

Bridge theory to practice: One-step full gradient can suffice for low-rank fine-tuning, provably and efficiently

Fanghui Liu

fanghui.liu@warwick.ac.uk

Department of Computer Science, University of Warwick, UK

Centre for Discrete Mathematics and its Applications (DIMAP), Warwick

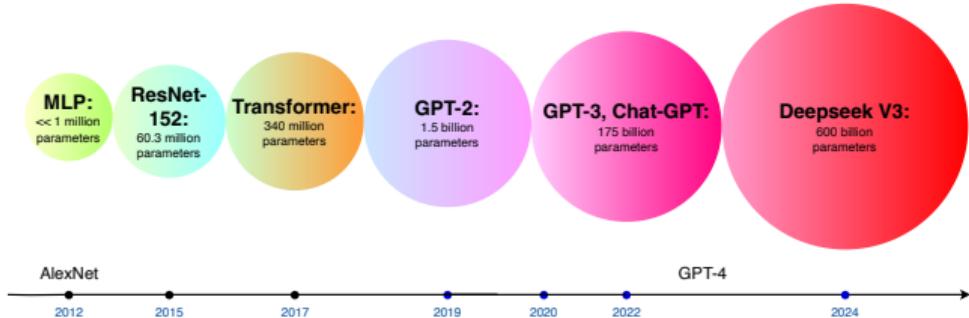
[joint work with Yuanhe Zhang (Warwick) and Yudong Chen (UW-Madison)]



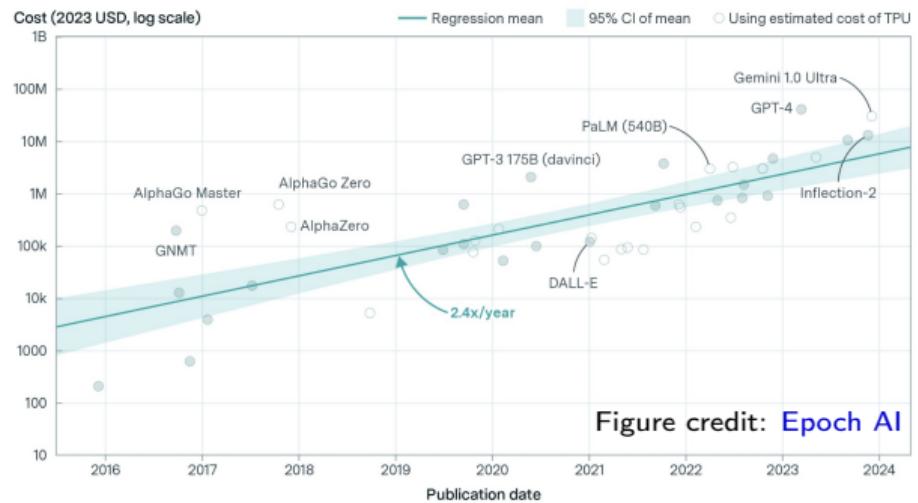
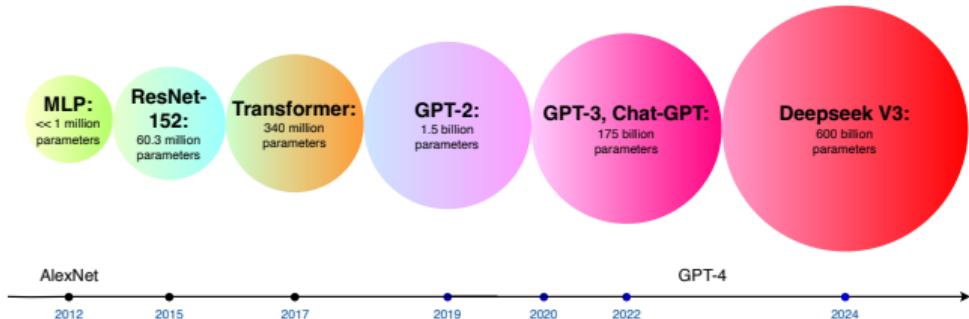
Outline

- ❑ Fine-tuning in LLMs
- ❑ How does theory contribute to practice?
 - understanding: subspace alignment
 - theory-grounded algorithm for efficiency improvement
 - clarify misconceptions
- ❑ Proofs

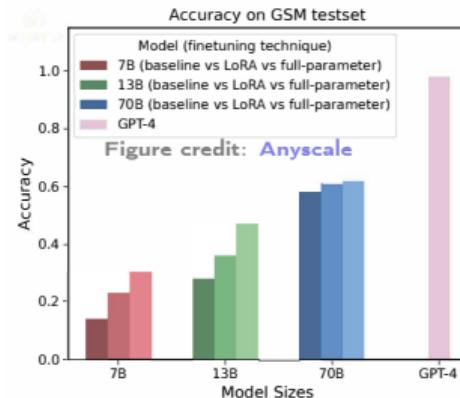
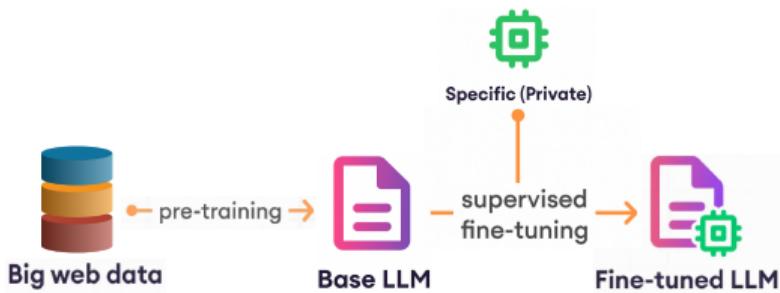
In the era of machine learning (Pre-training)



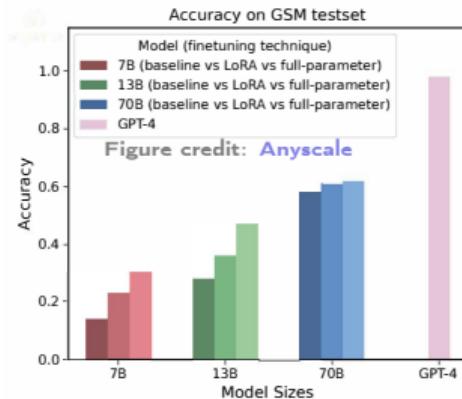
In the era of machine learning (Pre-training)



From pre-training to (parameter-efficient) fine-tuning

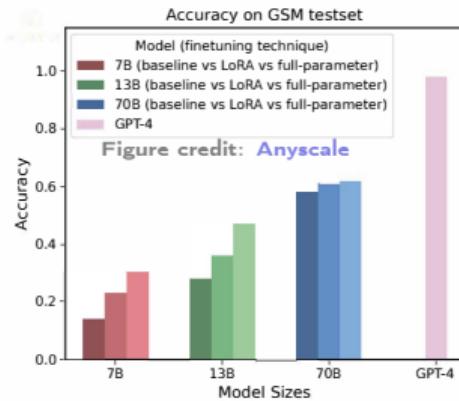


From pre-training to (parameter-efficient) fine-tuning



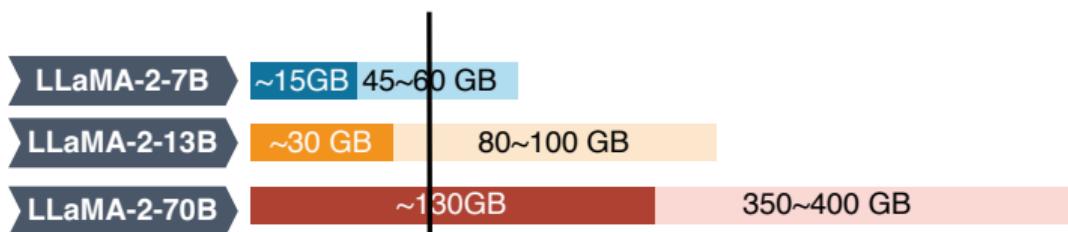
Full fine-tuning vs. parameter-efficient fine-tuning

