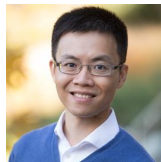
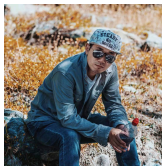


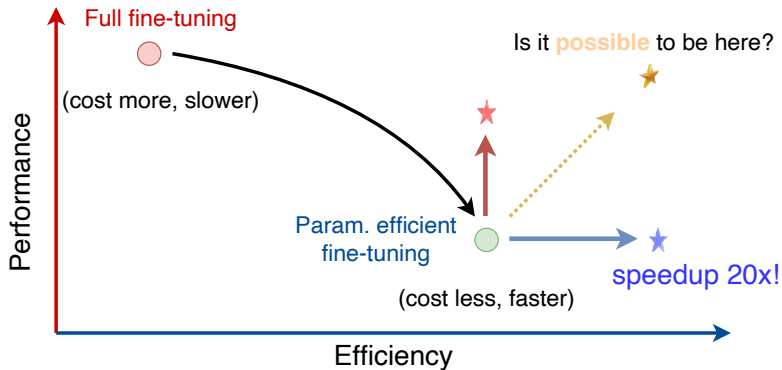
LoRA-One: One-step full gradient could suffice for fine-tuning large language models, provably and efficiently

Yuanhe Zhang (Warwick), **Fanghui Liu** (Warwick), **Yudong Chen** (UW-Madison)



at ICML'25, Vancouver

How can **theory** contribute to efficiency in LLMs?



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

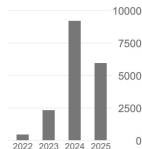
Weizhu Chen

Microsoft Corporation

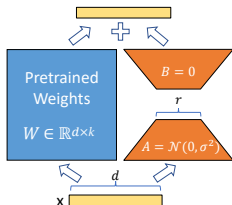
edward.hu@mila.quebec

{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:

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

Today's talk: Improve “sub-optimal” LoRA



How can theory guide practice

- understanding: training dynamics of $(\mathbf{A}_t, \mathbf{B}_t)$
- design new algorithm \rightarrow performance improvement
- 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 $\text{rank}(\mathbf{G}) = r^*$

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

Today's talk: Improve “sub-optimal” LoRA



How can theory guide practice

- understanding: training dynamics of $(\mathbf{A}_t, \mathbf{B}_t)$
- design new algorithm \rightarrow performance improvement
- 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 new algorithm \rightarrow performance improvement
- clarify some misconceptions in previous algorithm designs

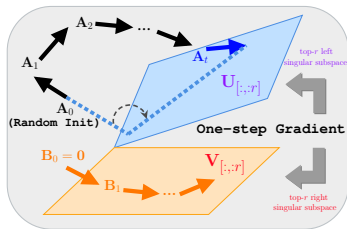
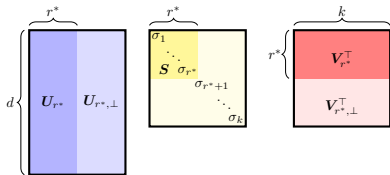
□ Even for linear model (pre-training and fine-tuning), **nonlinear dynamics...**

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

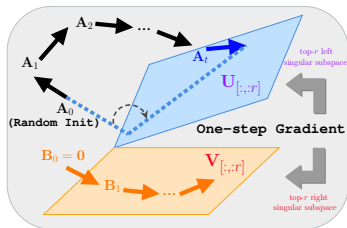
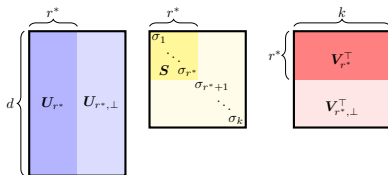
□ One-step full gradient: $\mathbf{G} \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}) = r^*$

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

Understanding: Alignment on B_t



Understanding: Alignment on B_t

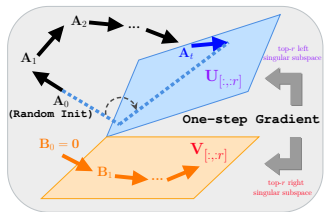


Theorem (Alignment between G and B_t , informal)

For the linear setting, LoRA via gradient descent yields

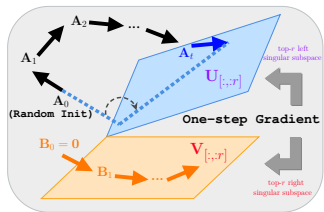
$$\angle(\mathbf{V}_{r^*}(\mathbf{B}_t), \mathbf{V}_{r^*}(\mathbf{G})) = 0, \quad \forall t \in \mathbb{N}_+.$$

