

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

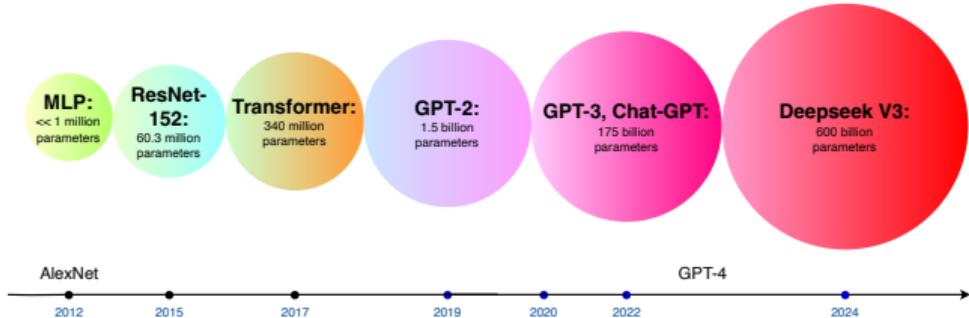
Department of Computer Science, University of Warwick, UK

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

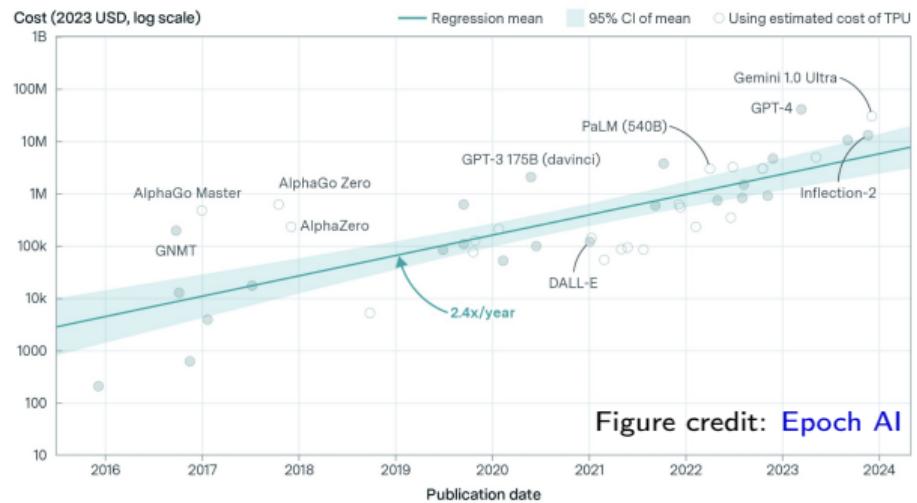
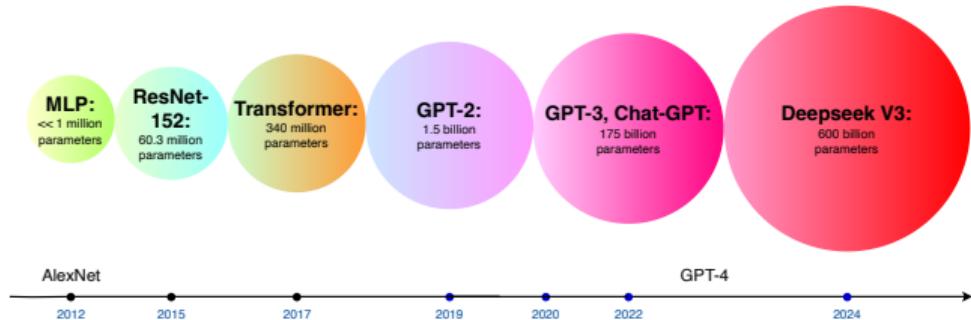
[joint work with Yuanhe Zhang (Warwick) and Yudong Chen (UW-Madison)]



In the era of machine learning (Pre-training)



In the era of machine learning (Pre-training)

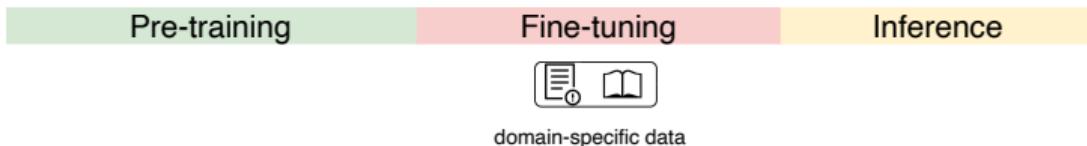


From pre-training to (parameter-efficient) fine-tuning

- GPT3: 175 billion parameters
- Llama3.1: > 400 billion parameters
- Gemini 1.5 Pro 300–500 billion parameters (**unconfirmed**)
- Deepseek-v3: > 600 billion parameters
- Llama 4 Behemoth: > 2,000 billion parameters

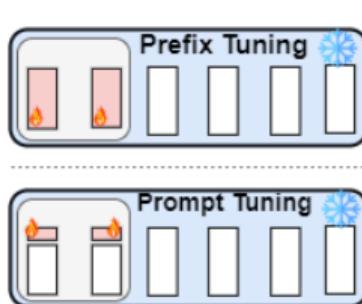
From pre-training to (parameter-efficient) fine-tuning

- GPT3: 175 billion parameters
- Llama3.1: > 400 billion parameters
- Gemini 1.5 Pro 300–500 billion parameters (**unconfirmed**)
- Deepseek-v3: > 600 billion parameters
- Llama 4 Behemoth: > 2,000 billion parameters

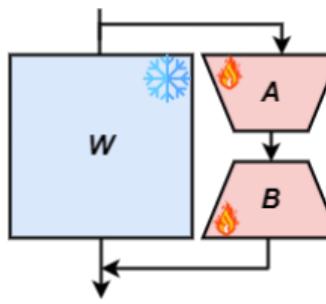


From pre-training to (parameter-efficient) fine-tuning

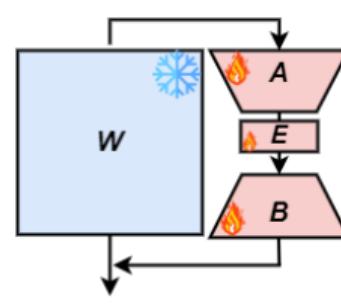
- GPT3: 175 billion parameters
- Llama3.1: > 400 billion parameters
- Gemini 1.5 Pro 300–500 billion parameters (**unconfirmed**)
- Deepseek-v3: > 600 billion parameters
- Llama 4 Behemoth: > 2,000 billion parameters



(a) Prefix & Prompt



(b) LoRA



(c) LoRA variants

Low-rank adaption (LoRA) for fine-tuning [2]