From pre-training to (parameter-efficient) fine-tuning

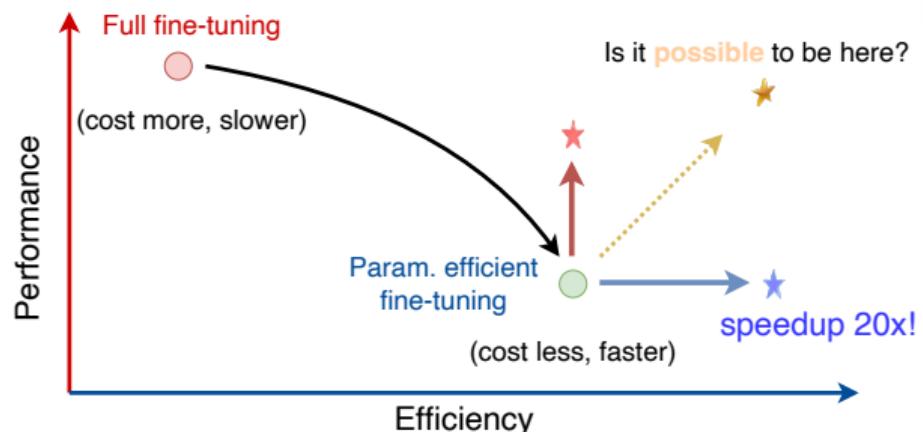
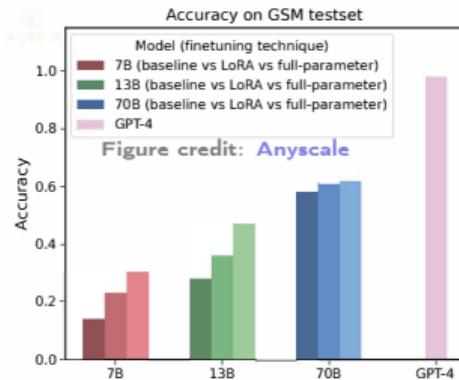


Full fine-tuning vs. parameter-efficient fine-tuning

A100 40GB GPU



From pre-training to (parameter-efficient) fine-tuning

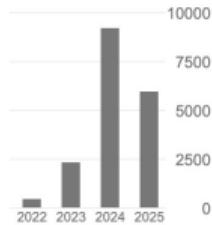


LoRA: Low-rank adaption

Published as a conference paper at ICLR 2022

LORA: LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS

Edward Hu* Yelong Shen* Phillip Wallis Zeyuan Allen-Zhu
Yuanzhi Li Shean Wang Lu Wang Weizhu Chen
Microsoft Corporation
edward.hu@mila.quebec
{yeshe, phwallis, zeyuana, swang, luw, wzchen}@microsoft.com
yuanzhil@andrew.cmu.edu

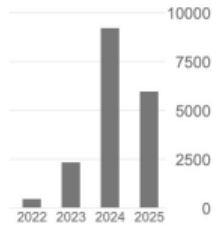


LoRA: Low-rank adaption

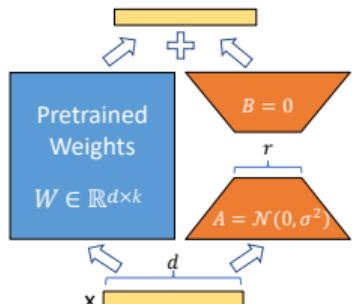
Published as a conference paper at ICLR 2022

LORA: LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS

Edward Hu* Yelong Shen* Phillip Wallis Zeyuan Allen-Zhu
Yuanzhi Li Shean Wang Lu Wang Weizhu Chen
Microsoft Corporation
edward.hu@mila.quebec
{yeshe, phwallis, zeyuana, swang, luw, wzchen}@microsoft.com
yuanzhil@andrew.cmu.edu



$$\mathbf{W}^{\text{FT}} = \mathbf{W}^{\text{pre}} + \Delta \in \mathbb{R}^{d \times k}$$



- Formulation:

$$\Delta \approx \mathbf{AB} \text{ with } \mathbf{A} \in \mathbb{R}^{d \times r} \text{ and } \mathbf{B} \in \mathbb{R}^{r \times k}$$

- Initialization:

$$[\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2) \quad \text{and} \quad [\mathbf{B}_0]_{ij} = 0. \quad (\text{LoRA-init.})$$



How can theory guide practice

- understanding: training dynamics of $(\mathbf{A}_t, \mathbf{B}_t)$
- design theory-grounded algorithm (LoRA-One)
- clarify some misconceptions in previous algorithm designs

□ Even for linear model (pre-training and fine-tuning), **nonlinear dynamics...**

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G} \\ \eta \mathbf{G}^\top & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term} \quad \begin{cases} [\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2) \\ [\mathbf{B}_0]_{ij} = 0. \end{cases}$$

□ One-step full gradient: $\mathbf{G} \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}) = r^*$

$$\mathbf{G} := -\nabla_{\mathbf{W}} L(\mathbf{W}^{\text{pre}}) = \frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \Delta.$$



How can theory guide practice

- understanding: training dynamics of $(\mathbf{A}_t, \mathbf{B}_t)$
- design theory-grounded algorithm (LoRA-One)
- clarify some misconceptions in previous algorithm designs

□ Even for linear model (pre-training and fine-tuning), **nonlinear dynamics...**

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G} \\ \eta \mathbf{G}^\top & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term} \quad \begin{cases} [\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2) \\ [\mathbf{B}_0]_{ij} = 0. \end{cases}$$

□ One-step full gradient: $\mathbf{G} \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}) = r^*$

$$\mathbf{G} := -\nabla_{\mathbf{W}} L(\mathbf{W}^{\text{pre}}) = \frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \Delta.$$



How can theory guide practice

- understanding: training dynamics of $(\mathbf{A}_t, \mathbf{B}_t)$
- design theory-grounded algorithm (LoRA-One)
- clarify some misconceptions in previous algorithm designs

□ Even for linear model (pre-training and fine-tuning), **nonlinear dynamics...**

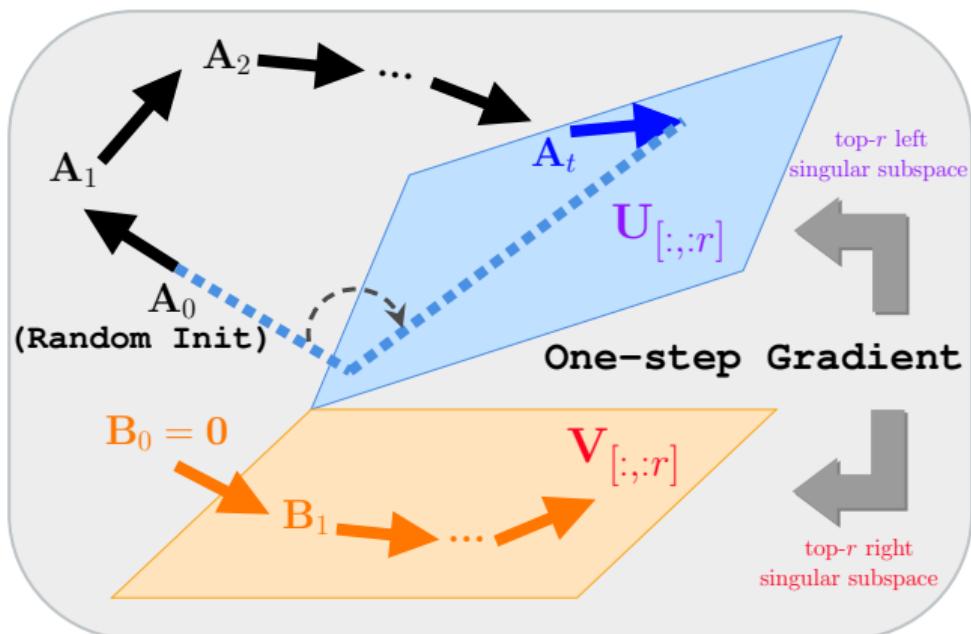
$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G} \\ \eta \mathbf{G}^\top & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term} \quad \begin{cases} [\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2) \\ [\mathbf{B}_0]_{ij} = 0. \end{cases}$$

□ One-step full gradient: $\mathbf{G} \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}) = r^*$

$$\mathbf{G} := -\nabla_{\mathbf{W}} L(\mathbf{W}^{\text{pre}}) = \frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \Delta.$$

Alignment and theory-grounded algorithm

Pipeline



Problem setting and assumptions

- (Downstream) data : with unknown low-rank feature shift $\text{rank}(\Delta) = r^*$

$$\tilde{\mathbf{y}} := \begin{cases} (\mathbf{W}^\natural + \Delta)^\top \tilde{\mathbf{x}}, & \tilde{\mathbf{x}} \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, \quad \text{linear} \\ \sigma[(\mathbf{W}^\natural + \Delta)^\top \tilde{\mathbf{x}}], & \tilde{\mathbf{x}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) \quad \text{nonlinear} \end{cases}.$$

Problem setting and assumptions

- (Downstream) data : with unknown low-rank feature shift $\text{rank}(\Delta) = r^*$

$$\tilde{\mathbf{y}} := \begin{cases} (\mathbf{W}^\natural + \Delta)^\top \tilde{\mathbf{x}}, & \tilde{\mathbf{x}} \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, \quad \text{linear} \\ \sigma[(\mathbf{W}^\natural + \Delta)^\top \tilde{\mathbf{x}}], & \tilde{\mathbf{x}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) \quad \text{nonlinear} \end{cases}.$$

- Model : LoRA starting from **known** pre-trained $\mathbf{W}^\natural \in \mathbb{R}^{d \times k}$

$$f(\mathbf{x}; \mathbf{A}, \mathbf{B}) := \begin{cases} (\mathbf{W}^\natural + \mathbf{AB})^\top \tilde{\mathbf{x}} & \text{linear} \\ \sigma[(\mathbf{W}^\natural + \mathbf{AB})^\top \tilde{\mathbf{x}}] & \text{nonlinear} \end{cases}.$$

