



# A Memory Efficient Subspace Optimization Method for Training Large Language Models

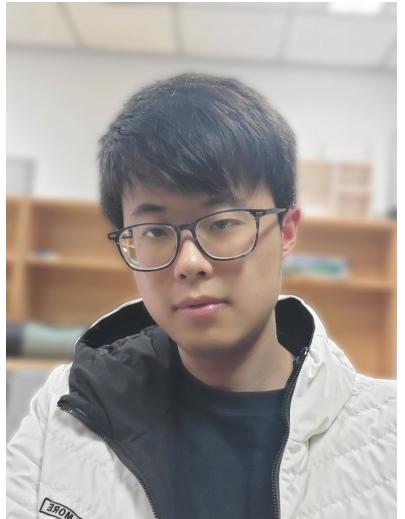
Kun Yuan (袁 坤)

Center for Machine Learning Research @ Peking University

June 24, 2025 Beijing

# Joint work with

---



Yiming Chen  
(PKU)



Yuan Zhang  
(PKU)



Yin Liu  
(PKU)



Zaiwen Wen  
(PKU)

Y. Chen, Y. Zhang, Y. Liu, K. Yuan, Z. Wen, *A Memory Efficient Randomized Subspace Optimization Method for Training Large Language Models*, ICML 2025



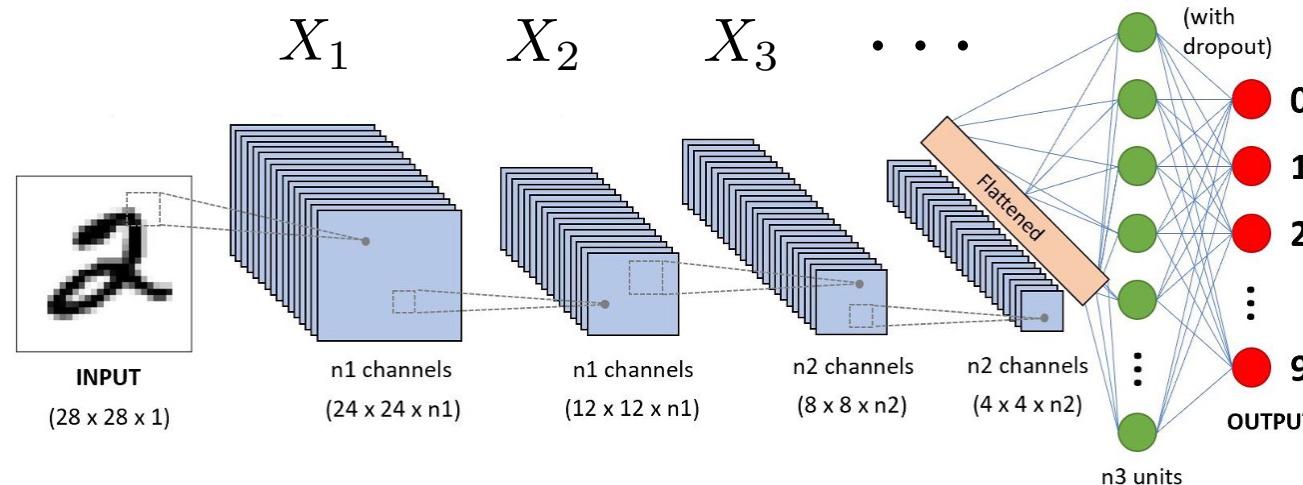
# PART 01

---

## Basics and Motivation

# LLM pretraining is essentially solving stochastic optimization

- The model weights in neural networks are a set of matrices  $\mathbf{X} = \{\mathbf{X}_\ell\}_{\ell=1}^L$



- Let  $h(\mathbf{X}; \xi)$  be the language model;  $\hat{y} = h(\mathbf{X}; \xi)$  is the predicted token

**cross entropy**

LLM cost function:

$$\mathbf{X}^* = \arg \min_{\mathbf{X}} \left\{ \mathbb{E}_{\xi \sim \mathcal{D}} \left[ L(h(\mathbf{X}; \xi), y) \right] \right\}$$

↑  
 data distribution    pred. token    real token

# LLM pretraining is essentially solving stochastic optimization

---

- If we define  $\xi = (\xi, y)$  and  $F(\mathbf{X}; \xi) = L(h(\mathbf{X}; \xi), y)$ , the LLM problem becomes

Stochastic optimization:  $\mathbf{X}^* = \arg \min_{\mathbf{X}} \left\{ \mathbb{E}_{\xi \sim \mathcal{D}} [F(\mathbf{X}; \xi)] \right\}$

- In other words, LLM pretraining is essentially solving a stochastic optimization problem
- Adam is the standard approach in LLM pretraining

