

# LEARNING DYNAMICS OF LLM FINETUNING

Yi Ren and Danica J. Sutherland  
ICLR 2025 Outstanding Paper Awards

PRESENTER: LILANG LIN

2025/04/27

# Outline

**1** / **Authors**

**2** / **Background**

**3** / **Method**

**4** / **Experiments**

**5** / **Discussion**

# Outline

1 / Authors

2 / **Background**

3 / Method

4 / Experiments

5 / Discussion

# Background

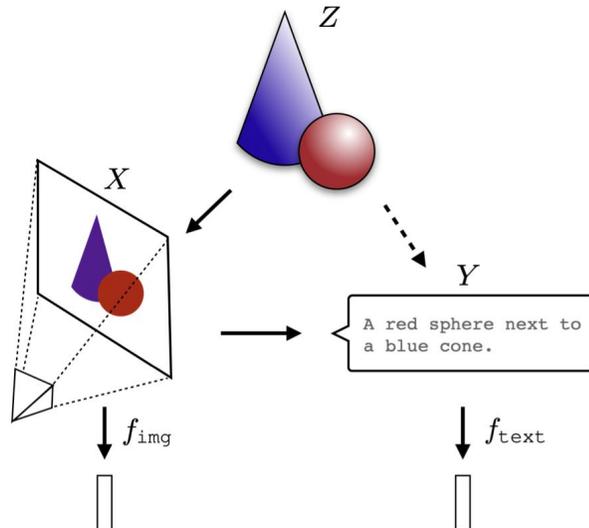
## The Platonic Representation Hypothesis

Minyoung Huh<sup>\*1</sup> Brian Cheung<sup>\*1</sup> Tongzhou Wang<sup>\*1</sup> Phillip Isola<sup>\*1</sup>

### ■ The Platonic Representation Hypothesis (ICML 2024)

#### The Platonic Representation Hypothesis

Neural networks, trained with different objectives on different data and modalities, are converging to a shared statistical model of reality in their representation spaces.



$$K_{\text{img}}(i, j) = \langle f_{\text{img}}(x_i), f_{\text{img}}(x_j) \rangle$$