Problem setting and assumptions

- (Downstream) data : with unknown low-rank feature shift $\text{rank}(\Delta) = r^*$

$$\tilde{\mathbf{y}} := \begin{cases} (\mathbf{W}^\natural + \Delta)^\top \tilde{\mathbf{x}}, & \tilde{\mathbf{x}} \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, \quad \text{linear} \\ \sigma[(\mathbf{W}^\natural + \Delta)^\top \tilde{\mathbf{x}}], & \tilde{\mathbf{x}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) \quad \text{nonlinear} \end{cases}.$$

- Model : LoRA starting from **known** pre-trained $\mathbf{W}^\natural \in \mathbb{R}^{d \times k}$

$$f(\mathbf{x}; \mathbf{A}, \mathbf{B}) := \begin{cases} (\mathbf{W}^\natural + \mathbf{AB})^\top \tilde{\mathbf{x}} & \text{linear} \\ \sigma[(\mathbf{W}^\natural + \mathbf{AB})^\top \tilde{\mathbf{x}}] & \text{nonlinear} \end{cases}.$$

- Algorithm : Squared loss, gradient descent

$$\begin{cases} L(\mathbf{W}) = \frac{1}{2N} \sum_{i=1}^N \|f(\mathbf{x}_i; \mathbf{W}) - \tilde{\mathbf{y}}_i\|_2^2, & \text{Full fine-tuning} \\ L(\mathbf{A}, \mathbf{B}) = \frac{1}{2N} \sum_{i=1}^N \|f(\mathbf{x}_i; \mathbf{A}, \mathbf{B}) - \tilde{\mathbf{y}}_i\|_2^2, & \text{LoRA} \end{cases}.$$

Problem setting and assumptions

- (Downstream) data : with unknown low-rank feature shift $\text{rank}(\Delta) = r^*$

$$\tilde{\mathbf{y}} := \begin{cases} (\mathbf{W}^\natural + \Delta)^\top \tilde{\mathbf{x}}, & \tilde{\mathbf{x}} \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, \quad \text{linear} \\ \sigma[(\mathbf{W}^\natural + \Delta)^\top \tilde{\mathbf{x}}], & \tilde{\mathbf{x}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) \quad \text{nonlinear} \end{cases}.$$

- Model : LoRA starting from **known** pre-trained $\mathbf{W}^\natural \in \mathbb{R}^{d \times k}$

$$f(\mathbf{x}; \mathbf{A}, \mathbf{B}) := \begin{cases} (\mathbf{W}^\natural + \mathbf{AB})^\top \tilde{\mathbf{x}} & \text{linear} \\ \sigma[(\mathbf{W}^\natural + \mathbf{AB})^\top \tilde{\mathbf{x}}] & \text{nonlinear} \end{cases}.$$

- Algorithm : Squared loss, gradient descent

$$\begin{cases} L(\mathbf{W}) = \frac{1}{2N} \sum_{i=1}^N \|f(\mathbf{x}_i; \mathbf{W}) - \tilde{\mathbf{y}}_i\|_2^2, & \text{Full fine-tuning} \\ L(\mathbf{A}, \mathbf{B}) = \frac{1}{2N} \sum_{i=1}^N \|f(\mathbf{x}_i; \mathbf{A}, \mathbf{B}) - \tilde{\mathbf{y}}_i\|_2^2, & \text{LoRA} \end{cases}.$$

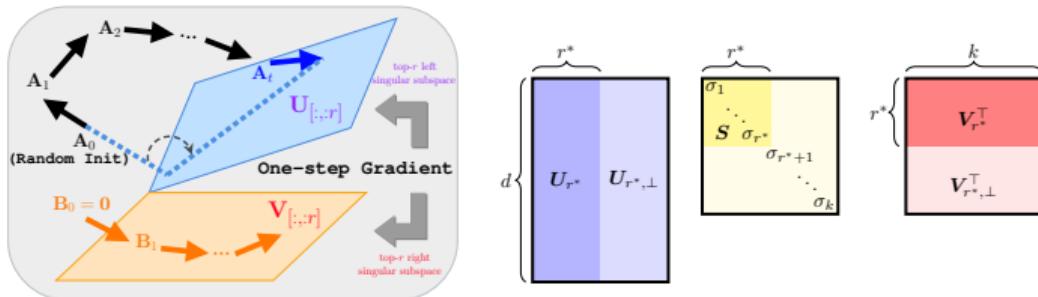
- Evaluate $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}$: optimization and generalization!

$$\mathbb{E}_{\tilde{\mathbf{x}}} \left\| \tilde{\mathbf{y}} - \sigma(\mathbf{W}^\natural + \mathbf{A}_t \mathbf{B}_t)^\top \tilde{\mathbf{x}} \right\|_2^2 \lesssim \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}^2$$

Our results: Alignment on B_t

- one-step full gradient: $\mathbf{G} \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}) = r^*$

$$\mathbf{G} := -\nabla_{\mathbf{W}} L(\mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \Delta.$$



Theorem (Alignment between G and B_t)

For the linear setting, LoRA trained by gradient descent admits

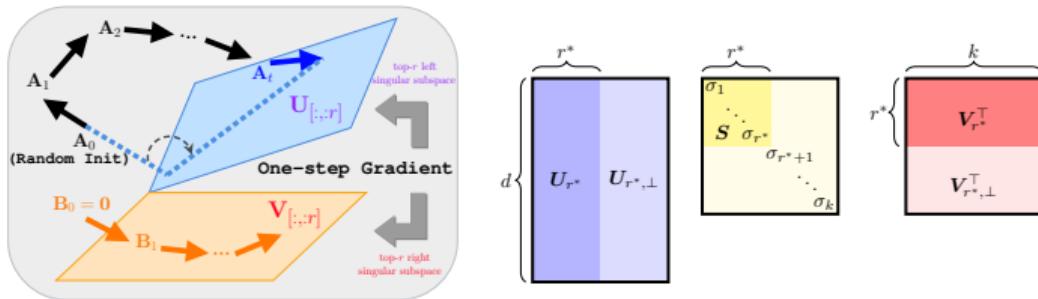
$$\angle(\mathbf{V}_{r^*}(\mathbf{B}_t), \mathbf{V}_{r^*}(\mathbf{G})) = 0, \quad \forall t \in \mathbb{N}_+.$$

Remark: $\mathbf{B}_1 = \eta_1 \mathbf{A}_0^\top \mathbf{G}$ with $\text{Rank}(\mathbf{B}_1) \leq r^*$

Our results: Alignment on B_t

□ one-step full gradient: $\mathbf{G} \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}) = r^*$

$$\mathbf{G} := -\nabla_{\mathbf{W}} L(\mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \Delta.$$



Theorem (Alignment between \mathbf{G} and B_t)

For the linear setting, LoRA trained by gradient descent admits

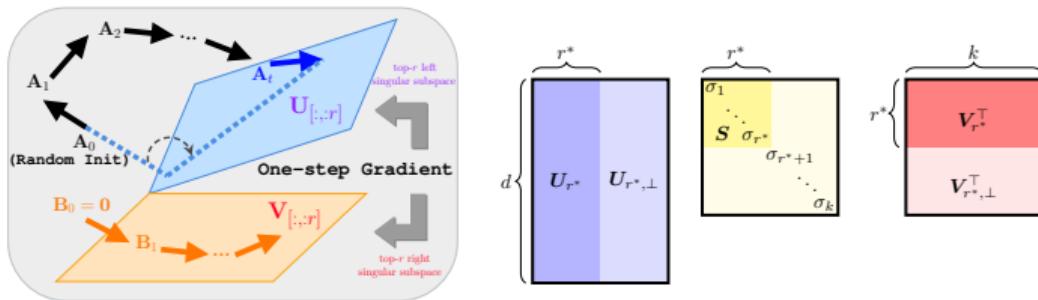
$$\angle(\mathbf{V}_{r^*}(B_t), \mathbf{V}_{r^*}(\mathbf{G})) = 0, \quad \forall t \in \mathbb{N}_+.$$

Remark: $B_1 = \eta_1 \mathbf{A}_0^\top \mathbf{G}$ with $\text{Rank}(B_1) \leq r^*$

Our results: Alignment on B_t

□ one-step full gradient: $\mathbf{G} \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}) = r^*$

$$\mathbf{G} := -\nabla_{\mathbf{W}} L(\mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \Delta.$$



Theorem (Alignment between \mathbf{G} and B_t)

For the linear setting, LoRA trained by gradient descent admits

$$\angle(\mathbf{V}_{r^*}(B_t), \mathbf{V}_{r^*}(\mathbf{G})) = 0, \quad \forall t \in \mathbb{N}_+.$$