**Optimizer states**

$$\begin{aligned}
 \mathbf{G}_t &= \nabla F(\mathbf{X}_t; \xi_t) && \text{(stochastic gradient)} \\
 \boxed{\mathbf{M}_t} &= (1 - \beta_1) \mathbf{M}_{t-1} + \beta_1 \mathbf{G}_t && \text{(first-order momentum)} \\
 \boxed{\mathbf{V}_t} &= (1 - \beta_2) \mathbf{V}_{t-1} + \beta_2 \mathbf{G}_t \odot \mathbf{G}_t && \text{(second-order momentum)} \\
 \mathbf{X}_{t+1} &= \mathbf{X}_t - \frac{\gamma}{\sqrt{\mathbf{V}_t} + \epsilon} \odot \mathbf{M}_t && \text{(adaptive SGD)}
 \end{aligned}$$

Memory = **Model** + **Gradient** + **Optimizer states** + **Activations**

- Given a **model** with **P** parameters, **gradient** will consume **P** parameters, and **optimizer states** will consume **2P** parameters; **4P parameters in total**.

$$\mathbf{P} \quad \mathbf{G}_t = \nabla F(\mathbf{X}_t; \boldsymbol{\xi}_t)$$

$$2\mathbf{P} \quad \left\{ \begin{array}{l} \mathbf{M}_t = (1 - \beta_1)\mathbf{M}_{t-1} + \beta_1 \mathbf{G}_t \\ \mathbf{V}_t = (1 - \beta_2)\mathbf{V}_{t-1} + \beta_2 \mathbf{G}_t \odot \mathbf{G}_t \end{array} \right.$$

$$\mathbf{P} \quad \mathbf{X}_{t+1} = \mathbf{X}_t - \frac{\gamma}{\sqrt{\mathbf{V}_t} + \epsilon} \odot \mathbf{M}_t$$

**Optimizer states introduces significant memory cost**

# Memory cost to pre-train LLMs



Memory = **Model** + **Gradient** + **Optimizer states** + **Activations**

- Activations are auxiliary variables to facilitate the gradient calculations

Consider a linear neural network

$$z_i = X_i z_{i-1}, \forall i = 1, \dots, L$$

$$f = \mathcal{L}(z_L; y)$$

The gradient is derived as follows

$$\frac{\partial f}{\partial X_i} = \frac{\partial f}{\partial z_i} z_{i-1}^\top$$

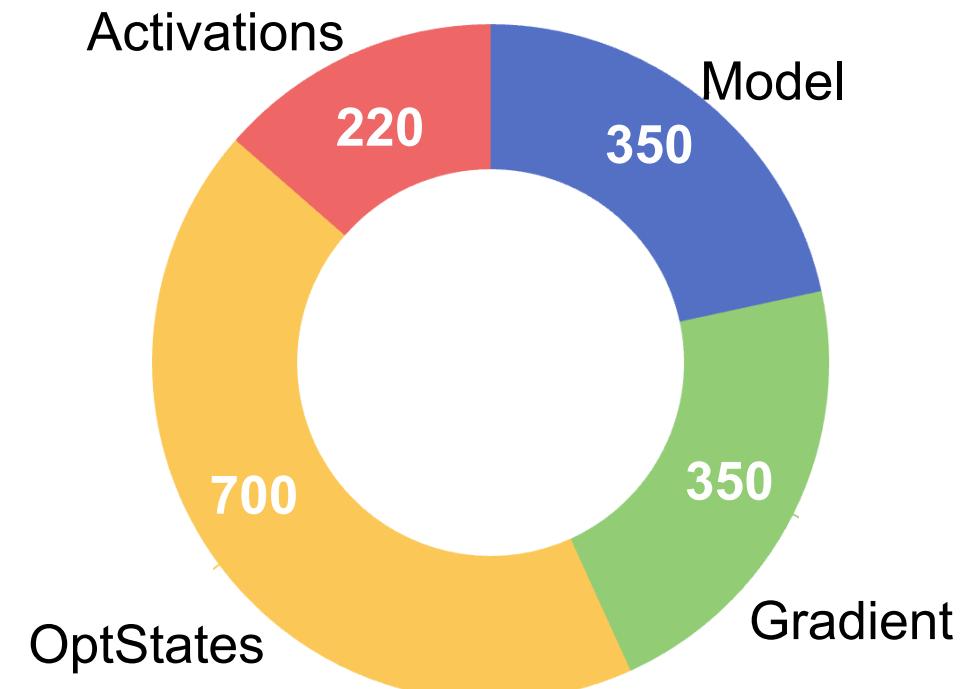
Need to store activations  $z_1, z_2, \dots, z_L$

