



Optimization for Deep Learning

Memory-Efficient Optimizers

Kun Yuan Peking University

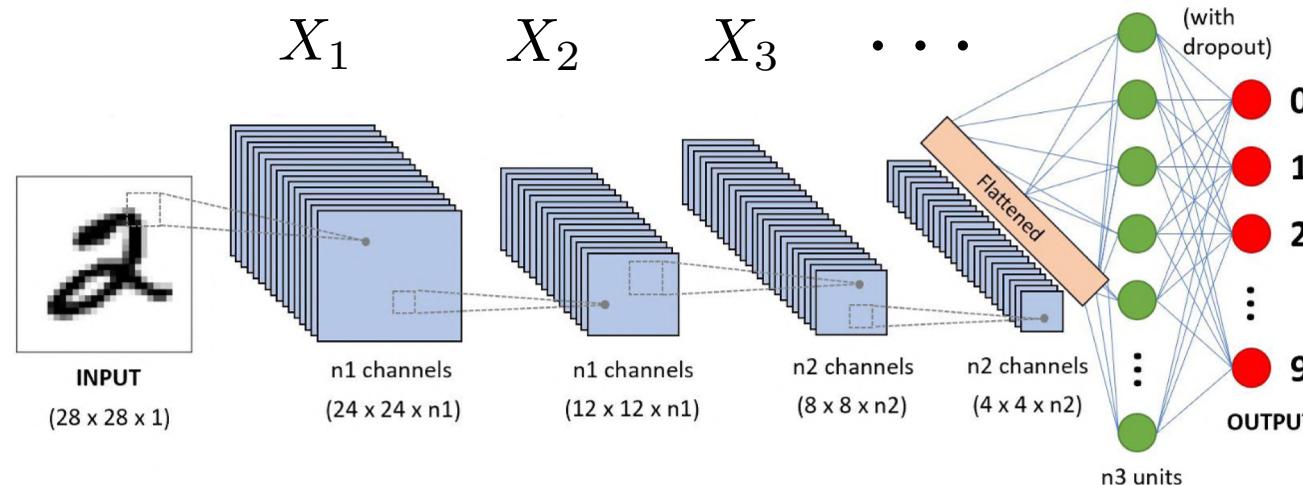


PART 01

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 = **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

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

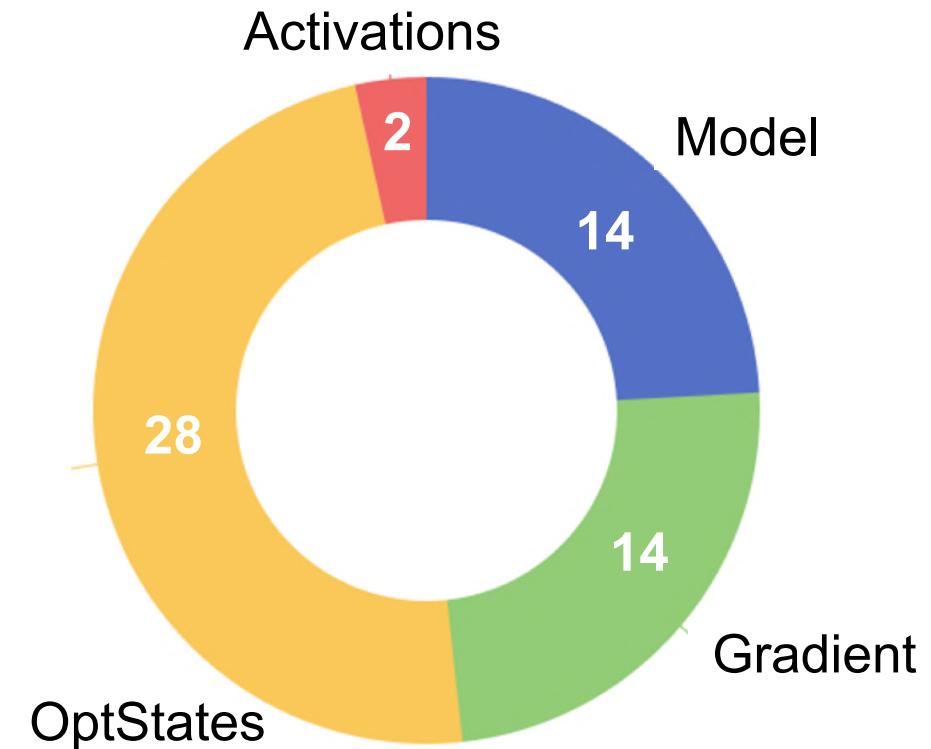
Need to store activations z_1, z_2, \dots, z_L

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

Minimum memory requirement: LLaMA 7B

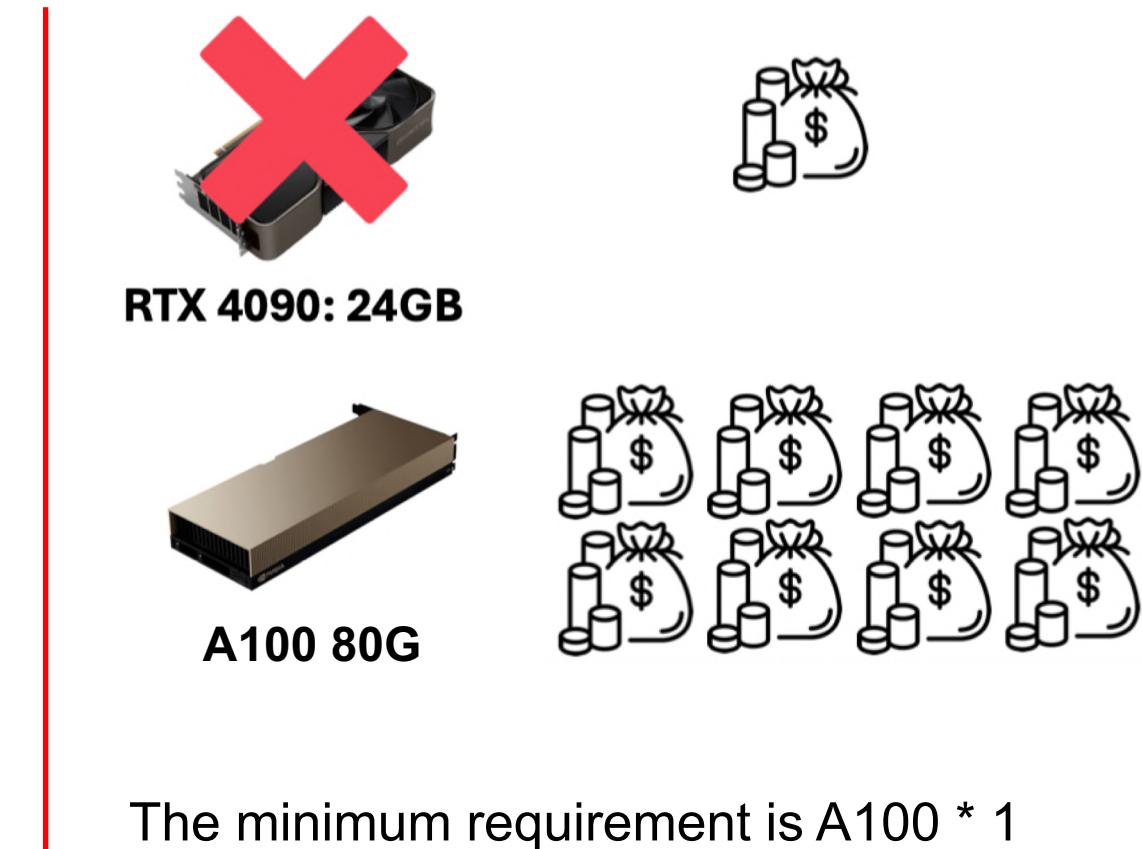
- Pretrain LLaMA 7B model (BF16) from scratch with **a single batch size** requires

- Parameters: 7B
- Model storage: $7B * 2 \text{ Bytes} = 14 \text{ GB}$
- Gradient storage: 14 GB
- Optimizer states: 28 GB (using Adam)
- Activation storage: 2 GB [Zhao et. al., 2024]
- In total: **58 GB**