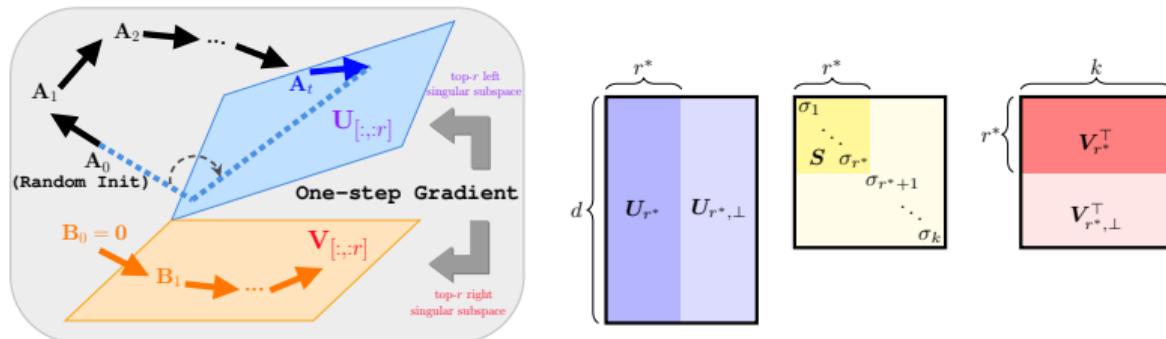
Remark: $B_1 = \eta_1 \mathbf{A}_0^\top \mathbf{G}$ with $\text{Rank}(B_1) \leq r^*$

Our results: Alignment on A_t

Theorem (Informal, LoRA initialization)

For $r \geq r^*$, $[A_0]_{ij} \sim \mathcal{N}(0, \alpha^2)$, for any $\epsilon \in (0, 1)$, choosing $\alpha = \mathcal{O}(\epsilon d^{-\frac{3}{4}\kappa^\natural - \frac{1}{2}})$, running GD with $t^* = \Theta(\ln d)$ steps, then we have

$$\angle(\mathbf{U}_{r^*}(A_{t^*}), \mathbf{U}_{r^*}(\mathbf{G})) \lesssim \epsilon, \text{ w.h.p.}$$



Our results: Alignment on A_t

Theorem (Informal, LoRA initialization)

For $r \geq r^*$, $[A_0]_{ij} \sim \mathcal{N}(0, \alpha^2)$, for any $\epsilon \in (0, 1)$, choosing $\alpha = \mathcal{O}(\epsilon d^{-\frac{3}{4}\kappa^\natural - \frac{1}{2}})$, running GD with $t^* = \Theta(\ln d)$ steps, then we have

$$\angle(\mathbf{U}_{r^*}(A_{t^*}), \mathbf{U}_{r^*}(\mathbf{G})) \lesssim \epsilon, \text{ w.h.p.}$$

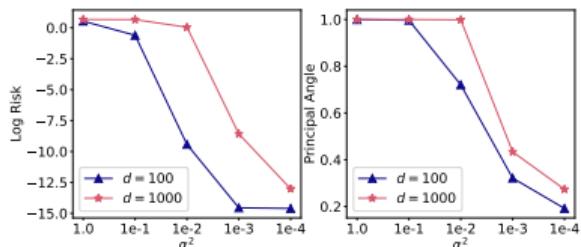


Figure 2: Left: the risk $\frac{1}{2} \|A_t B_t - \Delta\|_F^2$. Right: the principal angle is $\min_t \|\mathbf{U}_{r^*, \perp}(\mathbf{G}) \mathbf{U}_{r^*}(A_t)\|_{op}$.

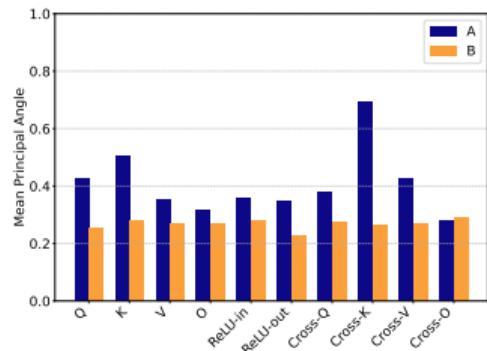


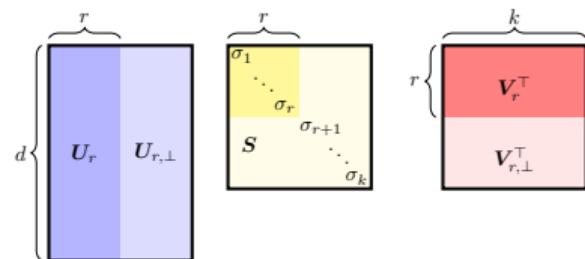
Figure 3: Principal angle of fine-tuning T5 on MRPC.

Algorithm design principle

□ SVD: $\mathbf{G} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$

$$\mathbf{A}_0 = \mathbf{U}_{[:,1:r]} \mathbf{S}_{[1:r]}^{\frac{1}{2}}. \quad (\text{Spec-init.})$$

$$\mathbf{B}_0 = \mathbf{S}_{[1:r]}^{\frac{1}{2}} \mathbf{V}_{[:,1:r]}^\top.$$

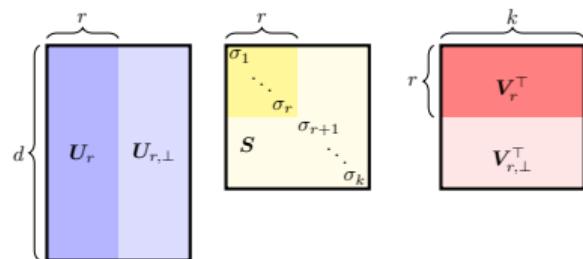


Algorithm design principle

□ SVD: $\mathbf{G} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$

$$\mathbf{A}_0 = \mathbf{U}_{[:,1:r]} \mathbf{S}_{[1:r]}^{\frac{1}{2}} . \quad (\text{Spec-init.})$$

$$\mathbf{B}_0 = \mathbf{S}_{[1:r]}^{\frac{1}{2}} \mathbf{V}_{[:,1:r]}^\top .$$



Key Message: we can “escape” the alignment stage

Under (Spec-init.), for both linear/nonlinear models, we can directly achieve the alignment at initialization.

$$\|\mathbf{A}_0\mathbf{B}_0 - \Delta\|_F \text{ is small, w.h.p.}$$

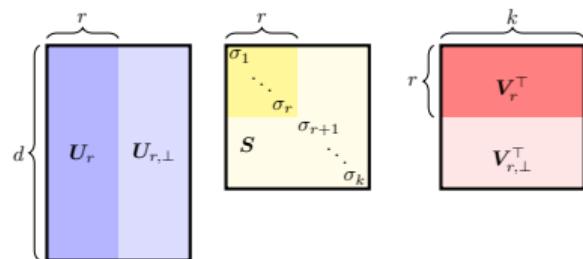
Algorithm design principle

□ SVD: $\mathbf{G} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$

$$\mathbf{A}_0 = \mathbf{U}_{[:,1:r]} \mathbf{S}_{[1:r]}^{\frac{1}{2}}.$$

(Spec-init.)

$$\mathbf{B}_0 = \mathbf{S}_{[1:r]}^{\frac{1}{2}} \mathbf{V}_{[:,1:r]}^\top.$$



Key Message: we can “escape” the alignment stage

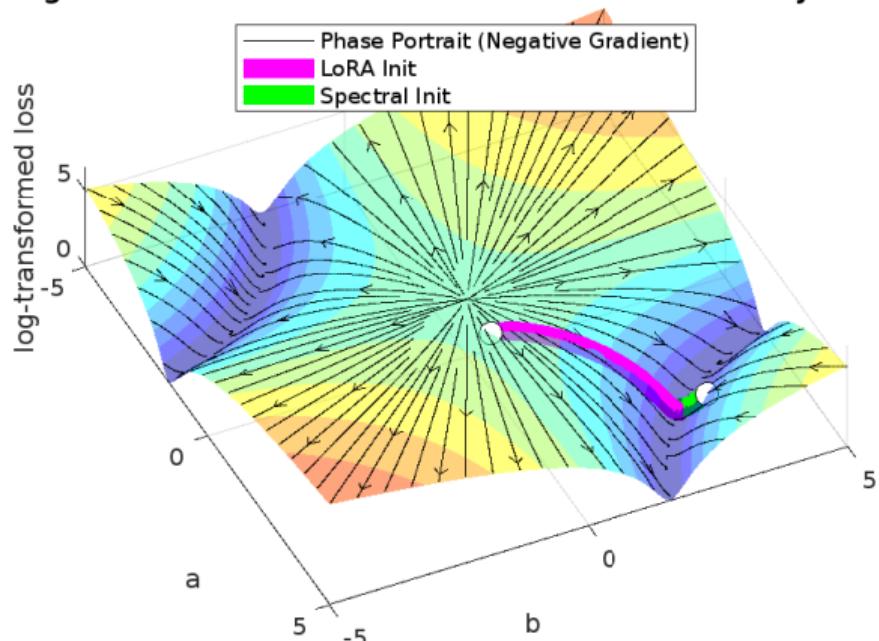
Under (Spec-init.), for both linear/nonlinear models, we can directly achieve the alignment at initialization.

$$\|\mathbf{A}_0\mathbf{B}_0 - \Delta\|_F \text{ is small, w.h.p.}$$

The “best” initialization strategy!

