

# Subspace Optimization for Large Language Models with Convergence Guarantees

Kun Yuan (袁坤)

Center for Machine Learning Research @ Peking University

National University of Singapore (NUS), Nov. 9, 2024

## Joint work with

---



**Yutong He**  
(Peking U.)



**Pengrui Li**  
(Beihang U.)



**Yipeng Hu**  
(Peking U.)



**Chuyan Chen**  
(Peking U.)

# PART 01

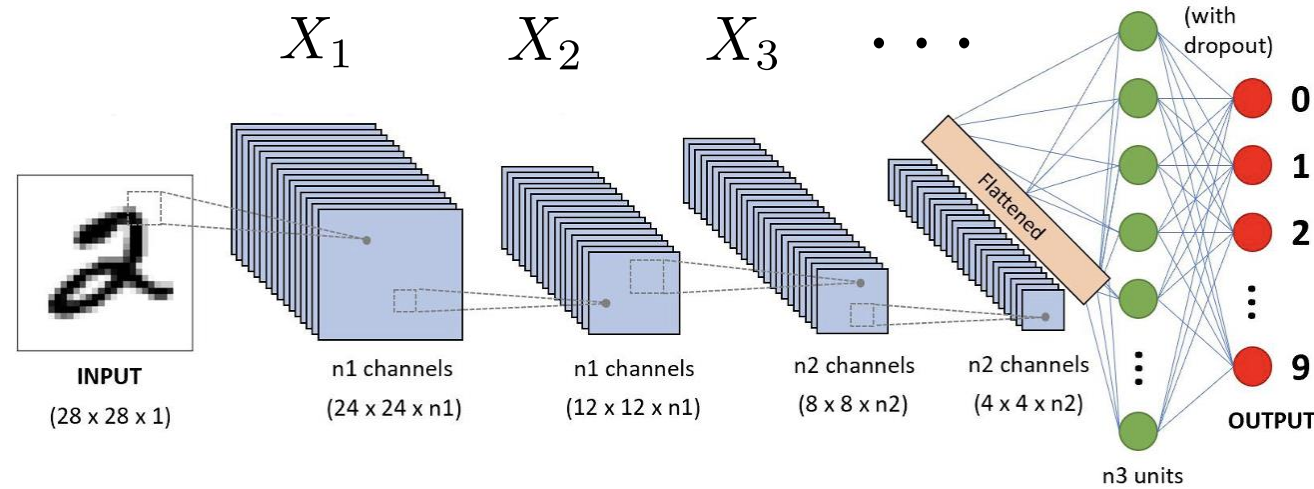
---

## 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

LLM cost function:

$$\mathbf{X}^* = \arg \min_{\mathbf{X}} \left\{ \mathbb{E}_{\xi \sim \mathcal{D}} \left[ \overset{\text{cross entropy}}{\underset{\substack{\text{data distribution} \quad \text{pred. token} \quad \text{real token}}}{L(h(\mathbf{X}; \xi), y)}} \right] \right\}$$

# 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

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

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

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

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**.

**P**  $G_t = \nabla F(\mathbf{X}_t; \xi_t)$

**2P**  $\left\{ \begin{array}{l} M_t = (1 - \beta_1)M_{t-1} + \beta_1 G_t \\ V_t = (1 - \beta_2)V_{t-1} + \beta_2 G_t \odot G_t \end{array} \right.$

**P**  $\mathbf{X}_{t+1} = \mathbf{X}_t - \frac{\gamma}{\sqrt{V_t} + \epsilon} \odot M_t$

**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

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

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

The gradient is derived as follows

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

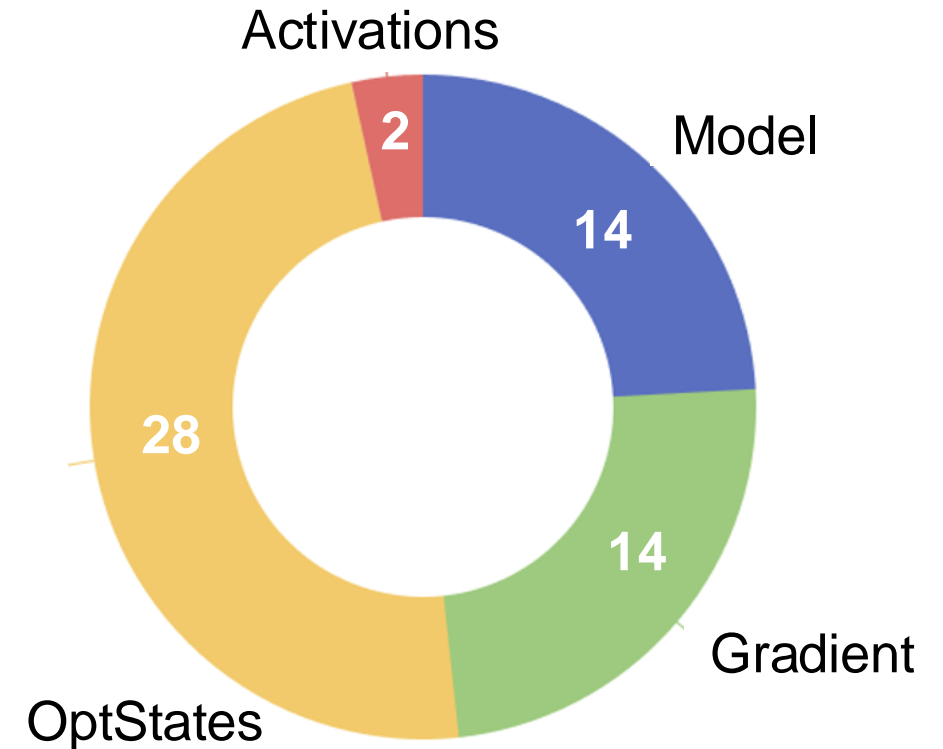
Need to store activations  $z_1, z_2, \dots, z_L$

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

# Minimum memory requirement

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

- Parameters: 7B
- Model storage:  $7\text{B} * 2 \text{ Bytes} = 14 \text{ GB}$
- 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**



