

# Low-Rank Training and Low-Rank Adapters

Timon Klein and Piotr Minakowski  
Otto von Guericke University Magdeburg

March 10th, 2026  
MLaftermath



OTTO VON GUERICKE  
UNIVERSITÄT  
MAGDEBURG

## A geometric framework for momentum-based optimizers for low-rank training

---

Steffen Schotthöfer\* Timon Klein<sup>†</sup> and Jonas Kusch<sup>‡</sup>

\*Computer Science and Mathematics Division; Oak Ridge National Laboratory; Oak Ridge, TN 37831 USA;  
Mail correspondence: [schotthofers@ornl.gov](mailto:schotthofers@ornl.gov)

<sup>†</sup>Department of Mathematics; Otto von Guericke University Magdeburg; 39106 Magdeburg; Germany

<sup>‡</sup>Scientific Computing; Norwegian University of Life Sciences; Drøbakveien 31, 1433 Ås; Norway

# Authors



Steffen Schotthöfer

Householder Fellow in the  
Mathematics in Computation Section  
at Oak Ridge National Laboratory

[schotthofers@ornl.gov](mailto:schotthofers@ornl.gov)



Timon Klein

PhD Student in Mathematical  
Optimization at Otto von Guericke  
University Magdeburg

[timon.klein@ovgu.de](mailto:timon.klein@ovgu.de)

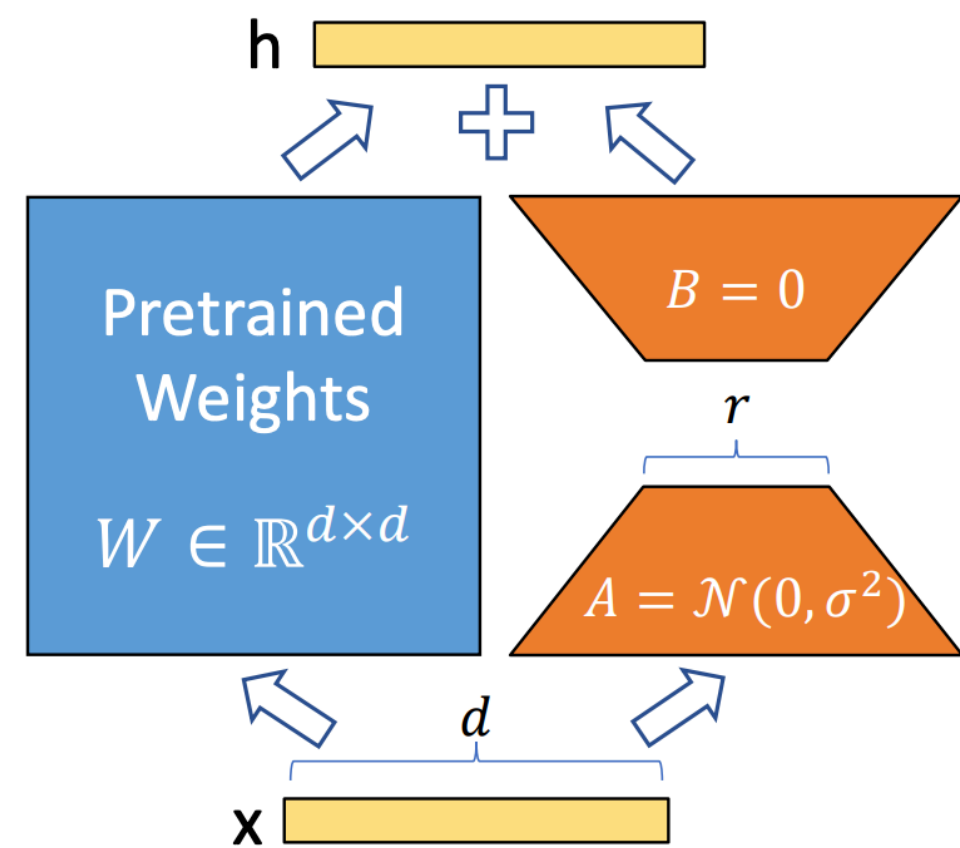


Jonas Kusch

Associate Professor in Scientific  
Computing at Norwegian University of  
Life Sciences

[jonas.kusch@nmbu.no](mailto:jonas.kusch@nmbu.no)

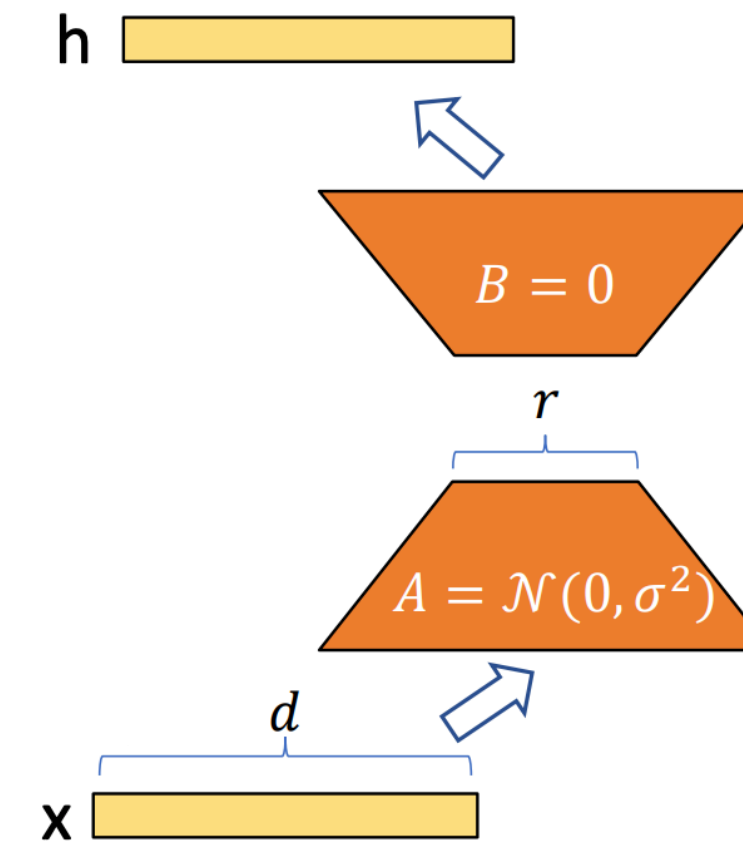
# Low Rank Compression and Fine-Tuning of Neural Networks



## Low Rank Finetuning

- Hu et al., *LoRA: Low-Rank Adaptation of Large Language Models*. 2021
- Zhang et al., *AdaLoRA: Adaptive Budget Allocation for Parameter-Efficient Fine-Tuning*. 2023
- Liu et al., *DoRA: Weight-Decomposed Low-Rank Adaptation*. 2024

$$z(x) = \sigma(Wx + AB^T x)$$



## Low Rank Compression

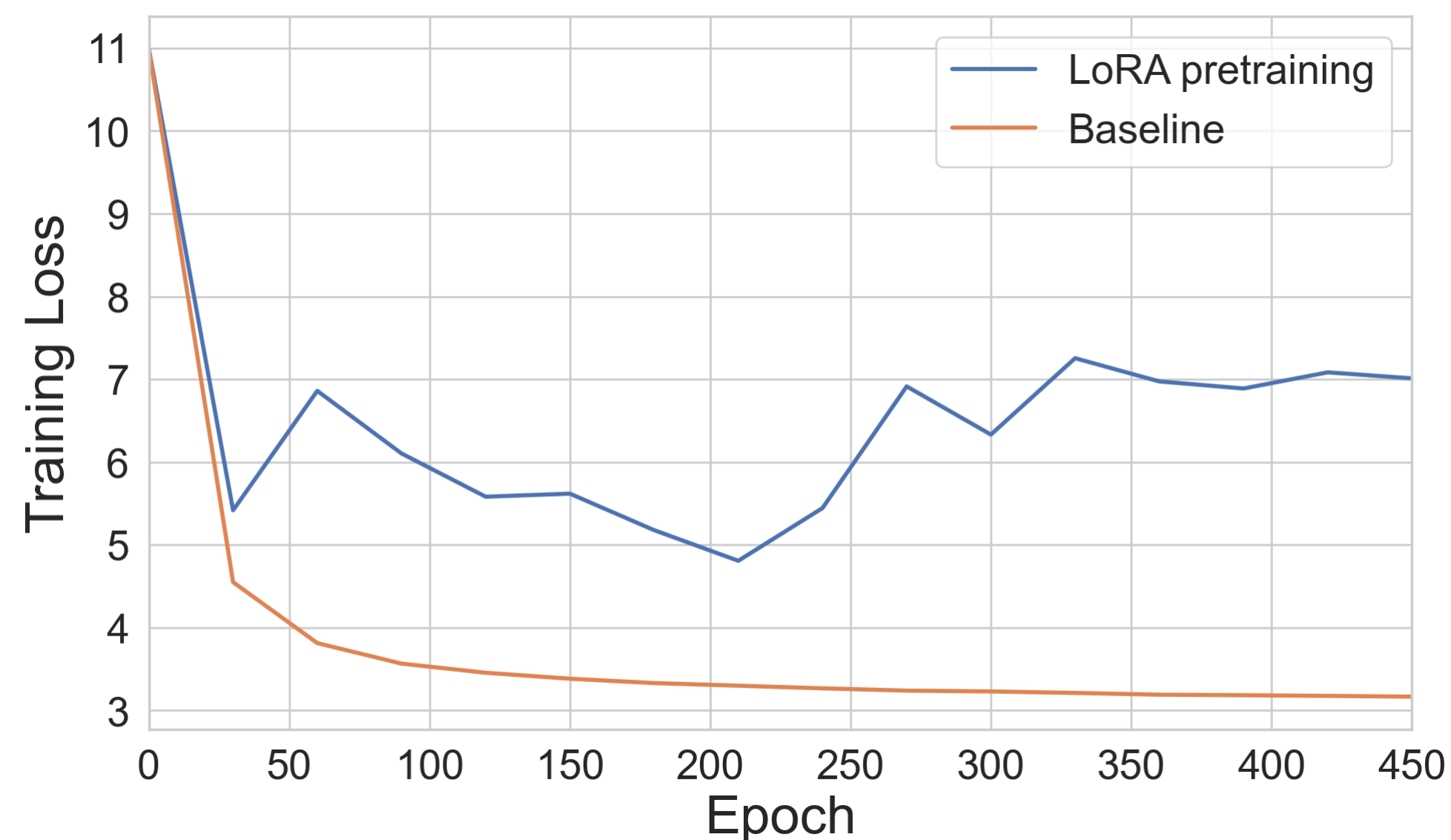
- Denton et al., *Exploiting Linear Structure Within Convolutional Networks*. 2014
- A. Novikov, et al., *Tensorizing neural networks*. 2015
- A. Tjandra, S. Sakti, and S. Nakamura. *Compressing recurrent neural network with tensor train*. 2017