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

LoRA

- 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:

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

How theory guides practice (not limited to understanding)

- design new algorithm -> performance improvement (accuracy, efficiency)
- clarify some misconceptions in algorithm design

Low-rank adaption (LoRA) for fine-tuning [2]

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

LoRA

- 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:

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

How theory guides practice (not limited to understanding)

- design new algorithm -> performance improvement (accuracy, efficiency)
- clarify some misconceptions in algorithm design

Low-rank adaption (LoRA) for fine-tuning [2]

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

LoRA

- 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:

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

How theory guides practice (not limited to understanding)

- design new algorithm -> performance improvement (accuracy, efficiency)
- clarify some misconceptions in algorithm design

Motivation: non-linear dynamics and subspace alignment

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

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G}^\natural \\ \eta \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term}.$$

- \mathbf{G}^\natural : one-step full gradient (from full fine-tuning)
- The dynamics $(\mathbf{A}_t, \mathbf{B}_t)$ heavily depends on \mathbf{G}^\natural !

Target

- Q1: How to characterize low-rank dynamics of LoRA and the associated subspace alignment in theory?
- Q2: How can our theoretical results contribute to algorithm design for LoRA in practice?

Motivation: non-linear dynamics and subspace alignment

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

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G}^\natural \\ \eta \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term}.$$

- \mathbf{G}^\natural : one-step full gradient (from full fine-tuning)
- The dynamics $(\mathbf{A}_t, \mathbf{B}_t)$ heavily depends on \mathbf{G}^\natural !

Target

- Q1: How to characterize low-rank dynamics of LoRA and the associated subspace alignment in theory?
- Q2: How can our theoretical results contribute to algorithm design for LoRA in practice?

Motivation: non-linear dynamics and subspace alignment

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

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G}^\natural \\ \eta \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term}.$$

- \mathbf{G}^\natural : one-step full gradient (from full fine-tuning)
- The dynamics $(\mathbf{A}_t, \mathbf{B}_t)$ heavily depends on \mathbf{G}^\natural !

Target

- Q1: *How to characterize low-rank dynamics of LoRA and the associated subspace alignment in theory?*
- Q2: *How can our theoretical results contribute to algorithm design for LoRA in practice?*

Motivation: non-linear dynamics and subspace alignment

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

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta \mathbf{G}^\natural \\ \eta \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term}.$$

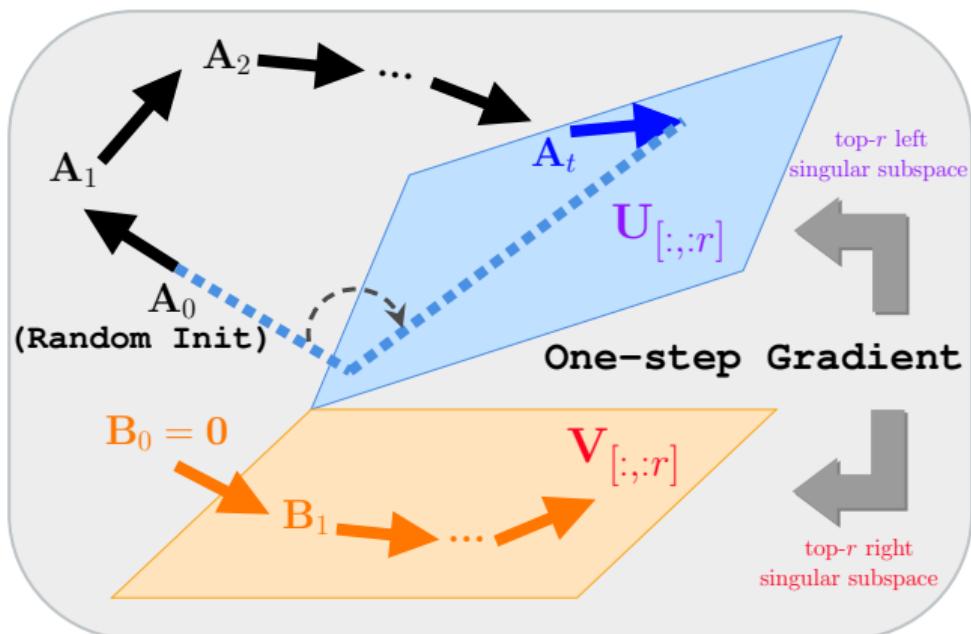
- \mathbf{G}^\natural : one-step full gradient (from full fine-tuning)
- The dynamics $(\mathbf{A}_t, \mathbf{B}_t)$ heavily depends on \mathbf{G}^\natural !

Target

- Q1: How to characterize low-rank dynamics of LoRA and the associated subspace alignment in theory?*
- Q2: How can our theoretical results contribute to algorithm design for LoRA in practice?*

Alignment and theory-grounded algorithm

Pipeline



Problem setting and assumptions

- Pre-trained model: known $\mathbf{W}^\natural \in \mathbb{R}^{d \times k}$ and the ReLU activation σ

$$f_{\text{pre}}(\mathbf{x}) := \begin{cases} (\mathbf{x}^\top \mathbf{W}^\natural)^\top \in \mathbb{R}^k & \text{linear} \\ \sigma[(\mathbf{x}^\top \mathbf{W}^\natural)^\top] \in \mathbb{R}^k & \text{nonlinear} \end{cases}.$$

- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^\natural := \mathbf{W}^\natural + \Delta$
- $\text{Rank}(\Delta) = r^* < \min\{d, k\}$ with unknown r^*
- Downstream well-behaved data $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^N$ for fine-tuning:
$$\tilde{\mathbf{y}} := \begin{cases} (\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\natural)^\top \in \mathbb{R}^k, \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, & \text{linear} \\ \sigma[(\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\natural)^\top], \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) & \text{nonlinear} \end{cases}$$
- We assume $N > d$, e.g., MetaMathQA, Code-Feedback, $d = 1,024$ and $N \sim 10^5$

Problem setting and assumptions

- Pre-trained model: known $\mathbf{W}^\natural \in \mathbb{R}^{d \times k}$ and the ReLU activation σ

$$f_{\text{pre}}(\mathbf{x}) := \begin{cases} (\mathbf{x}^\top \mathbf{W}^\natural)^\top \in \mathbb{R}^k & \text{linear} \\ \sigma[(\mathbf{x}^\top \mathbf{W}^\natural)^\top] \in \mathbb{R}^k & \text{nonlinear} \end{cases}.$$

- Unknown low-rank feature shift Δ : $\tilde{\mathbf{W}}^\natural := \mathbf{W}^\natural + \Delta$