- The size of activations depends on sequence length and batch size

## Minimum memory requirement: GPT-3

- Pretrain GPT-3 model (BF16) from scratch with **a single batch size** requires

- Parameters: 175B
- Model storage:  $175\text{B} * 2 \text{ Bytes} = 350 \text{ GB}$
- Gradient storage: 350 GB
- Optimizer states: 700 GB (using Adam)
- Activation storage: ~220 GB
- In total: **1620 GB**



## Minimum memory requirement: GPT-3

- Pretrain GPT-3 model (BF16) from scratch with **a single batch size** requires

- Parameters: 175B
- Model storage:  $175\text{B} * 2 \text{ Bytes} = 350 \text{ GB}$
- Gradient storage: 350 GB
- Optimizer states: 700 GB (using Adam)
- Activation storage: ~220 GB
- In total: **1620 GB**

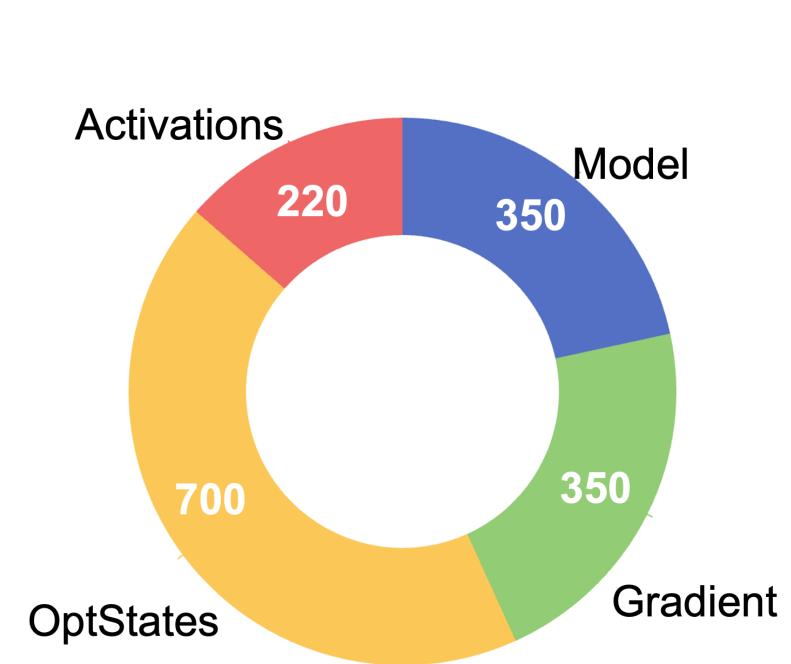


The minimum requirement is A100 \* 21  
**Very expensive!**

# Memory-efficient algorithm is in urgent need

Memory = **Model** + **Gradient** + **Optimizer states** + **Activations**

- With memory-efficient algorithms, we can
  - Train larger models** on limited computing resources
  - Use a larger training batch size to **improve throughput**
- The memory for model cannot be saved
- Is it possible to save memory for gradient, optimizer states, and activations?**





## PART 02

---

# Randomized Subspace Optimization

# Randomized Subspace Optimization (RSO)

- Consider the optimization problem for the LLM training:

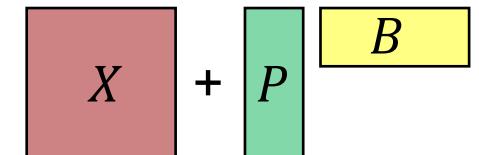
$$\min_{X \in \mathbb{R}^{m \times n}} f(X) := \mathbb{E}_{\xi}[F(X; \xi)].$$

- We solve small problems sequentially; optimizing over **one subspace** (spanned by  $P$ ) at a time:

Sample  $P$  for each iteration  $k$

$$\tilde{B}^k \approx \operatorname{argmin}_{B \in \mathbb{R}^{r \times n}} \left\{ f(X^k + P^k B) + \frac{1}{2\eta^k} \|B\|^2 \right\},$$

$$X^{k+1} = X^k + P^k \tilde{B}^k.$$



Here  $X \in R^{m \times n}$ ,  $P \in R^{m \times r}$  and  $B \in R^{r \times n}$ . Dimension  $r \ll \min\{m, n\}$ . The new variable  $B$  has much fewer parameters than  $X$