$$z(x) = \sigma(AB^T x)$$

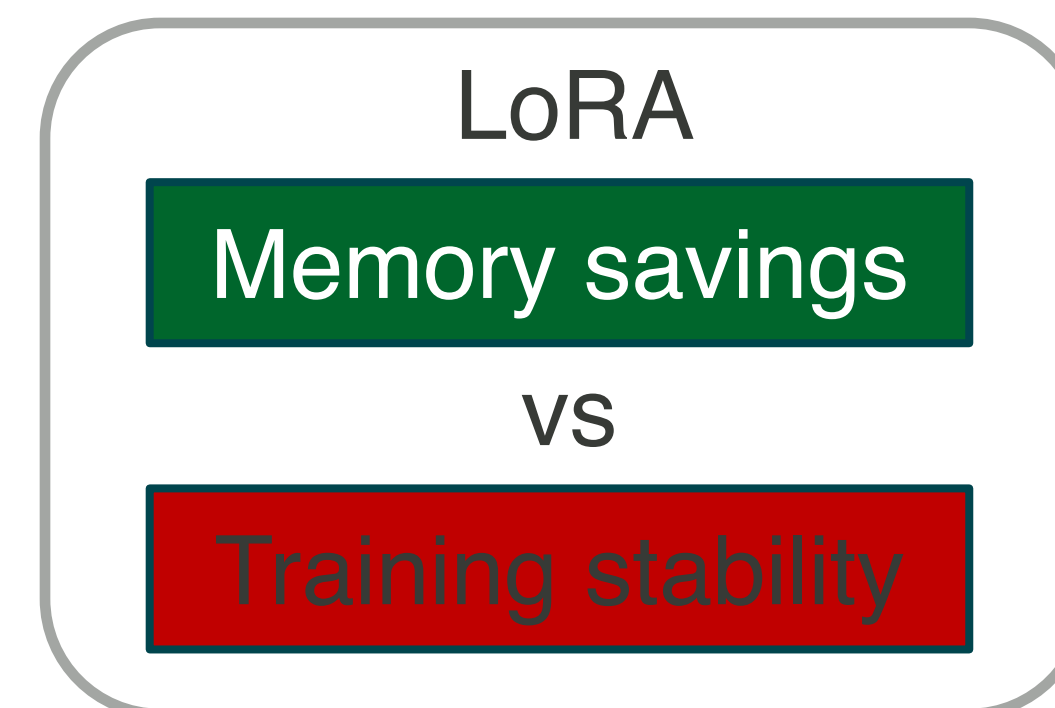
## Low Rank Attention

- DeepSeek-AI, "Multihhead Latent Attention". 2024
- Ainslie, et al., *GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints*. '23

$$z(x) = \sigma(xAB^T x^T)$$



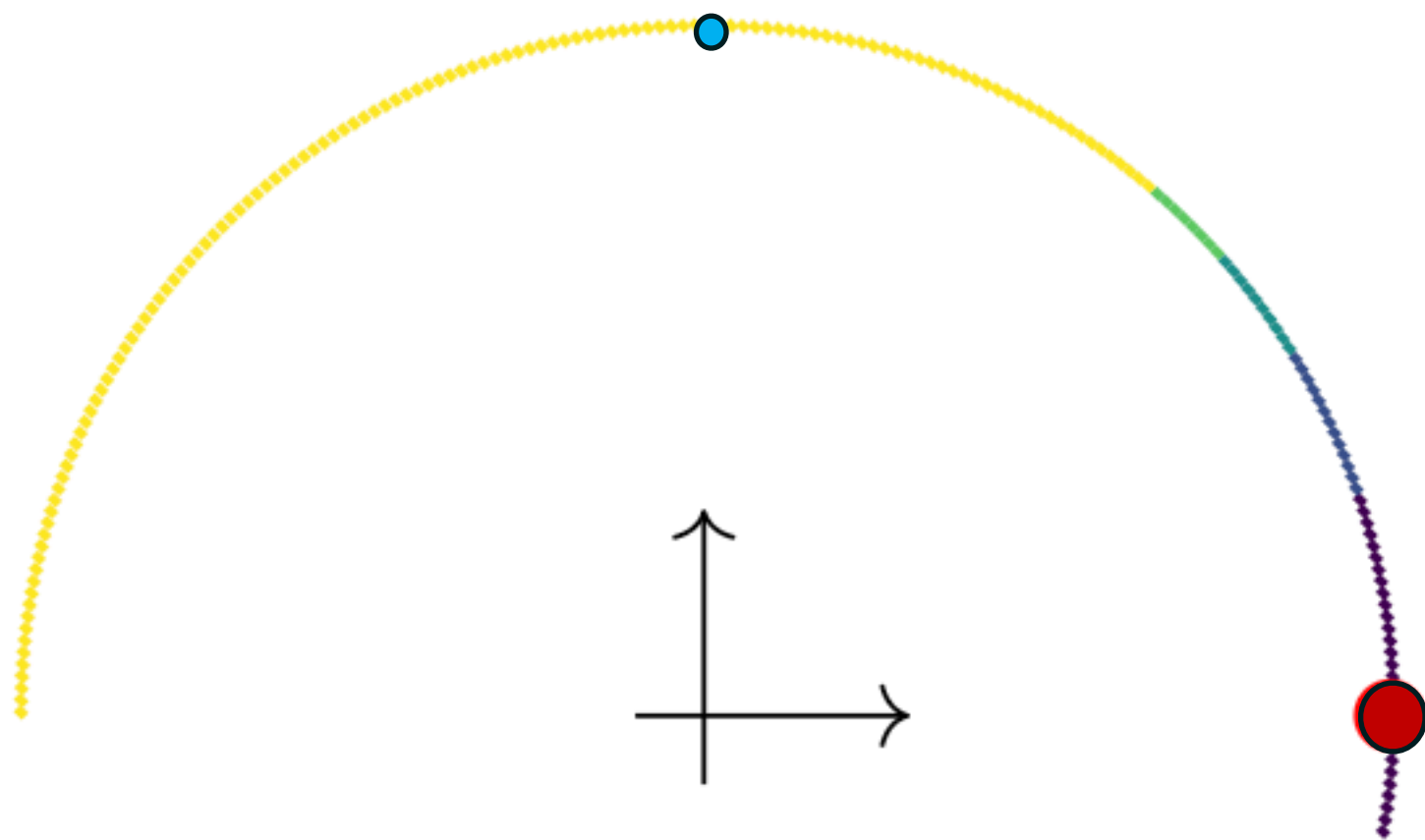
GPT-2 on OWT, low-rank MLP: Low rank training stalls. Why?



# Thought Experiment: Manifold Constrained Optimization

$$\min_{w \in \mathcal{M}} \mathcal{L}(w) = \frac{1}{2} \left\| [1,0] - w \right\|_2^2$$

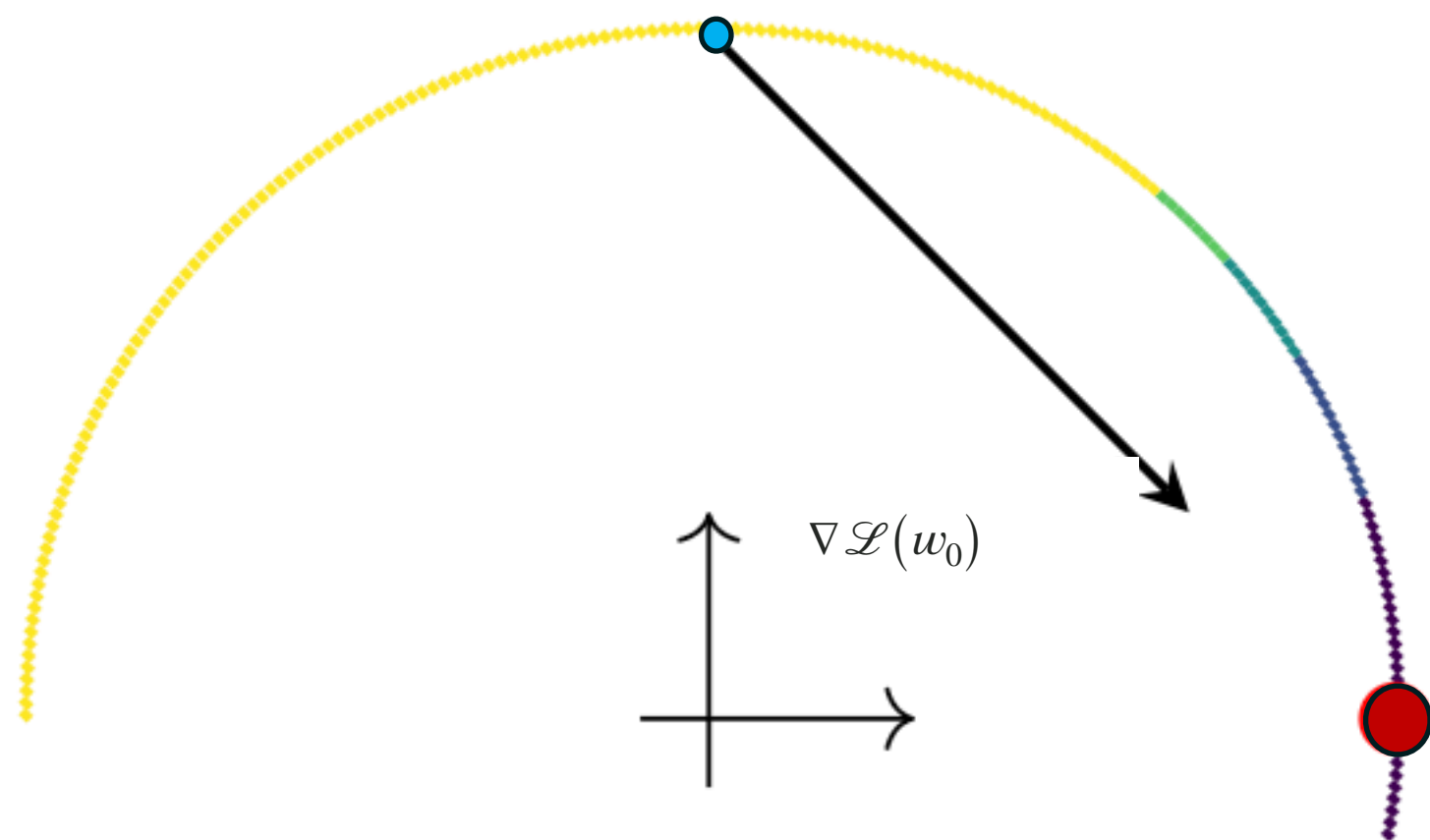
- Manifold: Unit circle  $\mathcal{M} = \{w \in \mathbb{R}^2 : \|w\| = 1\}$
- Initialization at  $w_0 = [0,1]$  • solution at  $w_* = [1,0]$  •