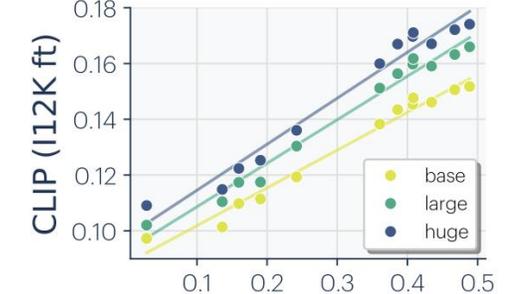
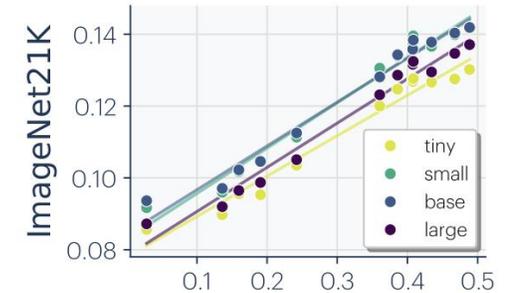
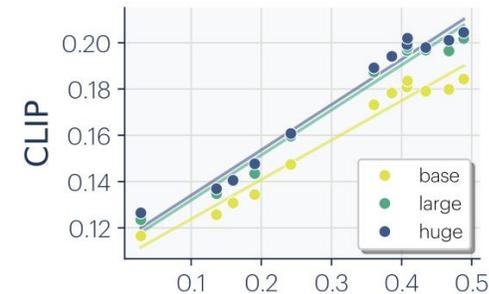
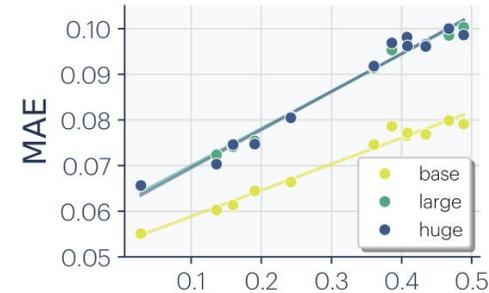
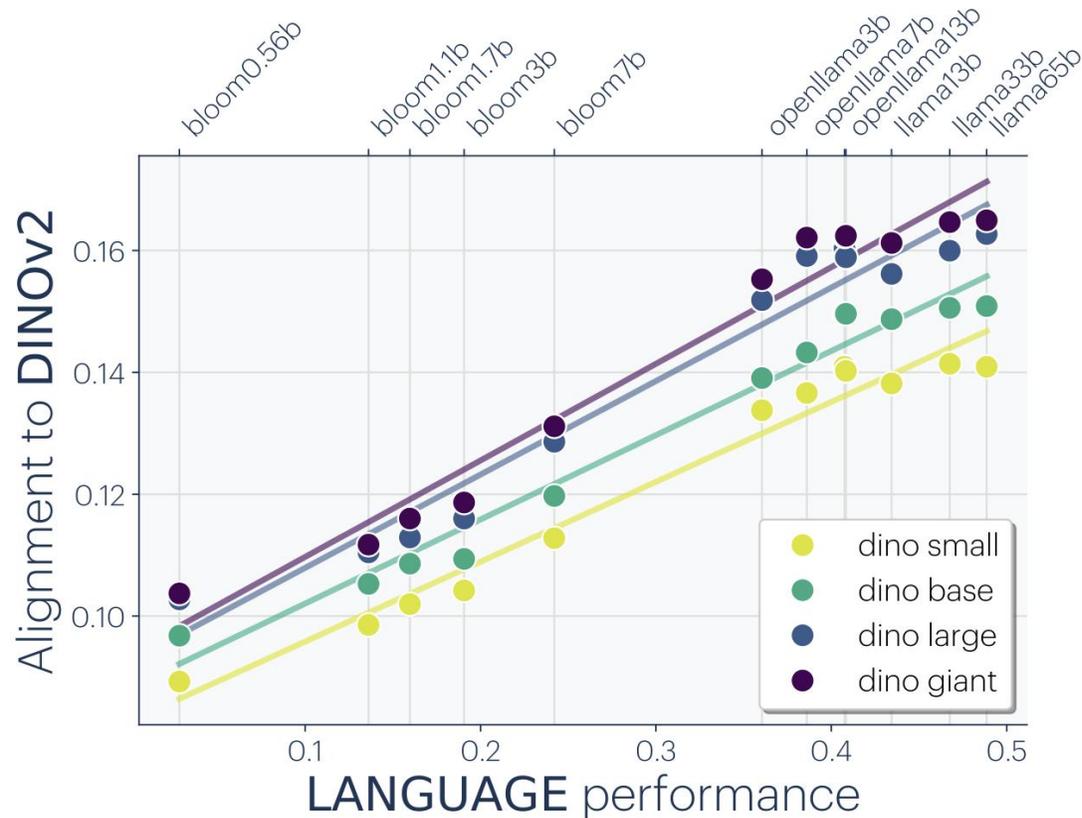
$$K_{\text{text}}(i, j) = \langle f_{\text{text}}(y_i), f_{\text{text}}(y_j) \rangle.$$

# Background

## The Platonic Representation Hypothesis

Minyoung Huh<sup>\*1</sup> Brian Cheung<sup>\*1</sup> Tongzhou Wang<sup>\*1</sup> Phillip Isola<sup>\*1</sup>