Minimum memory requirement: LLaMA 7B

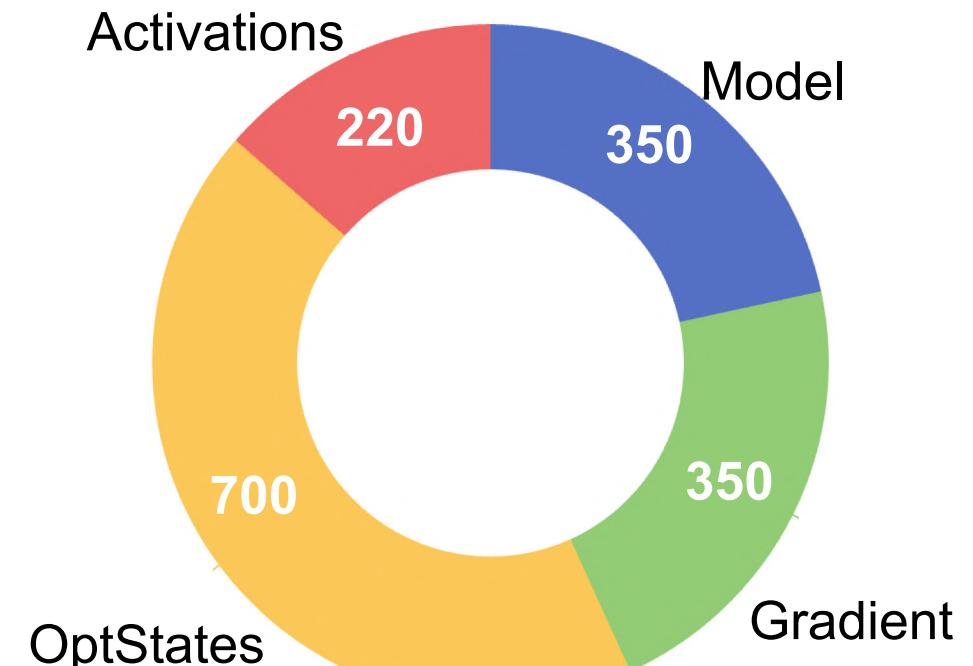
- Pretrain LLaMA 7B model (BF16) from scratch with **a single batch size** requires
 - Parameters: 7B
 - Model storage: $7B * 2 \text{ Bytes} = 14 \text{ GB}$
 - Gradient storage: 14 GB
 - Optimizer states: 28 GB (using Adam)
 - Activation storage: 2 GB [Zhao et. al., 2024]
 - In total: **58 GB**



Minimum memory requirement: GPT-3

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

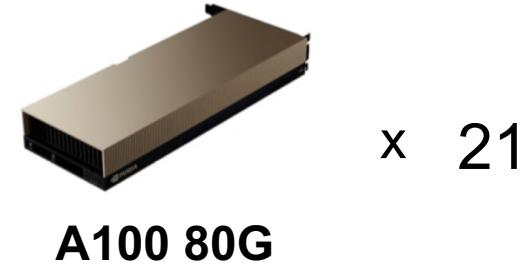
- Parameters: 175B
- Model storage: $175\text{B} * 2 \text{ Bytes} = 350 \text{ GB}$
- Gradient storage: 350 GB
- Optimizer states: 700 GB (using Adam)
- Activation storage: ~220 GB
- In total: **1620 GB**



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 optimizers are in urgent need!

With memory-efficient algorithms, we can

- **Train larger models** on limited computing resources
- Use a larger training batch size to **improve throughput**



PART 02

GaLore: Gradient Low-Rank Projection

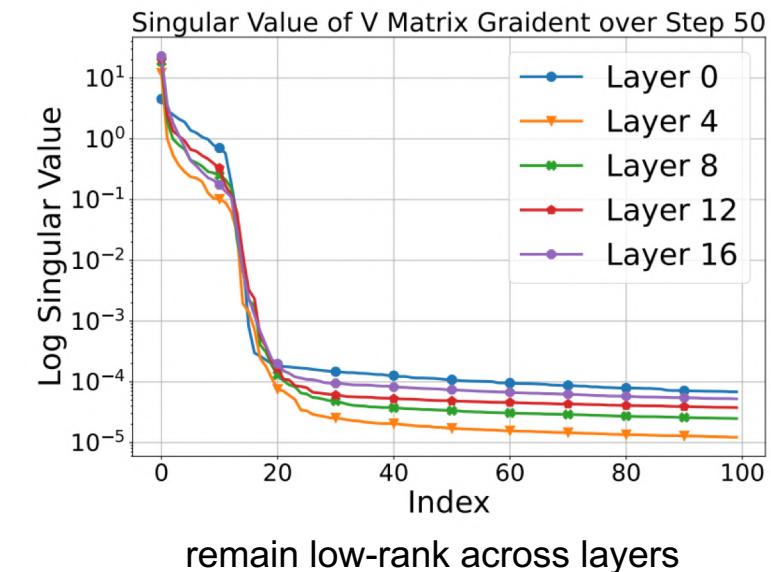
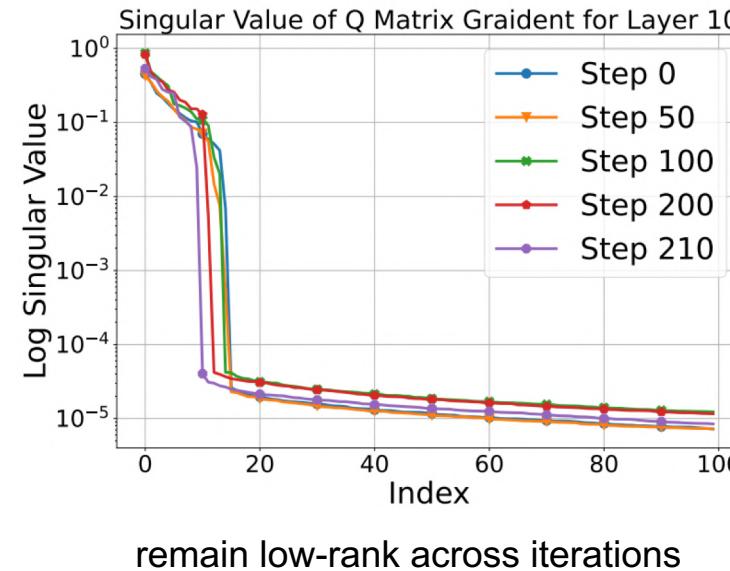
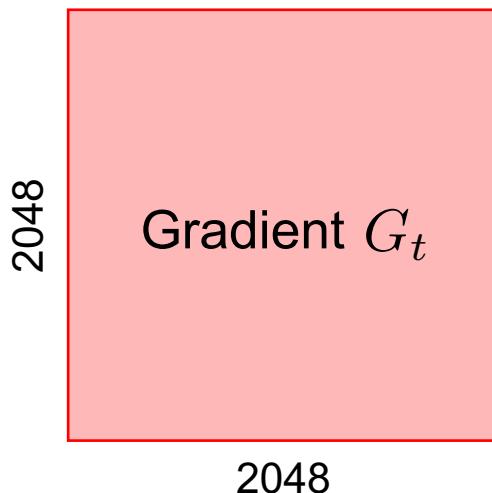
GaLore: Memory-Efficient LLM Training by Gradient Low-Rank Projection

Jiawei Zhao¹ Zhenyu Zhang³ Beidi Chen^{2,4} Zhangyang Wang³ Anima Anandkumar^{*1} Yuandong Tian^{*2}

- GaLore is a novel approach for reducing memory consumption from optimizer states
- Memory-efficient without severe performance degradation

GaLore: Gradient Low-Rank Projection

- Observation: gradient in LLMs becomes **low-rank** during training



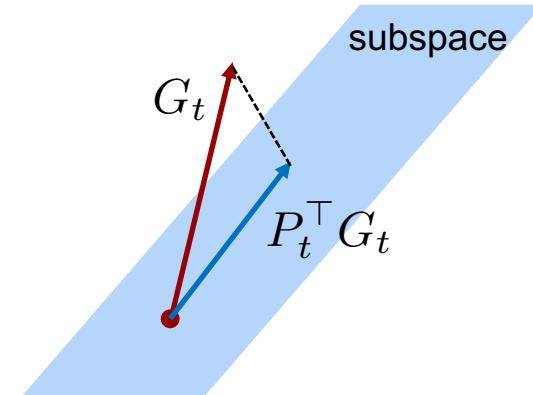
- Given a gradient matrix with dimensions 2048 by 2048, around **top 10** eigenvalues dominate
- How to utilize the low-rank structure in gradients?