# Randomized Subspace Optimization (RSO)

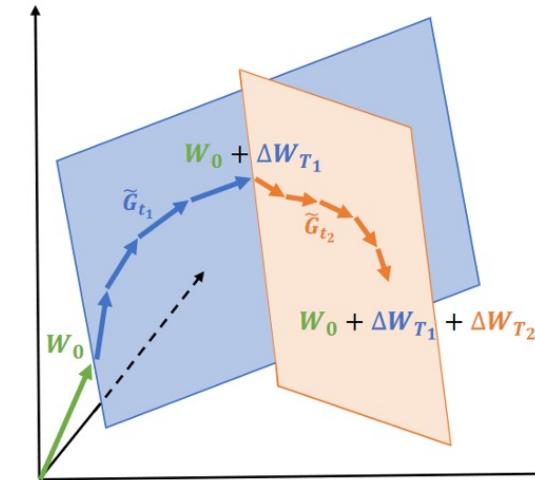
Can be solved with  
any optimizer

$$\tilde{\mathbf{B}}^k \approx \underset{\mathbf{B} \in \mathbb{R}^{r \times n}}{\operatorname{argmin}} \left\{ f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta^k} \|\mathbf{B}\|^2 \right\}, \quad (4a)$$

$$\mathbf{X}^{k+1} = \mathbf{X}^k + \mathbf{P}^k \tilde{\mathbf{B}}^k. \quad (4b)$$

Promote strong convexity

(Zhao et. al. ICML 2023)



---

## Algorithm 1: Randomized Subspace Optimization

---

**Input:** Initialization  $\mathbf{X}^0$

**Output:** Solution  $\mathbf{X}^K$

**for**  $k = 0, 1, \dots, K - 1$  **do**

Sample  $\mathbf{P}^k$  according to a given distribution;

Solve subproblem (4a) and obtain the approximate solution  $\tilde{\mathbf{B}}^k$  using a given optimizer such as Adam;

Update the weights by  $\mathbf{X}^{k+1} = \mathbf{X}^k + \mathbf{P}^k \tilde{\mathbf{B}}^k$ ;

---

## RSO is essentially a generalized method of coordinate minimization

- Recall the random **coordinate** minimization:

$$b_k^* = \arg \min_{b \in \mathbb{R}} \{f(x^k + b \cdot e_{i_k})\}, \text{ where } i_k \sim \mathcal{U}\{1, \dots, d\},$$

$$x^{k+1} = x^k + b_k^* \cdot e_{i_k}, \quad b \text{ is the increment in the } i\text{-th coordinate}$$

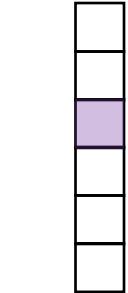
where the variable  $x$  is a vector in the above method

- Random subspace optimization can be regarded as a generalized method; not in the coordinate subspace but in the low-rank subspace

$$\tilde{\mathbf{B}}^k \approx \operatorname{argmin}_{\mathbf{B} \in \mathbb{R}^{r \times n}} \left\{ f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta^k} \|\mathbf{B}\|^2 \right\},$$

$$\mathbf{X}^{k+1} = \mathbf{X}^k + \mathbf{P}^k \tilde{\mathbf{B}}^k.$$

$\tilde{\mathbf{B}}$  is the increment in the  
subspace spanned by  $P^k$



$P$  B

Memory = **Model** + **Gradient** + **Optimizer states** + **Activations**

$$\tilde{\mathbf{B}}^k \approx \operatorname{argmin}_{\mathbf{B} \in \mathbb{R}^{r \times n}} \left\{ f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta_k} \|\mathbf{B}\|^2 \right\},$$

$$\mathbf{X}^{k+1} = \mathbf{X}^k + \mathbf{P}^k \tilde{\mathbf{B}}^k.$$

$$\text{Let } h(\mathbf{B}) = f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta_k} \|\mathbf{B}\|_F^2$$

We solve  $\tilde{\mathbf{B}}^k \approx \arg \min_{\mathbf{B}} \{h(\mathbf{B})\}$  using **Adam**

For  $t = 1, 2, \dots, \tau$

**Adam for solving  $h(\mathbf{B})$**

$$\mathbf{g}_t = \nabla_{\mathbf{B}} h(\mathbf{B}_t)$$

$$\mathbf{m}_t = (1 - \beta_1) \mathbf{m}_{t-1} + \beta_1 \mathbf{g}_t$$

$$\mathbf{v}_t = (1 - \beta_2) \mathbf{v}_{t-1} + \beta_2 \mathbf{g}_t \odot \mathbf{g}_t$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \frac{\gamma}{\sqrt{\mathbf{v}_t + \epsilon}} \odot \mathbf{m}_t$$

## RSO saves gradient and optimizer state

$$h(\mathbf{B}) = f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta_k} \|\mathbf{B}\|_F^2$$

Memory = **Model** + **Gradient** + **Optimizer states** + **Activations**

Adam for solving  $f(X)$

$$\mathbf{G}_t = \nabla F(\mathbf{X}_t; \boldsymbol{\xi}_t) \quad (\text{m, n})$$

$$\mathbf{M}_t = (1 - \beta_1)\mathbf{M}_{t-1} + \beta_1 \mathbf{G}_t \quad (\text{m, n})$$

$$\mathbf{V}_t = (1 - \beta_2)\mathbf{V}_{t-1} + \beta_2 \mathbf{G}_t \odot \mathbf{G}_t \quad (\text{m, n})$$