- Rank(Δ) = $r^* < \min\{d, k\}$ with unknown r^*

- Downstream well-behaved data $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^N$ for fine-tuning:

$$\tilde{\mathbf{y}} := \begin{cases} (\tilde{\mathbf{x}}^\top \tilde{\mathbf{W}}^\natural)^\top \in \mathbb{R}^k, \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, & \text{linear} \\ \sigma[(\tilde{\mathbf{x}}^\top \tilde{\mathbf{W}}^\natural)^\top], \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) & \text{nonlinear} \end{cases}.$$

- We assume $N > d$, e.g., MetaMathQA, Code-Feedback, $d = 1,024$ and $N \sim 10^5$

Problem setting and assumptions

- Pre-trained model: known $\mathbf{W}^\natural \in \mathbb{R}^{d \times k}$ and the ReLU activation σ

$$f_{\text{pre}}(\mathbf{x}) := \begin{cases} (\mathbf{x}^\top \mathbf{W}^\natural)^\top \in \mathbb{R}^k & \text{linear} \\ \sigma[(\mathbf{x}^\top \mathbf{W}^\natural)^\top] \in \mathbb{R}^k & \text{nonlinear} \end{cases}.$$

- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^\natural := \mathbf{W}^\natural + \Delta$
- $\text{Rank}(\Delta) = r^* < \min\{d, k\}$ with unknown r^*
- Downstream well-behaved data $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^N$ for fine-tuning:
$$\tilde{\mathbf{y}} := \begin{cases} (\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\natural)^\top \in \mathbb{R}^k, \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, & \text{linear} \\ \sigma[(\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\natural)^\top], \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) & \text{nonlinear} \end{cases}.$$
- We assume $N > d$, e.g., MetaMathQA, Code-Feedback, $d = 1,024$ and $N \sim 10^5$

Problem setting and assumptions

- Pre-trained model: known $\mathbf{W}^\natural \in \mathbb{R}^{d \times k}$ and the ReLU activation σ

$$f_{\text{pre}}(\mathbf{x}) := \begin{cases} (\mathbf{x}^\top \mathbf{W}^\natural)^\top \in \mathbb{R}^k & \text{linear} \\ \sigma[(\mathbf{x}^\top \mathbf{W}^\natural)^\top] \in \mathbb{R}^k & \text{nonlinear} \end{cases}.$$

- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^\natural := \mathbf{W}^\natural + \Delta$
- $\text{Rank}(\Delta) = r^* < \min\{d, k\}$ with unknown r^*
- Downstream well-behaved data $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^N$ for fine-tuning:
$$\tilde{\mathbf{y}} := \begin{cases} (\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\natural)^\top \in \mathbb{R}^k, \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, & \text{linear} \\ \sigma[(\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\natural)^\top], \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) & \text{nonlinear} \end{cases}.$$
- We assume $N > d$, e.g., MetaMathQA, Code-Feedback, $d = 1,024$ and $N \sim 10^5$

Full fine-tuning and LoRA updates

- full fine-tuning (initialized at $\mathbf{W}_0 := \mathbf{W}^\natural$)

$$L(\mathbf{W}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}}\mathbf{W} - \tilde{\mathbf{Y}} \right\|_F^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}}\mathbf{W}) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{nonlinear} \end{cases}$$

- LoRA update

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{AB}) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{AB})) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{nonlinear} \end{cases}$$

- Gradient descent with step-size η

$$\mathbf{A}_{t+1} = \mathbf{A}_t - \eta \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \eta \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

- Evaluation by $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F$: optimization and generalization!

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

Full fine-tuning and LoRA updates

- full fine-tuning (initialized at $\mathbf{W}_0 := \mathbf{W}^\natural$)

$$L(\mathbf{W}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}}\mathbf{W} - \tilde{\mathbf{Y}} \right\|_F^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}}\mathbf{W}) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{nonlinear} \end{cases}$$

- LoRA update

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{AB}) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{AB})) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{nonlinear} \end{cases}$$

- Gradient descent with step-size η

$$\mathbf{A}_{t+1} = \mathbf{A}_t - \eta \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \eta \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

- Evaluation by $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F$: optimization and generalization!

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

Full fine-tuning and LoRA updates

- full fine-tuning (initialized at $\mathbf{W}_0 := \mathbf{W}^\natural$)

$$L(\mathbf{W}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}}\mathbf{W} - \tilde{\mathbf{Y}} \right\|_F^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}}\mathbf{W}) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{nonlinear} \end{cases}$$

- LoRA update

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{AB}) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{AB})) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{nonlinear} \end{cases}$$

- Gradient descent with step-size η

$$\mathbf{A}_{t+1} = \mathbf{A}_t - \eta \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \eta \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

- Evaluation by $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F$: optimization and generalization!

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

Full fine-tuning and LoRA updates

- full fine-tuning (initialized at $\mathbf{W}_0 := \mathbf{W}^\natural$)

$$L(\mathbf{W}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}}\mathbf{W} - \tilde{\mathbf{Y}} \right\|_F^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}}\mathbf{W}) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{nonlinear} \end{cases}$$

- LoRA update

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{AB}) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}}(\mathbf{W}^\natural + \mathbf{AB})) - \tilde{\mathbf{Y}} \right\|_F^2 & \text{nonlinear} \end{cases}$$

- Gradient descent with step-size η

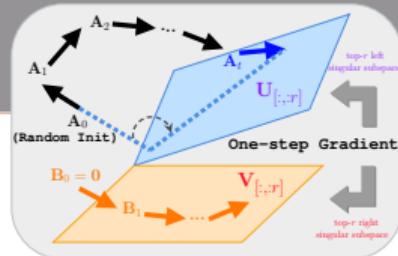
$$\mathbf{A}_{t+1} = \mathbf{A}_t - \eta \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \eta \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

- Evaluation by $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F$: optimization and generalization!