Understanding: Alignment on A_t



$$\begin{bmatrix} A_{t+1} \\ B_{t+1}^\top \end{bmatrix} = \begin{bmatrix} I_d & \eta G \\ \eta G^\top & I_k \end{bmatrix} \begin{bmatrix} A_t \\ B_t^\top \end{bmatrix} + \text{nonlinear term}$$

Understanding: Alignment on A_t



$$\begin{bmatrix} A_{t+1} \\ B_{t+1}^\top \end{bmatrix} = \begin{bmatrix} I_d & \eta G \\ \eta G^\top & I_k \end{bmatrix} \begin{bmatrix} A_t \\ B_t^\top \end{bmatrix} + \text{nonlinear term}$$

Theorem (Alignment between G and A_t , informal)

For small initialization over A_0 , after $t^* = \Theta(\ln d)$ steps, LoRA updates yield $\angle(U_{r^*}(A_{t^*}), U_{r^*}(G))$ is small, w.h.p.

Understanding: Alignment on A_t

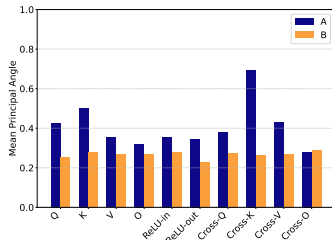
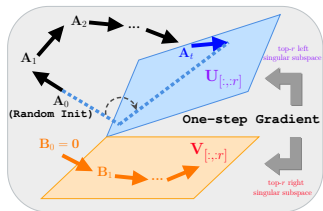


Figure 2: Principal angle of fine-tuning T5 on MRPC.

Theorem (Alignment between G and A_t , informal)

For small initialization over A_0 , after $t^* = \Theta(\ln d)$ steps, LoRA updates yield

$$\angle(U_{r^*}(A_{t^*}), U_{r^*}(G)) \text{ is small, w.h.p.}$$

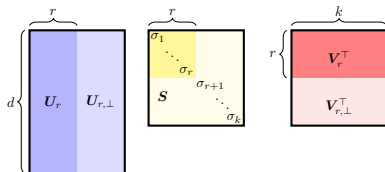
Algorithm design principle

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

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

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

(Spec-init.)



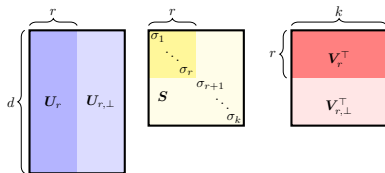
Algorithm design principle

□ SVD: $G = USV^\top$

$$A_0 = U_{[:,1:r]} S_{[1:r]}^{\frac{1}{2}}$$

$$B_0 = S_{[1:r]}^{\frac{1}{2}} V_{[:,1:r]}^\top$$

(Spec-init.)



Key Message: we can “escape” the alignment stage

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

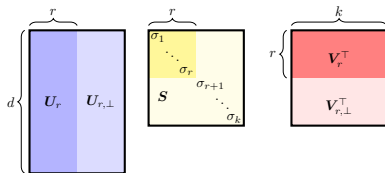
$$\|A_0 B_0 - \Delta\|_F \text{ is small, w.h.p.}$$

Algorithm design principle

□ SVD: $\mathbf{G} = \mathbf{USV}^\top$

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

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



Key Message: we can “escape” the alignment stage

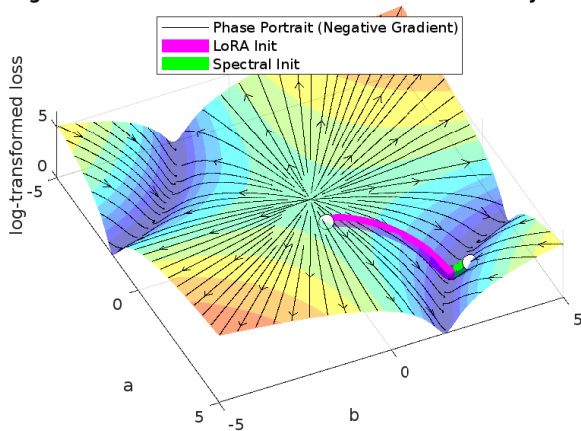
Under (**Spec-init.**), for both **linear/nonlinear** models, we can directly achieve the alignment at initialization.