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \frac{\gamma}{\sqrt{\mathbf{V}_t} + \epsilon} \odot \mathbf{M}_t \quad (\text{m, n})$$

Adam for solving  $h(\mathbf{B})$

$$\mathbf{g}_t = \nabla_{\mathbf{B}} h(\mathbf{B}_t) \quad (\text{r, n})$$

$$\mathbf{m}_t = (1 - \beta_1)\mathbf{m}_{t-1} + \beta_1 \mathbf{g}_t \quad (\text{r, n})$$

$$\mathbf{v}_t = (1 - \beta_2)\mathbf{v}_{t-1} + \beta_2 \mathbf{g}_t \odot \mathbf{g}_t \quad (\text{r, n})$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \frac{\gamma}{\sqrt{\mathbf{v}_t} + \epsilon} \odot \mathbf{m}_t \quad (\text{r, n})$$

- Reduce gradient and optimizer states from  $mn$  parameters to  $rn$  parameters;
- Incurs **tiny** gradient and optimizer when  $r \ll m$

Memory = **Model** + **Gradient** + **Optimizer states** + **Activations**

Forward:  $Z = YX, y = L(Z)$ .

Backward:  $\frac{\partial y}{\partial X} = Y^\top \frac{\partial y}{\partial Z}$

Adam stores (**s by m**) activations

Forward:  $Z = Y(X + PB), y = L(Z)$ .

Backward:  $\frac{\partial y}{\partial B} = (YP)^\top \frac{\partial y}{\partial Z}$

RSO stores (**s by r**) activations

- Reduce activations from  $sm$  parameters to  $sr$  parameters;
- Incurs **tiny** activations when  $r \ll m$

# Memory comparison with existing optimizer

$s$ : #tokens in a sequence     $n$ : hidden dimension     $b$ : batch size     $r$ : rank

Algorithm	Memory		
	Optimizer States	Activations	Gradients
<b>RSO</b>	$24nr$	$8bsn + 4bsr + 2bs^2$	$12nr$
GaLore	$24nr$	$15bsn + 2bs^2$	$12n^2$
LoRA	$48nr$	$15bsn + 2bs^2$	$24nr$
Adam	$24n^2$	$15bsn + 2bs^2$	$12n^2$

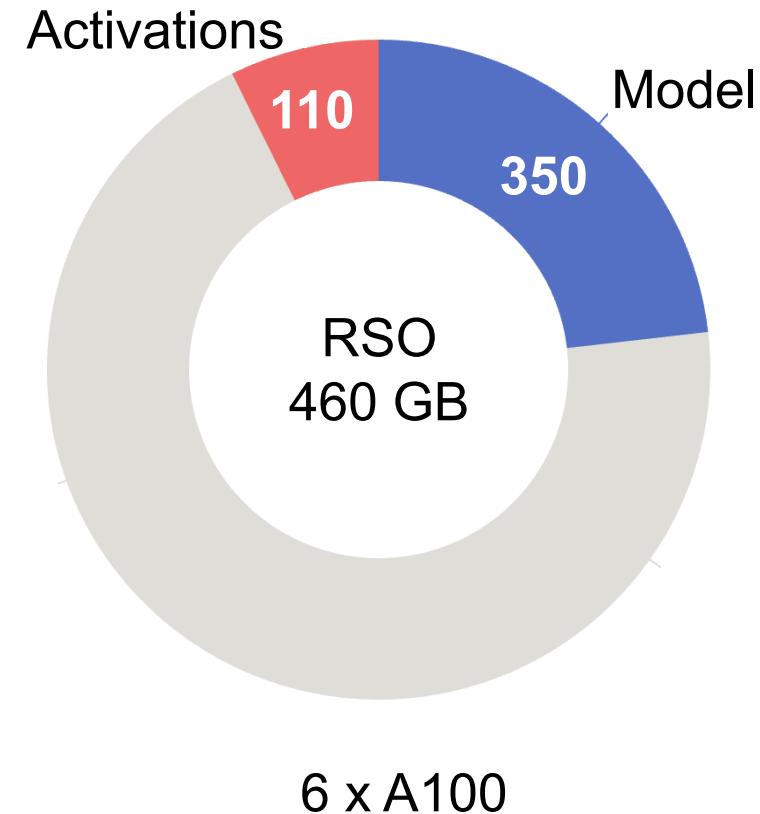
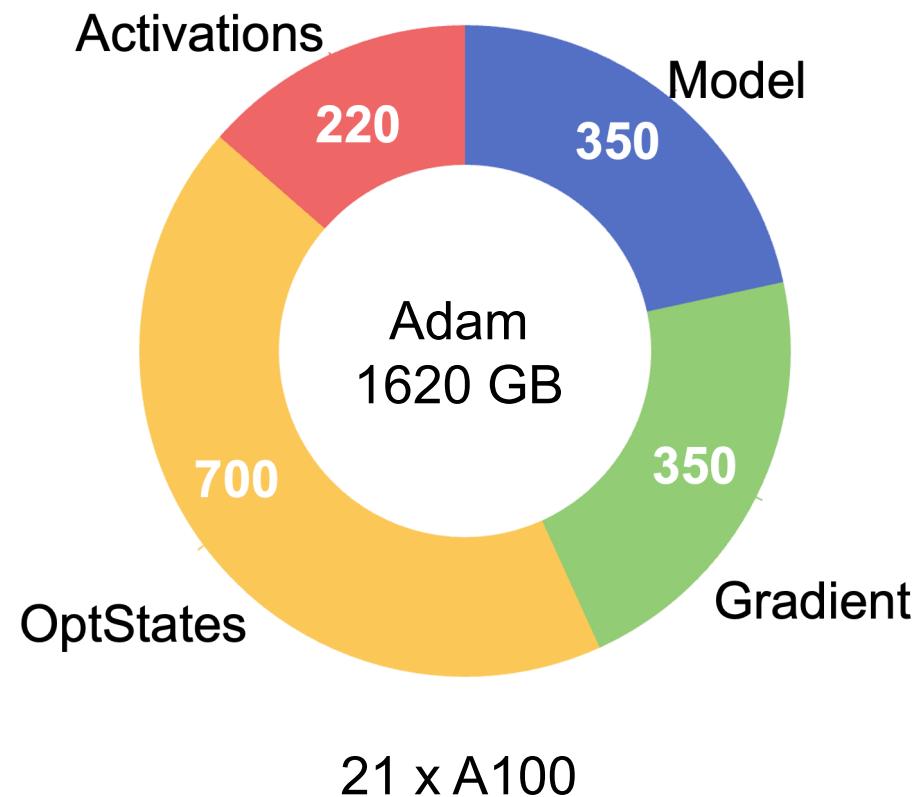
$$\tilde{\mathbf{B}}^k \approx \underset{\mathbf{B} \in \mathbb{R}^{r \times n}}{\operatorname{argmin}} \left\{ f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta^k} \|\mathbf{B}\|^2 \right\},$$

