $T$  updates one has

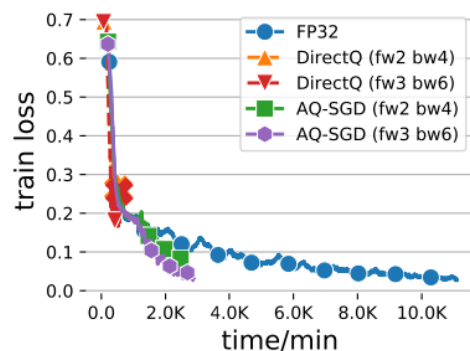
$$\frac{1}{T} \sum_{t \in [T]} \mathbb{E} \|\nabla f(x_t)\|^2 \lesssim \frac{(C + L_f)(f(x_1) - f^*)}{\sqrt{T}} + \frac{\sigma^2 + (c_Q C_a C_{f \circ b})^2}{\sqrt{T}}. \quad (3.1)$$

# AQ-SGD Results

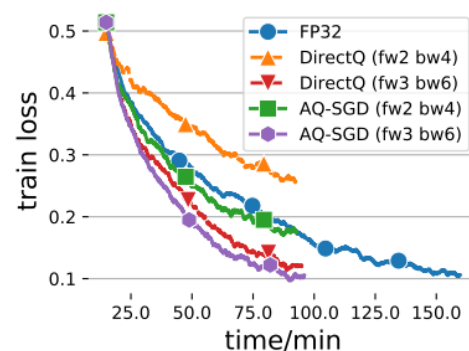
- End-to-end training performance over different networks. x represents divergence.