GaLore: Gradient Low-Rank Projection

- Main idea: Projecting gradient onto the low-rank subspace
- Given a gradient $G_t \in \mathbb{R}^{m \times n}$ and a projection $P_t \in \mathbb{R}^{m \times r}$, we project Gradient onto low-rank subspace

$$G_t \in \mathbb{R}^{m \times n} \quad P_t^\top G_t \quad g_t = P_t^\top G_t \in \mathbb{R}^{r \times n}$$

Low-rank projection



- Since $r \ll m$, the low-rank gradient g_t has much smaller parameters than G_t

GaLore: Gradient Low-Rank Projection

- Low-rank optimizer states:

$$\mathbf{g}_t = \mathbf{P}_t^\top \mathbf{G}_t$$

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

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

$$\boldsymbol{\delta}_t = \frac{\gamma}{\sqrt{\mathbf{v}_t} + \epsilon} \odot \mathbf{m}_t$$

▷ dims r x n



Simplified as

$$\mathbf{X}_{t+1} = \mathbf{X}_t + \mathbf{P}_t \boldsymbol{\rho}(\mathbf{P}_t^\top \mathbf{G}_t)$$

- Parameter updates:

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \mathbf{P}_t \boldsymbol{\delta}_t$$

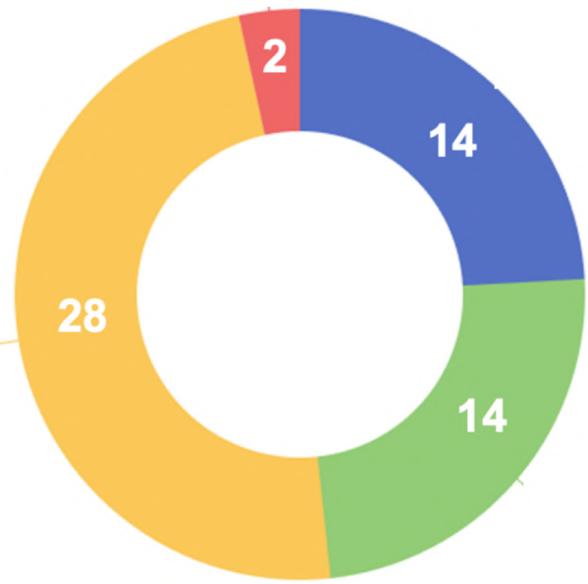
▷ dims m x n

- Memory cost: Model \mathbf{X} , Gradient \mathbf{G} , Projection \mathbf{P} , OptStates \mathbf{m}, \mathbf{v} and activations

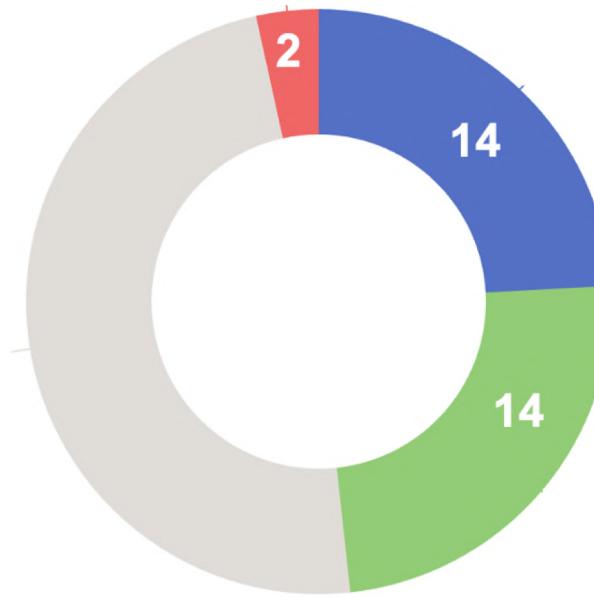
trivial memory cost

GaLore: Gradient Low-Rank Projection

LLaMA 7B



Adam



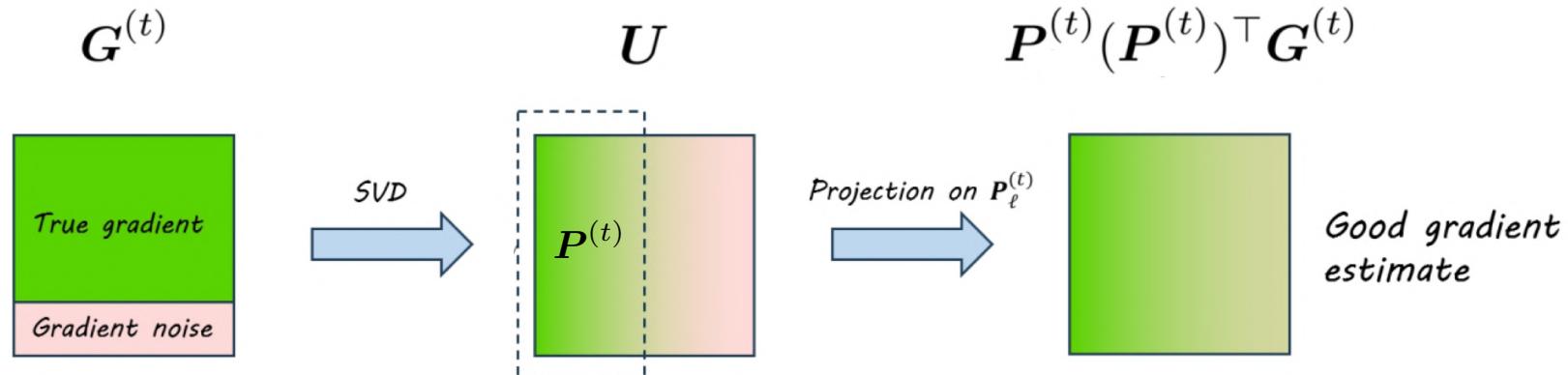
Adam + GaLore

GaLore: Gradient Low-Rank Projection

- Recall the GaLore update: $X_{t+1} = X_t + P_t \rho(P_t^\top G_t)$
- How to find the projection matrix? **SVD decomposition!**

$$G_t = U\Sigma V^\top \longrightarrow P_t = \boxed{U[:, :r]} \in \mathbb{R}^{m \times r}$$

Select the first r columns



[Subspace Optimization for Large Language Models with Convergence Guarantees]

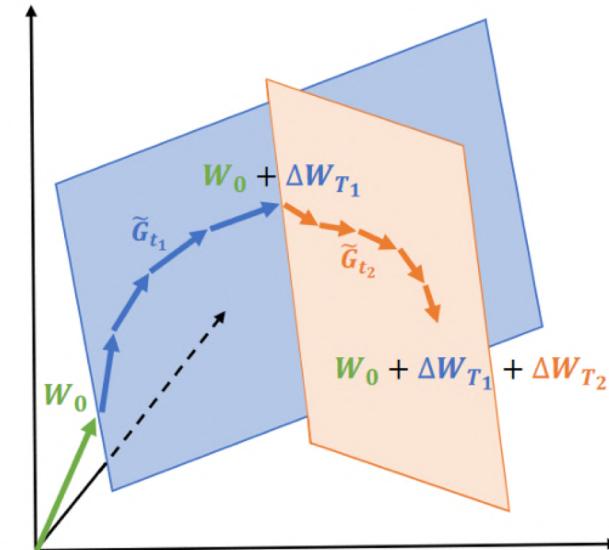
GaLore: Gradient Low-Rank Projection

- It is computationally expensive to perform SVD in each iteration
- Lazy SVD**: perform SVD every τ iterations

(The complete GaLore algorithm)

$$\begin{cases} \mathbf{P}_t \leftarrow \text{SVD}(\mathbf{G}_t) & \text{if } t \bmod \tau = 0 \\ \mathbf{P}_t \leftarrow \mathbf{P}_{t-1} & \text{otherwise} \end{cases}$$

$$\mathbf{X}_{t+1} = \mathbf{X}_t + \mathbf{P}_t \boldsymbol{\rho}(\mathbf{P}_t^\top \mathbf{G}_t)$$



[GaLore: Memory-Efficient LLM Training by Gradient Low-Rank Projection]

GaLore: Gradient Low-Rank Projection



Pretraining LLaMA on C4 dataset

	60M	130M	350M	1B
Full-Rank	34.06 (0.36G)	25.08 (0.76G)	18.80 (2.06G)	15.56 (7.80G)
GaLore	34.88 (0.24G)	25.36 (0.52G)	18.95 (1.22G)	15.64 (4.38G)
Low-Rank	78.18 (0.26G)	45.51 (0.54G)	37.41 (1.08G)	142.53 (3.57G)
LoRA	34.99 (0.36G)	33.92 (0.80G)	25.58 (1.76G)	19.21 (6.17G)
ReLoRA	37.04 (0.36G)	29.37 (0.80G)	29.08 (1.76G)	18.33 (6.17G)
r/d_{model}	128 / 256	256 / 768	256 / 1024	512 / 2048
Training Tokens	1.1B	2.2B	6.4B	13.1B

GaLore looks great

But does GaLore provably converge to the desired solution?