# Thought experiment: Manifold Constrained Optimization

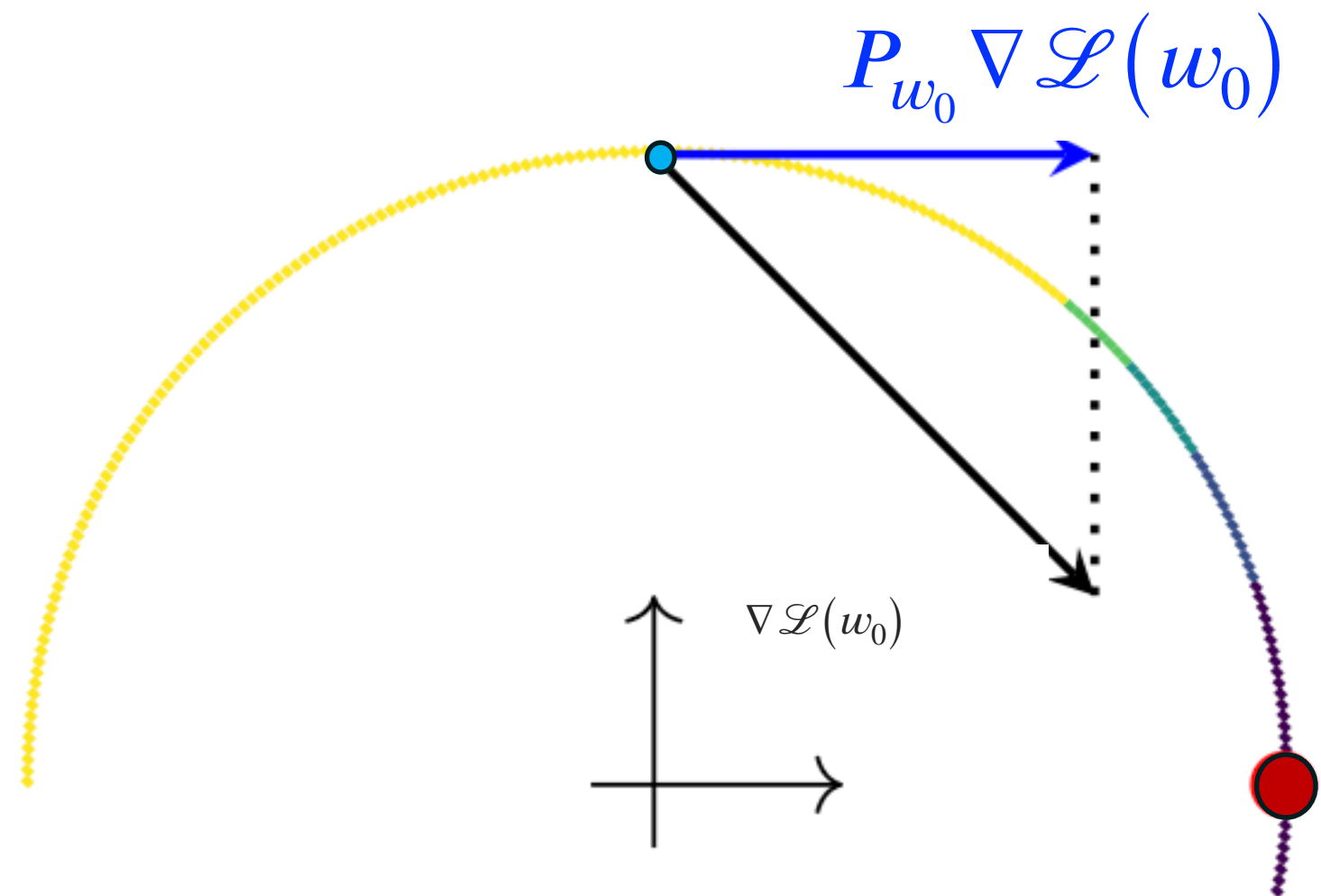
$$\min_{w \in \mathcal{M}} \mathcal{L}(w) = \frac{1}{2} \left\| [1, 0] - w \right\|_2^2$$

- Manifold: Unit circle  $\mathcal{M} = \{ w \in \mathbb{R}^2 : \|w\| = 1 \}$
- Initialization at  $w_0 = [0, 1]$  solution at  $w_* = [1, 0]$
- Gradient  $\nabla \mathcal{L}(w_0) = [-1, 1]$



# Thought experiment: Manifold Constrained Optimization

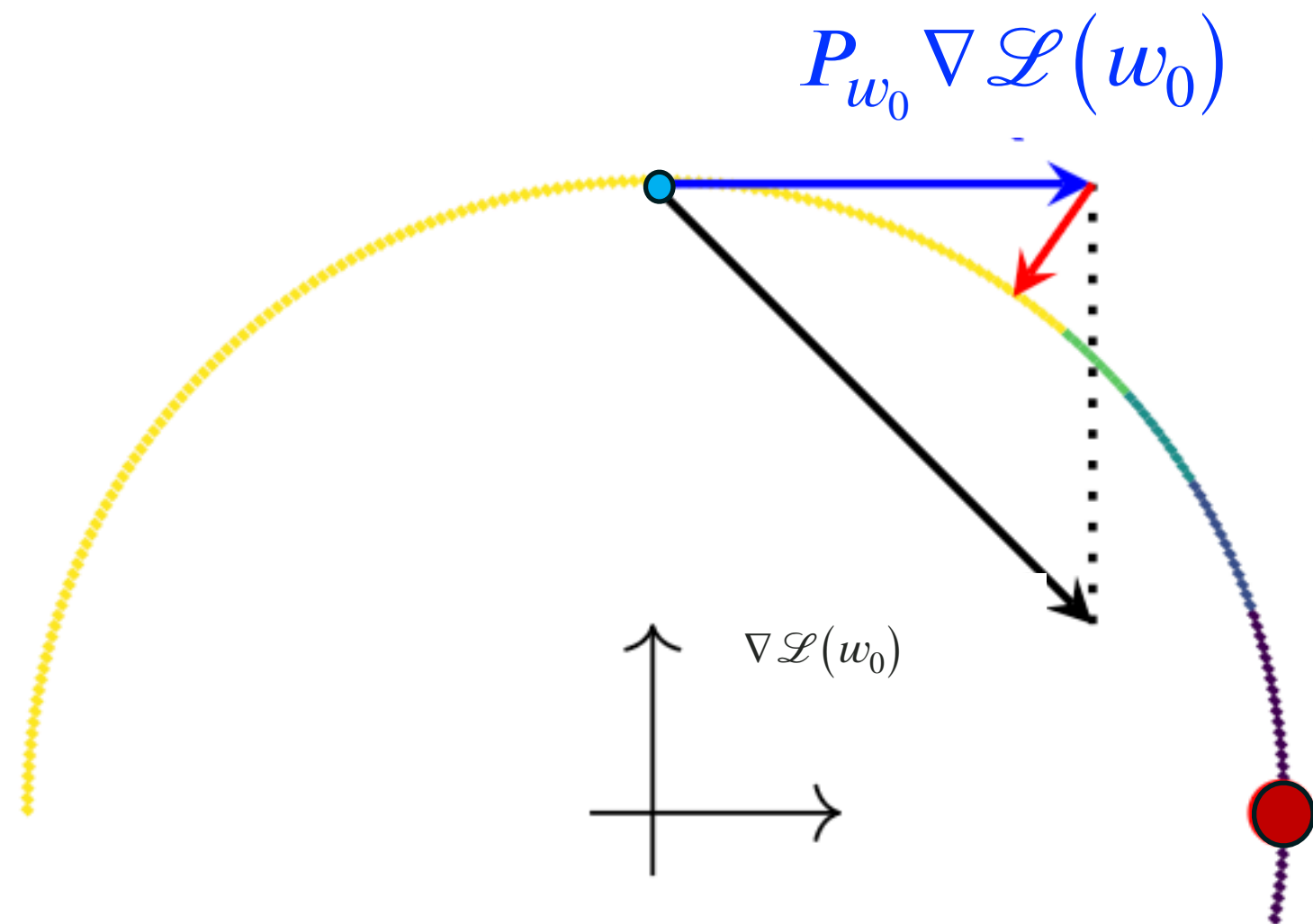
$$\min_{w \in \mathcal{M}} \mathcal{L}(w) = \frac{1}{2} \left\| [1,0] - w \right\|_2^2$$



- Manifold: Unit circle  $\mathcal{M} = \{w \in \mathbb{R}^2 : \|w\| = 1\}$
- Initialization at  $w_0 = [0,1]$  • solution at  $w_* = [1,0]$
- Gradient  $\nabla \mathcal{L}(w_0) = [-1,1]$
- Riemannian gradient  $P_{w_0} \nabla \mathcal{L}(w_0) = [0,1]$   
→ orthogonal projection  $P_{w_0}$