RTX 4090: 24GB



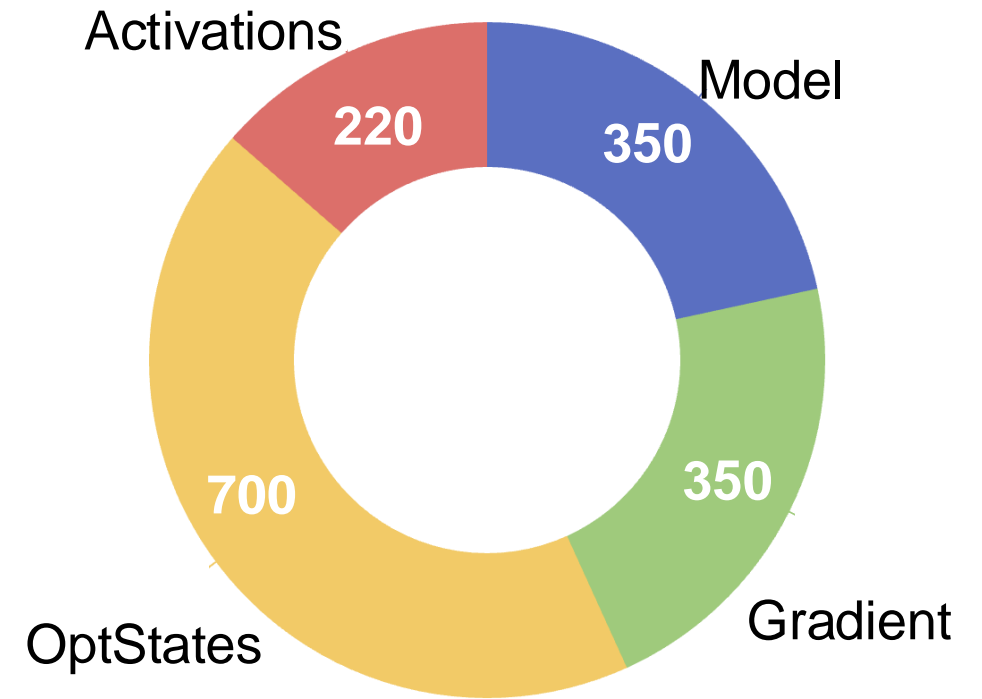
A100 80G



The minimum requirement is A100 \* 1

# 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**



x 21

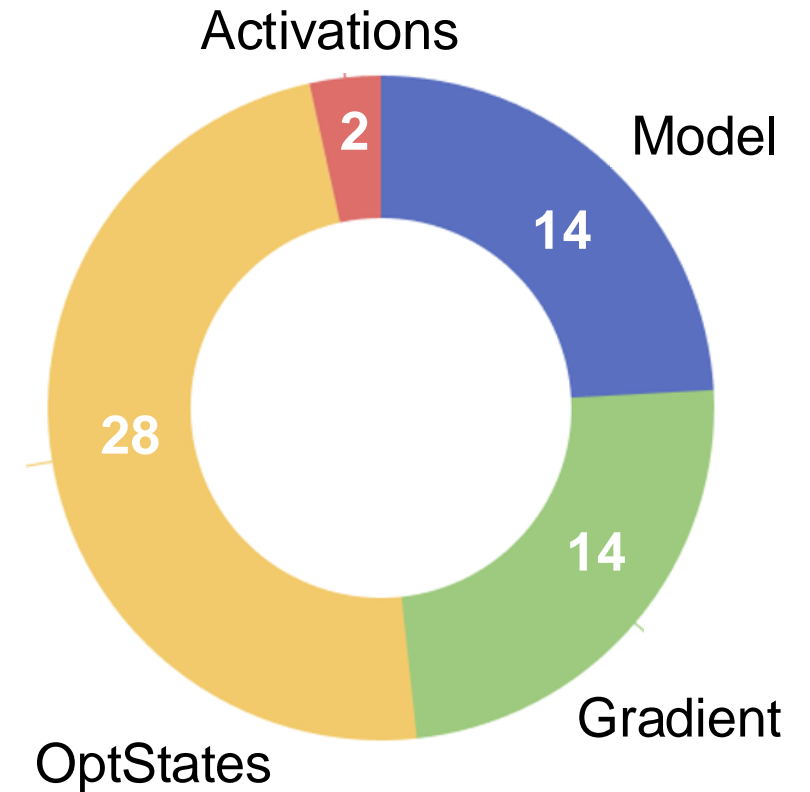
**A100 80G**

The minimum requirement is  $A100 * 21$

**Very expensive!**

# Memory-efficient algorithm is 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**
- Activation-incurred memory is **relatively minor** when using a single batch size
- Gradient-incurred memory can be **removed** by layer-wise calculation and dropping
- How to save memory caused by optimizer states?



## PART 02

---

# GaLore: Gradient Low-Rank Projection



---

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

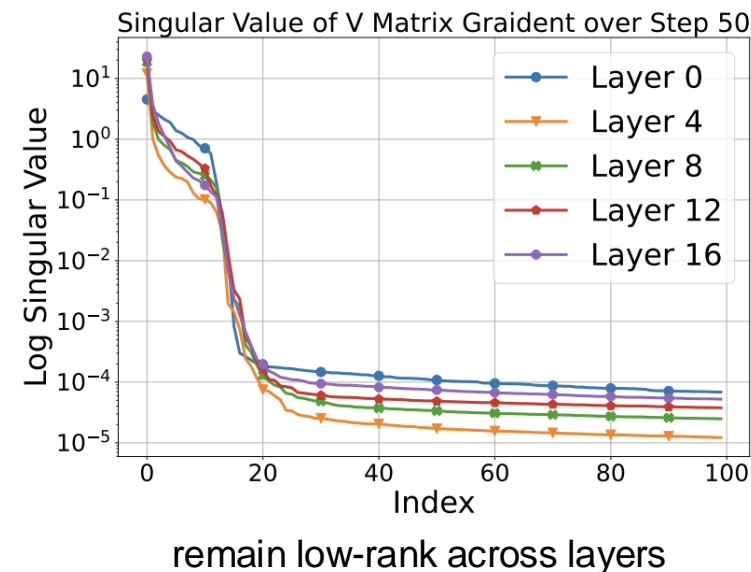
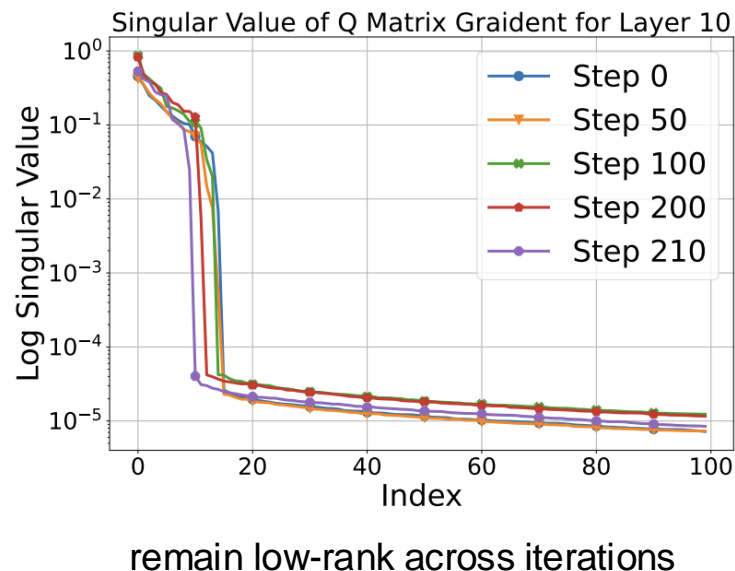
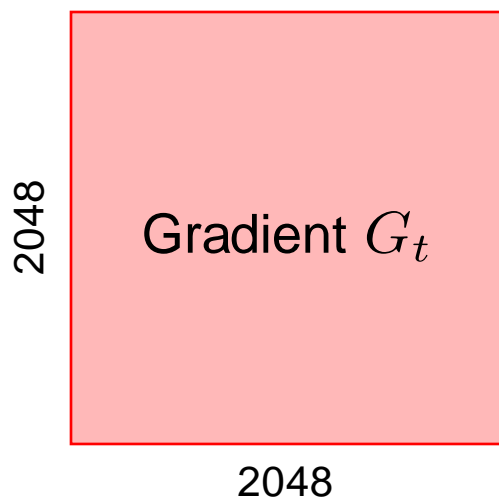
---

Jiawei Zhao<sup>1</sup> Zhenyu Zhang<sup>3</sup> Beidi Chen<sup>2,4</sup> Zhangyang Wang<sup>3</sup> Anima Anandkumar<sup>\*1</sup> Yuandong Tian<sup>\*2</sup>

- 
- GaLore is a novel approach for reducing memory consumption from optimizer states
  - The first algorithm that enables LLaMA-7B pre-training on a single 4090 GPU (24GB)
  - Memory-efficient without severe performance degradation

# GaLore: Gradient Low-Rank Projection

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



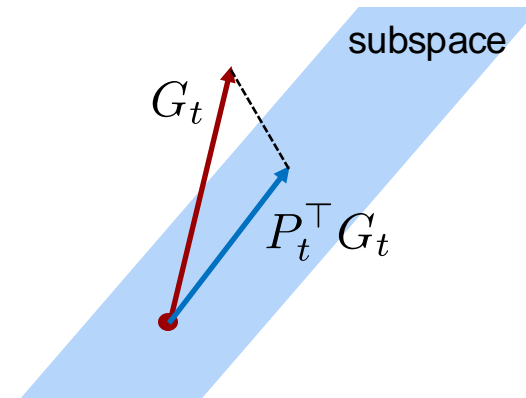
- Given a gradient matrix with dimensions 2048 by 2048, around **top 10** eigenvalues dominate
- How to utilize the low-rank structure in gradients?