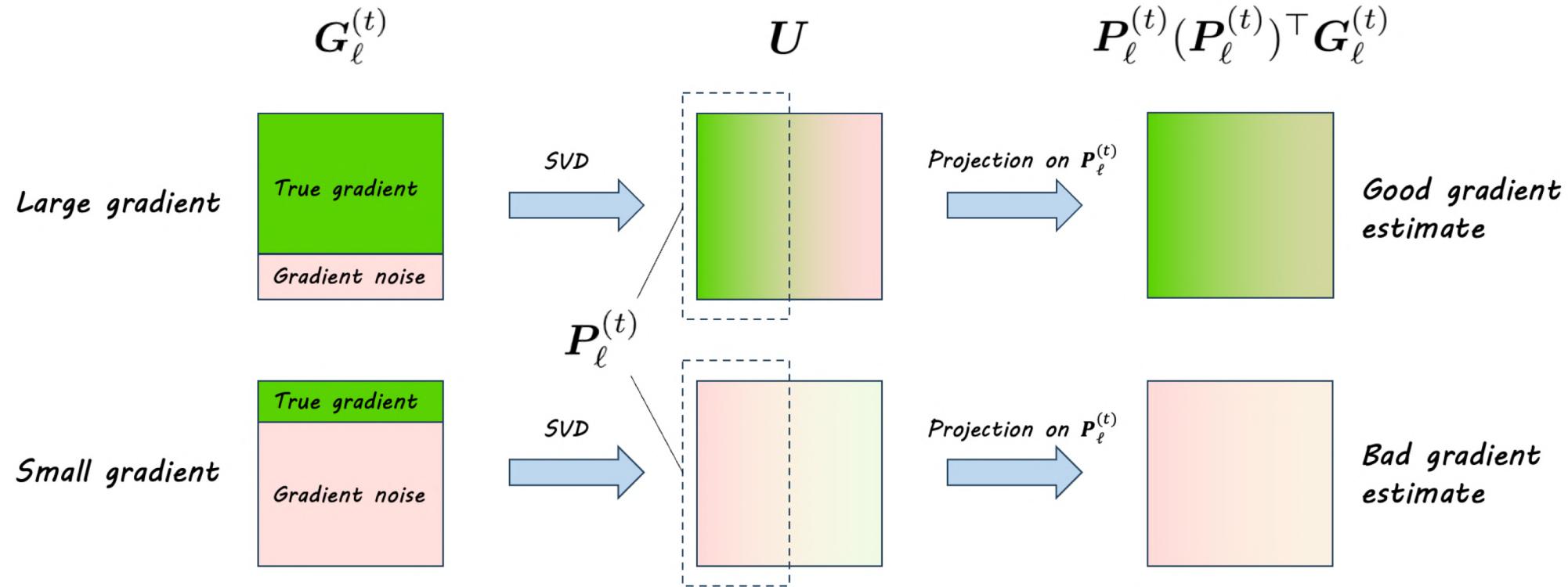
Not Necessarily True!



PART 03

GaLore is NOT GUARANTEED to converge

Intuition behind GaLore's non-convergence

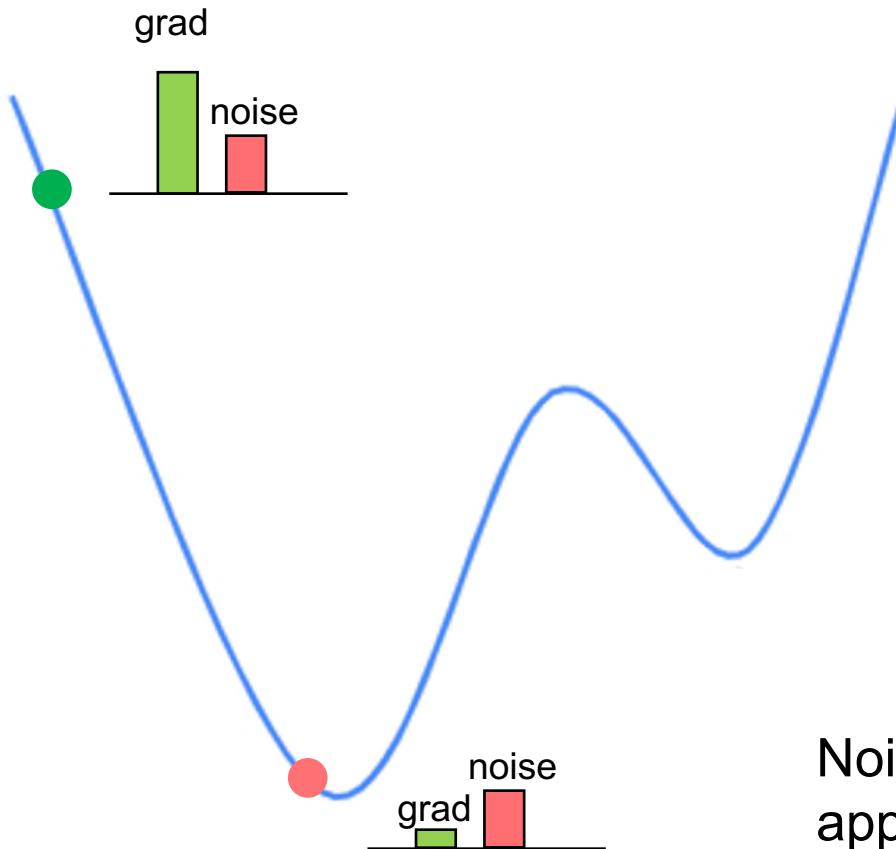


When gradient noise dominates the stochastic gradient, SVD captures **noise-dominated** subspace!

All gradient information is lost !

Is noise-dominant scenario common? Yes!

Gradient dominates
during the initial stages



Noise dominates when
approaching the local minimum

SVD 投影的理论缺陷：反例构造

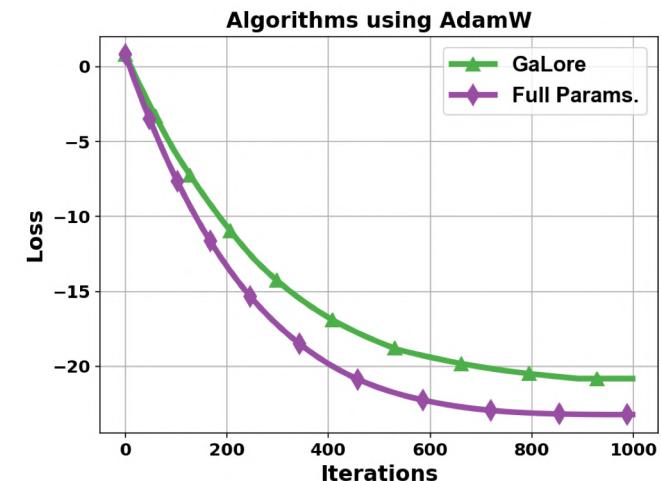
Counter-Example. We consider the following quadratic problem with gradient noise:

$$f(\mathbf{X}) = \frac{1}{2} \|\mathbf{AX}\|_F^2 + \langle \mathbf{B}, \mathbf{X} \rangle_F, \quad \nabla F(\mathbf{X}; \xi) = \nabla f(\mathbf{X}) + \xi \sigma \mathbf{C}, \quad (1)$$

where $\mathbf{A} = (\mathbf{I}_{n-r} \quad 0) \in \mathbb{R}^{(n-r) \times n}$, $\mathbf{B} = \begin{pmatrix} \mathbf{D} & 0 \\ 0 & 0 \end{pmatrix} \in \mathbb{R}^{n \times n}$ with $\mathbf{D} \in \mathbb{R}^{(n-r) \times (n-r)}$ generated randomly, $\mathbf{C} = \begin{pmatrix} 0 & 0 \\ 0 & \mathbf{I}_r \end{pmatrix} \in \mathbb{R}^{n \times n}$, ξ is a random variable uniformly sampled from $\{1, -1\}$ per iteration, and σ is used to control the gradient noise.