$$\mathbf{X}^{k+1} = \mathbf{X}^k + \mathbf{P}^k \tilde{\mathbf{B}}^k.$$

- Saves gradient, optimizer states, and activations simultaneously
- RSO is the most memory-efficient optimizer to our best knowledge

# Memory saving visualization

Suppose  $r \ll m$ , and  $s^2$  is not the dominant term



## PART 03

---

### Convergence Guarantees

# Assumptions

---

**Definition 5.1** (Expected  $\epsilon$ -inexact solution). A solution  $\tilde{\mathbf{B}}^k$  is said to be an expected  $\epsilon$ -inexact solution if it satisfies:

$$\mathbb{E}[g^k(\tilde{\mathbf{B}}^k)] - g^k(\mathbf{B}_\star^k) \leq \epsilon,$$

where  $g^k(\mathbf{B}) := f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta^k} \|\mathbf{B}\|^2$ , and  $\mathbf{B}_\star^k$  is the optimal solution define as  $\mathbf{B}_\star^k := \arg \min_{\mathbf{B}} g^k(\mathbf{B})$ .

**Assumption 5.2.** The objective function  $f(\mathbf{W})$  is  $L$ -smooth, i.e., it holds for any  $\mathbf{W}^1$  and  $\mathbf{W}^2$  that

$$\|\nabla f(\mathbf{W}^1) - \nabla f(\mathbf{W}^2)\| \leq L \|\mathbf{W}^1 - \mathbf{W}^2\|, \quad (\text{smoothness})$$

where  $\|\mathbf{W}\| := \sqrt{\sum_{\ell=1}^L \|W_\ell\|_F^2}$  for any  $\mathbf{W} = \{W_\ell\}_{\ell=1}^L$ .

**Assumption 5.3.** The random matrix  $\mathbf{P} = \{P_\ell\}_{\ell=1}^L$  is sampled from a distribution such that  $P_\ell^\top P_\ell = (m_\ell/r_\ell) I_{r_\ell}$  and  $\mathbb{E}[P_\ell P_\ell^\top] = I_{m_\ell}$  for each  $\ell$ . (random subspace)

## Theorem

Let each subproblem in RSO iteration be solved starting from the initial point  $\mathbf{B}^0 = \mathbf{0}$  to an expected  $\epsilon$ -inexact solution  $\tilde{\mathbf{B}}^k$  with a suitable choice of  $\eta^k$ . The sequence  $\{\mathbf{X}^k\}$  generated by the RSO method satisfies the following bound:

subproblem optimizer

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\mathbf{X}^k)\|^2 \leq \boxed{\frac{18\hat{L}\Delta_0}{K}} + \boxed{18\hat{L}\epsilon},$$

RSO framework

where  $\Delta_0 := f(\mathbf{X}^0) - f^*$  and  $\hat{L} := \max_\ell \{m_\ell/r_\ell\} L$ .