# 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:

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

▷ dims  $r \times n$

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

▷ dims  $r \times n$

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

▷ dims  $r \times n$

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

▷ dims  $r \times n$

**Simplified as**

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

- Parameter updates:

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

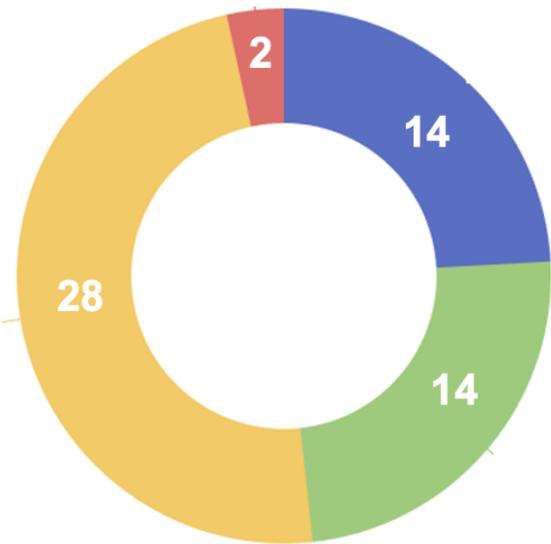
▷ dims  $m \times n$

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

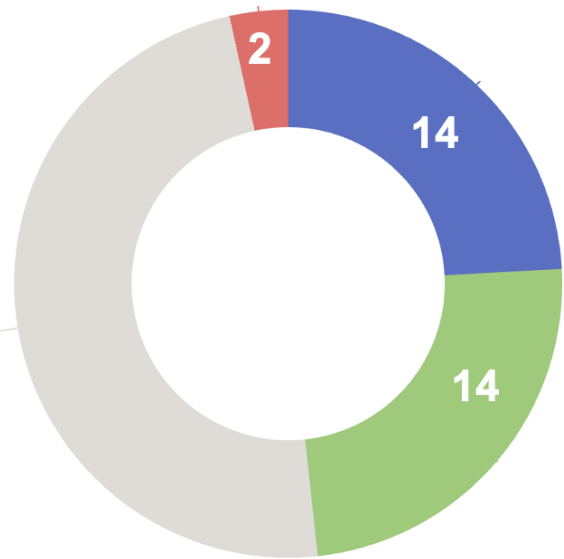
**trivial memory cost**

# GaLore: Gradient Low-Rank Projection

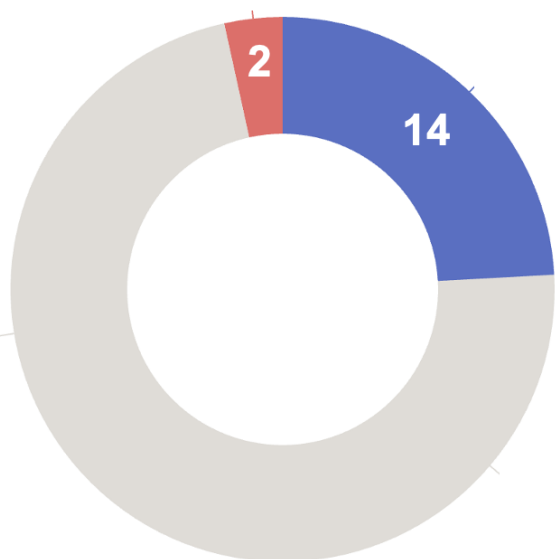
LLaMA 7B



Adam



Adam + GaLore



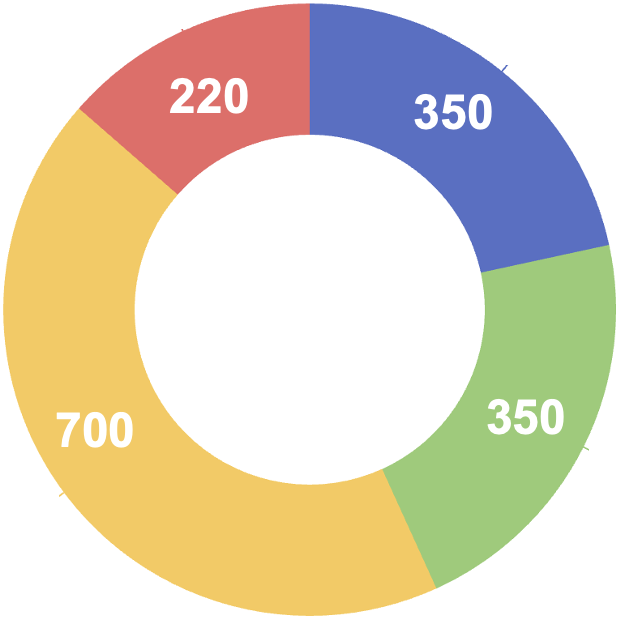
Adam + GaLore + Layerwise  
gradient dropping (LWGD)



**RTX 4090 affordable !!**

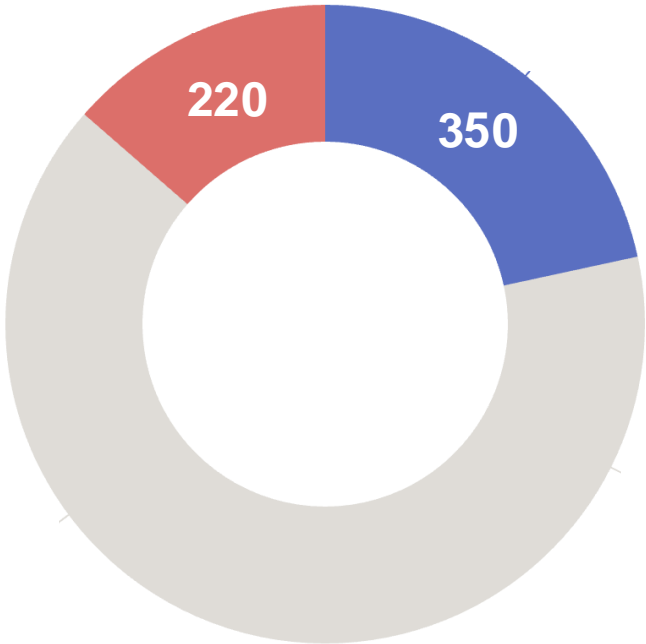
# GaLore: Gradient Low-Rank Projection

GPT3 175B



Adam

**A100 \* 21**



Adam + GaLore + LWGD

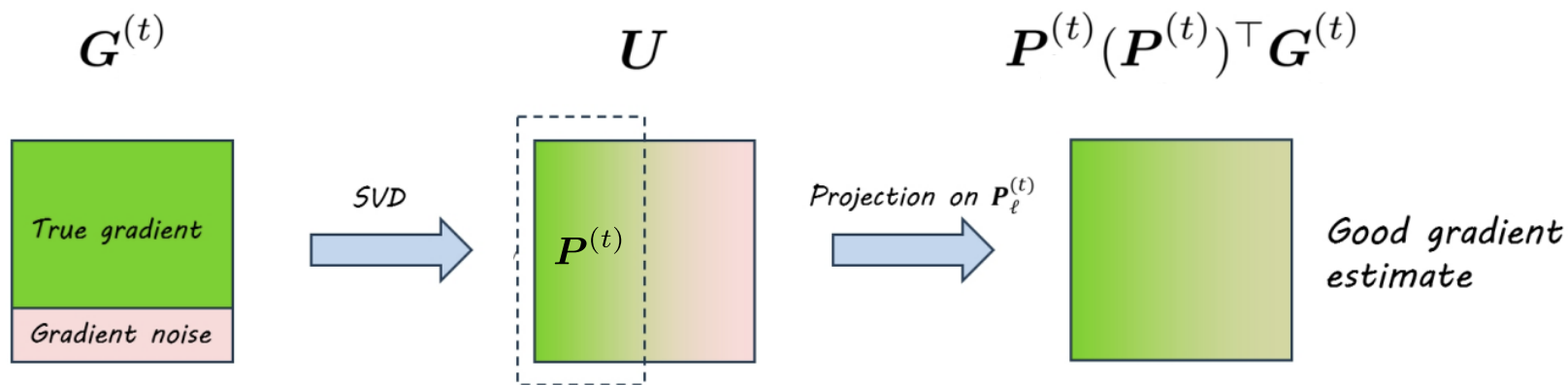
**A100 \* 8**

# GaLore: Gradient Low-Rank Projection

- Recall the GaLore update:  $\mathbf{X}_{t+1} = \mathbf{X}_t + \mathbf{P}_t \rho(\mathbf{P}_t^\top \mathbf{G}_t)$
- How to find the projection matrix? **SVD decomposition!**