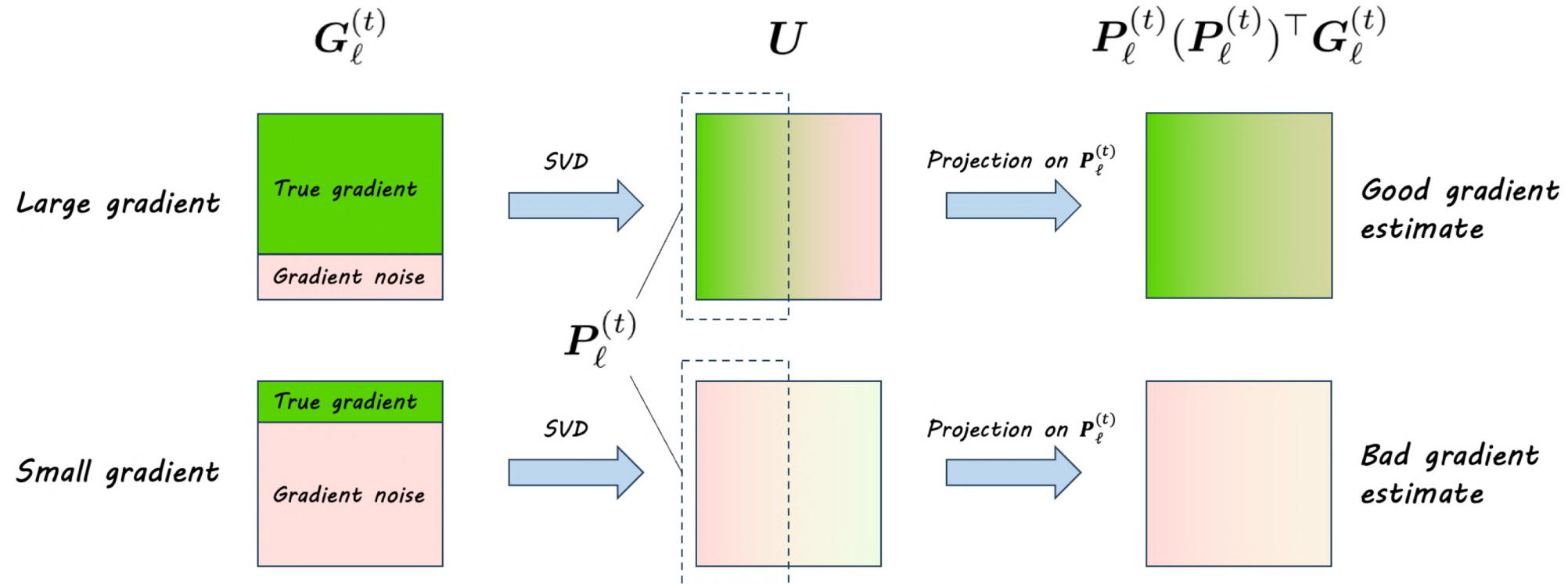
Theorem (Non-convergence of GaLore): There exists an objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying Assumptions 1 and 2, a stochastic gradient oracle (F, D) satisfying Assumption 3, an initial point $\mathbf{x}^{(0)} \in \mathbb{R}^d$, and a constant $\epsilon_0 > 0$ such that for any rank $r_\ell < \min\{m_\ell, n_\ell\}$, subspace changing frequency τ , any optimizer ρ that inputs a subspace gradient of shape $r_\ell \times n_\ell$ and outputs a subspace update direction of the same shape, and for any $t > 0$, it holds that

$$\|\nabla f(\mathbf{x}^{(t)})\|_2^2 \geq \epsilon_0.$$



GaLore does **NOT** converge to desired solutions

Under what conditions can GaLore converge?



GaLore can converge if we can avoid the noise-dominant scenarios

Condition I: noise-free

- Consider GaLore with deterministic gradient: $G_\ell^{(t)} = \nabla_\ell f(\mathbf{x}^{(t)})$

Theorem 2 (Convergence rate of deterministic GaLore). *Under Assumptions 1-2, if the number of iterations $T \geq 64/(3\underline{\delta})$ and we choose*

$$\beta_1 = 1, \quad \tau = \left\lceil \frac{64}{3\underline{\delta}\beta_1} \right\rceil, \quad \text{and} \quad \eta = \left(4L + \sqrt{\frac{80L^2}{3\underline{\delta}\beta_1^2}} + \sqrt{\frac{80\tau^2L^2}{3\underline{\delta}}} + \sqrt{\frac{16\tau L^2}{3\beta_1}} \right)^{-1},$$

GaLore using deterministic gradients and MSGD with MP converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O}\left(\frac{L\Delta}{\underline{\delta}^{5/2}T}\right),$$

where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$ and $\underline{\delta} := \min_\ell \frac{r_\ell}{\min\{m_\ell, n_\ell\}}$.

- Noise-free GaLore **converges** at rate $\mathcal{O}(1/T)$.

Condition II: Large batch-size

- Consider GaLore with large-batch stochastic gradient: $G_\ell^{(t)} = \frac{1}{\mathcal{B}} \sum_{b=1}^{\mathcal{B}} \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t,b)})$
- Batch-size \mathcal{B} increases with iteration T , e.g., $\mathcal{B} = \mathcal{O}(\sqrt{T})$

Theorem 3 (Convergence rate of large-batch GaLore). *Under Assumptions 1-3, if $T \geq 2 + 128/(3\underline{\delta}) + (128\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$ and we choose $\tau = \lceil 64/(3\underline{\delta}\beta_1) \rceil$, $\mathcal{B} = \lceil 1/(\underline{\delta}\beta_1) \rceil$,*

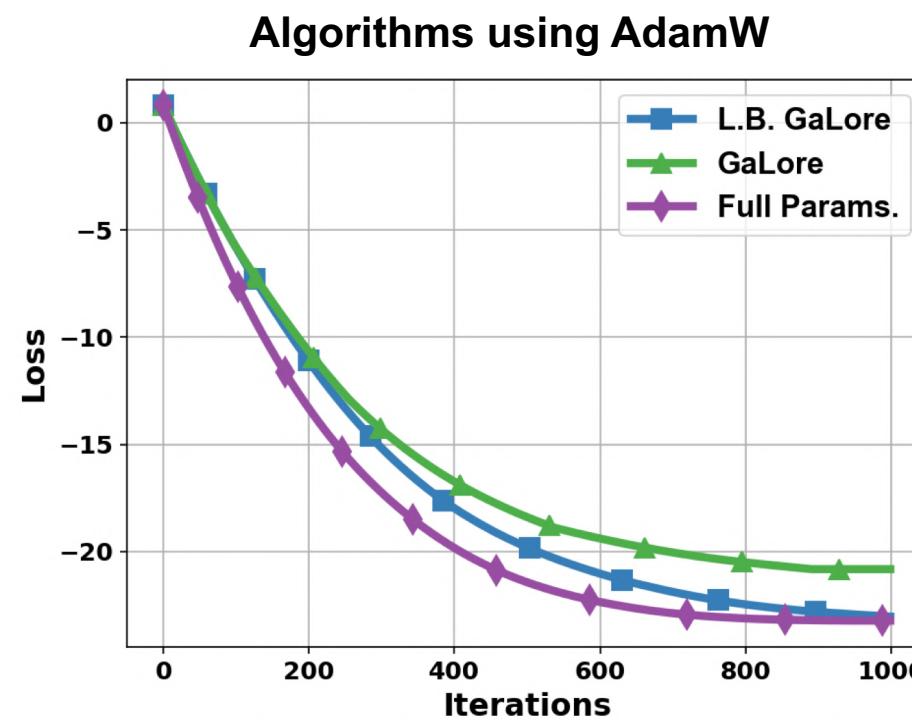
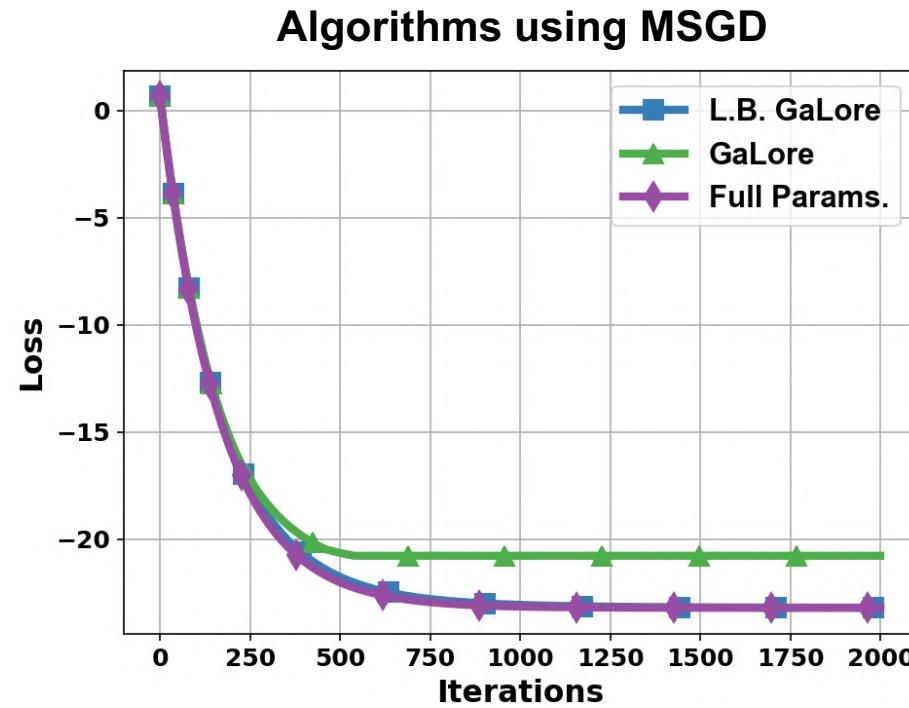
$$\beta_1 = \left(1 + \sqrt{\frac{\underline{\delta}^{3/2}\sigma^2 T}{L\Delta}} \right)^{-1}, \quad \text{and} \quad \eta = \left(4L + \sqrt{\frac{80L^2}{3\underline{\delta}\beta_1^2}} + \sqrt{\frac{40\tau^2 L^2}{\underline{\delta}}} + \sqrt{\frac{32\tau L^2}{3\beta_1}} \right)^{-1},$$