# Thought experiment: Manifold Constrained Optimization

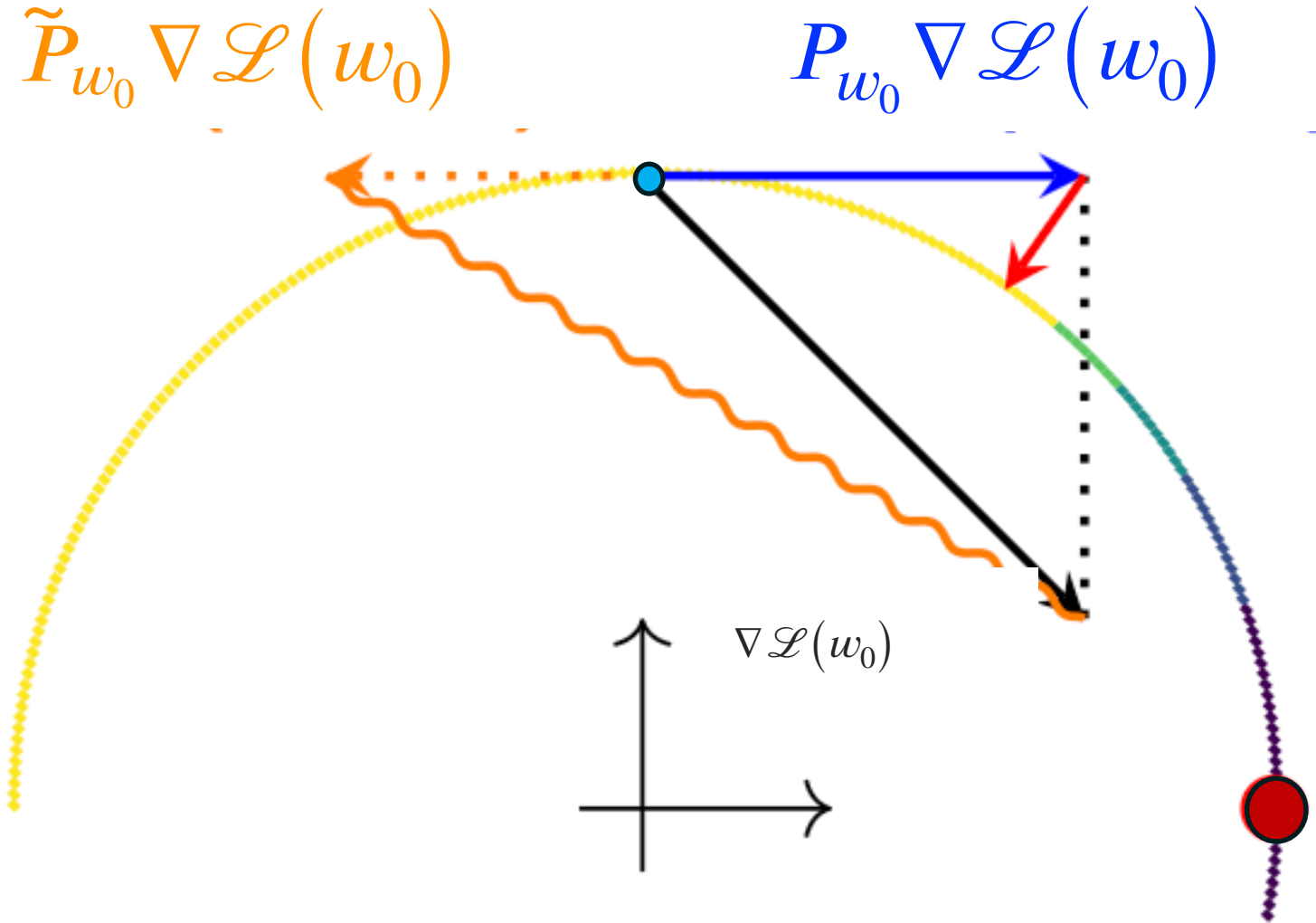
$$\min_{w \in \mathcal{M}} \mathcal{L}(w) = \frac{1}{2} \left\| [1,0] - w \right\|_2^2$$



- Manifold: Unit circle  $\mathcal{M} = \{w \in \mathbb{R}^2 : \|w\| = 1\}$
- Initialization at  $w_0 = [0,1]$  • solution at  $w_* = [1,0]$  •
- Gradient  $\nabla \mathcal{L}(w_0) = [-1,1]$
- Riemannian gradient  $P_{w_0} \nabla \mathcal{L}(w_0) = [0,1]$   
→ orthogonal projection  $P_{w_0}$
- Retraction onto unit circle  $\mathcal{M}$

# Thought experiment: Manifold Constrained Optimization

$$\min_{w \in \mathcal{M}} \mathcal{L}(w) = \frac{1}{2} \left\| [1,0] - w \right\|_2^2$$



- Manifold: Unit circle  $\mathcal{M} = \{w \in \mathbb{R}^2 : \|w\| = 1\}$
- Initialization at  $w_0 = [0,1]$  • solution at  $w_* = [1,0]$  •
- Gradient  $\nabla \mathcal{L}(w_0) = [-1, 1]$
- Riemannian gradient  $P_{w_0} \nabla \mathcal{L}(w_0) = [0, 1]$   
→ orthogonal projection  $P_{w_0}$
- **Retraction** onto unit circle  $\mathcal{M}$
- Non orthogonal  $\tilde{P}_{w_0}$  breaks the structure

$$\dot{w} = -\nabla \mathcal{L}$$

gradient flow

vs

$$\dot{w} = -P_w \nabla \mathcal{L}$$

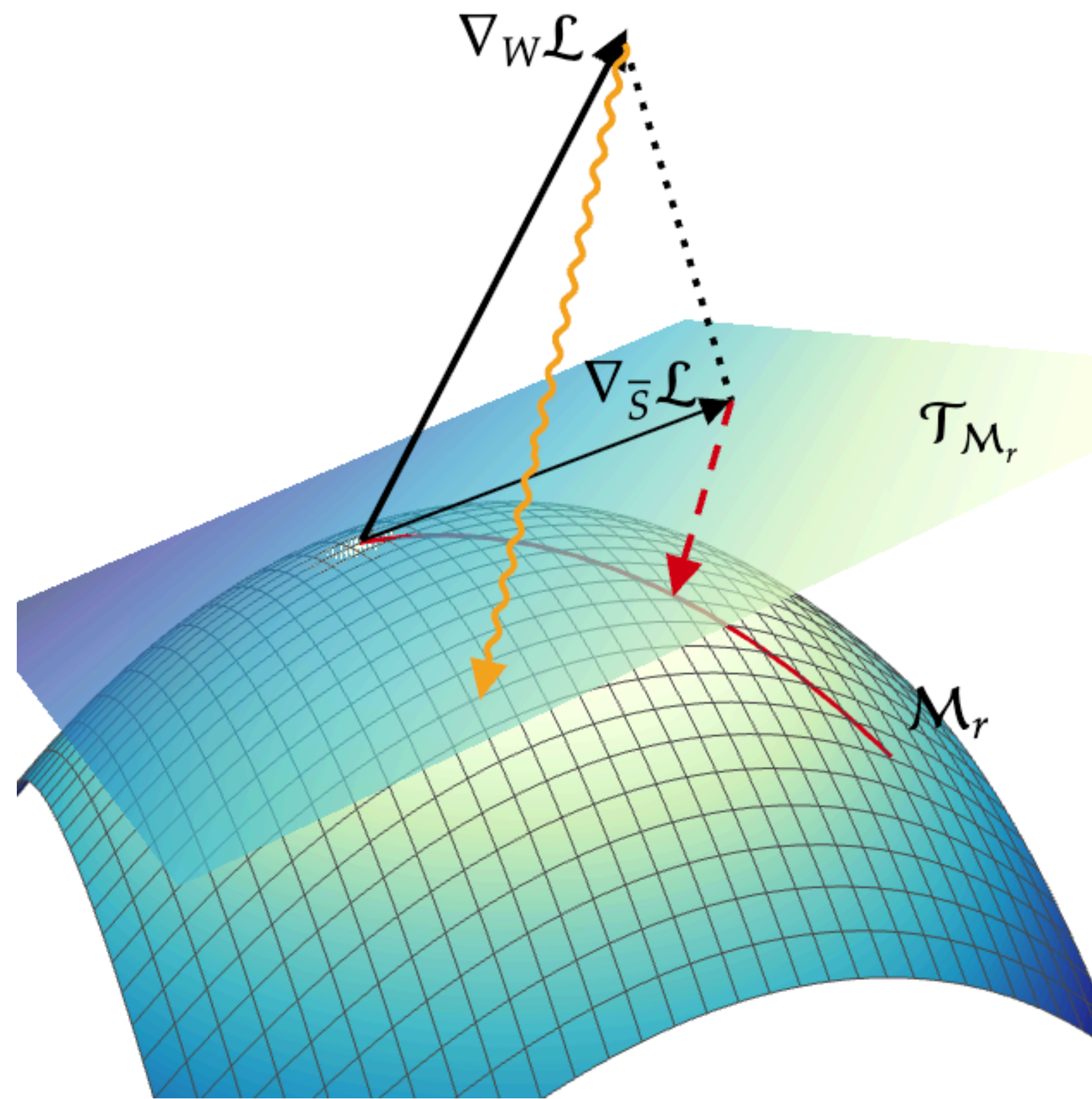
Riemannian gradient flow

vs

$$\dot{w} = -\tilde{P}_w \nabla \mathcal{L}$$

# Low Rank training is manifold constrained optimization

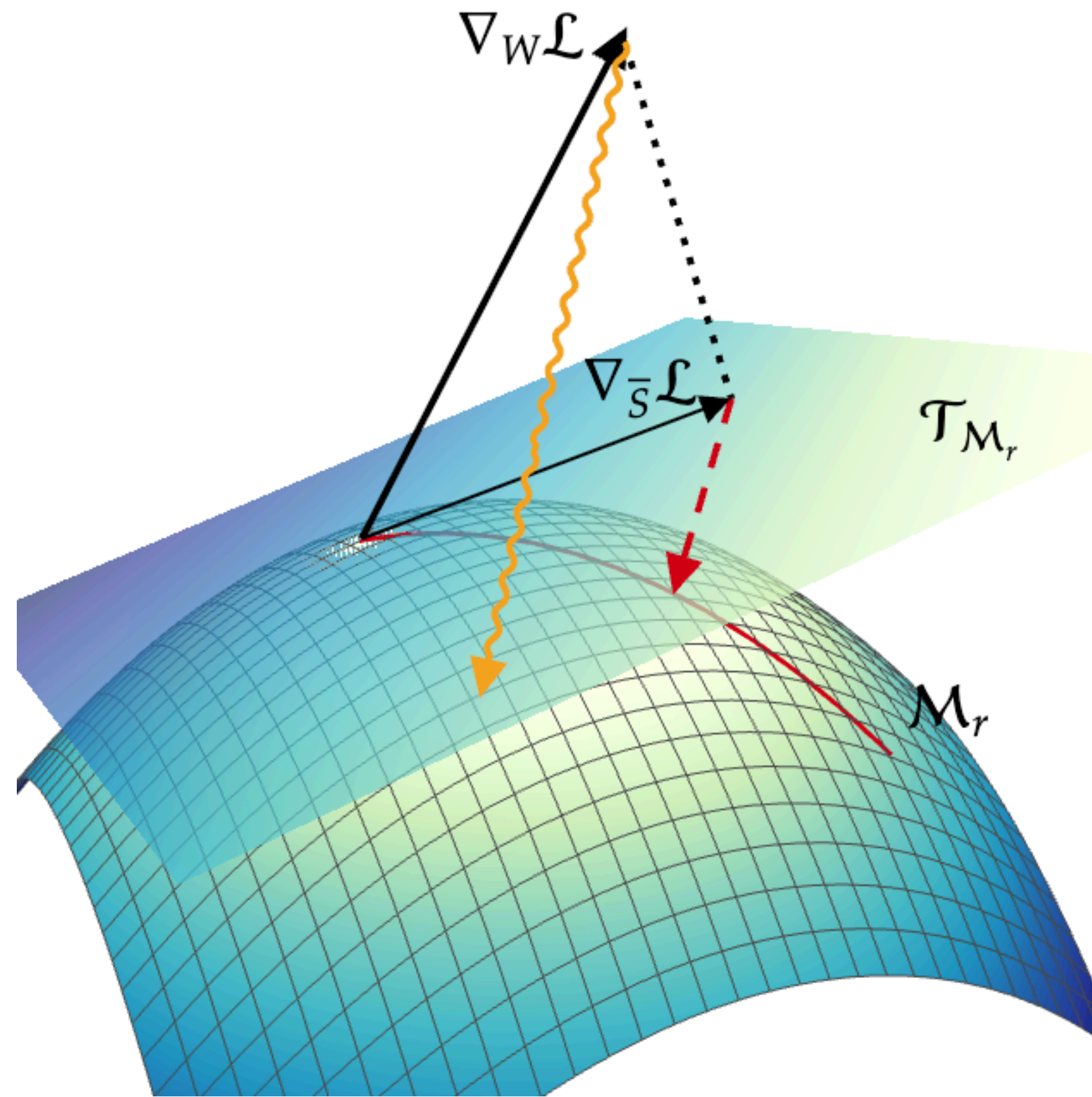
$$\min_{W \in \mathcal{M}} \mathcal{L}(X, Y; W)$$