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

Our results: Alignment on B_t



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

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

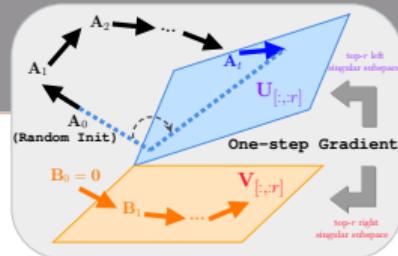
Theorem (Alignment between \mathbf{G}^\natural and B_t)

For the linear setting, consider the LoRA updates with (LoRA-init.). We have

$$\left\| \mathbf{V}_{r^*, \perp}^\top (\mathbf{G}^\natural) \mathbf{V}_{r^*}(B_t) \right\|_{op} = 0, \quad \forall t \in \mathbb{N}_+.$$

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

Our results: Alignment on B_t



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

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

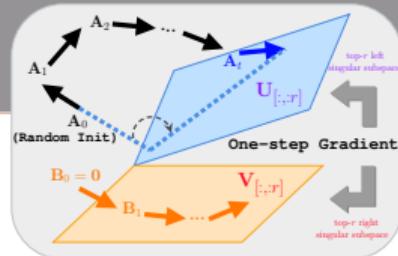
Theorem (Alignment between \mathbf{G}^\natural and \mathbf{B}_t)

For the linear setting, consider the LoRA updates with (LoRA-init.). We have

$$\left\| \mathbf{V}_{r^*, \perp}^\top (\mathbf{G}^\natural) \mathbf{V}_{r^*} (\mathbf{B}_t) \right\|_{op} = 0, \quad \forall t \in \mathbb{N}_+.$$

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

Our results: Alignment on B_t



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

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

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

For the linear setting, consider the LoRA updates with (LoRA-init.). We have

$$\left\| \mathbf{V}_{r^*, \perp}^\top (\mathbf{G}^\natural) \mathbf{V}_{r^*} (\mathbf{B}_t) \right\|_{op} = 0, \quad \forall t \in \mathbb{N}_+.$$

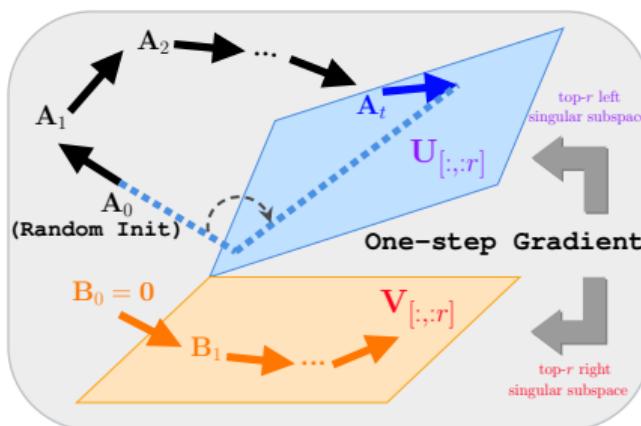
Remark: $\mathbf{B}_1 = \eta_1 \mathbf{A}_0^\top \mathbf{G}^\natural$ with $\text{Rank}(\mathbf{B}_1) \leq r^*$

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

$$\left\| U_{r^*, \perp}^\top(G^\natural) U_{r^*}(A_{t^*}) \right\|_{op} \lesssim \epsilon, \text{ w.h.p.}$$



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

$$\left\| U_{r^*, \perp}^\top(G^\natural) U_{r^*}(A_{t^*}) \right\|_{op} \lesssim \epsilon, \text{ w.h.p.}$$

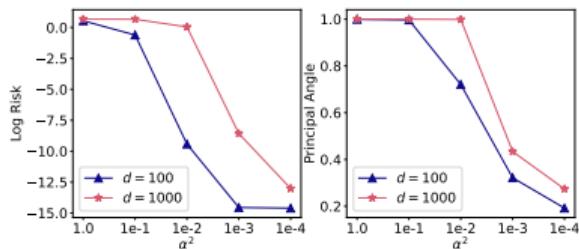


Figure 2: Left: the risk $\frac{1}{2} \|A_t B_t - \Delta\|_F^2$.
 Right: the principal angle is $\min_t \|U_{r^*, \perp}^\top(G^\natural) U_{r^*}(A_t)\|_{op}$.

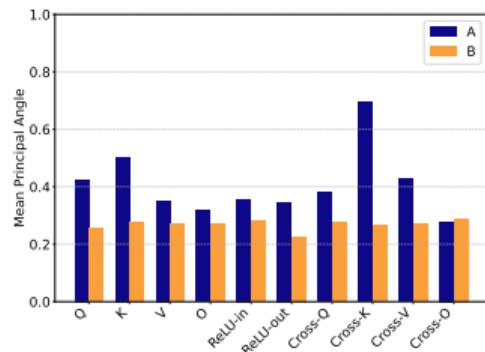


Figure 3: Principal angle of fine-tuning T5 on MRPC.

Key message: Algorithm design principle

Can we “escape” the alignment stage?

- Take the SVD of \mathbf{G}^\natural : $\mathbf{G}^\natural = \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$

$$\mathbf{A}_0 = \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]} \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} . \quad (\text{Spec-init.})$$

$$\mathbf{B}_0 = \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}^\top .$$

Message

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

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F \leq \epsilon \|\Delta\|_{op}, \quad w.p. \ 1 - \exp(-\epsilon^2 N)$$

The “best” initialization strategy!

Key message: Algorithm design principle

Can we “escape” the alignment stage?

- Take the SVD of \mathbf{G}^\natural : $\mathbf{G}^\natural = \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$

$$\mathbf{A}_0 = \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]} \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} . \quad (\text{Spec-init.})$$

$$\mathbf{B}_0 = \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}^\top .$$

Message

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

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F \leq \epsilon \|\Delta\|_{op}, \quad w.p. 1 - \exp(-\epsilon^2 N)$$

The “best” initialization strategy!

Key message: Algorithm design principle

Can we “escape” the alignment stage?

- Take the SVD of \mathbf{G}^\natural : $\mathbf{G}^\natural = \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$

$$\mathbf{A}_0 = \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]} \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} . \quad (\text{Spec-init.})$$

$$\mathbf{B}_0 = \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}^\top .$$

Message

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

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F \leq \epsilon \|\Delta\|_{op}, \quad w.p. 1 - \exp(-\epsilon^2 N)$$

The “best” initialization strategy!

Key message: Algorithm design principle

Can we “escape” the alignment stage?

- Take the SVD of \mathbf{G}^\natural : $\mathbf{G}^\natural = \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$

$$\mathbf{A}_0 = \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]} \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} . \quad (\text{Spec-init.})$$

$$\mathbf{B}_0 = \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}^\top .$$

Message

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

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F \leq \epsilon \|\Delta\|_{op}, \quad w.p. \ 1 - \exp(-\epsilon^2 N)$$

The “best” initialization strategy!

Key message: Algorithm design principle

Can we “escape” the alignment stage?