$$\mathbf{G}_t = \mathbf{U} \Sigma \mathbf{V}^\top \longrightarrow \mathbf{P}_t = \boxed{\mathbf{U}[:, : r]} \in \mathbb{R}^{m \times r}$$

|  
Select the first  $r$  columns



[Subspace Optimization for Large Language Models with Convergence Guarantees]

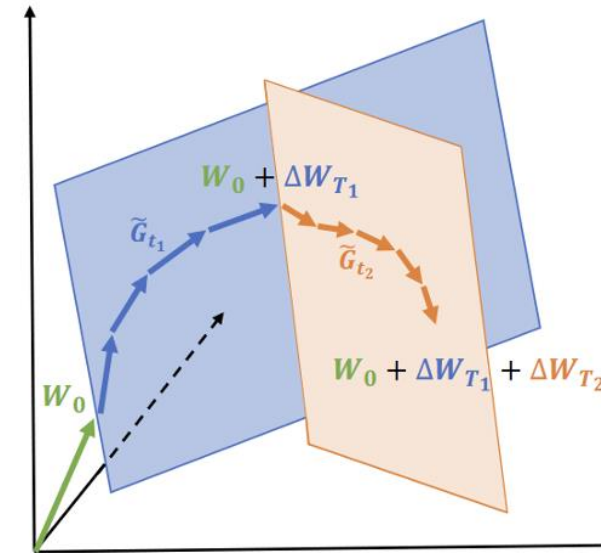
# GaLore: Gradient Low-Rank Projection

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

(The complete GaLore algorithm)

$$\begin{cases} P_t \leftarrow \text{SVD}(G_t) & \text{if } t \bmod \tau = 0 \\ P_t \leftarrow P_{t-1} & \text{otherwise} \end{cases}$$

$$X_{t+1} = X_t + P_t \rho(P_t^\top G_t)$$



[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)
<b>GaLore</b>	<b>34.88</b> (0.24G)	<b>25.36</b> (0.52G)	<b>18.95</b> (1.22G)	<b>15.64</b> (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

Fine-tuning RoBERTa-Base on GLUE

	Memory	CoLA	STS-B	MRPC	RTE	SST2	MNLI	QNLI	QQP	Avg
Full Fine-Tuning	747M	62.24	90.92	91.30	79.42	94.57	87.18	92.33	92.28	86.28
<b>GaLore (rank=4)</b>	253M	60.35	<b>90.73</b>	<b>92.25</b>	<b>79.42</b>	<b>94.04</b>	<b>87.00</b>	<b>92.24</b>	91.06	<b>85.89</b>
LoRA (rank=4)	257M	<b>61.38</b>	90.57	91.07	78.70	92.89	86.82	92.18	<b>91.29</b>	85.61
<b>GaLore (rank=8)</b>	257M	60.06	<b>90.82</b>	<b>92.01</b>	<b>79.78</b>	<b>94.38</b>	<b>87.17</b>	92.20	91.11	<b>85.94</b>
LoRA (rank=8)	264M	<b>61.83</b>	90.80	91.90	79.06	93.46	86.94	<b>92.25</b>	<b>91.22</b>	85.93

GaLore looks great

But does GaLore provably converge to the desired solution?

**Not Necessarily True!**

Y. He, P. Li, Y. Hu, C. Chen, and K. Yuan, *Subspace Optimization for Large Language Models with Convergence Guarantees*, 2024