- Manifold  $\mathcal{M} = \{W \in \mathbb{R}^{n \times n} : \text{rank}(W) = r\}$
- LoRA ansatz:  $W = AB^T$  with  $A, B \in \mathbb{R}^{n \times r}$
- Gradient flow:  $\dot{W} = \dot{A}B^T + A\dot{B}^T = \tilde{P}_W \nabla_W \mathcal{L}$  (chain rule)
- $\tilde{P}_W$  is defined by  $[A, \dot{A}]$ , and  $[B, \dot{B}]$ 
  - not orthogonal
  - no steepest descent on  $\mathcal{M}$

# Low Rank training is manifold constrained optimization

$$\min_{W \in \mathcal{M}} \mathcal{L}(X, Y; W)$$



- Manifold  $\mathcal{M} = \{W \in \mathbb{R}^{n \times n} : \text{rank}(W) = r\}$
- AdaLoRA ansatz:  $W = USV^T$  with  $U, V \in \mathbb{R}^{n \times r}, S \in \mathbb{R}^{r \times r}$   
chain rule
- Gradient flow:  $\dot{W} = \dot{U}SV^T + U\dot{S}V^T + US\dot{V}^T = \tilde{P}_W \nabla_W \mathcal{L}$
- $\tilde{P}_W$  is defined by  $[U, \dot{U}]$ , and  $[V, \dot{V}]$ 
  - not orthogonal
  - no steepest descent on  $\mathcal{M}$

# DLRT: Dynamical Low Rank Training

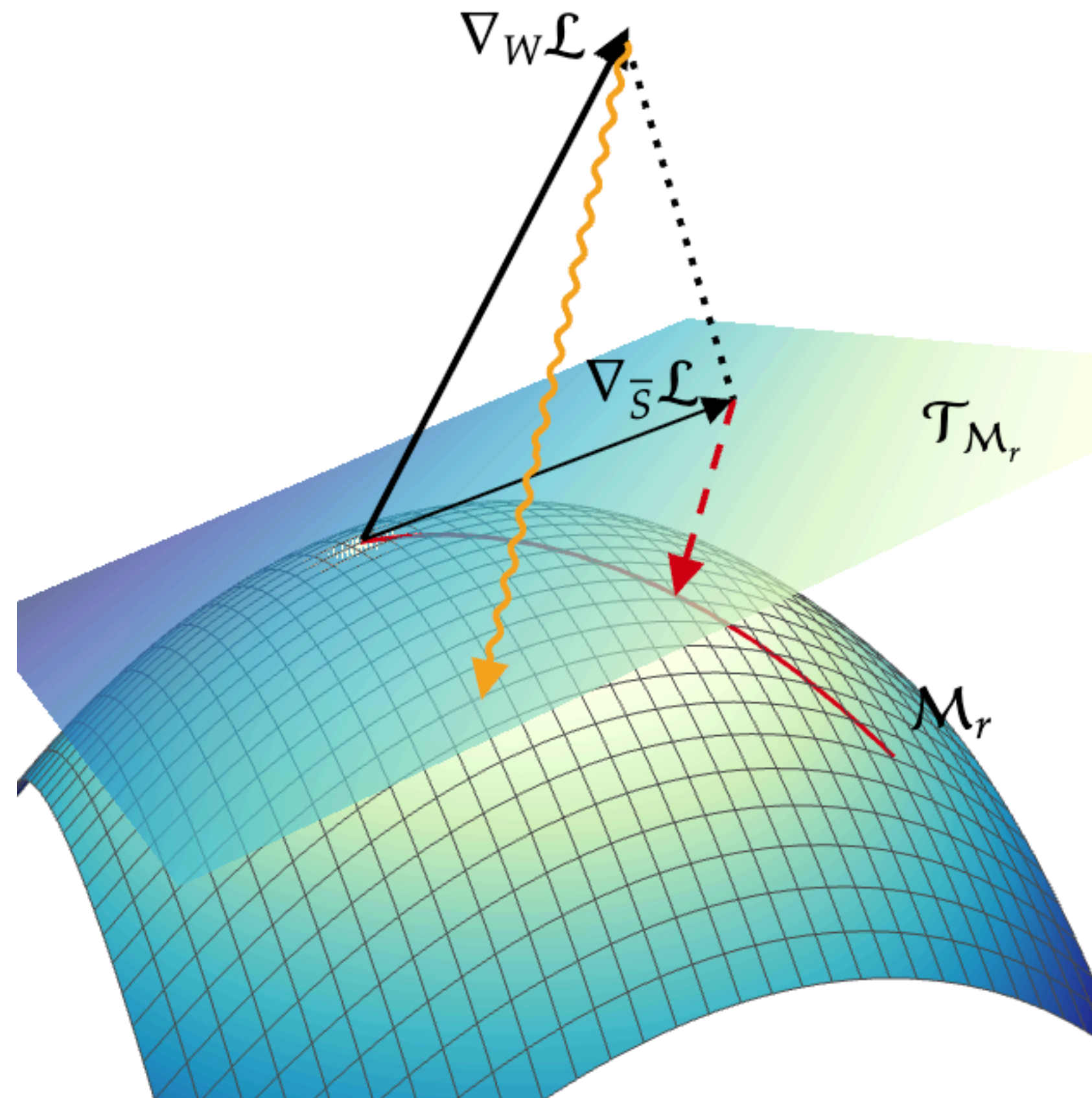
Efficient evolution of projected gradient flow

S.S., Zangrando, Kusch, Ceruti, Tudisco; *Low-Rank Lottery Tickets ...*; NeurIPS 2022

$$\dot{W} = P_W \nabla \mathcal{L}$$

$$\min_{W \in \mathcal{M}} \mathcal{L}(X, Y; W)$$

- Manifold  $\mathcal{M} = \{W \in \mathbb{R}^{n \times n} : \text{rank}(W) = r\}$
  - DLRT ansatz:  $W = USV^T$  with  $U, V \in \mathbb{R}^{n \times r}, S \in \mathbb{R}^{r \times r}$
  - Gradient flow:  $\dot{W} = \dot{U}S V^T + U\dot{S}V^T + US\dot{V}^T = P_W \nabla_W \mathcal{L}$
  - Construct  $P_W$  with orthogonal bases  
 $\bar{U} = \text{ortho}\{[U, \dot{U}]\}$ , and  $\bar{V} = \text{ortho}\{[V, \dot{V}]\}$
- Basis for tangent space  $\mathcal{T}_{\mathcal{M}}$   
 → enables steepest descent on  $\mathcal{M}$



Memory cost – slightly better than LoRA

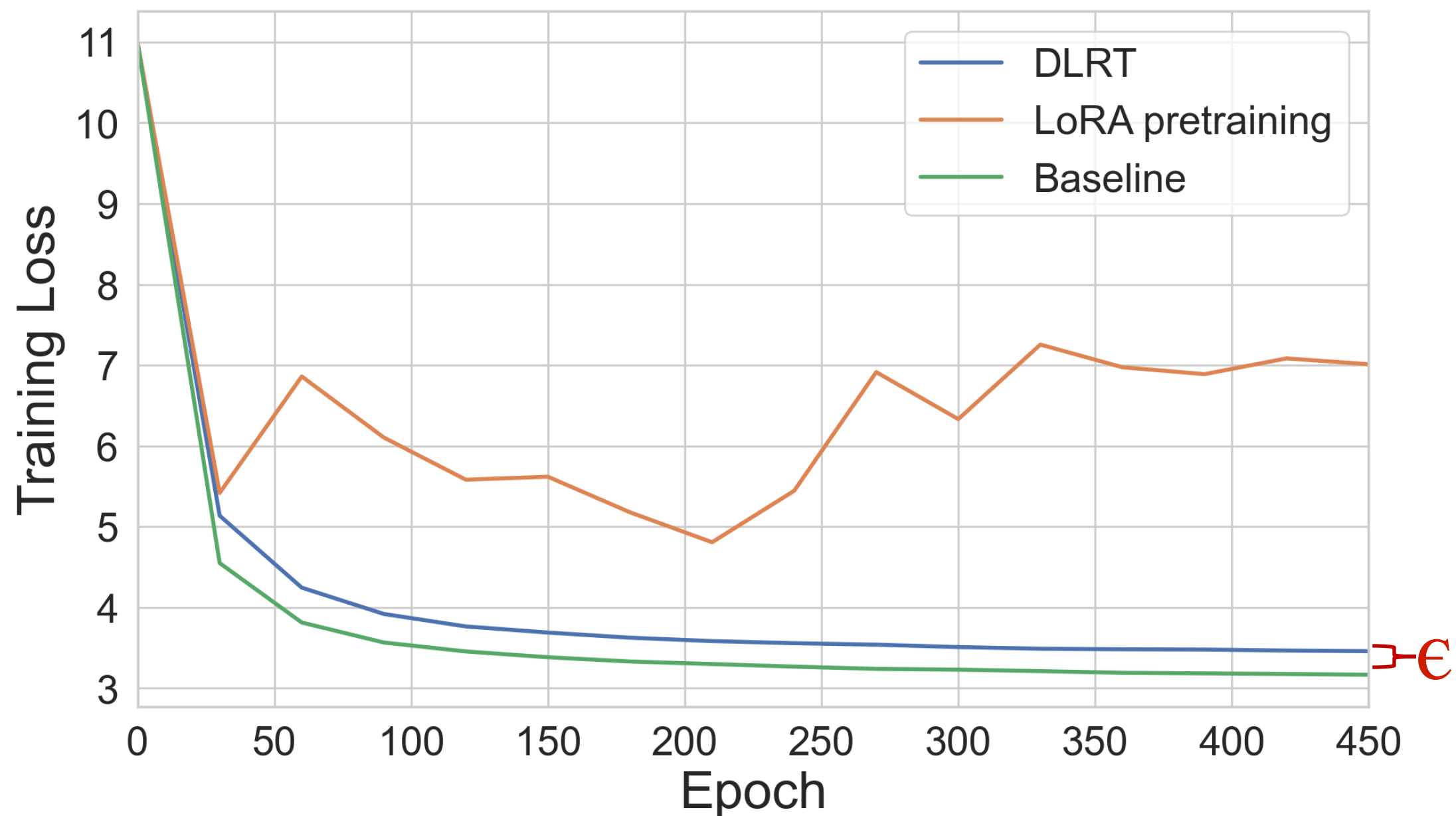
- Weights:  $\mathcal{O}(2nr + r^2)$
- Gradients:  $\mathcal{O}(2nr)$  for basis update  $\mathcal{O}(r^2)$  for optimization
- Optimizer states:  $\mathcal{O}(r^2)$