- Take the SVD of \mathbf{G}^\natural : $\mathbf{G}^\natural = \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$

$$\mathbf{A}_0 = \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]} \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} . \quad (\text{Spec-init.})$$

$$\mathbf{B}_0 = \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}^\top .$$

Message

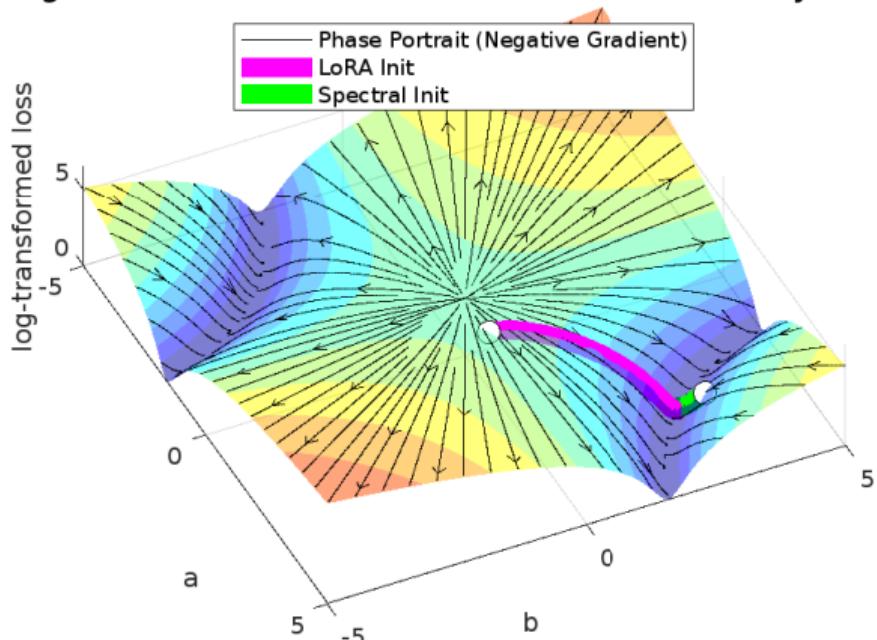
If we choose (Spec-init.), for both linear/nonlinear models, we can directly achieve the alignment at initialization.

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F \leq \epsilon \|\Delta\|_{op}, \quad w.p. \ 1 - \exp(-\epsilon^2 N)$$

The “best” initialization strategy!

How to understand “best” initialization? - Phase portrait

Log-Transformed Surface with Phase Portrait and Trajectories



Toy example (I)

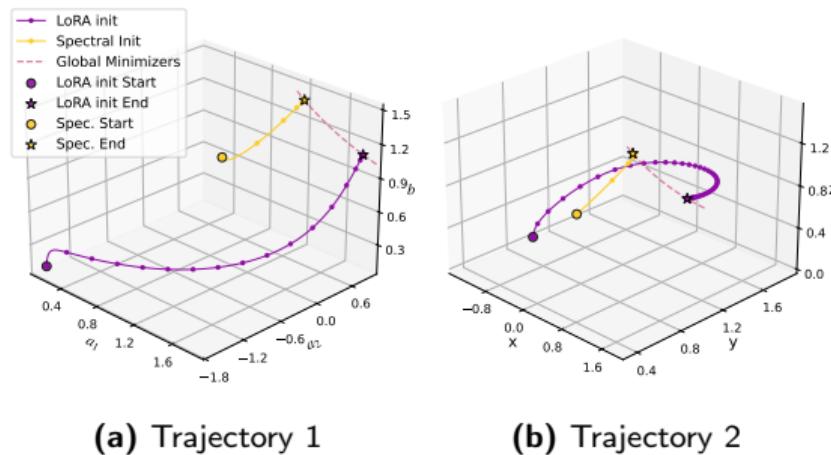


Figure 4: Comparison of the GD trajectories between LoRA and ours. (a) and (b): $\mathbf{A} \in \mathbb{R}^2$ and $B \in \mathbb{R}$ with different initializations. The set of global minimizers is $\{a_1^* = 2/t, a_2^* = 1/t, b^* = t \mid t \in \mathbb{R}\}$.

Toy example (II)

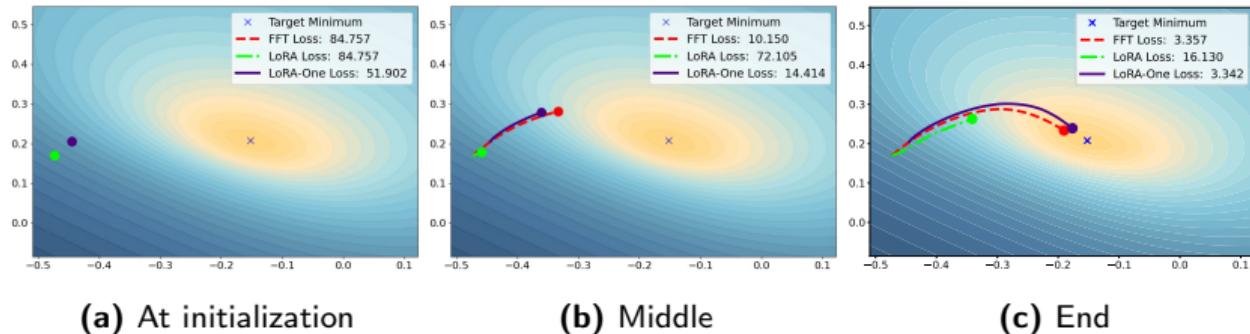


Figure 5: Comparison of the GD trajectories between LoRA and ours. We use two-layer neural networks pre-trained on odd-labeled data and fine-tuned on even-labeled data on MNIST, see [GIF illustration](#).

One-step full gradient may suffice for low-rank fine-tuning!

Table 1: Fine-tuning T5 model across NLP tasks from GLUE.

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
One-step GD	-	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

- **CoLA** LoRA: 47s, one-step: <1s
- **MRPC** LoRA: 25s, one-step: <1s

One-step full gradient may suffice for low-rank fine-tuning!

Table 1: Fine-tuning T5 model across NLP tasks from GLUE.

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
One-step GD	-	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

- **CoLA** LoRA: 47s, one-step: <1s
- **MRPC** LoRA: 25s, one-step: <1s

Implementation remarks

- For $\nabla_{\mathbf{W}} L(\mathbf{W}^\dagger)$, we employ a memory-efficient approach [3, 5], the usage remains at $\mathcal{O}(1)$ instead of $\mathcal{O}(\text{number of layers})$
- For SVD $(-\nabla_{\mathbf{W}} L(\mathbf{W}^\dagger))$, we use randomized algorithm [1] with $20\times$ speedup.
- The initialization across layers is parallelizable.