“Best” initialization: phase portrait

Log-Transformed Surface with Phase Portrait and Trajectories



Toy example (I)

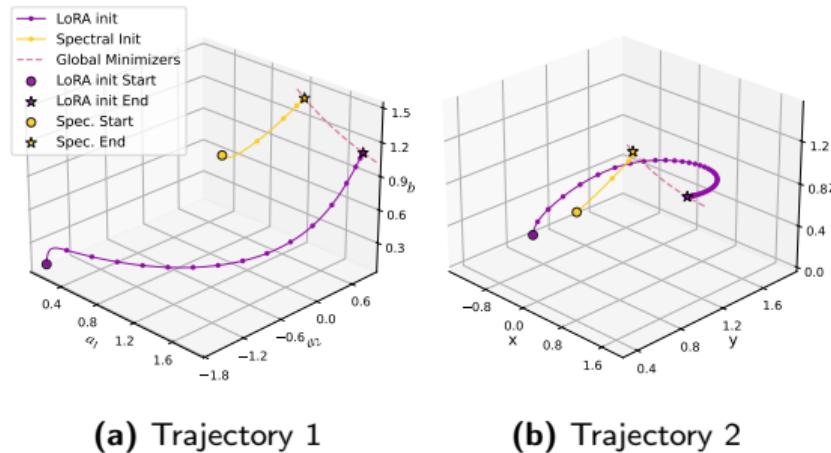


Figure 4: Comparison of the GD trajectories between LoRA and ours. (a) and (b): $\mathbf{A} \in \mathbb{R}^2$ and $B \in \mathbb{R}$ with different initializations. The set of global minimizers is $\{a_1^* = 2/t, a_2^* = 1/t, b^* = t \mid t \in \mathbb{R}\}$.

Toy example (II)

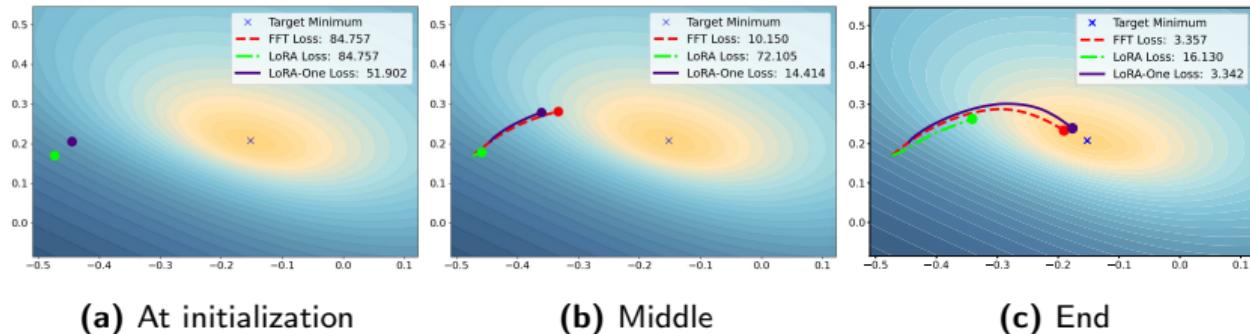


Figure 5: Comparison of the GD trajectories between LoRA and ours. We use two-layer neural networks pre-trained on odd-labeled data and fine-tuned on even-labeled data on MNIST, see [GIF illustration](#).

One-step gradient can suffice on small-scale datasets!

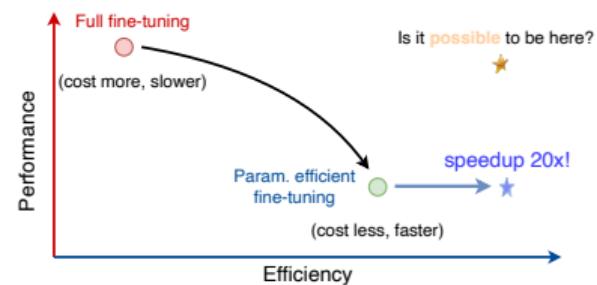
Dataset	MNLI	SST-2	CoLA	QNLI	MRPC
Size	393k	67k	8.5k	105k	3.7k
Pre-trained	-	89.79	59.03	49.28	63.48
Spectral init.	-	90.48	73.00	76.64	68.38
LoRA ₈	85.30 _{±0.04}	94.04 _{±0.09}	72.84 _{±1.25}	93.02 _{±0.07}	68.38 _{±0.01}

One-step gradient can suffice on small-scale datasets!

Dataset	MNLI	SST-2	CoLA	QNLI	MRPC
Size	393k	67k	8.5k	105k	3.7k
Pre-trained	-	89.79	59.03	49.28	63.48
Spectral init.	-	90.48	73.00	76.64	68.38
LoRA ₈	85.30 ± 0.04	94.04 ± 0.09	72.84 ± 1.25	93.02 ± 0.07	68.38 ± 0.01

Time cost (sec.)	LoRA	Spectral init.
CoLA	47s	<1s
MRPC	25s	<1s

memory-efficient [1] + randomized SVD + parallel initialization



Key features in our LoRA-One algorithm

Algorithm 1 LoRA-One training for a specific layer

Input: Pre-trained weight \mathbf{W}^\natural , batched data $\{\mathcal{D}_t\}_{t=1}^T$, LoRA rank r , LoRA alpha α , loss function L

Output: $\mathbf{W}^\natural + \frac{\alpha}{\sqrt{r}} \mathbf{A}_T \mathbf{B}_T$

Compute $\nabla_{\mathbf{W}} L(\mathbf{W}^\natural)$ and $\mathbf{U}, \mathbf{S}, \mathbf{V} \leftarrow \text{SVD}(\nabla_{\mathbf{W}} L(\mathbf{W}^\natural))$

$$\mathbf{A}_0 \leftarrow \sqrt{\gamma} \cdot \mathbf{U}_{[:,1:r]} \mathbf{S}_{[:,r,:r]}^{1/2}$$

$$\mathbf{B}_0 \leftarrow \sqrt{\gamma} \cdot \mathbf{S}_{[:,r,:r]}^{1/2} \mathbf{V}_{[:,1:r]}^\top$$

Clear $\nabla_{\mathbf{W}} L(\mathbf{W}^\natural)$

for $t = 1, \dots, T$ **do**

$$\mathbf{G}_t^A \leftarrow \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_{t-1}, \mathbf{B}_{t-1}) \left(\mathbf{B}_{t-1} \mathbf{B}_{t-1}^\top + \lambda \mathbf{I}_r \right)^{-1}$$

$$\mathbf{G}_t^B \leftarrow \left(\mathbf{A}_{t-1}^\top \mathbf{A}_{t-1} + \lambda \mathbf{I}_r \right)^{-1} \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_{t-1}, \mathbf{B}_{t-1})$$

$$\text{Update } \mathbf{A}_t, \mathbf{B}_t \leftarrow \text{AdamW}(\mathbf{G}_t^A, \mathbf{G}_t^B)$$

end

Results on LLaMA 2-7B (for one epoch)

(r = 8)	GSM8K		MMLU Avg.	HumanEval PASS@1
	Direct	8s-CoT		
LoRA	59.26 \pm 0.76	53.36 \pm 0.77	45.73 \pm 0.30	25.85 \pm 1.75
LoRA-GA	56.44 \pm 1.37	46.07 \pm 1.01	45.70 \pm 0.77	26.95 \pm 1.30
LoRA-One	60.44 \pm 0.17	55.88 \pm 0.44	47.12 \pm 0.12	28.66 \pm 0.39

- One epoch, rank 8, three runs
- Hyperparameter optimized over learning rate, batch size
- Train: 100k subset from MetaMathQA
- Test: GSM8K, Accuracy (%)

Results on LLaMA 2-7B (for one epoch)

$(r = 8)$	GSM8K		MMLU	HumanEval
	Direct	8s-CoT	Avg.	PASS@1
LoRA	59.26 ± 0.76	53.36 ± 0.77	45.73 ± 0.30	25.85 ± 1.75
LoRA-GA	56.44 ± 1.37	46.07 ± 1.01	45.70 ± 0.77	26.95 ± 1.30
LoRA-One	60.44 ± 0.17	55.88 ± 0.44	47.12 ± 0.12	28.66 ± 0.39

- One epoch, rank 8, three runs
- Hyperparameter optimized over learning rate, batch size
- Train: 100k subset from MetaMathQA
- Test: GSM8K, Accuracy (%)



Results on LLaMA 2-7B (for one epoch)