# DLRT: Dynamical Low Rank Training

Schotthöfer., Zangrando, Kusch, Ceruti, Tudisco; *Low-Rank Lottery Tickets ...* ; NeurIPS 2022

Zangrando , S.S., Ceruti, Kusch, Tudisco; *Geometry-aware training [...] in tensor Tucker format* ; NeurIPS '24

S.S., Zangrando Ceruti, Kusch, Tudisco; *GeoLoRA [...]*; ICLR '25



GPT-2 on OWT, low-rank MLP: DLRT beats LoRA

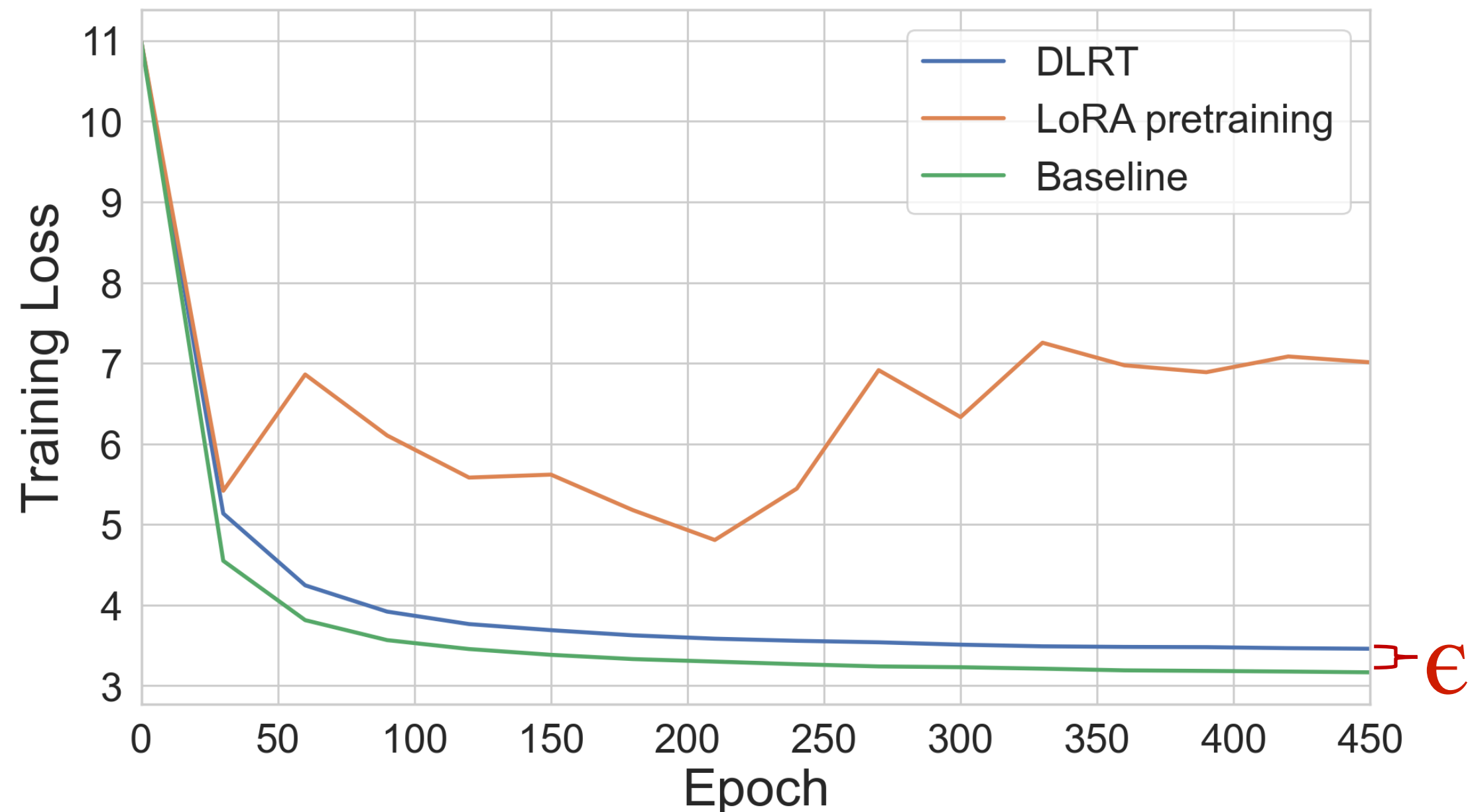
- DLRT inherits training robustness from full training
  - Provable: Optimality, loss descent, convergence
  - Hyperparameter can be transferred from full training
- Provable error bound to full rank training
$$\|W_{\text{full rank}}(t) - W_{\text{DLRT}}(t)\| < \epsilon(\lambda, \vartheta)$$
- Automatic rank selection (like AdaLoRA)

# DLRT: Dynamical Low Rank Training

Schotthöfer, Zangrando, Kusch, Ceruti, Tudisco; *Low-Rank Lottery Tickets ...*; NeurIPS 2022

Zangrando, Schotthöfer, Ceruti, Kusch, Tudisco; *Geometry-aware training [...] in tensor Tucker format*; NeurIPS '24

Schotthöfer, Zangrando, Ceruti, Tudisco, Kusch; *GeoLoRA [...]*; ICLR '25



GPT-2 on OWT, low-rank MLP: DLRT beats LoRA

- DLRT inherits training robustness from full training
  - Provable: Optimality, loss descent, convergence
  - Hyperparameter can be transferred from full training
- Provable error bound to full rank training
 
$$\|W_{\text{full rank}}(t) - W_{\text{DLRT}}(t)\| < \epsilon(\lambda, \vartheta)$$
- Automatic rank selection (like AdaLoRA)

## What about momentum methods/Adam?

Schotthöfer, Klein, Kusch; *A geometric framework for momentum-based optimizers for low-rank training*; NeurIPS '25

Momentum gradient flow

$$\begin{aligned} \dot{W} &= \mathcal{V} \\ \dot{\mathcal{V}} + \gamma \mathcal{V} &= -\nabla_W \mathcal{L} \end{aligned}$$

DLRT Momentum gradient flow

$$\begin{aligned} \dot{W} &= P_W \mathcal{V} \\ \dot{\mathcal{V}} + \gamma \mathcal{V} &= -P_W \nabla_W \mathcal{L} \end{aligned}$$

- Bases  $U, V$  for  $W$  can be re-used for momentum terms
- Extendable for Adam, AdamW

# Low-Rank Adapters

---

## Mitigating Subject Dependency in EEG Decoding with Subject-Specific Low-Rank Adapters

---

Timon Klein<sup>1</sup> Piotr Minakowski<sup>1</sup> Sebastian Sager<sup>1,2</sup> Steffen Schotthöfer<sup>3</sup>

---

<sup>1</sup>Department of Mathematics; Otto von Guericke University Magdeburg; 39106 Magdeburg; Germany. <sup>2</sup>Max Planck Institute for Dynamics of Complex Technical Systems; 39106 Magdeburg; Germany. <sup>3</sup>Computer Science and Mathematics Division; Oak Ridge National Laboratory; Oak Ridge, TN 37831; USA. Correspondence to: Timon Klein <timon.klein@ovgu.de>.



# Authors



Timon Klein

PhD Student in Mathematical  
Optimization at Otto von Guericke  
University Magdeburg

timon.klein@ovgu.de



Piotr Minakowski

AI Expert for Energy Load  
Forecasting and Balancing

pminak@gmail.com



Sebastian Sager

Professor in Mathematical  
Optimization at Otto von  
Guericke University  
Magdeburg