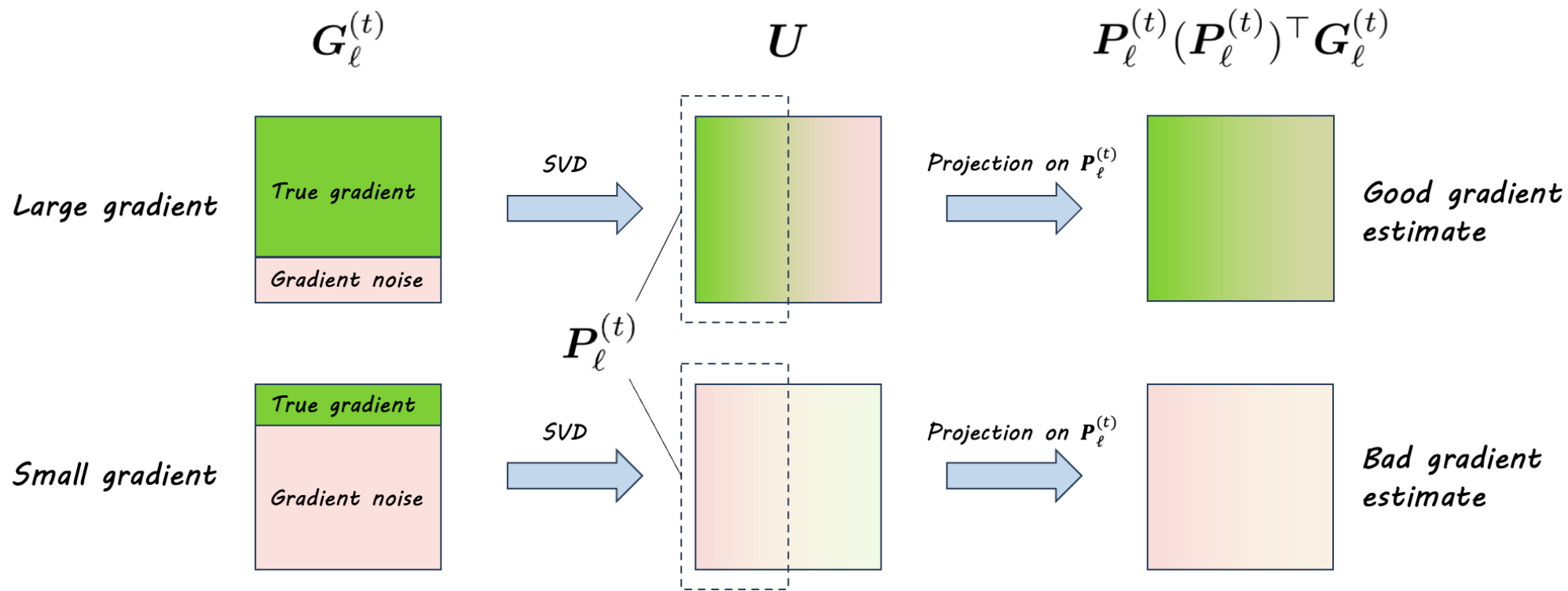
## PART 03

---

# Non-convergence and convergence in GaLore



# Intuition behind GaLore's non-convergence

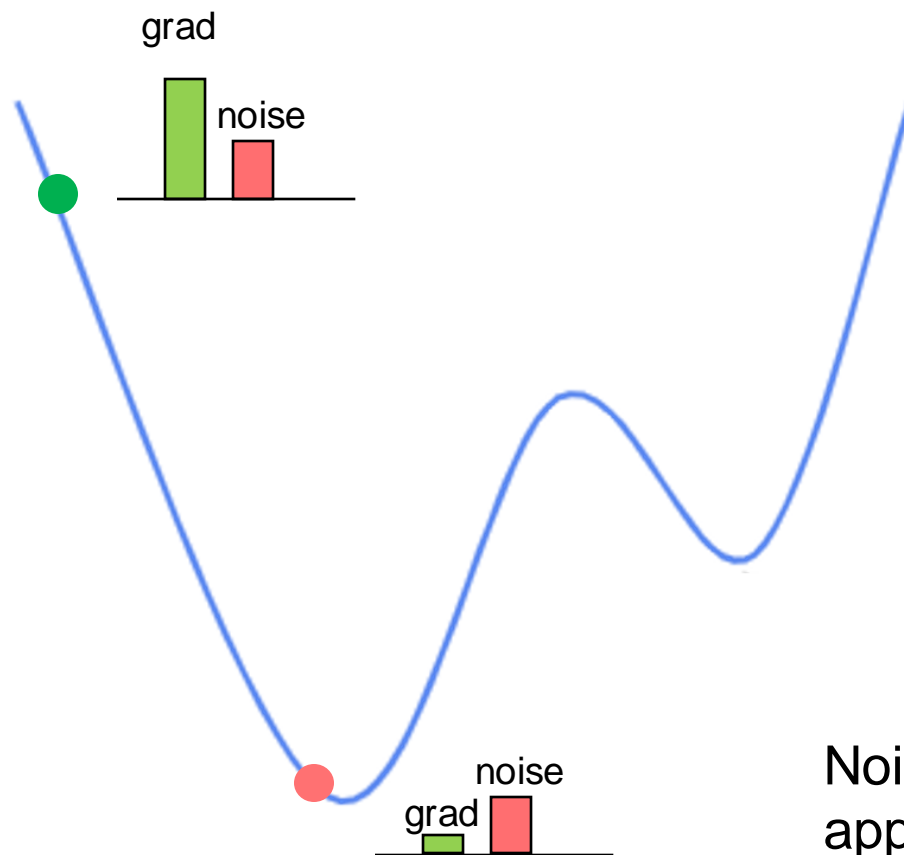


When gradient noise dominates the stochastic gradient, SVD captures **noise-dominated** subspace!

**All gradient information is lost !**

# Is non-convergence common? Yes!

Gradient dominates  
during the initial stages

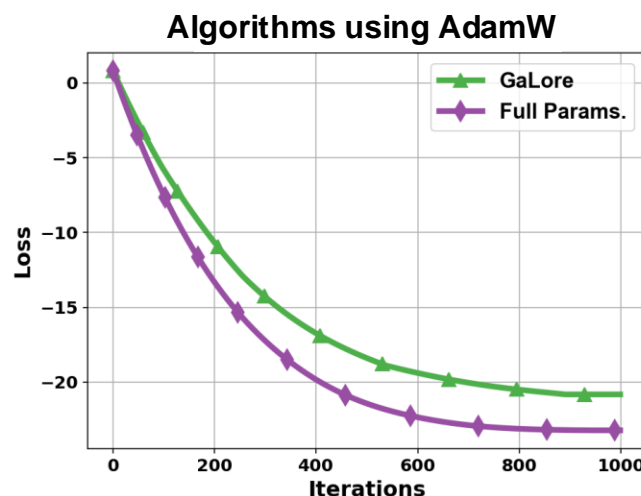
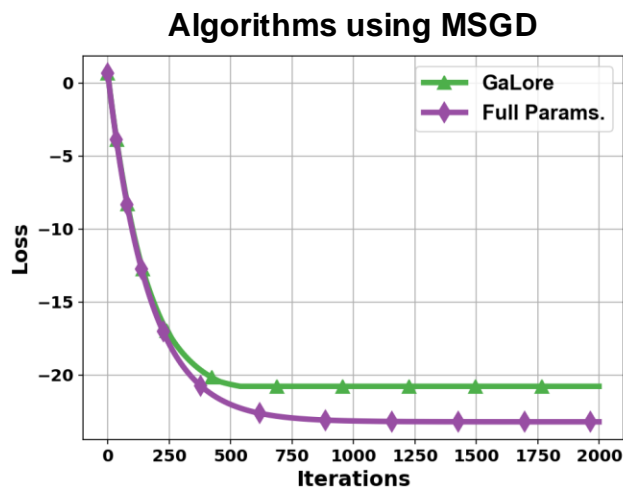


Noise dominates when  
approaching the local minimum

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

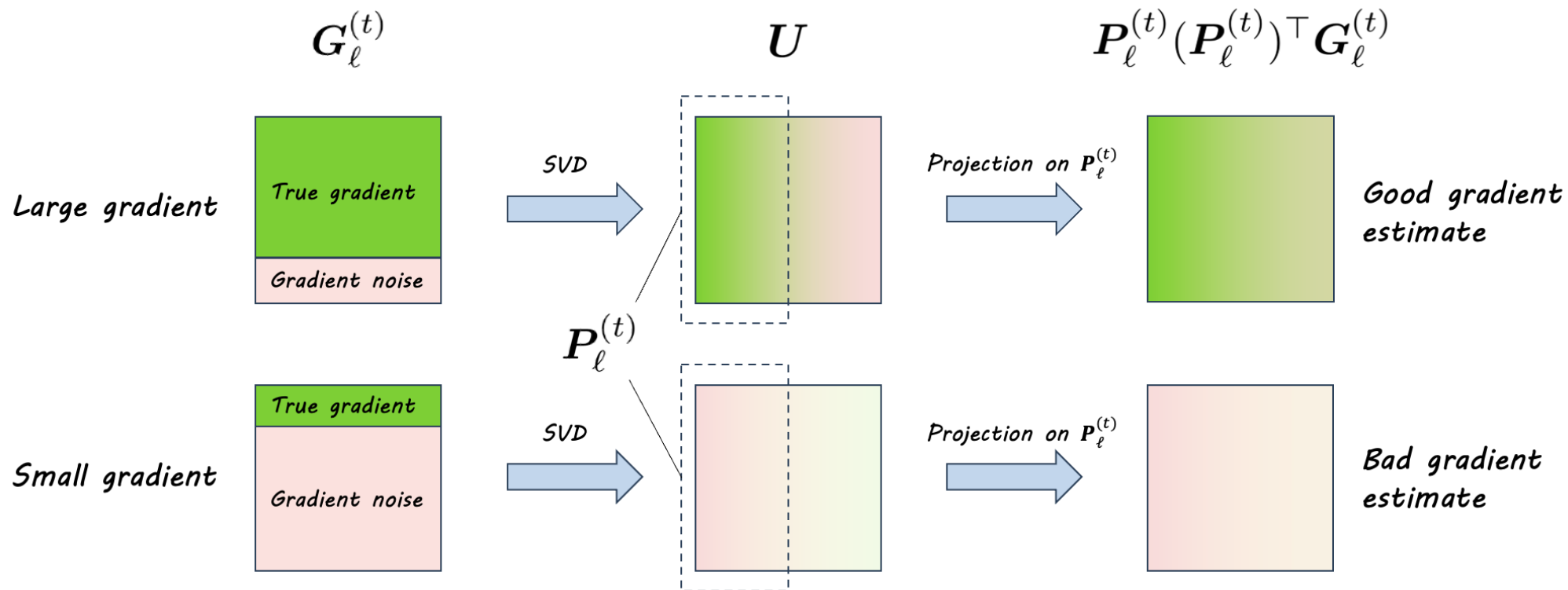
$$f(\mathbf{X}) = \frac{1}{2} \|\mathbf{A}\mathbf{X}\|_F^2 + \langle \mathbf{B}, \mathbf{X} \rangle_F, \quad \nabla F(\mathbf{X}; \xi) = \nabla f(\mathbf{X}) + \xi \sigma \mathbf{C}, \quad (1)$$