Implementation remarks

- For $\nabla_{\mathbf{W}} L(\mathbf{W}^\natural)$, we employ a memory-efficient approach [3, 5], the usage remains at $\mathcal{O}(1)$ instead of $\mathcal{O}(\text{number of layers})$
- For $\text{SVD}\left(-\nabla_{\mathbf{W}} L(\mathbf{W}^\natural)\right)$, we use randomized algorithm [1] with $20\times$ speedup.
- The initialization across layers is parallelizable.

Implementation remarks

- For $\nabla_{\mathbf{W}} L(\mathbf{W}^\natural)$, we employ a memory-efficient approach [3, 5], the usage remains at $\mathcal{O}(1)$ instead of $\mathcal{O}(\text{number of layers})$
- For $\text{SVD}\left(-\nabla_{\mathbf{W}} L(\mathbf{W}^\natural)\right)$, we use randomized algorithm [1] with $20\times$ speedup.
- The initialization across layers is parallelizable.

Clarification on gradient alignment based work

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

Motivation of LoRA-GA [5, 6]

make LoRA's gradients align to full fine-tuning!

Clarification on gradient alignment based work

- best- $2r$ approximation: $\text{rank}(\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) + \text{rank}(\nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) \leq 2r$

$$\mathbf{A}_0 \leftarrow \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}, \mathbf{B}_0 \leftarrow \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,r+1:2r]}^\top. \quad (\text{LoRA-GA})$$

- But! \mathbf{B}_t will align to the right-side rank- r^* singular subspace of \mathbf{G}^\natural .

- biased estimation: $\mathbf{W}^\natural - \mathbf{A}_0 \mathbf{B}_0$

- small learning rate for Taylor expansion

Clarification on gradient alignment based work

- best- $2r$ approximation: $\text{rank}(\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) + \text{rank}(\nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) \leq 2r$

$$\mathbf{A}_0 \leftarrow \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}, \mathbf{B}_0 \leftarrow \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,r+1:2r]}^\top. \quad (\text{LoRA-GA})$$

- But! \mathbf{B}_t will align to the right-side rank- r^* singular subspace of \mathbf{G}^\natural .

- biased estimation: $\mathbf{W}^\natural - \mathbf{A}_0 \mathbf{B}_0$

- small learning rate for Taylor expansion

Clarification on gradient alignment based work

- best- $2r$ approximation: $\text{rank}(\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) + \text{rank}(\nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) \leq 2r$

$$\mathbf{A}_0 \leftarrow \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}, \mathbf{B}_0 \leftarrow \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,r+1:2r]}^\top. \quad (\text{LoRA-GA})$$

- But! \mathbf{B}_t will align to the right-side rank- r^* singular subspace of \mathbf{G}^\natural .

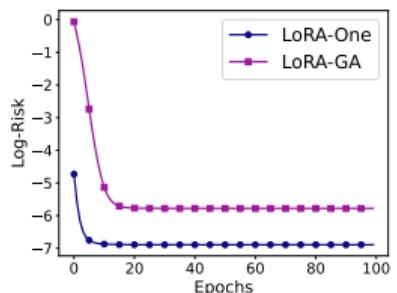
- biased estimation: $\mathbf{W}^\natural - \mathbf{A}_0 \mathbf{B}_0$
- small learning rate for Taylor expansion

Clarification on gradient alignment based work

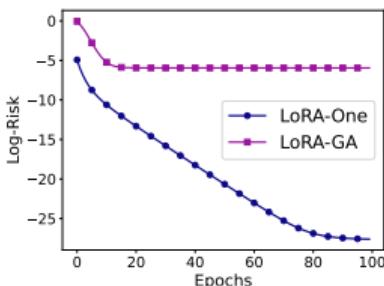
- best- $2r$ approximation: $\text{rank}(\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) + \text{rank}(\nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) \leq 2r$

$$\mathbf{A}_0 \leftarrow \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}, \mathbf{B}_0 \leftarrow \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,r+1:2r]}^\top. \quad (\text{LoRA-GA})$$

- But! \mathbf{B}_t will align to the right-side rank- r^* singular subspace of \mathbf{G}^\natural .



(a) $r < r^*$



(b) $r > r^*$

- biased estimation: $\mathbf{W}^\natural - \mathbf{A}_0 \mathbf{B}_0$
- small learning rate for Taylor expansion

Experiments

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

Experiments on NLP tasks from GLUE

Method	MNLI	SST-2	CoLA	QNLI	MRPC
LoRA	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
LoRA+	85.81 \pm 0.09	93.85 \pm 0.24	77.53 \pm 0.20	93.14 \pm 0.03	74.43 \pm 1.39
P-LoRA	85.28 \pm 0.15	93.88 \pm 0.11	79.58 \pm 0.67	93.00 \pm 0.07	83.91 \pm 1.16
PiSSA	85.75 \pm 0.07	94.07 \pm 0.06	74.27 \pm 0.39	93.15 \pm 0.14	76.31 \pm 0.51
LoRA-GA	85.70 \pm 0.09	94.11 \pm 0.18	80.57 \pm 0.20	93.18 \pm 0.06	85.29 \pm 0.24
LoRA-Pro	86.03 \pm 0.19	94.19 \pm 0.13	81.94 \pm 0.24	93.42 \pm 0.05	86.60 \pm 0.14
LoRA-One	85.89 \pm 0.08	94.53 \pm 0.13	82.04 \pm 0.22	93.37 \pm 0.02	87.83 \pm 0.37

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

(r = 8)	GSM8K		MMLU	HumanEval
	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)

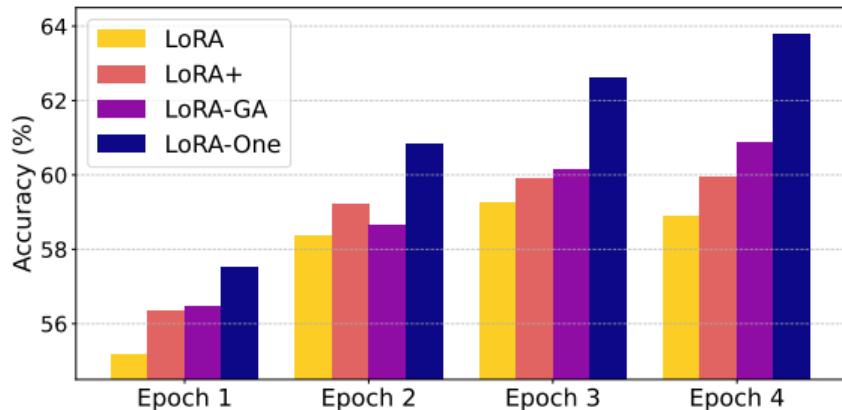


Figure 7: Accuracy comparison across different methods over epochs on GSM8K.

