

# Do You Want to Fine-tune a 175B LLM with a 24GB GPU?

Liangyu Wang, Jie Ren, Hang Xu, Junxiao Wang, David E. Keyes, Di Wang

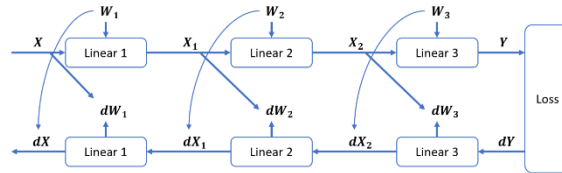
## ZO (Zeroth-order Optimizer) & Motivation

Given a loss function  $f(\cdot)$  and a model  $x$  with parameters in  $d$  dimensions, the gradient  $\nabla f(x)$  can be estimated by:

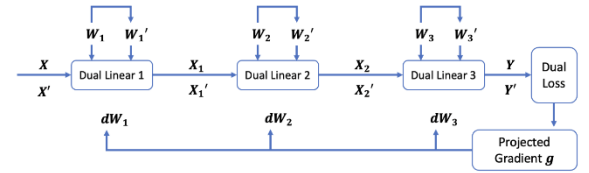
$$\hat{\nabla} f(x) = gz \in \mathbb{R}^d,$$

$$g = \frac{f(x + \epsilon z) - f(x - \epsilon z)}{2\epsilon} \in \mathbb{R}^1,$$

where  $z$  is a random direction vector drawn from the standard Gaussian Distribution  $N(0, I)$ , and  $\epsilon$  is a small perturbation step size.