GaLore using large-batch gradients and MSGD with MP converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O} \left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}} \right),$$

where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$ and $\underline{\delta} := \min_\ell \frac{r_\ell}{\min\{m_\ell, n_\ell\}}$.

Numerical simulations



However, neither noise-free nor large-batch is practical for LLMs settings

PART 04

GoLore: Gradient RANDOM Low-rank projection

- In LLM settings, gradient noise exists and batch-size does not increase with iterations
- The root reason that GaLore has convergence issues is the SVD-incurred subspace
- Random projection can possibly capture gradient information when noise dominates

(Stiefel manifold) $\text{St}_{m,r} = \{P \in \mathbb{R}^{m \times r} \mid P^\top P = I_r\}.$

Proposition 1 (Chikuse (2012), Theorem 2.2.1). *A random matrix X uniformly distributed on $\text{St}_{m,r}$ is expressed as $X = Z(Z^\top Z)^{-1/2}$, where the elements of an $m \times r$ random matrix Z are independent and identically distributed as normal $\mathcal{N}(0, 1)$.*

Following Prop. 1, we can sample random projections from Stiefel manifold

GoLore:基于随机均匀投影的子空间训练算法

- Gradient random Low-rank projection (GoLore) 将梯度投影至随机

$$P_t \sim \mathcal{U}(\text{St}_{m,r})$$

Lemma 5 (Error of GoLore's projection). *Let $P \sim \mathcal{U}(\text{St}_{m,r})$, $Q \sim \mathcal{U}(\text{St}_{n,r})$, it holds for all $G \in \mathbb{R}^{m \times n}$ that*

$$\mathbb{E}[PP^\top] = \frac{r}{m} \cdot I, \quad \mathbb{E}[QQ^\top] = \frac{r}{n} \cdot I,$$

and

$$\mathbb{E}[\|PP^\top G - G\|_F^2] = \left(1 - \frac{r}{m}\right) \|G\|_F^2, \quad \mathbb{E}[\|GQQ^\top - G\|_F^2] = \left(1 - \frac{r}{n}\right) \|G\|_F^2.$$

$$\mathbb{E}[PP^\top G] = \mathbb{E}[PP^\top] \cdot \mathbb{E}[G] = \frac{r}{m} \nabla F(X)$$

The low-rank randomized projected gradient is an unbiased estimate of the true gradient

GoLore converges to desired solutions



Theorem 4 (Convergence rate of GoLore). *Under Assumptions 1-3, for any $T \geq 2 + 128/(3\underline{\delta}) + (128\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$, if we choose $\tau = \lceil 64/(3\underline{\delta}\beta_1) \rceil$,*

$$\beta_1 = \left(1 + \sqrt{\frac{\underline{\delta}^{3/2}\sigma^2 T}{L\Delta}}\right)^{-1}, \quad \text{and} \quad \eta = \left(4L + \sqrt{\frac{80L^2}{3\underline{\delta}\beta_1^2}} + \sqrt{\frac{80\tau^2 L^2}{3\underline{\delta}}} + \sqrt{\frac{16\tau L^2}{3\beta_1}}\right)^{-1},$$

GoLore using small-batch stochastic gradients and MSGD with MP converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O}\left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}}\right),$$

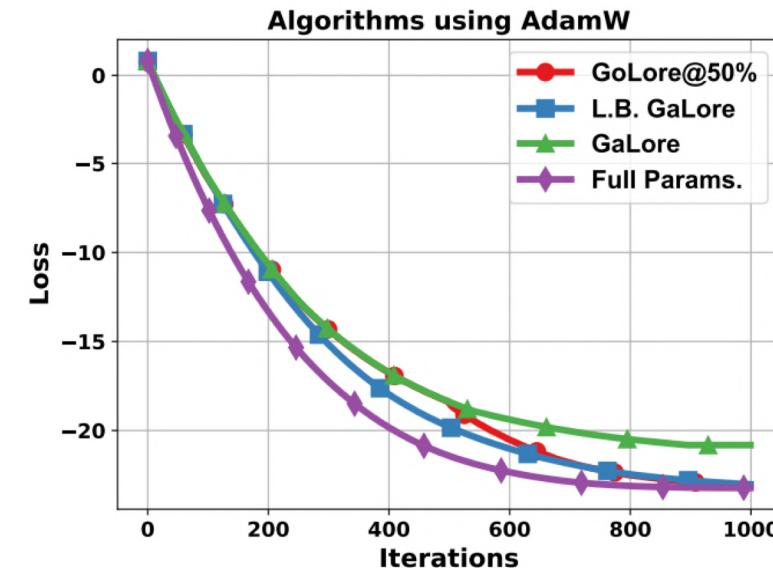
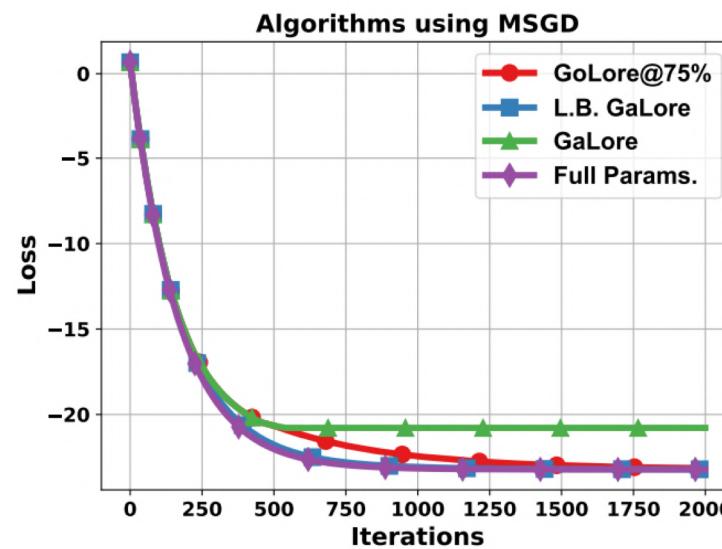
where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$ and $\underline{\delta} := \min_{\ell} \frac{r_{\ell}}{\min\{m_{\ell}, n_{\ell}\}}$.

- Theoretically, GoLore **converges** at rate $\mathcal{O}(1/\sqrt{T})$.