$(r = 8)$	GSM8K		MMLU	HumanEval
	Direct	8s-CoT	Avg.	PASS@1
LoRA	59.26 ± 0.76	53.36 ± 0.77	45.73 ± 0.30	25.85 ± 1.75
LoRA-GA	56.44 ± 1.37	46.07 ± 1.01	45.70 ± 0.77	26.95 ± 1.30
LoRA-One	60.44 ± 0.17	55.88 ± 0.44	47.12 ± 0.12	28.66 ± 0.39

- One epoch, rank 8, three runs
- Hyperparameter optimized over learning rate, batch size
- Train: 100k subset from Code-Feedback
- Test: Humaneval, Pass@1

Results on LLaMA 2-7B (for one epoch)

$(r = 8)$	GSM8K		MMLU	HumanEval
	Direct	8s-CoT	Avg.	PASS@1
LoRA	59.26 ± 0.76	53.36 ± 0.77	45.73 ± 0.30	25.85 ± 1.75
LoRA-GA	56.44 ± 1.37	46.07 ± 1.01	45.70 ± 0.77	26.95 ± 1.30
LoRA-One	60.44 ± 0.17	55.88 ± 0.44	47.12 ± 0.12	28.66 ± 0.39

- One epoch, rank 8, three runs
- Hyperparameter optimized over learning rate, batch size
- Train: 100k subset from Code-Foodback
- Test: Humaneval, Pass@1

Time cost

LoRA: 6h 24min

+ 2min

Memory

LoRA: 22.6 GB

- 1.1GB

Results on LLaMA 2-7B (for more epochs)

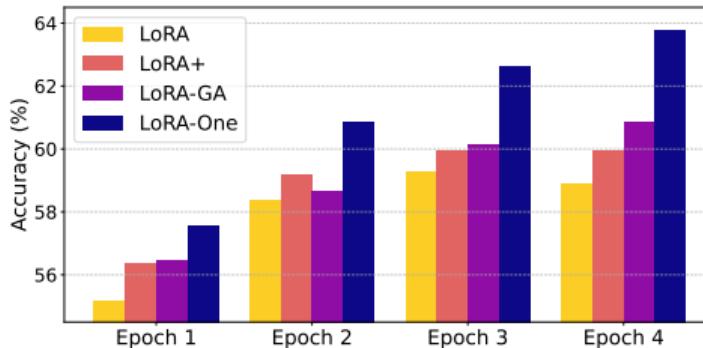


Figure 6: Accuracy comparison across different methods over epochs on GSM8K.

Results on LLaMA 2-7B (for more epochs)

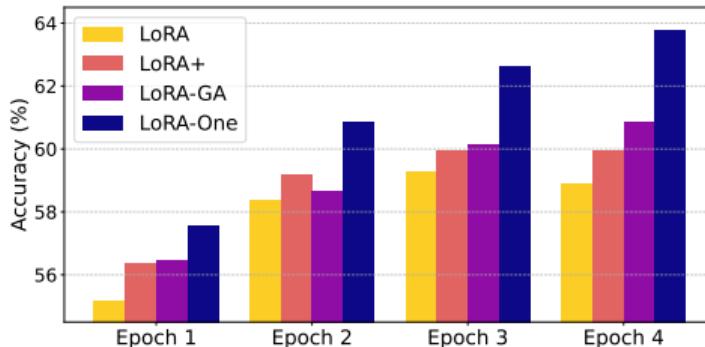
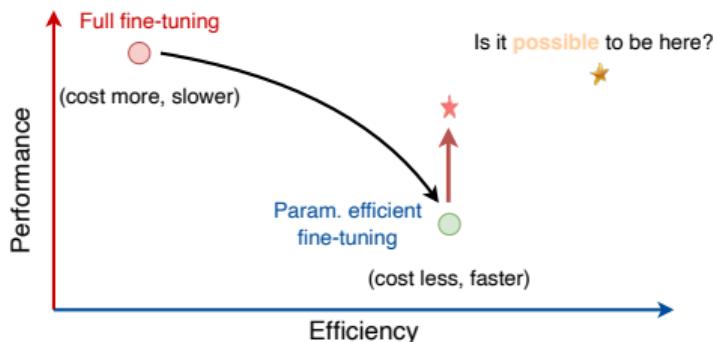


Figure 6: Accuracy comparison across different methods over epochs on GSM8K.



Clarification on gradient alignment based work

LoRA-GA ([Wang et al, 2024](#)): make LoRA's gradients align to full fine-tuning!

Clarification on gradient alignment based work

LoRA-GA ([Wang et al, 2024](#)): make LoRA's gradients align to full fine-tuning!

- ☐ best $2r$ approximation

$$\begin{cases} \text{rank}(\nabla_{\mathbf{A}} L(\mathbf{A}_t, \mathbf{B}_t)) \leq r \\ \text{rank}(\nabla_{\mathbf{B}} L(\mathbf{A}_t, \mathbf{B}_t)) \leq r \end{cases}$$

Clarification on gradient alignment based work

LoRA-GA ([Wang et al, 2024](#)): make LoRA's gradients align to full fine-tuning!

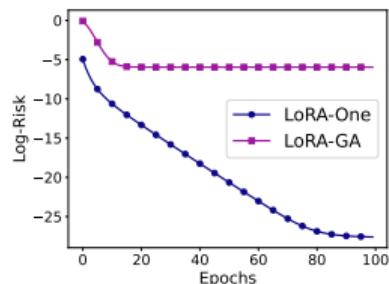
□ best $2r$ approximation

$$\begin{cases} \text{rank}(\nabla_{\mathbf{A}} L(\mathbf{A}_t, \mathbf{B}_t)) \leq r \\ \text{rank}(\nabla_{\mathbf{B}} L(\mathbf{A}_t, \mathbf{B}_t)) \leq r \end{cases}$$

Method	Init. on \mathbf{A}	Init. on \mathbf{B}	Calibration
LoRA	$\mathcal{N}(0, \alpha^2)$	0	-
LoRA-GA	$\mathbf{U}_{[:,1:r]}$	$\mathbf{V}_{[:,r+1:2r]}^\top$	$\mathbf{W}^{\text{pre}} - \mathbf{A}_0 \mathbf{B}_0$
LoRA-One	$\mathbf{U}_{[:,1:r]} \mathbf{S}_{[1:r]}^{1/2}$	$\mathbf{S}_{[1:r]}^{1/2} \mathbf{V}_{[:,1:r]}^\top$	-

Clarification on gradient alignment based work

LoRA-GA ([Wang et al, 2024](#)): make LoRA's gradients align to full fine-tuning!



Method	Init. on A	Init. on B	Calibration
LoRA	$\mathcal{N}(0, \alpha^2)$	0	-
LoRA-GA	$U_{[:,1:r]}$	$V_{[:,r+1:2r]}^\top$	$W^{\text{pre}} - A_0 B_0$
LoRA-One	$U_{[:,1:r]} S_{[1:r]}^{1/2}$	$S_{[1:r]}^{1/2} V_{[:,1:r]}^\top$	-

Theory and proof...

Model	Algorithm	Initialization	Results
Linear	GD	(LoRA-init.)	Subspace alignment of \mathbf{B}_t
	GD	(LoRA-init.)	Subspace alignment of \mathbf{A}_t
	GD	(Spec-init.)	$\ \mathbf{A}_0 \mathbf{B}_0 - \Delta\ _F$ is small
	GD	(Spec-init.)	Linear convergence of $\ \mathbf{A}_t \mathbf{B}_t - \Delta\ _F$
	Precondition GD	(Spec-init.)	Linear convergence rate independent of $\kappa(\Delta)$
Nonlinear	GD	(Spec-init.)	$\ \mathbf{A}_0 \mathbf{B}_0 - \Delta\ _F$ is small
	GD	(Spec-init.)	Linear convergence of $\ \mathbf{A}_t \mathbf{B}_t - \Delta\ _F$
	Precondition GD	(Spec-init.)	Linear convergence rate independent of $\kappa(\Delta)$

- subspace alignment
- global convergence

Proof of sketch: Control the dynamics for alignment

$$\underbrace{\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix}}_{:= \mathbf{Z}_{t+1}} = \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G} \\ \eta \mathbf{G}^\top & \mathbf{I}_k \end{bmatrix}}_{:= \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}}_{:= \mathbf{Z}_t} - \frac{1}{N} \begin{bmatrix} 0 & \eta \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \eta \mathbf{B}_t^\top \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}.$$

□ Approximated linear dynamical system $\mathbf{Z}_t^{\text{lin}} := \mathbf{H}^t \mathbf{Z}_0$

- Schur decomposition of \mathbf{H}
- obtain the dynamics of $\mathbf{Z}_t^{\text{lin}}$
- Define $\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}}$, control $\|\mathbf{E}_t\|_{op}$ in $\mathcal{O}(\log d)$ steps

□ Transfer the alignment from $\mathbf{A}_t^{\text{lin}}$ to \mathbf{A}_t [2] (Stöger & Soltanolkotabi, 2021)

$\|\mathbf{U}_{r^*, \perp}^\top (\mathbf{G}) \mathbf{U}_{r^*} (\mathbf{A}_{t^*})\|_{op} \lesssim \|\mathbf{U}_{r^*, \perp}^\top (\mathbf{P}_t^A) \mathbf{U}_{r^*} (\mathbf{P}_t^A \mathbf{A}_0 + \mathbf{E}_t)\|_{op}$ is small, w.h.p.