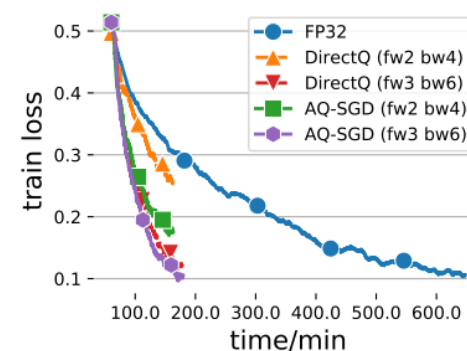
(a) QNLI, 500Mbps



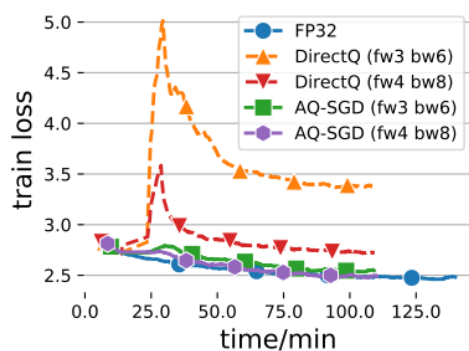
(b) QNLI, 100Mbps



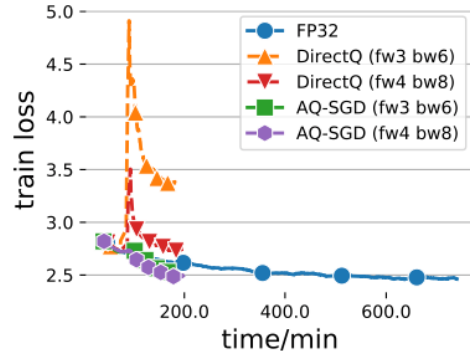
(c) CoLA, 500Mbps



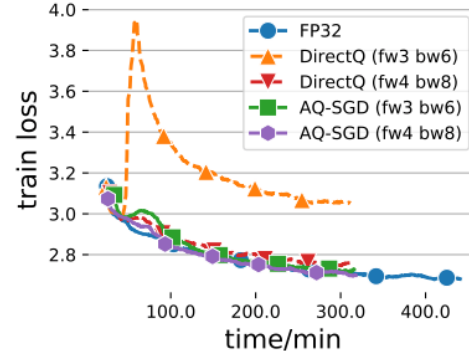
(d) CoLA, 100Mbps



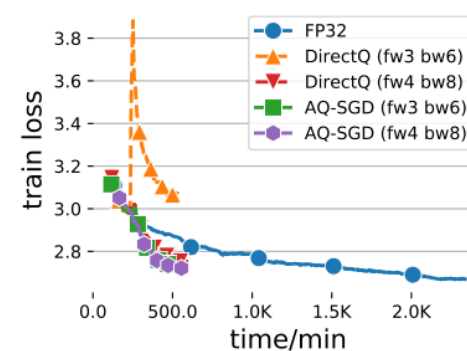
(e) WikiText2, 500Mbps



(f) WikiText2, 100Mbps



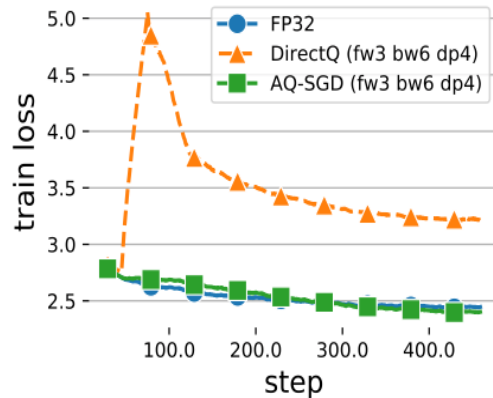
(g) arXiv, 500Mbps



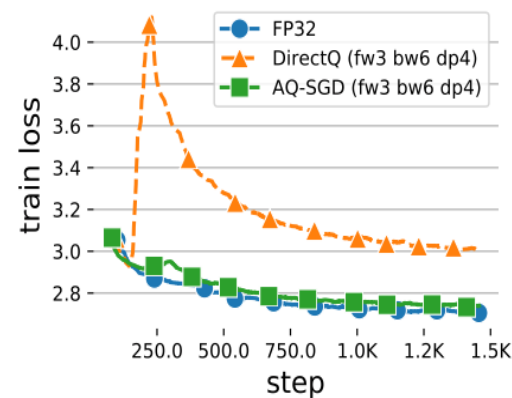
(h) arXiv, 100Mbps

# AQ-SGD Results

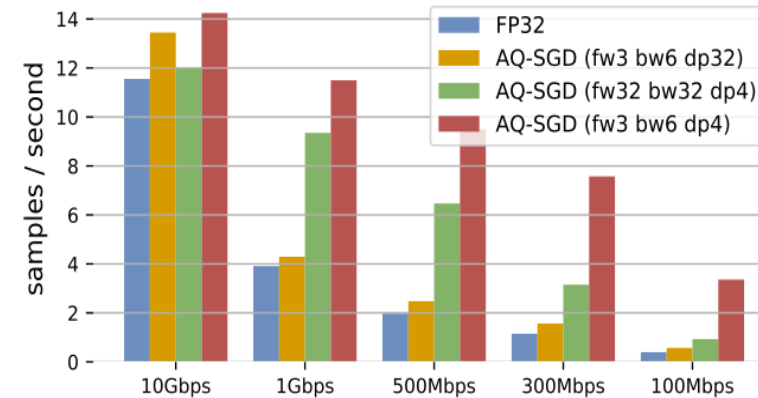
- Convergence and Throughput of AQ-SGD combined with gradient compression.