A hybrid strategy: GaLore + GoLore

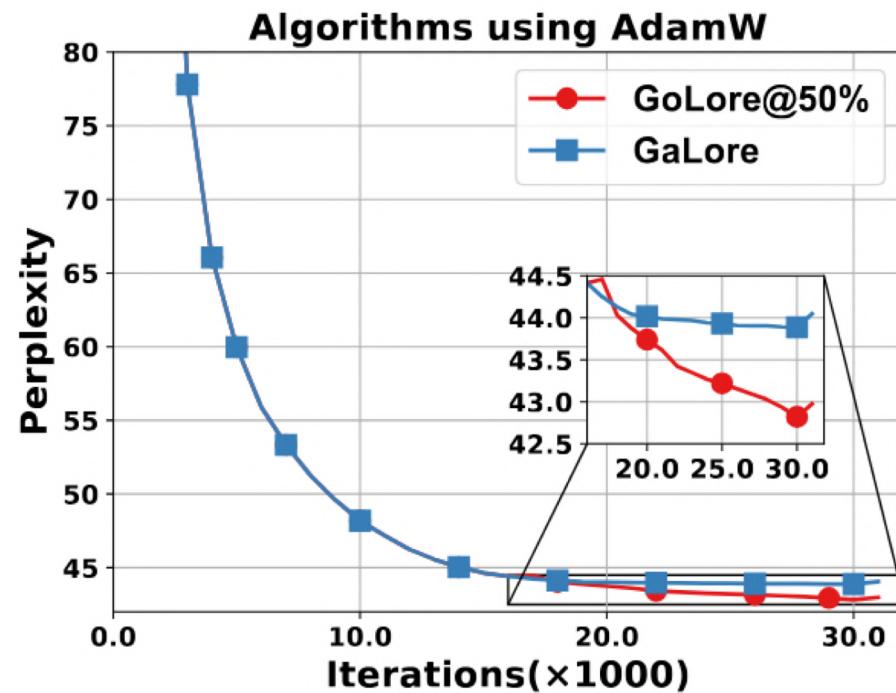
- SVD projection is preferred in initial stages: effectively capture gradient information
- Random projection is preferred when approaching solutions: avoid losing gradient information

$\text{GoLore}@x\% = \text{GaLore} (\text{first } (100-x)\% \text{ iters}) + \text{GoLore} (\text{last } x\% \text{ iters})$

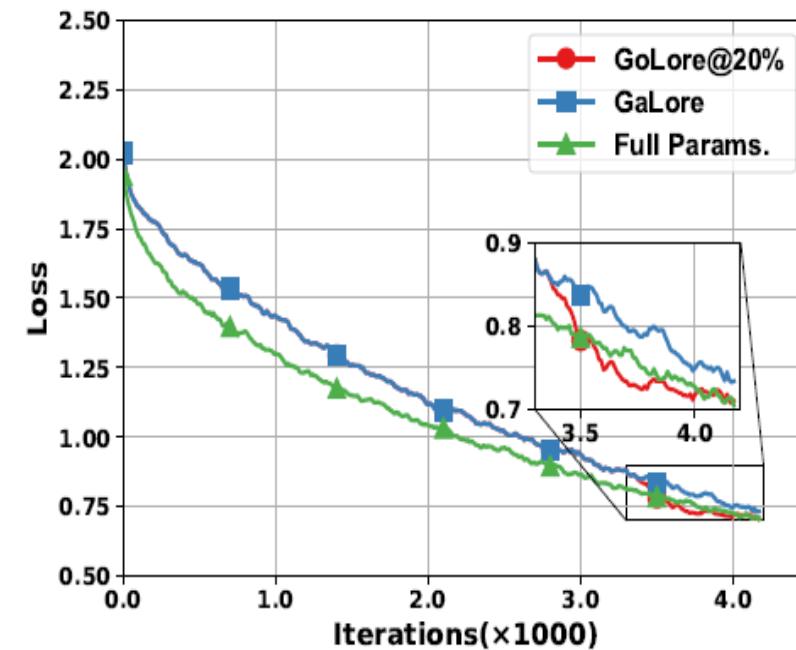


Experimental results

Pre-train LLaM2-60M on C4



Fine-tuning LLaMA2-7B on WinoGrande:



Experimental results

- Fine-tuning RoBERTa-BASE on GLUE benchmark:

Algorithm	CoLA	STS-B	MRPC	RTE	SST2	MNLI	QNLI	QQP	Avg
Full Params.	62.07	90.18	92.25	78.34	94.38	87.59	92.46	91.90	86.15
GaLore	61.32	90.24	92.55	77.62	94.61	86.92	92.06	90.84	85.77
GoLore@20%	61.66	90.55	92.93	78.34	94.61	87.02	92.20	90.91	86.03

- GoLore shows **superior** performance than GaLore in the above experiments.

Summary

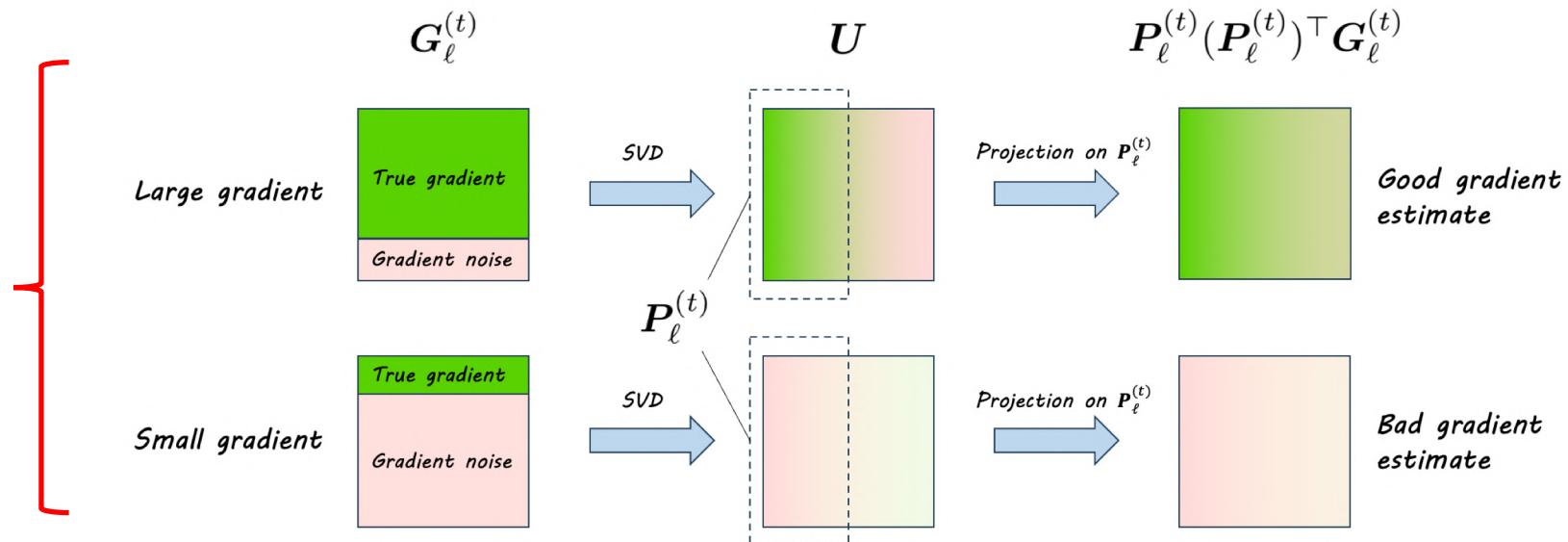
- Gradient low-rank projection can effectively save optimizer states
- GaLore cannot converge to desired solutions due to SVD projections
- Random projections enable GaLore to converge



paper

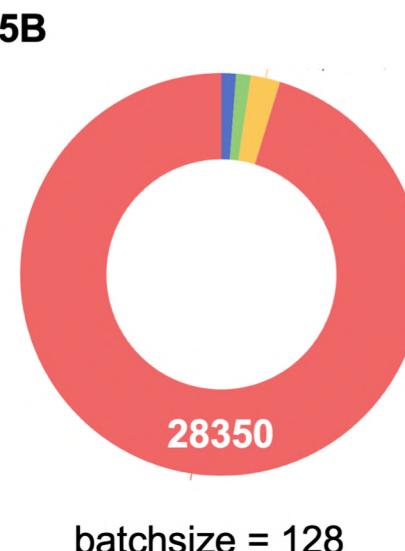
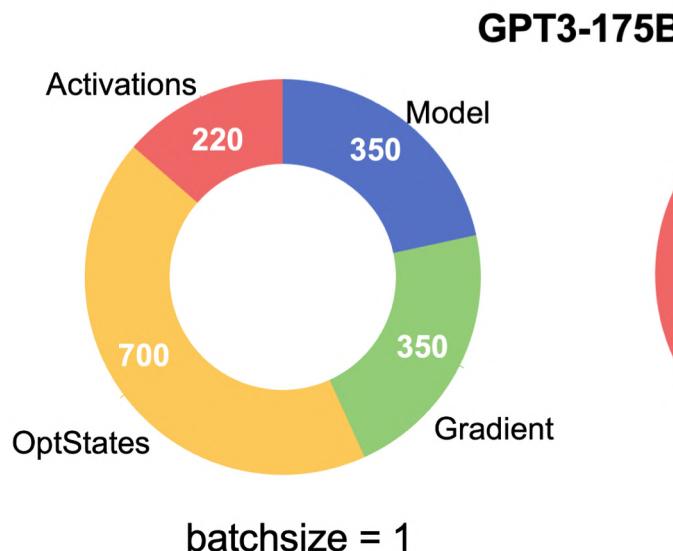
Github