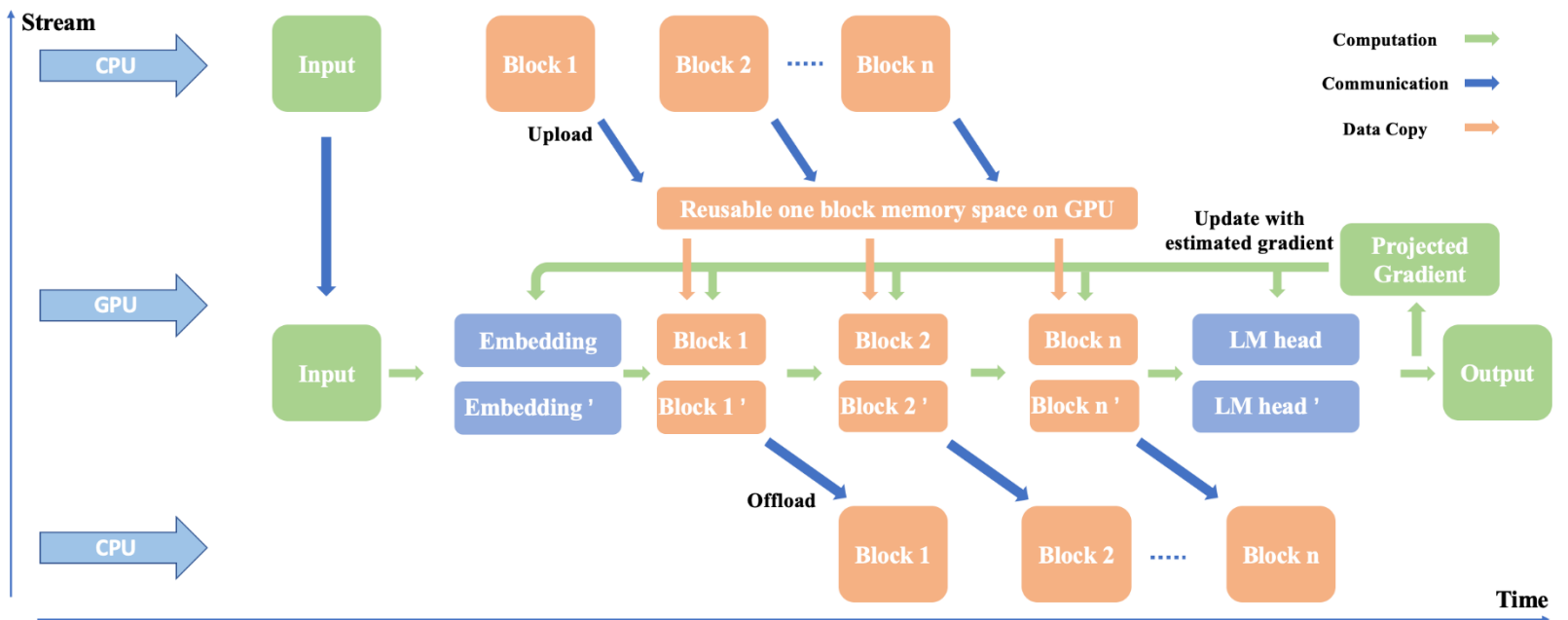


(a) Model using first-order optimizer with forward-backward passes workflow

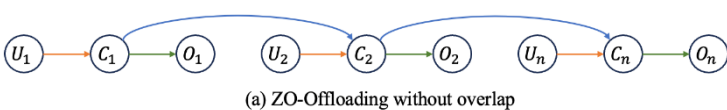


(b) Model using zero-order optimizer with only forward passes workflow

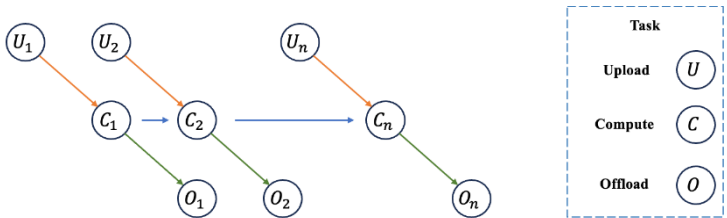
## ZO-Offloading Framework



## Dynamic Scheduler for ZO-Offloading



(a) ZO-Offloading without overlap



(b) ZO-Offloading with overlap

### Algorithm 1 ZO-Offloading Dynamic Scheduler

**Require:** Transformer blocks  $\{W_i\}_{i=1}^N$  with number of transformer blocks  $N$ , embedding parameters  $Embedding$ , and LM head  $LMhead$ .

- 1: Initialize a dynamic scheduler  $S\{\cdot\}$  to control dual forward computation  $C(\cdot)$ , uploading  $U(\cdot)$ , and offloading  $O(\cdot)$  operations.
- 2: Asynchronously launch  $S\{U(W_1), C(Embedding)\}$ .
- 3: **for**  $i = 1$  to  $N - 1$  **do**
- 4:     Synchronously wait until  $U(W_i)$  finished.
- 5:     **if**  $i = 1$  **then**
- 6:         Asynchronously launch  $S\{U(W_{i+1}), C(W_i)\}$ .
- 7:     **else**
- 8:         Synchronously wait until  $C(W_{i-1})$  finished.
- 9:         Asynchronously launch  $S\{U(W_{i+1}), C(W_i), O(W_{i-1})\}$ .
- 10:     **end if**
- 11: **end for**
- 12: Synchronously wait until  $U(W_N)$  and  $C(W_{N-1})$  finished.
- 13: Asynchronously launch  $S\{C(W_N), O(W_{N-1})\}$ .
- 14: Synchronously wait until  $C(W_N)$  finished.
- 15: Asynchronously launch  $S\{C(LMhead), O(W_N)\}$ .

## Experiment Results Compared with MeZO<sup>[1]</sup>

Model	GPU Memory Usage (MB) ↓				Throughput (tokens/sec) ↑			
	MeZO(32)	ZO-Offload(32)	MeZO(16)	ZO-Offload(16)	MeZO(32)	ZO-Offload(32)	MeZO(16)	ZO-Offload(16)
OPT-125M	3091	2941(x0.95)	1801(x0.58)	<b>1661(x0.54)</b>	14889	13074(x0.89)	<b>31058(x2.09)</b>	<b>31058(x2.09)</b>
OPT-350M	4219	3393(x0.81)	2389(x0.57)	<b>1643(x0.39)</b>	5274	5099(x0.97)	<b>13508(x2.56)</b>	12284(x2.32)
OPT-1.3B	9117	4413(x0.48)	4887(x0.54)	<b>2651(x0.29)</b>	1954	1954(x1.00)	<b>6788(x3.47)</b>	<b>6788(x3.47)</b>
OPT-2.7B	15277	5261(x0.34)	7933(x0.52)	<b>3111(x0.20)</b>	1087	1087(x1.00)	<b>4227(x3.89)</b>	<b>4227(x3.89)</b>
OPT-6.7B	32083	8329(x0.26)	16311(x0.51)	<b>4539(x0.14)</b>	499	499(x1.00)	<b>2455(x4.92)</b>	<b>2455(x4.92)</b>
OPT-13B	58251	12113(x0.21)	29411(x0.50)	<b>6445(x0.11)</b>	270	270(x1.00)	<b>1406(x5.21)</b>	1340(x4.96)
OPT-30B	-	18879	63953	<b>10369</b>	-	122	<b>651</b>	597
OPT-66B	-	29937	-	<b>14143</b>	-	40	-	<b>273</b>
OPT-175B	-	49203	-	<b>24667</b>	-	14	-	<b>37</b>

[1] Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. 2023. Fine-tuning language models with just forward passes. *Advances in Neural Information Processing Systems* 36 (2023), 53038–53075.