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



GaLore does **NOT** converge to desired solutions

# Under what conditions can GaLore converge?



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

- 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*

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

*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  $\tau = \lceil 64/(3\underline{\delta}\beta_1) \rceil$ ,  $\mathcal{B} = \lceil 1/(\underline{\delta}\beta_1) \rceil$ ,*

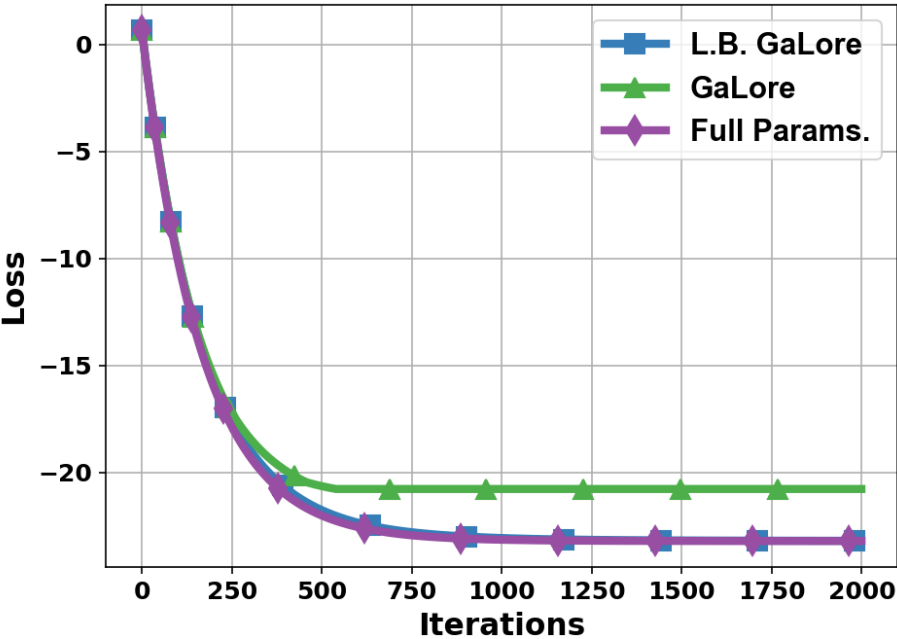
$$\beta_1 = \left(1 + \sqrt{\frac{\underline{\delta}^{3/2}\sigma^2 T}{L\Delta}}\right)^{-1}, \quad \text{and} \quad \eta = \left(4L + \sqrt{\frac{80L^2}{3\underline{\delta}\beta_1^2}} + \sqrt{\frac{40\tau^2 L^2}{\underline{\delta}}} + \sqrt{\frac{32\tau L^2}{3\beta_1}}\right)^{-1},$$

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

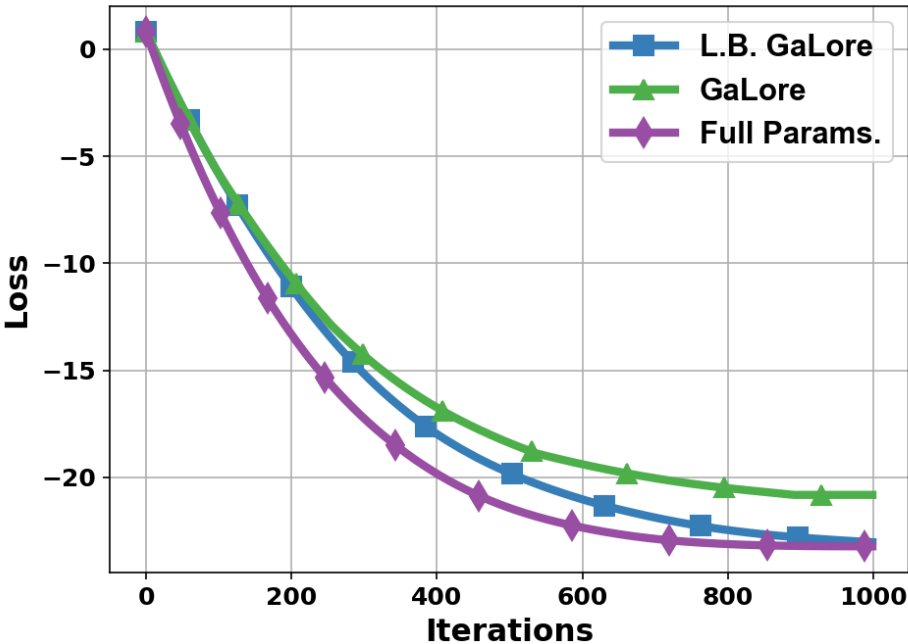
$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O}\left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}}\right),$$

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

Algorithms using MSGD



Algorithms using AdamW



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)  $\text{St}_{m,r} = \{\mathbf{P} \in \mathbb{R}^{m \times r} \mid \mathbf{P}^\top \mathbf{P} = \mathbf{I}_r\}.$

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

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

# Stiefel random projection injects contractive error

- Instead of SVD, GoLore samples the projection matrix uniformly on the Stiefel manifold:

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

- The following Lemma illustrates the projection error in GoLore:

**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*

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

**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*

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

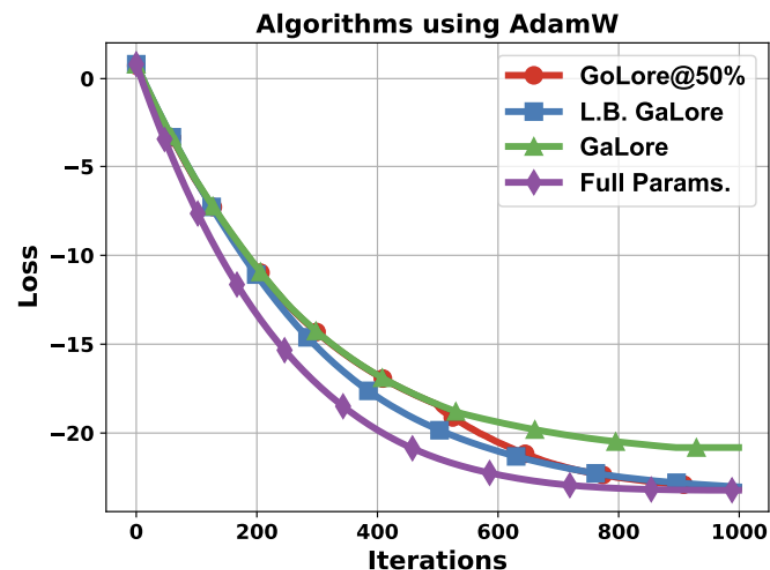
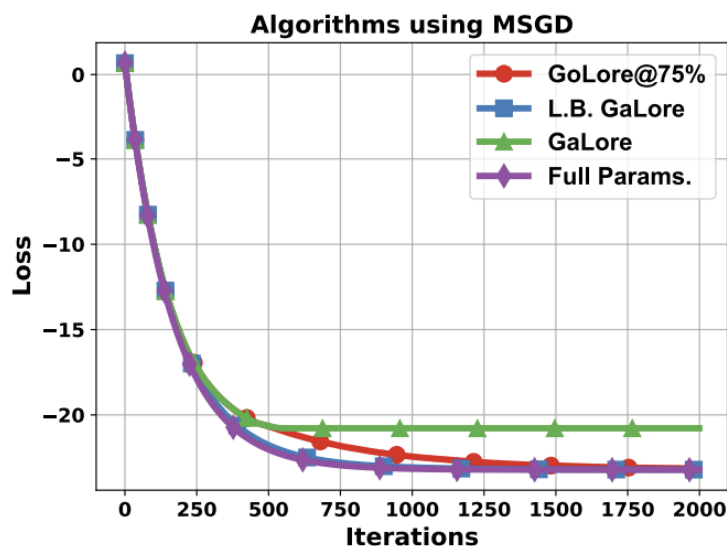
where  $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$  and  $\underline{\delta} := \min_{\ell} \frac{r_{\ell}}{\min\{m_{\ell}, n_{\ell}\}}$ .