(a) WikiText2, GPT2-1.5B



(b) arXiv, GPT2-1.5B



(c) Training Throughput

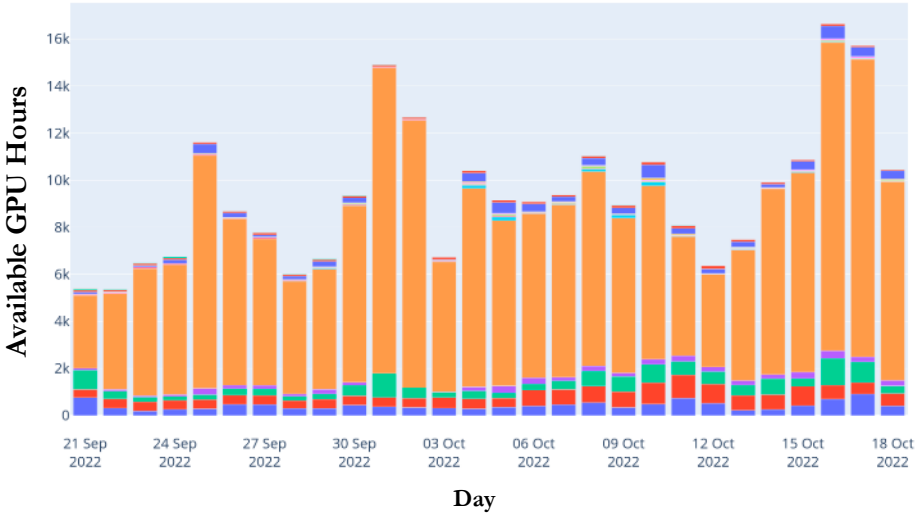


## Some Small Steps Towards *Decentralized ML*.



# Open Research on the Together Decentralized Cloud

Connecting idle compute across academic institutions.



**HELM** 11 billion tokens  
60K GPU Hours  
10 Open Models


BLOOM	176B	July 2022
T0pp	11B	October 2021
GPT-J	6B	July 2021
GPT-NeoX	20B	February 2022
GLM	130B	August 2022
UL2	20B	October 2022
T5	11B	February 2020
OPT	175B	June 2022
OPT	66B	June 2022
YaLM	100B	June 2022

# 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:
  - *Data Parallel*: {asynchronous, local training, compression, quantization, decentralized topology} & their combinations.
  - *Model Parallel*: Careful error compensation.
- Innovation of **Systems** is need to unleash the full potential **Algorithms**:
  - *Bagua*: Automatic optimization framework.
  - *System Scheduling of communication in decentralized environments*.



## *Ongoing Research.*



Large language model training goes  
*beyond* data & pipeline parallelism.

# Scheduling in a Bigger Scope

- Heterogeneous hardware:
  - A100, A800, A40, 3090, 4090, etc.
- Heterogeneous connections:
  - NVSwitch, NVLink, RDMA, etc.
  - Cross data-center/ @home Network.
- Parallel schemas:
  - Data parallelism;
  - Pipeline parallelism;
  - Tensor model parallelism;
  - Optimizer parallelism (FSDP).



# Communication Efficient Parallel LLM Training Algorithms

- Communication compression for different parallel schemas:
  - Data parallelism;

- Pipeline parallelism;
- Tensor model parallelism (?)
- Optimizer parallelism (FSDP) (?)

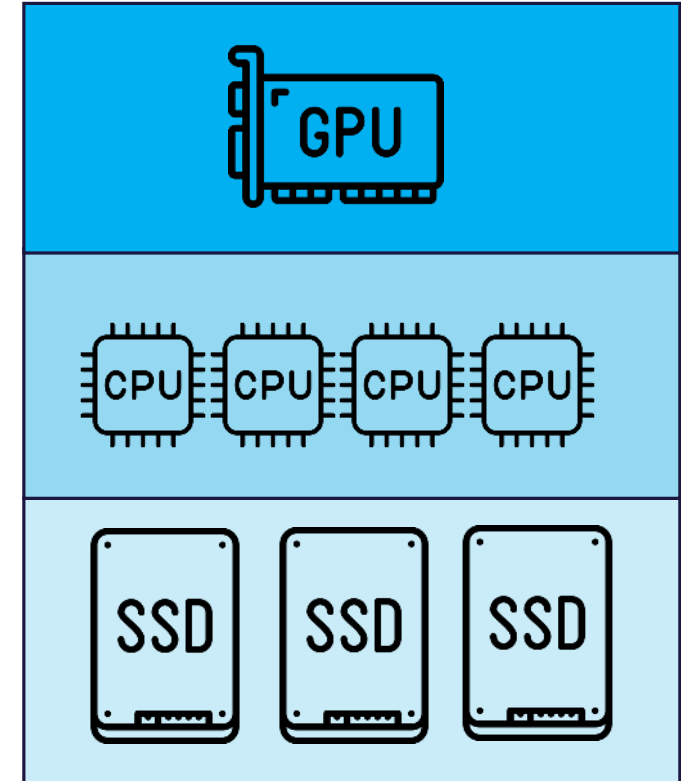


[ICML 2023]



[NeurIPS 2022]

*Enhancing High-Throughput LLM  
Inference through Off-loading  
Systems.*





# FlexGen

- OPT-175B Scale Inference

on a single GPU:

- 6.5K ★ on Github;
- Top discussion on Hacker News;
- High throughput scenario: *1 token/s*.

news.ycombinator.com

**Hacker News** login

new | past | comments | ask | show | jobs | submit

1. ▲ Running large language models like ChatGPT on a single GPU (github.com/ying1123)  
341 points by \_nhynes 3 hours ago | hide | 118 comments
2. ▲ CheatGPT (humphd.org)  
64 points by jicea 1 hour ago | hide | 53 comments
3. ▲ Show HN: Turn Your Pandas Dataframe into a Tableau-Style UI for Visual Analysis (github.com/kanaries)  
245 points by AwsDef 4 hours ago | hide | 24 comments
4. ▲ Opossum: Cross-platform web browser written in Golang, optimized for Plan 9 (github.com/psilva261)  
207 points by euclaise 5 hours ago | hide | 53 comments
5. ▲ California AG and state senators introduce bill to ban hidden fees (ca.gov)  
267 points by miguelazo 2 hours ago | hide | 261 comments
6. ▲ Creation happens in silence (josem.co)  
212 points by josem 6 hours ago | hide | 94 comments
7. ▲ The Philosophy of Computer Science (stanford.edu)  
20 points by lucidguppy 1 hour ago | hide | 16 comments
8. ▲ Photographer captures image of rare fish that walks on its 'hands' (cnn.com)  
93 points by bryan0 5 hours ago | hide | 27 comments
9. ▲ The Eleven Laws of Showrunning [pdf] (okbjgm.weebly.com)  
105 points by luu 6 hours ago | hide | 38 comments
10. ▲ Low-Level Software Security for Compiler Developers (lloftsec.github.io)  
17 points by signa11 2 hours ago | hide | 1 comment

Search or jump to... Pulls Issues Codespaces Marketplace Explore

FMInference / FlexGen Public Edit Pins Unwatch 77 Fork 332 Starred 6.5k

Code Issues 21 Pull requests 3 Actions Projects Wiki Security Insights

main Go to file Add file <> Code About

Running large language models on a single GPU for throughput-oriented scenarios.

machine-learning deep-learning offloading high-throughput opt gpt-3 large-language-models

Readme Apache-2.0 license 6.5k stars 77 watching 332 forks

**Releases**

No releases published  
[Create a new release](#)

**Packages**

No packages published  
[Publish your first package](#)

**Used by** 5

**Contributors** 18  
[+ 7 contributors](#)

BinhangYuan Data wrangle benchmark (#95)	yesterday	94
benchmark	Update Petals setup details	2 days ago
docs	Update gcp_setup.md	last week
flexgen	Data wrangle benchmark (#95)	yesterday
scripts	Simplify API (#68)	2 weeks ago
.gitignore	FlexGen for Data wrangle Tasks. (#91)	2 days ago
LICENSE	Release and merge commits	2 weeks ago
README.md	Update README.md	yesterday
pyproject.toml	Add more HELM examples (#82)	last week

**README.md**

## FlexGen

FlexGen is a high-throughput generation engine for running large language models with limited GPU memory. FlexGen allows **high-throughput** generation by IO-efficient offloading, compression, and **large effective batch sizes**.

## Throughput-Oriented Inference for Large Language Models

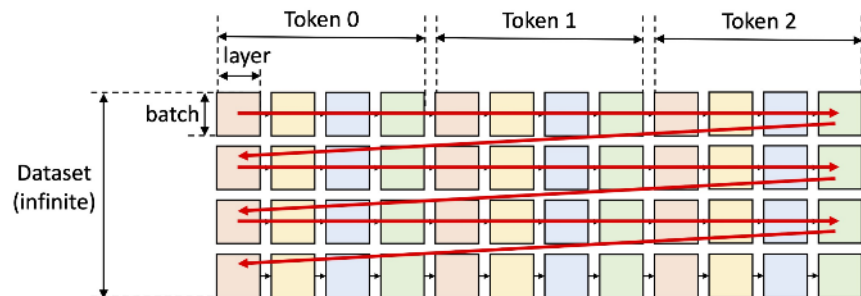
In recent years, large language models (LLMs) have shown great performance across a wide range of tasks. Increasingly, LLMs have been applied not only to interactive applications (such as chat), but also to many "back-of-house" tasks. These tasks include benchmarking, information extraction, data wrangling, and form processing.