- RSO converges to  $\epsilon$ -accurate solution at rate  $1/K$
- The accuracy is determined by the subproblem optimizer

# RSO convergence theorem with concrete subproblem optimizer



$$h(\mathbf{B}) = f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta_k} \|\mathbf{B}\|_F^2$$

Subproblem Solver	Subproblem Complexity	Total Complexity
Stochastic ZO Method	$O((\sum_{\ell=1}^{\mathcal{L}} n_\ell r_\ell)^2 \epsilon^{-2})$	$O((\sum_{\ell=1}^{\mathcal{L}} n_\ell r_\ell)^2 \epsilon^{-3})$
GD	$O(\log \epsilon^{-1})$	$\tilde{O}(\epsilon^{-1})$
Accelerated GD	$O(\log \epsilon^{-1})$	$\tilde{O}(\epsilon^{-1})$
SGD	$O(\epsilon^{-1})$	$O(\epsilon^{-2})$
Momentum SGD	$O(\epsilon^{-1})$	$O(\epsilon^{-2})$
Adam-family	$O(\epsilon^{-1})$	$O(\epsilon^{-2})$
Newton's method	$O(\log(\log \epsilon^{-1}))$	$\tilde{O}(\epsilon^{-1})$
Stochastic Quasi-Newton method	$O(\epsilon^{-1})$	$O(\epsilon^{-2})$

It is observed that RSO with Adam to solve subproblem has sample complexity  $O(\epsilon^{-2})$ , which is on the same order as vanilla Adam without subspace projection.

## PART 04

---

### Experiments

# RSO performance: Memory comparison in Pre-train



Algorithm	60M	130M	350M	1B
Adam	34.06 (0.22G)	25.08 (0.50G)	18.80 (1.37G)	15.56 (4.99G)
GaLore	34.88 (0.14G)	25.36 (0.27G)	18.95 (0.49G)	<b>15.64</b> (1.46G)
LoRA	34.99 (0.16G)	33.92 (0.35G)	25.58 (0.69G)	19.21 (2.27G)
ReLoRA	37.04 (0.16G)	29.37 (0.35G)	29.08 (0.69G)	18.33 (2.27G)
<b>RSO</b>	<b>34.55</b> (0.14G)	<b>25.34</b> (0.27G)	<b>18.86</b> (0.49G)	15.86 (1.46G)
$r/d_{model}$	128 / 256	256 / 768	256 / 1024	512 / 2048
Training Tokens (B)	1.1	2.2	6.4	13.1

Table: Comparison of validation perplexity and estimated memory usage for optimizer states across different algorithms during the pre-training of LLaMA models of various sizes on the C4 dataset.

- Saves significant optimizer states memory; performance degradation is within 2%
- Run for the same number of iterations/samples; performs well for small models (<350M, ~1%)
- Performance deteriorates for relatively large models ( $\geq 1B$ , 2%)