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

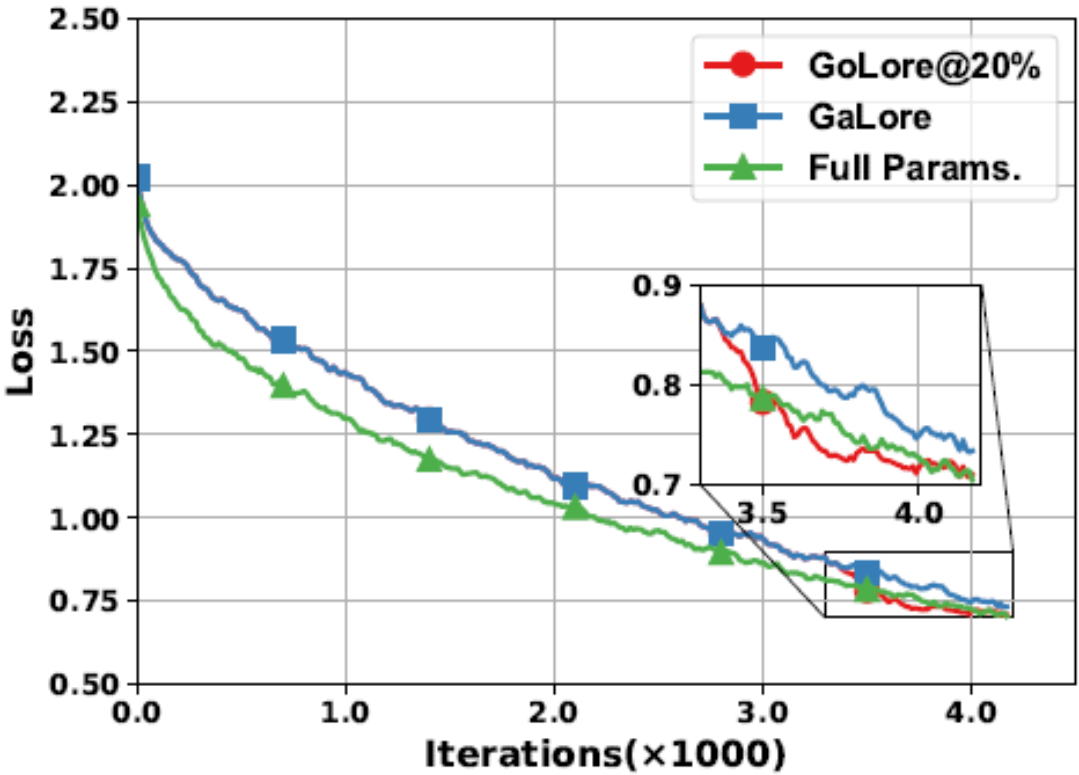
# A hybrid strategy: GaLore + GoLore

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

GoLore@x% = GaLore (first (100-x)% iters) + GoLore (last x% iters)



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	<b>94.61</b>	86.92	92.06	90.84	85.77
GoLore@20%	<b>61.66</b>	<b>90.55</b>	<b>92.93</b>	<b>78.34</b>	<b>94.61</b>	<b>87.02</b>	<b>92.20</b>	<b>90.91</b>	<b>86.03</b>

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

# Improving the computational efficiency

- GaLore/GoLore always computes the full gradient before compressing them into subspaces.
- Can we compute the compressed gradient directly, without computing the full gradient?

Original Implementation

$$y = Wx$$

$$\nabla_W \mathcal{L} = (\nabla_y \mathcal{L})x^\top \quad \text{Backpropagated gradient}$$

$$W \leftarrow W + B\rho(B^\top(\nabla_W \mathcal{L}))$$

$$W = W_0 + BA$$
$$\longrightarrow$$

New Implementation

$$y = W_0x + BAx$$

$$\nabla_A \mathcal{L} = B^\top(\nabla_y \mathcal{L})x^\top$$

$$A \leftarrow A + \rho(\nabla_A \mathcal{L})$$





# Improving the computational efficiency

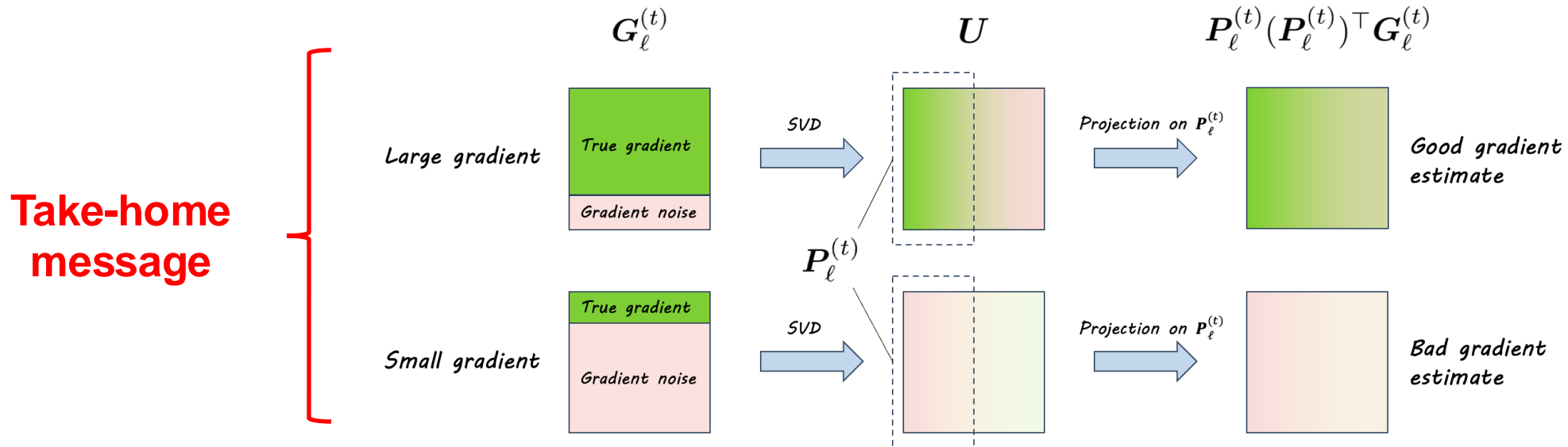
- Comparing the computational complexities:

GaLore Implementation	Memory	Computation
(Zhao et al., 2024)	$mn + rm + rn + bm$	$6bmn + 4rmn + 2mn + 3rn$
Our ReLoRA-like version	$mn + rm + 2rn + bm + br$	$4bmn + 4brm + 6brn + 5rn$

- When  $r \ll \min\{m, n\}$ , this new version reduces the computation complexity from  $(6b + 4r + 2)mn$  to  $4bmn$ , with minimal memory overhead.

# 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



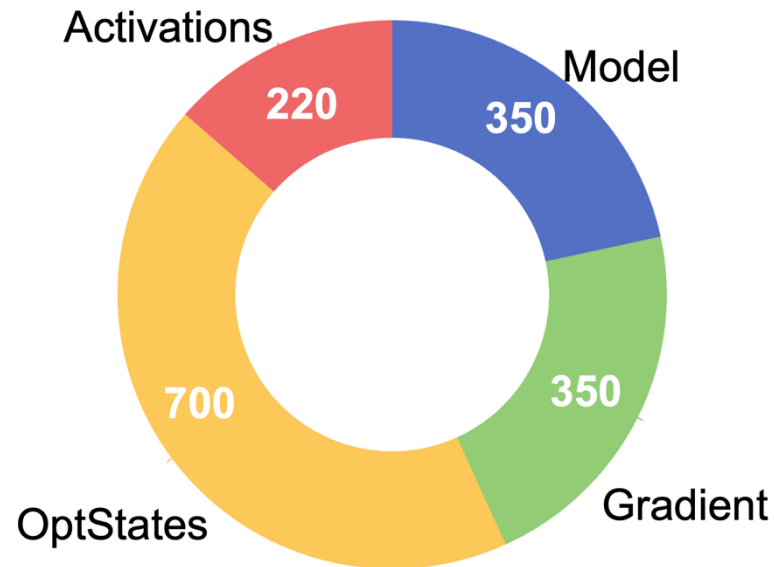
# Activation dominates when batch size is large

$$\text{Memory} = \text{Model} + \text{Gradient} + \text{Optimizer states} + \text{Activations}$$