One key characteristic of these applications is that they are **throughput-oriented**: they require running LLM inferences over millions of tokens in batches, e.g., all the private documents in a company's corpus, or all the tasks in the [HELM](#) benchmark. These workloads are less sensitive to latency - the user starts up a job and lets it

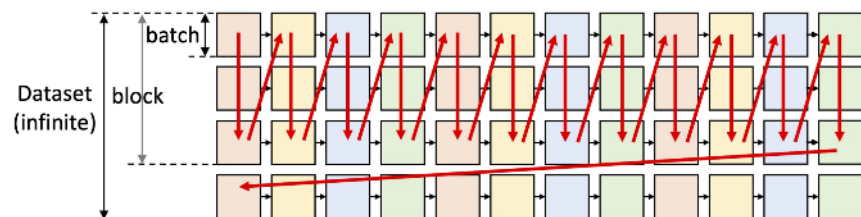


# FlexGen

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

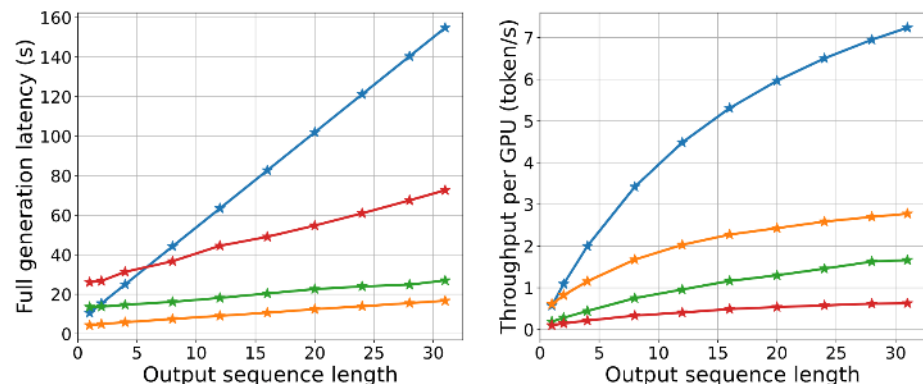


(a) Row-by-row schedule



(b) Zig-zag block schedule

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



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

Ying Sheng<sup>1</sup>, Liannan Zheng<sup>2</sup>, Binhang Yuan<sup>3</sup>, Zhuohan Li<sup>2</sup>, Max Ryabinin<sup>1\*</sup>,  
Beidi Chen<sup>4\*</sup>, Percy Liang<sup>1</sup>, Christopher Ré<sup>1</sup>, Ion Stoica<sup>2</sup>, Ce Zhang<sup>5</sup>

#### 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>.

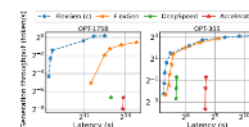


Figure 1. The total latency for a block and throughput trade-offs of three offloading-based systems for OPT-175B (left) and OPT-30B (right) on a single NVIDIA T4 (16 GB) GPU with 208 GB CPU DRAM and 1.5 TB SSD. FlexGen achieves a new Pareto-optimal frontier with 100x higher maximum throughput for OPT-175B.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

<sup>1</sup>Stanford University <sup>2</sup>UC Berkeley <sup>3</sup>ETH Zurich <sup>4</sup>Yandex <sup>5</sup>ISLE University \*Meta \*Carnegie Mellon University. Correspondence to: Ying Sheng <yingsheng@stanford.edu>.

Proceedings of the 40<sup>th</sup> International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

[ICML 2023 Oral]

## FlexGen 2.0

- Efficient support of multiple weak GPUs;
- Easy support for general transformer models;
- Integrate multiple relaxed computations (sparsified, quantized, etc.).

# 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:
  - *Data Parallel*: {asynchronous, local training, compression, quantization, decentralized topology} & their combinations.
  - *Model Parallel*: Careful error compensation.
- Innovation of **Systems** is need to unleash the full potential *Algorithms*:
  - *Bagua*: Automatic optimization framework.
  - *System Scheduling of communication in decentralized environments*.



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

Thank you!