Proof of sketch: Control the dynamics for alignment

$$\underbrace{\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix}}_{:= \mathbf{Z}_{t+1}} = \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G} \\ \eta \mathbf{G}^\top & \mathbf{I}_k \end{bmatrix}}_{:= \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}}_{:= \mathbf{Z}_t} - \frac{1}{N} \begin{bmatrix} 0 & \eta \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \eta \mathbf{B}_t^\top \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}.$$

□ Approximated linear dynamical system $\mathbf{Z}_t^{\text{lin}} := \mathbf{H}^t \mathbf{Z}_0$

- Schur decomposition of \mathbf{H}
- obtain the dynamics of $\mathbf{Z}_t^{\text{lin}}$
- Define $\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}}$, control $\|\mathbf{E}_t\|_{op}$ in $\mathcal{O}(\log d)$ steps

□ Transfer the alignment from $\mathbf{A}_t^{\text{lin}}$ to \mathbf{A}_t [2] (Stöger & Soltanolkotabi, 2021)

$\|\mathbf{U}_{r^*, \perp}^\top(\mathbf{G}) \mathbf{U}_{r^*}(\mathbf{A}_{t^*})\|_{op} \lesssim \|\mathbf{U}_{r^*, \perp}^\top(\mathbf{P}_t^A) \mathbf{U}_{r^*}(\mathbf{P}_t^A \mathbf{A}_0 + \mathbf{E}_t)\|_{op}$ is small, w.h.p.

Proof of sketch: Control the dynamics for alignment

$$\underbrace{\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix}}_{:= \mathbf{Z}_{t+1}} = \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G} \\ \eta \mathbf{G}^\top & \mathbf{I}_k \end{bmatrix}}_{:= \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}}_{:= \mathbf{Z}_t} - \frac{1}{N} \begin{bmatrix} 0 & \eta \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \eta \mathbf{B}_t^\top \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}.$$

□ Approximated linear dynamical system $\mathbf{Z}_t^{\text{lin}} := \mathbf{H}^t \mathbf{Z}_0$

- Schur decomposition of \mathbf{H}
- obtain the dynamics of $\mathbf{Z}_t^{\text{lin}}$
- Define $\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}}$, control $\|\mathbf{E}_t\|_{op}$ in $\mathcal{O}(\log d)$ steps

□ Transfer the alignment from $\mathbf{A}_t^{\text{lin}}$ to \mathbf{A}_t [2] (Stöger & Soltanolkotabi, 2021)

$\|\mathbf{U}_{r^*, \perp}^\top(\mathbf{G}) \mathbf{U}_{r^*}(\mathbf{A}_{t^*})\|_{op} \lesssim \|\mathbf{U}_{r^*, \perp}^\top(\mathbf{P}_t^A) \mathbf{U}_{r^*}(\mathbf{P}_t^A \mathbf{A}_0 + \mathbf{E}_t)\|_{op}$ is small, w.h.p.

Global convergence of nonlinear models

Theorem (Informal, linear convergence rate)

For nonlinear model with $r = r^*$ and gradient descent (with preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), \text{ w.h.p}$$

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{op} \leq \|\mathbf{A}_0 \mathbf{B}_0 - 2\mathbf{G}\|_{op} + 2\|\mathbf{G} - \mathbb{E}_{\tilde{x}}[\mathbf{G}]\|_{op} + \|2\mathbb{E}_{\tilde{x}}[\mathbf{G}] - \Delta\|_{op}$$

$$\mathbf{J}_{\mathbf{W}_t} := \frac{1}{N} \tilde{\mathbf{X}}^\top \left[\sigma(\tilde{\mathbf{X}} \tilde{\mathbf{W}}^\natural) - \frac{1}{N} \tilde{\mathbf{X}}^\top \sigma(\tilde{\mathbf{X}} \mathbf{W}_t) \right] \odot \sigma'(\tilde{\mathbf{X}} \mathbf{W}_t).$$

- low-rank approximation error $\leq 2\lambda_{r^*+1}(\mathbf{G})$
- population error: using $\mathbb{E}_{\tilde{x}}[-\mathbf{J}_{\mathbf{W}_t}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta) + \mathcal{O}(\frac{1}{\kappa r^*})$
- concentration error

$$\left\| \mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{x}}[\mathbf{J}_{\mathbf{W}_t}] \right\|_{\text{F}} \lesssim \sqrt{d}\epsilon \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, \text{ w.h.p.}$$

Global convergence of nonlinear models

Theorem (Informal, linear convergence rate)

For nonlinear model with $r = r^*$ and gradient descent (with preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), \text{ w.h.p}$$

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{op} \leq \|\mathbf{A}_0 \mathbf{B}_0 - 2\mathbf{G}\|_{op} + 2\|\mathbf{G} - \mathbb{E}_{\tilde{x}}[\mathbf{G}]\|_{op} + \|2\mathbb{E}_{\tilde{x}}[\mathbf{G}] - \Delta\|_{op}$$

$$\mathbf{J}_{W_t} := \frac{1}{N} \tilde{\mathbf{X}}^\top \left[\sigma(\tilde{\mathbf{X}} \tilde{\mathbf{W}}^\natural) - \frac{1}{N} \tilde{\mathbf{X}}^\top \sigma(\tilde{\mathbf{X}} W_t) \right] \odot \sigma'(\tilde{\mathbf{X}} W_t).$$

- low-rank approximation error $\leq 2\lambda_{r^*+1}(\mathbf{G})$
- population error: using $\mathbb{E}_{\tilde{x}}[-\mathbf{J}_{W_t}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta) + \mathcal{O}(\frac{1}{\kappa r^*})$
- concentration error

$$\left\| \mathbf{J}_{W_t} - \mathbb{E}_{\tilde{x}}[\mathbf{J}_{W_t}] \right\|_{\text{F}} \lesssim \sqrt{d\epsilon} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, \text{ w.h.p.}$$

Global convergence of nonlinear models

Theorem (Informal, linear convergence rate)

For nonlinear model with $r = r^*$ and gradient descent (with preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), \text{ w.h.p}$$

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{op} \leq \|\mathbf{A}_0 \mathbf{B}_0 - 2\mathbf{G}\|_{op} + 2\|\mathbf{G} - \mathbb{E}_{\tilde{x}}[\mathbf{G}]\|_{op} + \|2\mathbb{E}_{\tilde{x}}[\mathbf{G}] - \Delta\|_{op}$$

$$\mathbf{J}_{\mathbf{W}_t} := \frac{1}{N} \tilde{\mathbf{X}}^\top \left[\sigma(\tilde{\mathbf{X}} \tilde{\mathbf{W}}^\natural) - \frac{1}{N} \tilde{\mathbf{X}}^\top \sigma(\tilde{\mathbf{X}} \mathbf{W}_t) \right] \odot \sigma'(\tilde{\mathbf{X}} \mathbf{W}_t).$$

- low-rank approximation error $\leq 2\lambda_{r^*+1}(\mathbf{G})$
- population error: using $\mathbb{E}_{\tilde{x}}[-\mathbf{J}_{\mathbf{W}_t}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta) + \mathcal{O}(\frac{1}{\kappa r^*})$
- concentration error

$$\left\| \mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{x}}[\mathbf{J}_{\mathbf{W}_t}] \right\|_{\text{F}} \lesssim \sqrt{d\epsilon} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, \text{ w.h.p.}$$

Global convergence of nonlinear models

Theorem (Informal, linear convergence rate)

For nonlinear model with $r = r^*$ and gradient descent (with preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), \text{ w.h.p}$$

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{op} \leq \|\mathbf{A}_0 \mathbf{B}_0 - 2\mathbf{G}\|_{op} + 2\|\mathbf{G} - \mathbb{E}_{\tilde{x}}[\mathbf{G}]\|_{op} + \|2\mathbb{E}_{\tilde{x}}[\mathbf{G}] - \Delta\|_{op}$$

$$\mathbf{J}_{\mathbf{W}_t} := \frac{1}{N} \tilde{\mathbf{X}}^\top \left[\sigma(\tilde{\mathbf{X}} \tilde{\mathbf{W}}^\natural) - \frac{1}{N} \tilde{\mathbf{X}}^\top \sigma(\tilde{\mathbf{X}} \mathbf{W}_t) \right] \odot \sigma'(\tilde{\mathbf{X}} \mathbf{W}_t).$$

- low-rank approximation error $\leq 2\lambda_{r^*+1}(\mathbf{G})$
- population error: using $\mathbb{E}_{\tilde{x}}[-\mathbf{J}_{\mathbf{W}_t}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta) + \mathcal{O}(\frac{1}{\kappa r^*})$
- concentration error

$$\left\| \mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{x}}[\mathbf{J}_{\mathbf{W}_t}] \right\|_{\text{F}} \lesssim \sqrt{d\epsilon} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, \text{ w.h.p.}$$