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	Precondition GD	(Spec-init.)	Linear convergence rate independent of $\kappa(\Delta)$

- subspace alignment
- global convergence

Proof of sketch: Control the dynamics for alignment

$$\underbrace{\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix}}_{:= \mathbf{Z}_{t+1}} = \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^\natural \\ \eta_2 \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix}}_{:= \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}}_{:= \mathbf{Z}_t} - \frac{1}{N} \begin{bmatrix} 0 & \eta_1 \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \eta_2 \mathbf{B}_t^\top \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}.$$

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

- Schur decomposition of \mathbf{H}
- obtain the dynamics of $\mathbf{Z}_t^{\text{lin}}$ (decouple $\mathbf{A}_t^{\text{lin}}$ and $\mathbf{B}_t^{\text{lin}}$ and obtain the alignment to \mathbf{G}^\natural)
- Define the residual term $\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}}$, control $\|\mathbf{E}_t\|_{op}$ in early stage
 $t < T_1 \sim \ln \left(\frac{\|\mathbf{G}^\natural\|_{op}}{\|\mathbf{A}_0\|_{op}^2} \right)$

◦ Transfer the alignment from $\mathbf{A}_t^{\text{lin}}$ to \mathbf{A}_t [4] (Stöger & Soltanolkotabi)

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

Proof of sketch: Control the dynamics for alignment

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} := \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^\natural \\ \eta_2 \mathbf{G}^\natural \top & \mathbf{I}_k \end{bmatrix}}_{:= \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}}_{:= \mathbf{Z}_t} - \frac{1}{N} \begin{bmatrix} 0 & \eta_1 \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \eta_2 \mathbf{B}_t^\top \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}.$$

- Approximated linear dynamical system $\mathbf{Z}_t^{\text{lin}} := \mathbf{H}^t \mathbf{Z}_0$
- Schur decomposition of \mathbf{H}
- obtain the dynamics of $\mathbf{Z}_t^{\text{lin}}$ (decouple $\mathbf{A}_t^{\text{lin}}$ and $\mathbf{B}_t^{\text{lin}}$ and obtain the alignment to \mathbf{G}^\natural)
- Define the residual term $\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}}$, control $\|\mathbf{E}_t\|_{op}$ in early stage
 $t < T_1 \sim \ln \left(\frac{\|\mathbf{G}^\natural\|_{op}}{\|\mathbf{A}_0\|_{op}^2} \right)$

◦ Transfer the alignment from $\mathbf{A}_t^{\text{lin}}$ to \mathbf{A}_t [4] (Stöger & Soltanolkotabi)

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

Proof of sketch: Control the dynamics for alignment

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} := \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^\natural \\ \eta_2 \mathbf{G}^\natural \top & \mathbf{I}_k \end{bmatrix}}_{:= \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}}_{:= \mathbf{Z}_t} - \frac{1}{N} \begin{bmatrix} 0 & \eta_1 \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \eta_2 \mathbf{B}_t^\top \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}.$$

- Approximated linear dynamical system $\mathbf{Z}_t^{\text{lin}} := \mathbf{H}^t \mathbf{Z}_0$
- Schur decomposition of \mathbf{H}
- obtain the dynamics of $\mathbf{Z}_t^{\text{lin}}$ (decouple $\mathbf{A}_t^{\text{lin}}$ and $\mathbf{B}_t^{\text{lin}}$ and obtain the alignment to \mathbf{G}^\natural)
- Define the residual term $\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}}$, control $\|\mathbf{E}_t\|_{op}$ in early stage
 $t < T_1 \sim \ln \left(\frac{\|\mathbf{G}^\natural\|_{op}}{\|\mathbf{A}_0\|_{op}^2} \right)$
- Transfer the alignment from $\mathbf{A}_t^{\text{lin}}$ to \mathbf{A}_t [4] (Stöger & Soltanolkotabi)
 $\|\mathbf{U}_{r^*, \perp}^\top(\mathbf{G}^\natural) \mathbf{U}_{r^*}(\mathbf{A}_t)\|_{op} \lesssim \|\mathbf{U}_{r^*, \perp}^\top(\mathbf{P}_t^A) \mathbf{U}_{r^*}(\mathbf{P}_t^A \mathbf{A}_0 + \mathbf{E}_t)\|_{op}$ is small, w.h.p.

Global convergence on nonlinear models

Recall problem setting and assumptions for nonlinear model

- Pre-trained model $f_{\text{pre}}(\mathbf{x}) = \sigma[(\mathbf{x}^\top \mathbf{W}^\natural)^\top] \in \mathbb{R}^k$
- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^\natural := \mathbf{W}^\natural + \Delta$ with $\text{Rank}(\Delta) = r^*$
- We assume $r = r^*$.
- Downstream well-behaved data $\tilde{\mathbf{y}} = \sigma[(\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\natural)^\top]$, $\{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$
- training loss

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \left\| \sigma \left(\tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{AB}) \right) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2.$$

- gradient updates

$$\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top, \quad \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t},$$

where we define

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

- additional assumptions on $\widetilde{\mathbf{W}}^\natural$, e.g., adapted weight has smaller signal than pre-trained model

Recall problem setting and assumptions for nonlinear model

- Pre-trained model $f_{\text{pre}}(\mathbf{x}) = \sigma[(\mathbf{x}^\top \mathbf{W}^\natural)^\top] \in \mathbb{R}^k$
- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^\natural := \mathbf{W}^\natural + \Delta$ with $\text{Rank}(\Delta) = r^*$
- We assume $r = r^*$.
- Downstream well-behaved data $\widetilde{\mathbf{y}} = \sigma[(\widetilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\natural)^\top]$, $\{\widetilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$
- training loss

$$\widetilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \left\| \sigma \left(\widetilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{AB}) \right) - \widetilde{\mathbf{Y}} \right\|_{\text{F}}^2.$$

- gradient updates

$$\nabla_{\mathbf{A}} \widetilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top, \quad \nabla_{\mathbf{B}} \widetilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t},$$

where we define

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

- additional assumptions on $\widetilde{\mathbf{W}}^\natural$, e.g., adapted weight has smaller signal than pre-trained model

Recall problem setting and assumptions for nonlinear model

- Pre-trained model $f_{\text{pre}}(\mathbf{x}) = \sigma[(\mathbf{x}^\top \mathbf{W}^\natural)^\top] \in \mathbb{R}^k$
- Unknown low-rank feature shift Δ : $\tilde{\mathbf{W}}^\natural := \mathbf{W}^\natural + \Delta$ with $\text{Rank}(\Delta) = r^*$
- We assume $r = r^*$.
- Downstream well-behaved data $\tilde{\mathbf{y}} = \sigma[(\tilde{\mathbf{x}}^\top \tilde{\mathbf{W}}^\natural)^\top]$, $\{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$
- training loss