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

The “best” initialization strategy!

“Best” initialization: phase portrait

Log-Transformed Surface with Phase Portrait and Trajectories



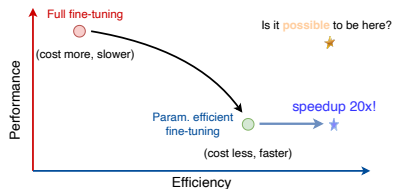
One-step gradient can suffice on small-scale datasets!

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 \pm 0.04	94.04 \pm 0.09	72.84 \pm 1.25	93.02 \pm 0.07	68.38 \pm 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 \pm 0.04	94.04 \pm 0.09	72.84 \pm 1.25	93.02 \pm 0.07	68.38 \pm 0.01

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



Results on LLaMA 2-7B (continue to run)

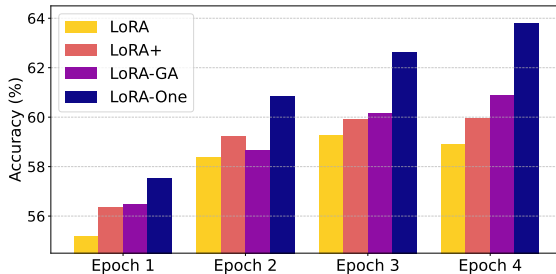


Figure 3: Accuracy comparison across different methods over epochs on GSM8K.

Time cost	LoRA: 6h 20min	+ 3min
Memory	LoRA: 21.6 GB	+ 0.1GB

Results on LLaMA 2-7B (continue to run)

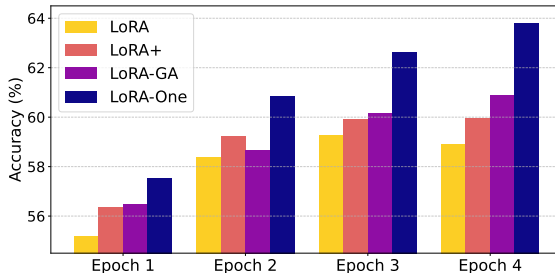
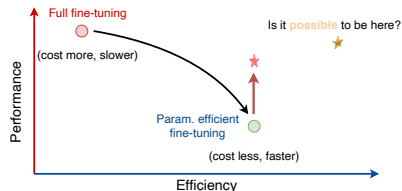


Figure 3: Accuracy comparison across different methods over epochs on GSM8K.

Time cost	LoRA: 6h 20min	+ 3min
Memory	LoRA: 21.6 GB	+ 0.1GB



Clarification on gradient alignment based work

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

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!

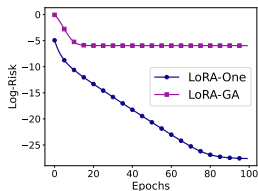
□ best $2r$ approximation

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

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!

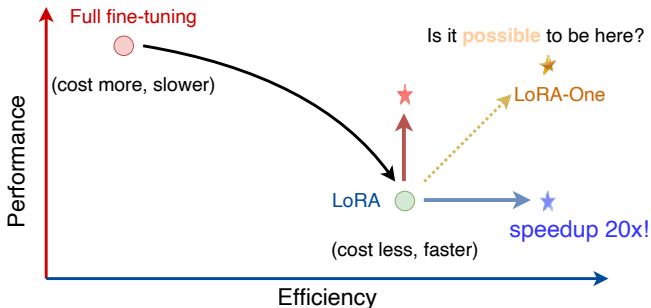


Method	Init. on A	Init. on B	Calibration
LoRA	$\mathcal{N}(0, \alpha^2)$	0	-
LoRA-GA	$U_{[:,1:r]}$	$V_{[:,r+1:2r]}^\top$	$W^{\text{pre}} - A_0 B_0$
LoRA-One	$U_{[:,1:r]} S_{[1:r]}^{1/2}$	$S_{[1:r]}^{1/2} V_{[:,1:r]}^\top$	-

Takeaway messages: speedup via spectral initialization



LoRA-One



- **subspace alignment:** \mathbf{G} and $(\mathbf{A}_t, \mathbf{B}_t) \Rightarrow$ theory-grounded algorithm design
- “optimal” non-zero initialization strategy
- clarification on gradient alignment based algorithms
- spectral initialization enables feature learning...
- global convergence on nonlinear models, scaled GD...