Global convergence of nonlinear models

Theorem (Informal, linear convergence rate)

For nonlinear model with $r = r^*$ and gradient descent (with preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), \text{ w.h.p}$$

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{op} \leq \|\mathbf{A}_0 \mathbf{B}_0 - 2\mathbf{G}\|_{op} + 2\|\mathbf{G} - \mathbb{E}_{\tilde{x}}[\mathbf{G}]\|_{op} + \|2\mathbb{E}_{\tilde{x}}[\mathbf{G}] - \Delta\|_{op}$$

$$\mathbf{J}_{\mathbf{W}_t} := \frac{1}{N} \tilde{\mathbf{X}}^\top \left[\sigma(\tilde{\mathbf{X}} \tilde{\mathbf{W}}^\natural) - \frac{1}{N} \tilde{\mathbf{X}}^\top \sigma(\tilde{\mathbf{X}} \mathbf{W}_t) \right] \odot \sigma'(\tilde{\mathbf{X}} \mathbf{W}_t).$$

- low-rank approximation error $\leq 2\lambda_{r^*+1}(\mathbf{G})$
- population error: using $\mathbb{E}_{\tilde{x}}[-\mathbf{J}_{\mathbf{W}_t}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta) + \mathcal{O}(\frac{1}{\kappa r^*})$
- concentration error

$$\left\| \mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{x}}[\mathbf{J}_{\mathbf{W}_t}] \right\|_{\text{F}} \lesssim \sqrt{d}\epsilon \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, \text{ w.h.p.}$$

Proof of sketch on $\mathbf{A}_t \mathbf{B}_t - \Delta$

$$\begin{aligned}\|\mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta\|_{\text{F}} &\lesssim \|\mathbf{J}_{\mathbf{w}_t} - \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta)\|_{\text{F}} \quad [\text{concentration+population}] \\ &+ (1 - \eta) \left\| \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^{\top} (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^{\top} \right\|_{\text{F}} \\ &+ \left\| (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^{\top}) (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^{\top}) \right\|_{\text{F}} \\ &+ \text{cross terms}\end{aligned}$$

□ projection

$$\mathbf{L} = \begin{bmatrix} \mathbf{U}_{\mathbf{A}_t} & \mathbf{0}_{d \times r} \\ \mathbf{0}_{k \times r} & \mathbf{V}_{\mathbf{B}_t} \end{bmatrix} \in \mathbb{R}^{(d+k) \times 2r},$$

then $\mathbf{L} \mathbf{L}^{\top}$ is a projection matrix, $\mathbf{I}_{d+k} - \mathbf{L} \mathbf{L}^{\top} = \mathbf{L}_{\perp} \mathbf{L}_{\perp}^{\top}$

□ lower bound $\left\| \mathbf{L}_{\perp}^{\top} \Delta \mathbf{L} \right\|_{\text{F}}^2$, upper bound $\left\| \mathbf{L}_{\perp}^{\top} \mathbf{U} \right\|_{\text{op}} < 1$

Proof of sketch on $\mathbf{A}_t \mathbf{B}_t - \Delta$

$$\begin{aligned}
\|\mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta\|_{\text{F}} &\lesssim \|\mathbf{J}_{\mathbf{w}_t} - \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta)\|_{\text{F}} \quad [\text{concentration+population}] \\
&+ (1 - \eta) \left\| \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^{\top} (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^{\top} \right\|_{\text{F}} \\
&+ \left\| (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^{\top}) (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^{\top}) \right\|_{\text{F}} \\
&+ \text{cross terms}
\end{aligned}$$

□ projection

$$\mathbf{L} = \begin{bmatrix} \mathbf{U}_{\mathbf{A}_t} & \mathbf{0}_{d \times r} \\ \mathbf{0}_{k \times r} & \mathbf{V}_{\mathbf{B}_t} \end{bmatrix} \in \mathbb{R}^{(d+k) \times 2r},$$

then $\mathbf{L}\mathbf{L}^{\top}$ is a projection matrix, $\mathbf{I}_{d+k} - \mathbf{L}\mathbf{L}^{\top} = \mathbf{L}_{\perp} \mathbf{L}_{\perp}^{\top}$

□ lower bound $\left\| \mathbf{L}_{\perp}^{\top} \Delta \mathbf{L} \right\|_{\text{F}}^2$, upper bound $\left\| \mathbf{L}_{\perp}^{\top} \mathbf{U} \right\|_{\text{op}} < 1$

Proof of sketch on $\mathbf{A}_t \mathbf{B}_t - \Delta$

$$\begin{aligned}
\|\mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta\|_{\text{F}} &\lesssim \|\mathbf{J}_{\mathbf{w}_t} - \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta)\|_{\text{F}} \quad [\text{concentration+population}] \\
&+ (1 - \eta) \left\| \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^{\top} (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^{\top} \right\|_{\text{F}} \\
&+ \left\| (\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^{\top}) (\mathbf{A}_t \mathbf{B}_t - \Delta) (\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^{\top}) \right\|_{\text{F}} \\
&+ \text{cross terms}
\end{aligned}$$

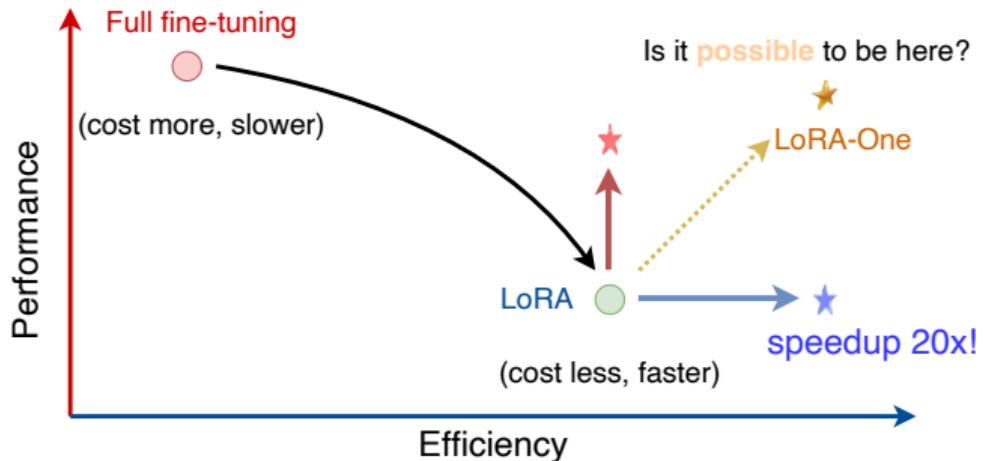
□ projection

$$\mathbf{L} = \begin{bmatrix} \mathbf{U}_{\mathbf{A}_t} & \mathbf{0}_{d \times r} \\ \mathbf{0}_{k \times r} & \mathbf{V}_{\mathbf{B}_t} \end{bmatrix} \in \mathbb{R}^{(d+k) \times 2r},$$

then $\mathbf{L}\mathbf{L}^{\top}$ is a projection matrix, $\mathbf{I}_{d+k} - \mathbf{L}\mathbf{L}^{\top} = \mathbf{L}_{\perp} \mathbf{L}_{\perp}^{\top}$

□ lower bound $\left\| \mathbf{L}_{\perp}^{\top} \Delta \mathbf{L} \right\|_{\text{F}}^2$, upper bound $\left\| \mathbf{L}_{\perp}^{\top} \mathbf{U} \right\|_{op} < 1$

Takeaway messages: speedup via spectral initialization



- *LoRA-One: One-step full gradient could suffice for fine-tuning large language models, provably and efficiently. ICML'25 Oral*

- **subspace alignment:** \mathbf{G} and $(\mathbf{A}_t, \mathbf{B}_t)$ \Rightarrow theory-grounded algorithm design
- “optimal” non-zero initialization strategy
- enables feature learning, precondition helps convergence...

Target

- How to handle **nonlinearity** at a theoretical level (e.g., training dynamics)
- How to precisely and efficiently approximate **nonlinearity** at a practical level under theoretical guidelines

Thank you!

fanghui.liu@warwick.ac.uk

www.lfhsgre.org

Target

- How to handle **nonlinearity** at a theoretical level (e.g., training dynamics)
- How to precisely and efficiently approximate **nonlinearity** at a practical level under theoretical guidelines

Thank you!

fanghui.liu@warwick.ac.uk

www.lfhsgre.org

-  Kai Lv, Yuqing Yang, Tengxiao Liu, Qipeng Guo, and Xipeng Qiu.
Full Parameter Fine-tuning for Large Language Models with Limited Resources.
In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8187–8198, 2024.
-  Dominik Stöger and Mahdi Soltanolkotabi.
Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction.
In *Advances in Neural Information Processing Systems*, pages 23831–23843, 2021.