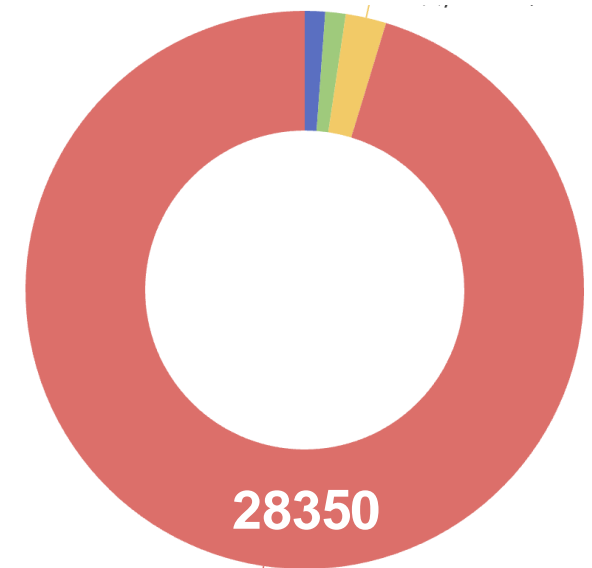
GPT3-175B

Activation dominates when  
batch-size is large

GaLore/GoLore cannot save  
activations



batchsize = 1



batchsize = 128

# Enhancing Zeroth-Order Fine-Tuning for Language Models with Low-Rank Structure

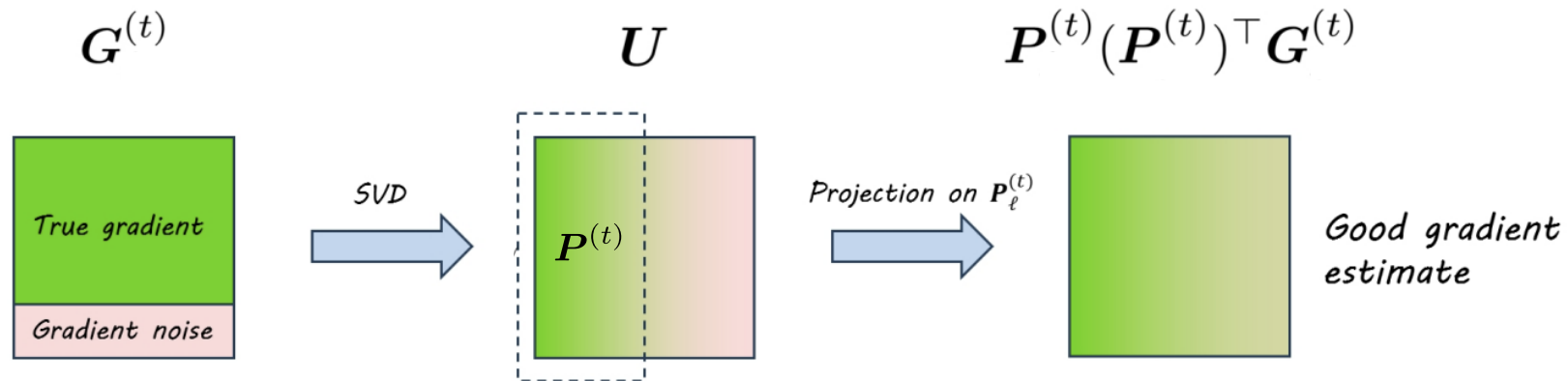
<https://arxiv.org/abs/2410.07698>

**Our new work saves the activation cost!**



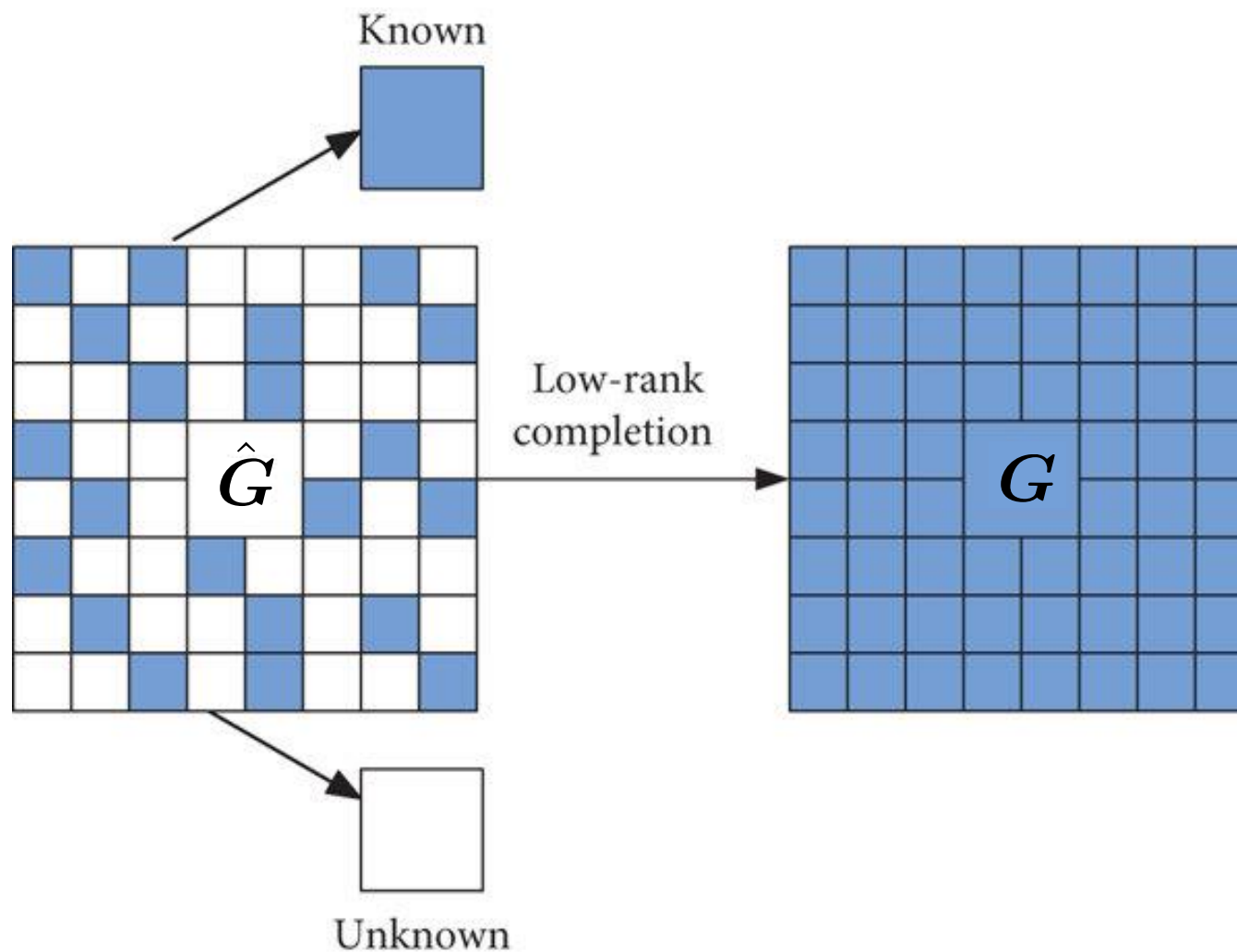
## Open question

- GaLore/GoLore always computes the gradient  $G$  before compressing them into subspaces.



- Can we avoid calculating the gradient  $G$  to save computations? Low-rank matrix completion?

# Low-rank matrix completion



# Thank you!

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