# RSO performance: Memory comparison in Pre-train

## Practical **Total Memory** Profiling (not just optimizer)

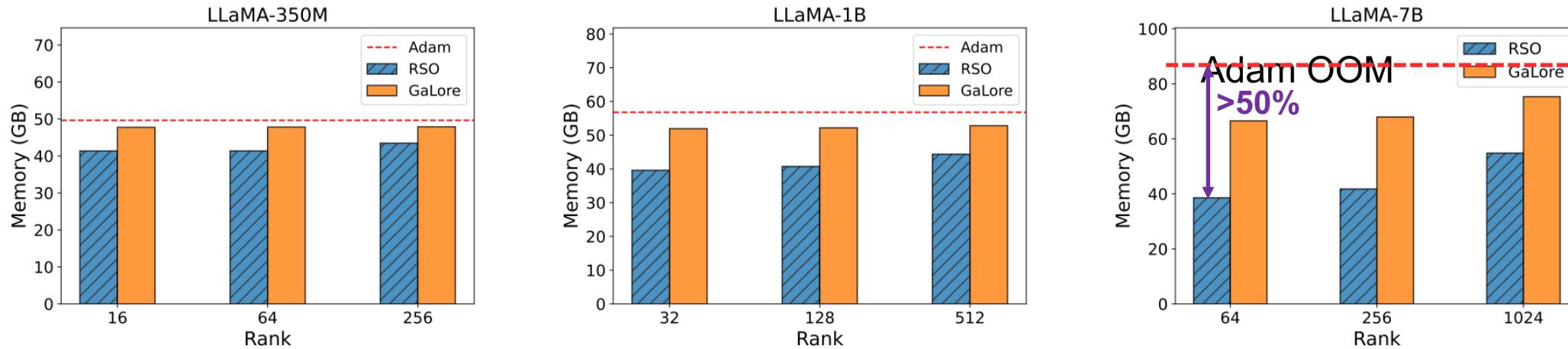


Figure 2: Comparison of peak memory usage (in GB) per device for RSO and GaLore during LLaMA training with varying ranks. All hyperparameters, except rank, are consistent with [Zhao et al., 2024a]. Adam's memory usage is reported for LLaMA-350M and LLaMA-1B but excluded for LLaMA-7B due to an out-of-memory (OOM) error.

- RSO saves gradient, optimizer, and activations, and hence it is more memory-efficient
- As rank increases, memory saving in RSO gets decreased

# RSO performance: Computation comparison in Pre-train



Method	LLaMA-1B (Seconds)			LLaMA-7B (Seconds)		
	Seq 64	Seq 128	Seq 256	Seq 64	Seq 128	Seq 256
RSO	0.94	1.70	3.29	2.40	2.94	4.60
GaLore	1.12	1.84	3.35	7.86	8.26	9.12
Adam	1.11	1.81	3.32	7.84	8.23	OOM

Table 5: Comparison of iteration time (in seconds) for different methods in LLaMA training across various sequence lengths. All hyperparameters, except sequence length, follow [Zhao et al., 2024a]. LLaMA-1B runs on 4× A800 GPUs, while LLaMA-7B uses 8× A800 GPUs. SVD decomposition time in GaLore is excluded. Additionally, for LLaMA-7B with a sequence length of 256, the Adam optimizer encounters an out-of-memory (OOM) error.

- RSO achieves much faster training speed due to efficient solving for small subproblems

# RSO performance: end-to-end pre-train

---

Run for the same number of iterations

Method	Memory (GB)	Training Time (h)	Perplexity
Adam	78.92	216	15.43
GaLore	75.33	134	15.59
<b>RSO</b>	<b>54.81</b>	<b>64</b>	15.99

Table 6: Comparison of various pre-training methods for the LLaMA-7B model on the C4 dataset. Perplexity is reported at 50K steps. The training is conducted on  $8 \times$  A800 GPUs. The actual memory cost per device and the total training time are also reported. RSO and GaLore are configured with a batch size of 16, while Adam uses a batch size of 8.

- With  $\sim 3.6\%$  performance degradation, RSO achieves **3.7x** speedup and **30%** memory saving
- It is conjectured that, given the same runtime, RSO outperforms Adam; however, we have not conducted experiments to verify this.

## Memory = Model + Gradient + Optimizer states + Activations

- Random subspace optimization (RSO) saves **gradient, optimizer, and activations**

$$\tilde{\mathbf{B}}^k \approx \underset{\mathbf{B} \in \mathbb{R}^{r \times n}}{\operatorname{argmin}} \left\{ f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta^k} \|\mathbf{B}\|^2 \right\},$$

$$\mathbf{X}^{k+1} = \mathbf{X}^k + \mathbf{P}^k \tilde{\mathbf{B}}^k.$$

- RSO has strong convergence guarantees
- Open questions:** How to improve RSO performance, especially for large LLMs?



paper



Github Code

# Ongoing work: Explore better subspace than random



$$\tilde{\mathbf{B}}^k \approx \operatorname{argmin}_{\mathbf{B} \in \mathbb{R}^{r \times n}} \left\{ f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta^k} \|\mathbf{B}\|^2 \right\},$$
$$\mathbf{X}^{k+1} = \mathbf{X}^k + \mathbf{P}^k \tilde{\mathbf{B}}^k.$$

- Random subspace  $P$  does not capture any problem structures
- Recall in coordinate gradient descent, there exists various greedy coordinate sampling strategies such as the *Gauss-Southwell* selection rules
- How to find better subspace  $P$  to enhance the performance of subspace methods?
- What structure shall the subspace  $P$  capture?

- Y. He, P. Li, Y. Hu, C. Chen, and **K. Yuan**, *Subspace Optimization for Large Language Models with Convergence Guarantees*, ICML 2025
- Y. Chen, Y. Zhang, Y. Liu, **K. Yuan**, Z. Wen, *A Memory Efficient Randomized Subspace Optimization Method for Training Large Language Models*, ICML 2025
- Y. Chen, Y. Zhang, L. Cao, **K. Yuan**, Z. Wen, *Enhancing Zeroth-Order Fine-tuning for Language Models with Low-Rank Structures*, ICLR 2025

# Thank you!