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \left\| \sigma \left(\tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{AB}) \right) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2.$$

- gradient updates

$$\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top, \quad \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t},$$

where we define

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

- additional assumptions on $\tilde{\mathbf{W}}^\natural$, e.g., adapted weight has smaller signal than pre-trained model

Recall problem setting and assumptions for nonlinear model

- Pre-trained model $f_{\text{pre}}(\mathbf{x}) = \sigma[(\mathbf{x}^\top \mathbf{W}^\natural)^\top] \in \mathbb{R}^k$
- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^\natural := \mathbf{W}^\natural + \Delta$ with $\text{Rank}(\Delta) = r^*$
- We assume $r = r^*$.
- Downstream well-behaved data $\widetilde{\mathbf{y}} = \sigma[(\widetilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\natural)^\top]$, $\{\widetilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$
- training loss

$$\widetilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \left\| \sigma \left(\widetilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{AB}) \right) - \widetilde{\mathbf{Y}} \right\|_{\text{F}}^2.$$

- gradient updates

$$\nabla_{\mathbf{A}} \widetilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top, \quad \nabla_{\mathbf{B}} \widetilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t},$$

where we define

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

- additional assumptions on $\widetilde{\mathbf{W}}^\natural$, e.g., adapted weight has smaller signal than pre-trained model

Global convergence

Theorem (Linear convergence rate)

Under (Spec-init.) and J_{W_t} for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

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

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

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

$$\left\| J_{W_t} - \mathbb{E}_{\tilde{x}}[J_{W_t}] \right\|_F \lesssim \sqrt{d\epsilon} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F, \text{ w.h.p.} \Rightarrow \text{control } \mathbf{G}^\natural$$

Global convergence

Theorem (Linear convergence rate)

Under (Spec-init.) and $\mathbf{J}_{\mathbf{W}_t}$ for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

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

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

- low-rank approximation error $\leq 2\lambda_{r^*+1}(\mathbf{G}^\natural)$
- 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

$$\|\mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{x}}[\mathbf{J}_{\mathbf{W}_t}]\|_F \lesssim \sqrt{d\epsilon} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F, \text{ w.h.p.} \Rightarrow \text{control } \mathbf{G}^\natural$$

Global convergence

Theorem (Linear convergence rate)

Under (Spec-init.) and J_{W_t} for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

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

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

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

$$\left\| J_{W_t} - \mathbb{E}_{\tilde{x}}[J_{W_t}] \right\|_F \lesssim \sqrt{d\epsilon} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F, \text{ w.h.p.} \Rightarrow \text{control } \mathbf{G}^\natural$$

Global convergence

Theorem (Linear convergence rate)

Under (Spec-init.) and $\mathbf{J}_{\mathbf{W}_t}$ for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

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

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

- low-rank approximation error $\leq 2\lambda_{r^*+1}(\mathbf{G}^\natural)$
- 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

$$\|\mathbf{J}_{\mathbf{W}_t} - \mathbb{E}_{\tilde{x}}[\mathbf{J}_{\mathbf{W}_t}]\|_F \lesssim \sqrt{d\epsilon} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F, \text{ w.h.p.} \Rightarrow \text{control } \mathbf{G}^\natural$$

Global convergence

Theorem (Linear convergence rate)

Under (Spec-init.) and $\mathbf{J}_{\mathbf{W}_t}$ for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

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

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

- low-rank approximation error $\leq 2\lambda_{r^*+1}(\mathbf{G}^\natural)$
- 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

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

Global convergence

Theorem (Linear convergence rate)

Under (Spec-init.) and $\mathbf{J}_{\mathbf{W}_t}$ for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

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

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

- low-rank approximation error $\leq 2\lambda_{r^*+1}(\mathbf{G}^\natural)$
- 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

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

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

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

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

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

- transformed to lower bound $\|\mathbf{L}_{\perp}^{\top} \Delta \mathbf{L}\|_{\text{F}}^2$
- upper bound $\|\mathbf{L}_{\perp}^{\top} \mathbf{U}\|_{\text{op}} < 1$ by Wedin's sin- θ theorem

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

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

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

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

- transformed to lower bound $\left\| \mathbf{L}_{\perp}^{\top} \Delta \mathbf{L} \right\|_{\text{F}}^2$
- upper bound $\left\| \mathbf{L}_{\perp}^{\top} \mathbf{U} \right\|_{op} < 1$ by Wedin's sin- θ theorem

Takeaway messages

- *LoRA-One: One-step full gradient could suffice for fine-tuning large language models, provably and efficiently.* ICML'25 Oral. [code](#)
- subspace alignment: \mathbf{G}^\natural and $(\mathbf{A}_t, \mathbf{B}_t)$ \Rightarrow theory-grounded algorithm design
- “optimal” non-zero initialization strategy
- clarification on gradient alignment based algorithms

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

Takeaway messages

- *LoRA-One: One-step full gradient could suffice for fine-tuning large language models, provably and efficiently.* ICML'25 Oral. [code](#)
- subspace alignment: \mathbf{G}^\natural and $(\mathbf{A}_t, \mathbf{B}_t)$ \Rightarrow theory-grounded algorithm design
- “optimal” non-zero initialization strategy
- clarification on gradient alignment based algorithms

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

Takeaway messages

- *LoRA-One: One-step full gradient could suffice for fine-tuning large language models, provably and efficiently.* ICML'25 Oral. [code](#)
- subspace alignment: \mathbf{G}^\natural and $(\mathbf{A}_t, \mathbf{B}_t)$ \Rightarrow theory-grounded algorithm design
- “optimal” non-zero initialization strategy
- clarification on gradient alignment based algorithms

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

-  Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp.
Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions.
SIAM review, 53(2):217–288, 2011.
-  Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.
LoRA: Low-rank adaptation of large language models.
In *International Conference on Learning Representations*, 2022.

-  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.

-  Shaowen Wang, Linxi Yu, and Jian Li.
LoRA-GA: Low-rank adaptation with gradient approximation.
In *Advances in Neural Information Processing Systems*, 2024.
-  Zhengbo Wang, Jian Liang, Ran He, Zilei Wang, and Tieniu Tan.
LoRA-Pro: Are Low-Rank Adapters Properly Optimized?
In *The Twelfth International Conference on Learning Representations*, 2025.