Take-home message



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

- When batch-size is large, activations will dominate the memory



GaLore cannot save activations

$$\mathbf{g}_t = \mathbf{P}_t^\top \boxed{\mathbf{G}_t}$$

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

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

$$\delta_t = \frac{\gamma}{\sqrt{\mathbf{v}_t} + \epsilon} \odot \mathbf{m}_t$$

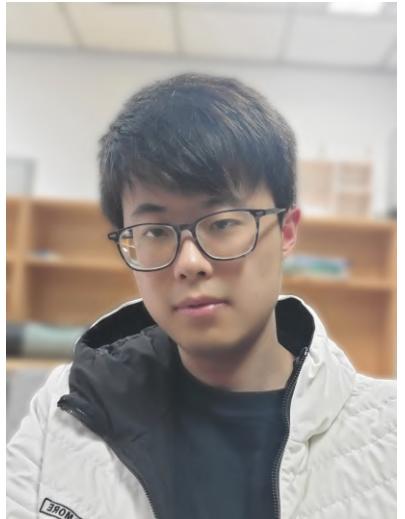


A Memory Efficient Randomized Subspace Optimization Method for Training Large Language Models

Kun Yuan (袁 坤)

Center for Machine Learning Research @ Peking University

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

Randomized Subspace Optimization

Randomized Subspace Optimization (RSO)

- Consider the optimization problem for the LLM training:

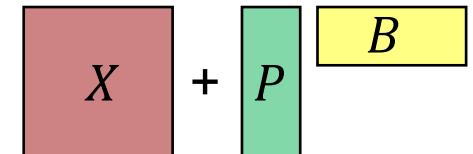
$$\min_{X \in \mathbb{R}^{m \times n}} f(X) := \mathbb{E}_{\xi}[F(X; \xi)].$$

- 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)

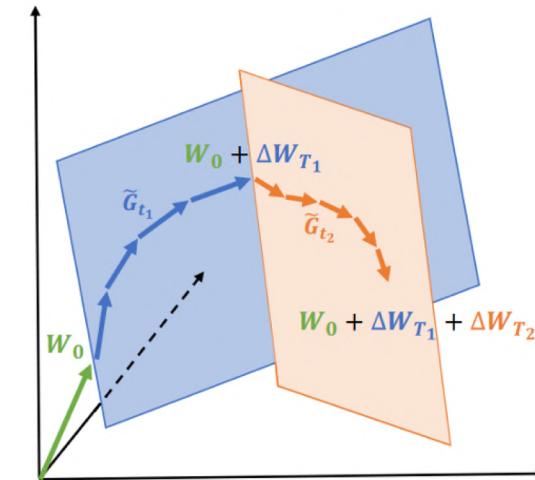
Can be solved with
any optimizer

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

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

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:

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

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

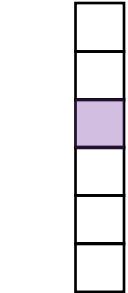
where the variable x is a vector in the above method

- 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 $f(\mathbf{X})$

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

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

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

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

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
RSO	$24nr$	$8bsn + 4bsr + 2bs^2$	$12nr$
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

PART 02

Convergence Guarantees

Assumptions

Definition 5.1 (Expected ϵ -inexact solution). A solution $\tilde{\mathbf{B}}^k$ is said to be an expected ϵ -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 ℓ . (random subspace)

RSO convergence theorem

Theorem

Let each subproblem in RSO iteration be solved starting from the initial point $\mathbf{B}^0 = \mathbf{0}$ to an expected ϵ -inexact solution $\tilde{\mathbf{B}}^k$ with a suitable choice of η^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 ϵ -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	$O((\sum_{\ell=1}^L n_\ell r_\ell)^2 \epsilon^{-2})$	$O((\sum_{\ell=1}^L n_\ell r_\ell)^2 \epsilon^{-3})$
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 03

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)	15.64 (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)
RSO	34.55 (0.14G)	25.34 (0.27G)	18.86 (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)

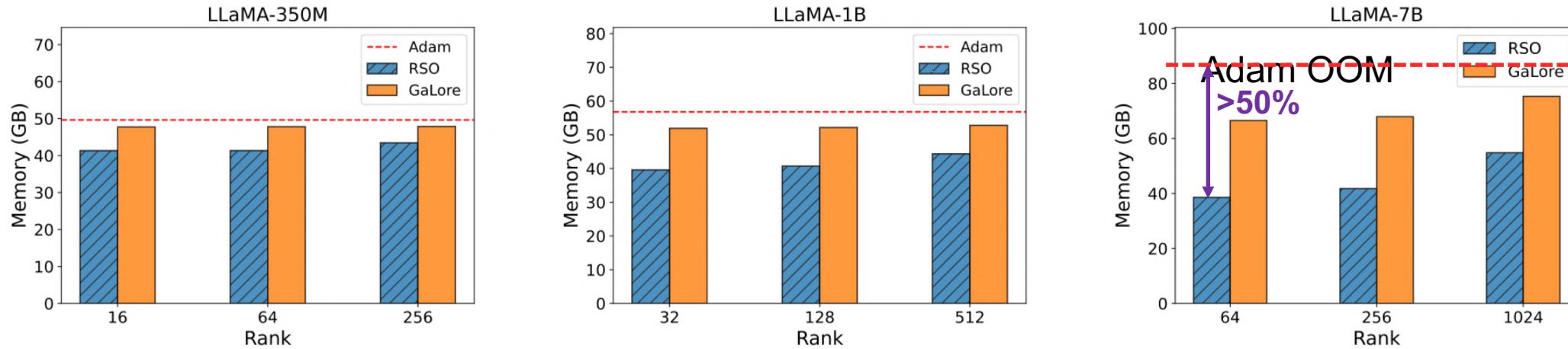


Figure 2: Comparison of peak memory usage (in GB) per device for RSO and GaLore during LLaMA training with varying ranks. All hyperparameters, except rank, are consistent with [Zhao et al., 2024a]. Adam’s memory usage is reported for LLaMA-350M and LLaMA-1B but excluded for LLaMA-7B due to an out-of-memory (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
RSO	54.81	64	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.

RSO performance : Fine-tune



	CoLA	STS-B	MRPC	RTE	SST2	MNLI	QNLI	QQP	Avg
Adam	62.24	90.92	91.30	79.42	94.57	87.18	92.33	92.28	86.28
GaLore (rank=4)	60.35	90.73	92.25	79.42	94.04	87.00	92.24	91.06	85.89
LoRA (rank=4)	61.38	90.57	91.07	78.70	92.89	86.82	92.18	91.29	85.61
RSO (rank=4)	62.47	90.62	92.25	78.70	94.84	86.67	92.29	90.94	86.10
GaLore (rank=8)	60.06	90.82	92.01	79.78	94.38	87.17	92.20	91.11	85.94
LoRA (rank=8)	61.83	90.80	91.90	79.06	93.46	86.94	92.25	91.22	85.93
RSO (rank=8)	64.62	90.71	93.56	79.42	95.18	86.96	92.44	91.26	86.77

Table 4: Evaluation of various fine-tuning methods on the GLUE benchmark using the pre-trained RoBERTa-Base model. The average score across all tasks is provided.

Memory = Model + Gradient + Optimizer states + Activations

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

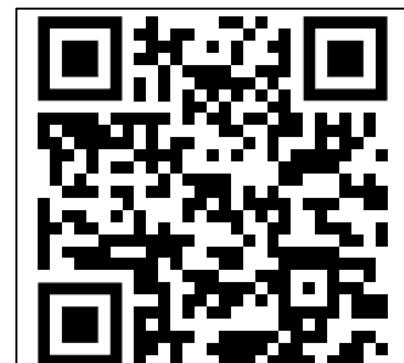
$$\tilde{\mathbf{B}}^k \approx \underset{\mathbf{B} \in \mathbb{R}^{r \times n}}{\operatorname{argmin}} \left\{ f(\mathbf{X}^k + \mathbf{P}^k \mathbf{B}) + \frac{1}{2\eta^k} \|\mathbf{B}\|^2 \right\},$$

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

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



paper



Github Code

References



- 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!

We're hiring PostDocs in optimization, machine learning, and LLMs

Kun Yuan homepage: <https://kunyuan827.github.io/>