Max Planck Fellow at MPI  
Magdeburg

sager@ovgu.de

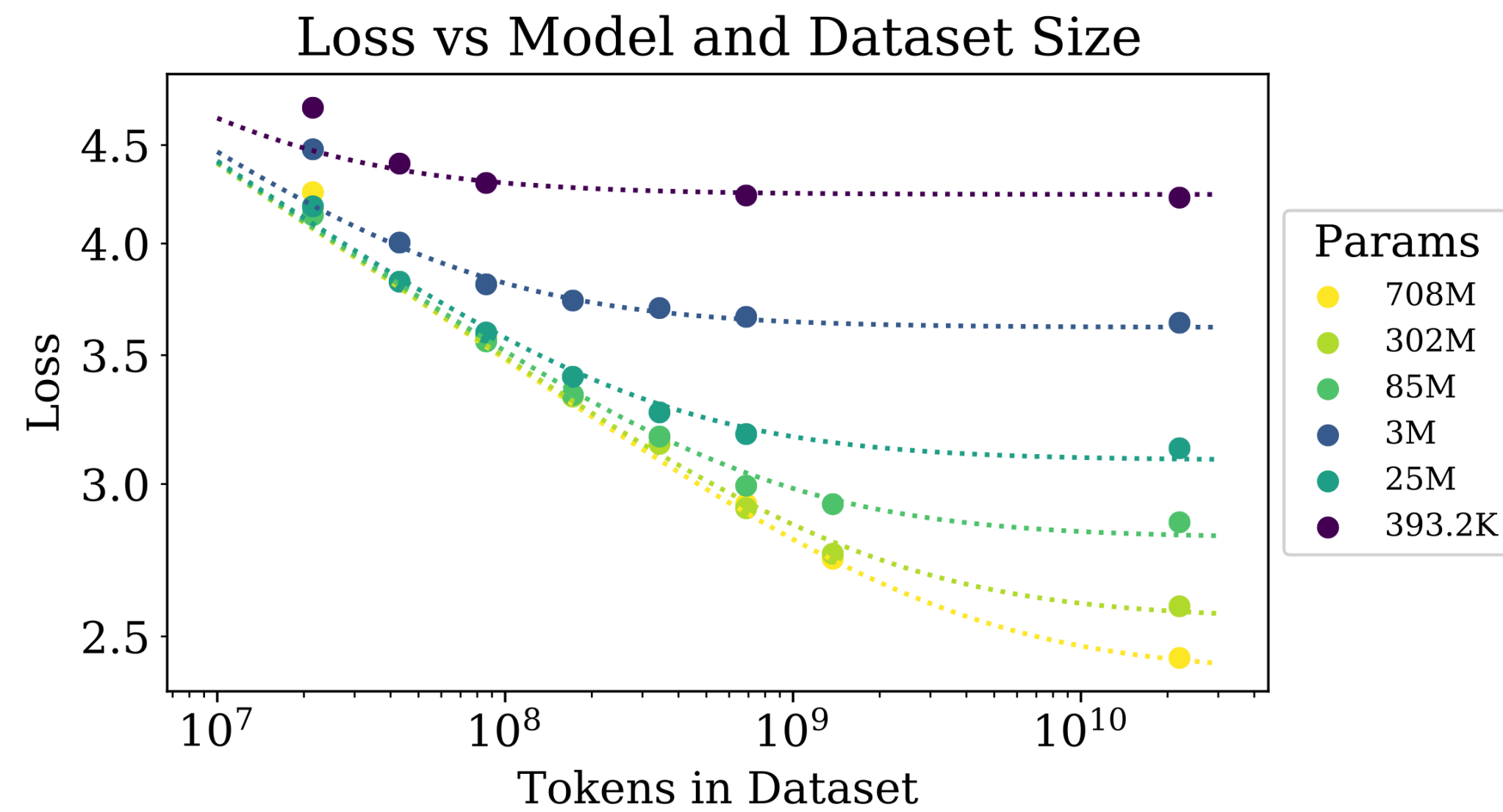


Steffen Schotthöfer

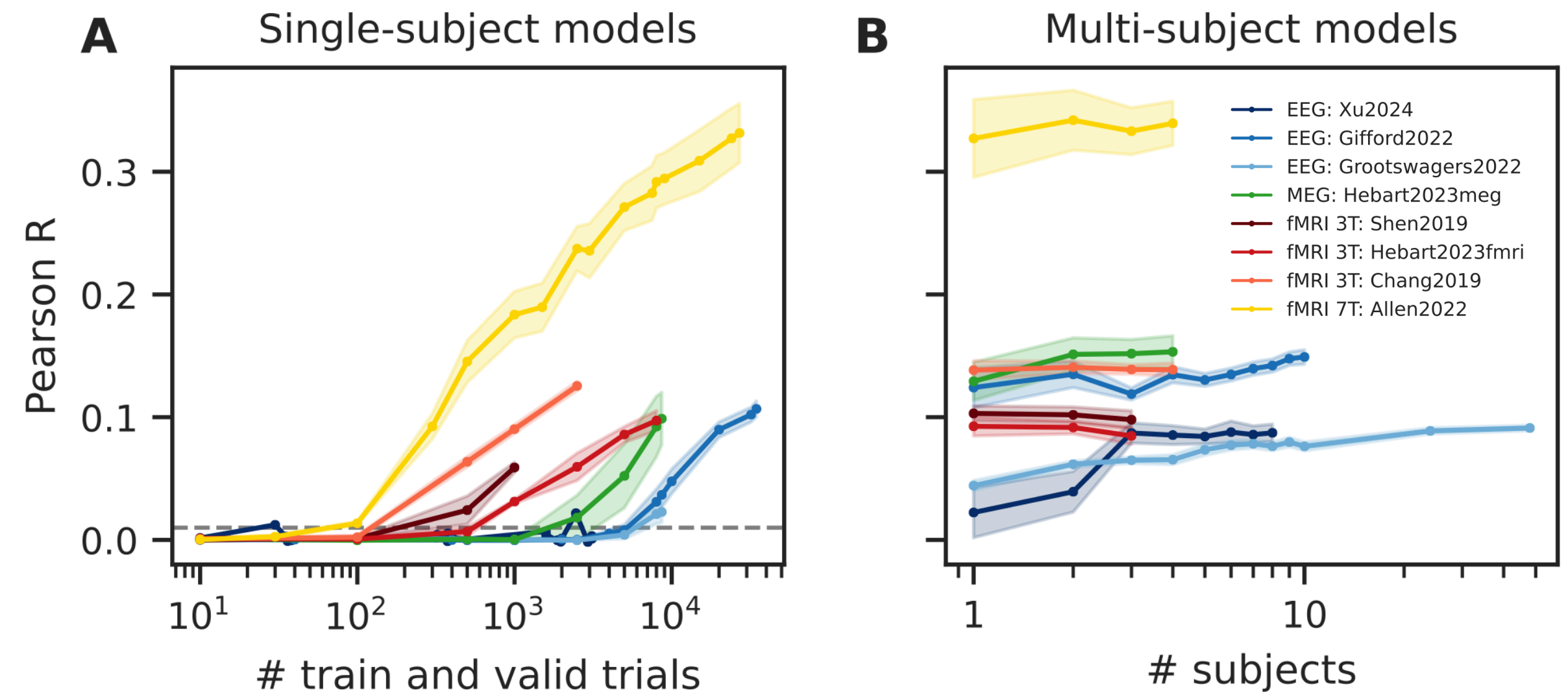
Householder Fellow in the  
Mathematics in  
Computation Section at Oak  
Ridge National Laboratory

schotthofers@ornl.gov

# Motivation: Scaling Laws

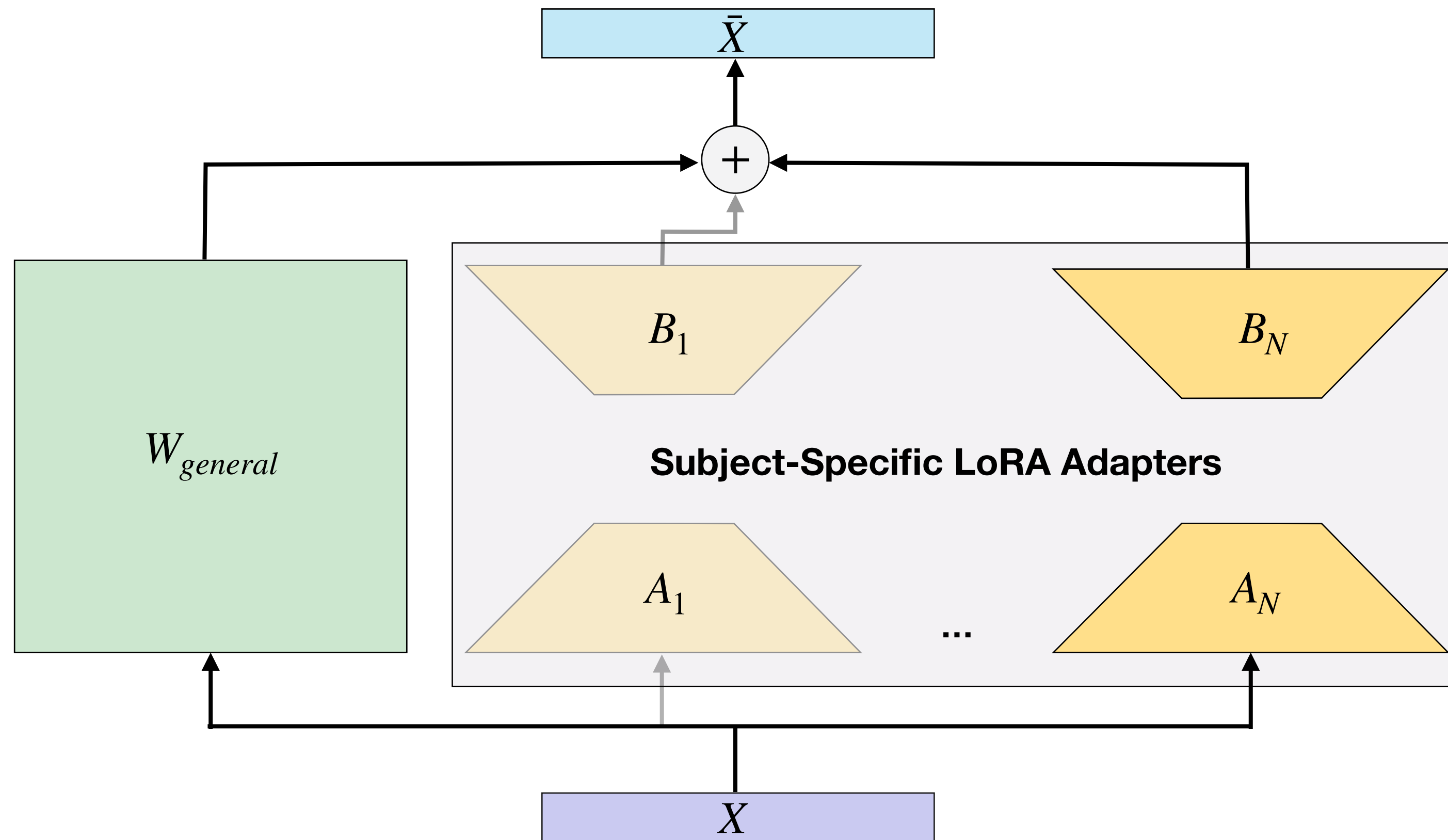


Scaling Laws for Neural Language Models  
J. Kaplan et al.



Scaling laws for decoding images from brain activity  
H. Banville et al.

# Subject-Conditioned Layer



$$\bar{X} = \sigma \left( XW_{general}^T + \sum_{s=1}^N (M_s \cdot X)A_sB_s \right)$$

$W_{general}$  shared Weights

$A_sB_s$  subject-specific Weights

$A_s \in \mathbb{R}^{n \times r}$ ,  $B_s \in \mathbb{R}^{r \times m}$

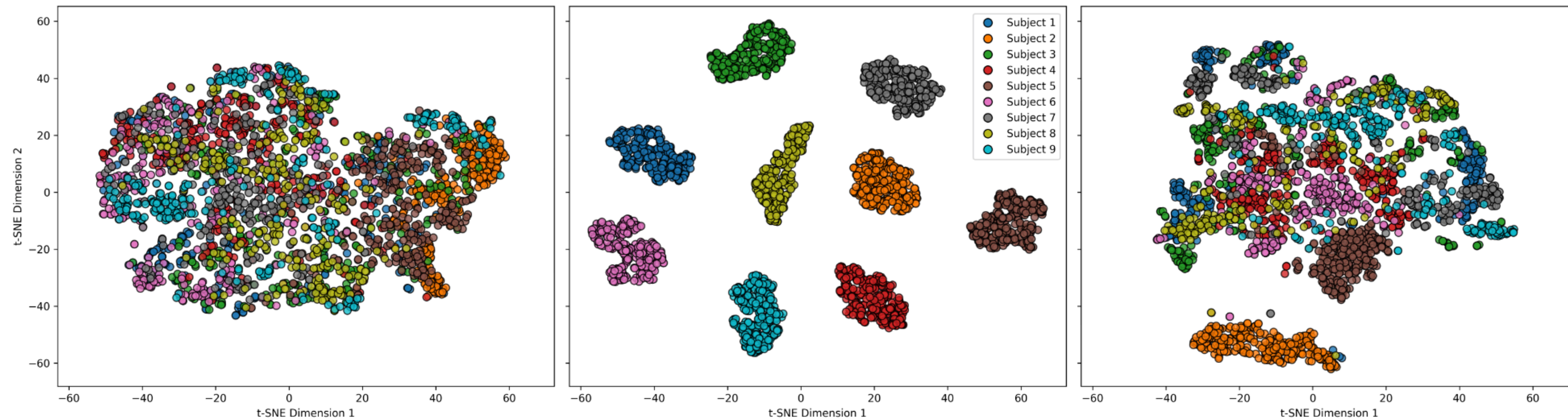
$M_s :=$  binary Mask

Subject  $i \in \{1, \dots, N\}$

Rank of Adapter:  $r$

Initialization:  $A_s = \mathcal{N}(0, \sigma^2)$ ,  $B_s = 0$

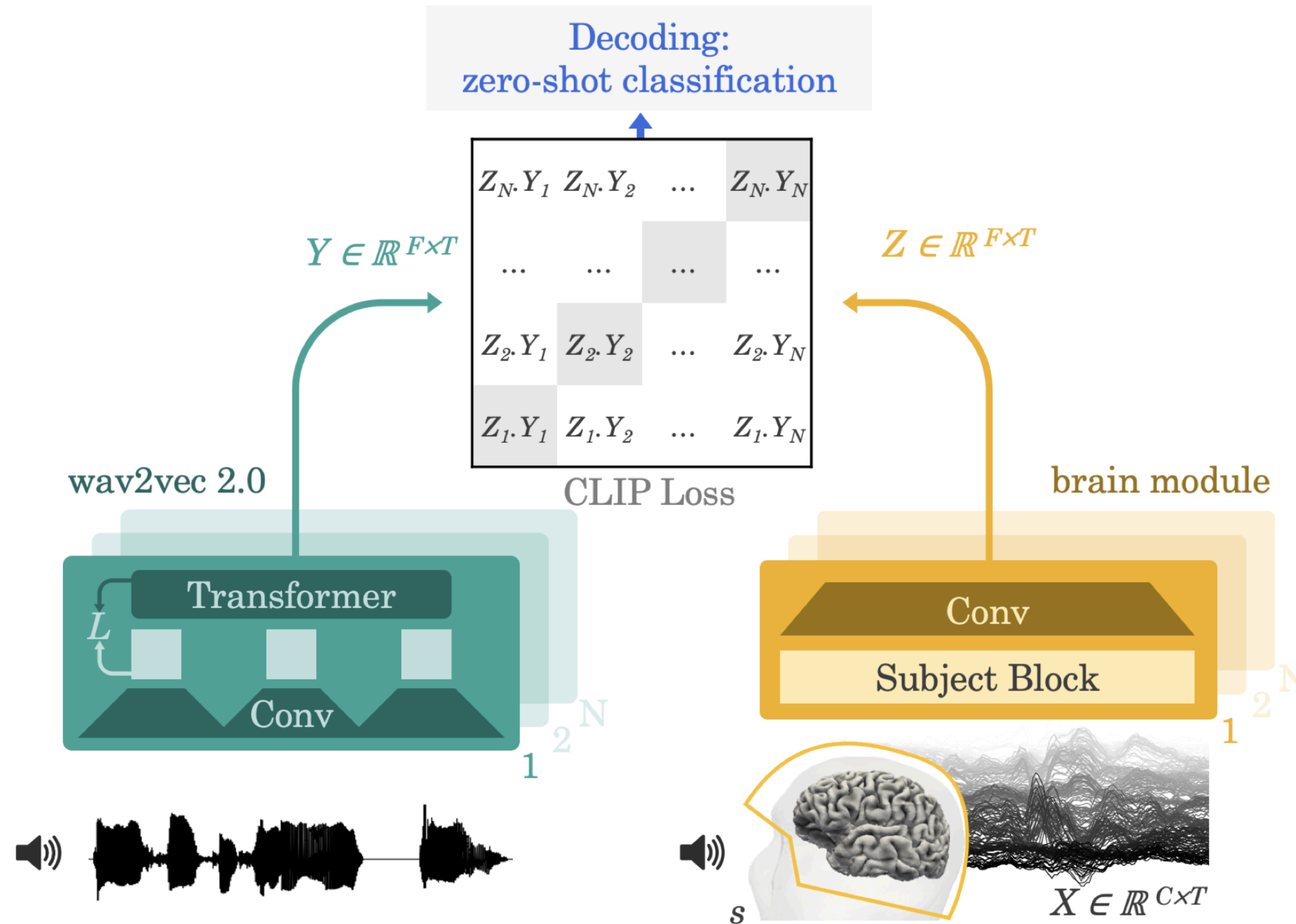
# Numerical Results



Method	Total # Params	Active # Params	BCI Comp. IV2a (4 Classes)	BCI Comp. IV2b (2 Classes)
EEGNeX Subject-Agnostic [1]	55 972	55 972	55.00% $\pm$ 0.66%	74.28% $\pm$ 0.39%
EEGNeX Subject-Specific	55 972	55 972	59.36% $\pm$ 14.42%	75.92% $\pm$ 14.32%
EEGNeX Subject-Specific LoRA	81 124	10 980	53.60% $\pm$ 2.91%	69.64% $\pm$ 0.27%
<b>EEGNeX Subject-Conditioned Layer</b>	134 884	64 740	<b>64.72% <math>\pm</math> 0.56%</b>	<b>76.48% <math>\pm</math> 0.25%</b>
PBT Subject-Agnostic	867 460	867 460	50.82% $\pm$ 0.85%	75.49% $\pm$ 1.10%
PBT Subject-Specific	867 460	867 460	44.16% $\pm$ 12.19%	75.63% $\pm$ 14.23%
PBT Subject-Specific LoRA	680 736	135 712	25.81% $\pm$ 0.52%	53.70% $\pm$ 0.80%
<b>PBT Subject-Conditioned Layer</b>	1 480 612	935 588	<b>54.51% <math>\pm</math> 0.30%</b>	<b>76.36% <math>\pm</math> 0.23%</b>

[1] Toward reliable signals decoding for electroencephalogram:  
A benchmark study to EEGNeX  
X. Chen et al.

# BrainMagick



Decoding speech perception from non-invasive brain recordings  
Defossez, A. et al.

# BrainMagick

	Method	subject-specific #Params	Accuracy [%] $\pm$ Std. [%]		
			Top 1	Top 5	Top 10
SuLoRA	r=4, $\alpha$ =4	131k	37.56 $\pm$ 0.67	58.72 $\pm$ 0.52	67.06 $\pm$ 0.54
	r=4, $\alpha$ =8	131k	38.85 $\pm$ 0.33	60.17 $\pm$ 0.45	68.52 $\pm$ 0.39
	r=4, $\alpha$ =12	131k	39.04 $\pm$ 0.20	60.36 $\pm$ 0.21	68.77 $\pm$ 0.36
	r=8, $\alpha$ =8	190k	39.82 $\pm$ 0.72	61.03 $\pm$ 0.68	69.39 $\pm$ 0.62
	r=8, $\alpha$ =16	190k	39.61 $\pm$ 0.48	61.02 $\pm$ 0.27	69.40 $\pm$ 0.37
	r=8, $\alpha$ =24	190k	40.32 $\pm$ 0.50	61.5 $\pm$ 0.23	69.79 $\pm$ 0.29
	r=16, $\alpha$ =16	306k	40.50 $\pm$ 0.27	61.77 $\pm$ 0.19	69.97 $\pm$ 0.22
	r=16, $\alpha$ =32	306k	40.50 $\pm$ 0.27	61.86 $\pm$ 0.32	70.09 $\pm$ 0.29
	r=16, $\alpha$ =48	306k	40.68 $\pm$ 0.51	61.94 $\pm$ 0.39	70.20 $\pm$ 0.36
	r=32, $\alpha$ =32	540k	40.95 $\pm$ 0.54	62.22 $\pm$ 0.64	70.34 $\pm$ 0.47
	r=32, $\alpha$ =64	540k	41.65 $\pm$ 0.76	62.63 $\pm$ 0.84	70.76 $\pm$ 0.76
	r=32, $\alpha$ =96	540k	41.44 $\pm$ 0.32	62.67 $\pm$ 0.31	70.60 $\pm$ 0.18
	r=64, $\alpha$ =64	1,006k	40.97 $\pm$ 0.44	62.15 $\pm$ 0.49	70.30 $\pm$ 0.49
	r=64, $\alpha$ =128	1,006k	41.74 $\pm$ 0.49	62.90 $\pm$ 0.43	70.98 $\pm$ 0.45
	r=64, $\alpha$ =192	1,006k	<b>41.94 <math>\pm</math> 0.92</b>	<b>63.20 <math>\pm</math> 0.65</b>	<b>71.18 <math>\pm</math> 0.61</b>
	r=64, $\alpha$ =256	1,006k	41.81 $\pm$ 0.21	63.00 $\pm$ 0.11	70.95 $\pm$ 0.21
	r=96, $\alpha$ =192	1,473k	41.62 $\pm$ 0.13	62.73 $\pm$ 0.17	70.85 $\pm$ 0.16
		Baseline	1,968k	40.60 $\pm$ 0.18	61.84 $\pm$ 0.21
	BL w/o Sub.-Layer	73k	19.23 $\pm$ 0.41	36.69 $\pm$ 0.54	45.48 $\pm$ 0.69



# Institute of Mathematical Optimization