### ■ The Platonic Representation Hypothesis (ICML 2024)

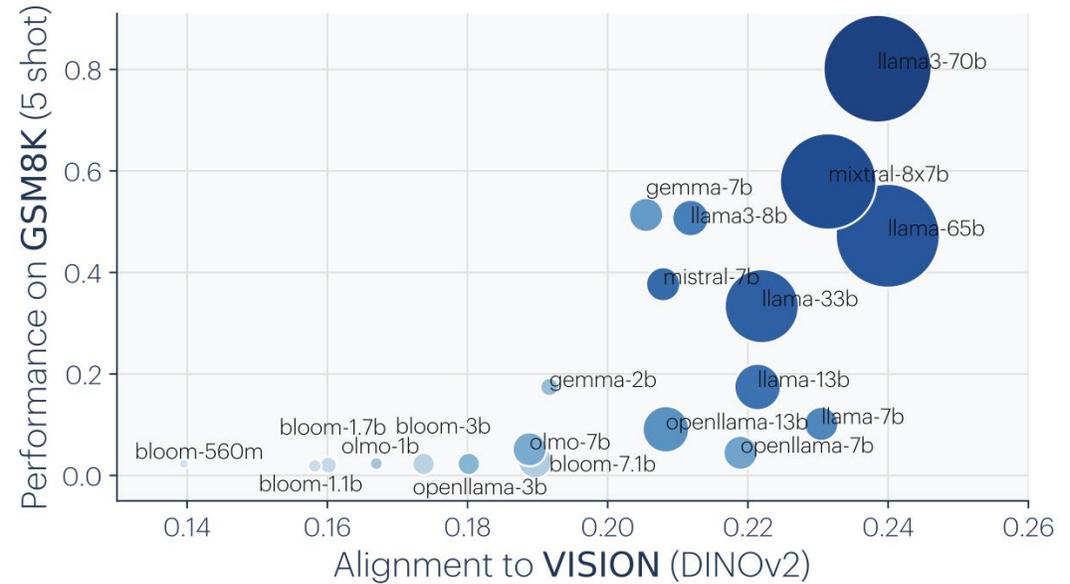
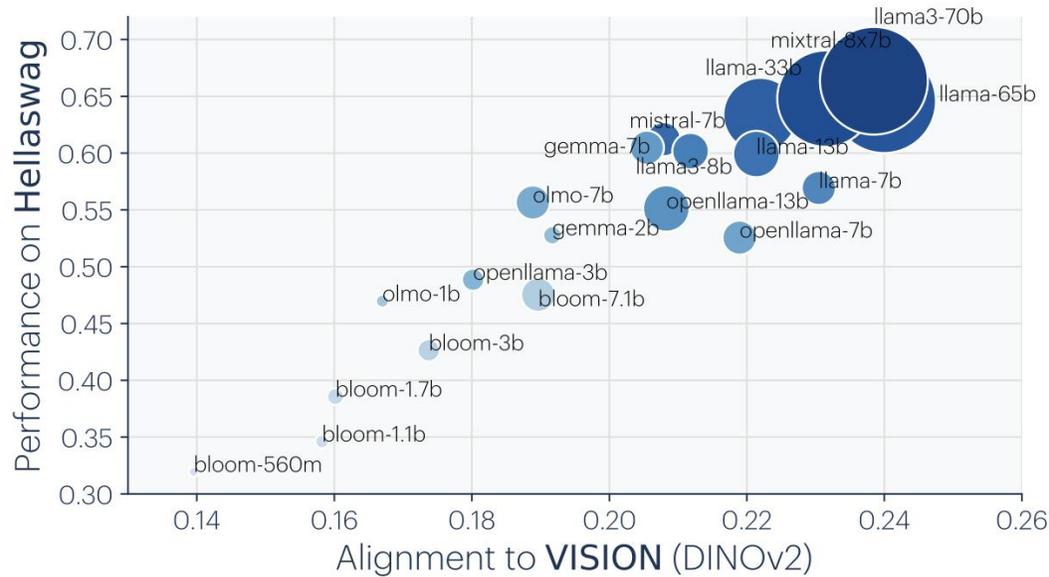


# Background

## The Platonic Representation Hypothesis

Minyoung Huh<sup>\*1</sup> Brian Cheung<sup>\*1</sup> Tongzhou Wang<sup>\*1</sup> Phillip Isola<sup>\*1</sup>

### ■ The Platonic Representation Hypothesis (ICML 2024)



# Background

## The Platonic Representation Hypothesis

Minyoung Huh<sup>\*1</sup> Brian Cheung<sup>\*1</sup> Tongzhou Wang<sup>\*1</sup> Phillip Isola<sup>\*1</sup>

### ■ The Platonic Representation Hypothesis (ICML 2024)

$$\overbrace{f^*}^{\text{trained model}} = \underset{\underbrace{f \in \mathcal{F}}_{\text{function class}}}{\text{arg min}} \mathbb{E}_{x \sim \text{dataset}} [\overbrace{\mathcal{L}(f, x)}^{\text{training objective}}] + \underbrace{\mathcal{R}(f)}_{\text{regularization}}$$

# Background

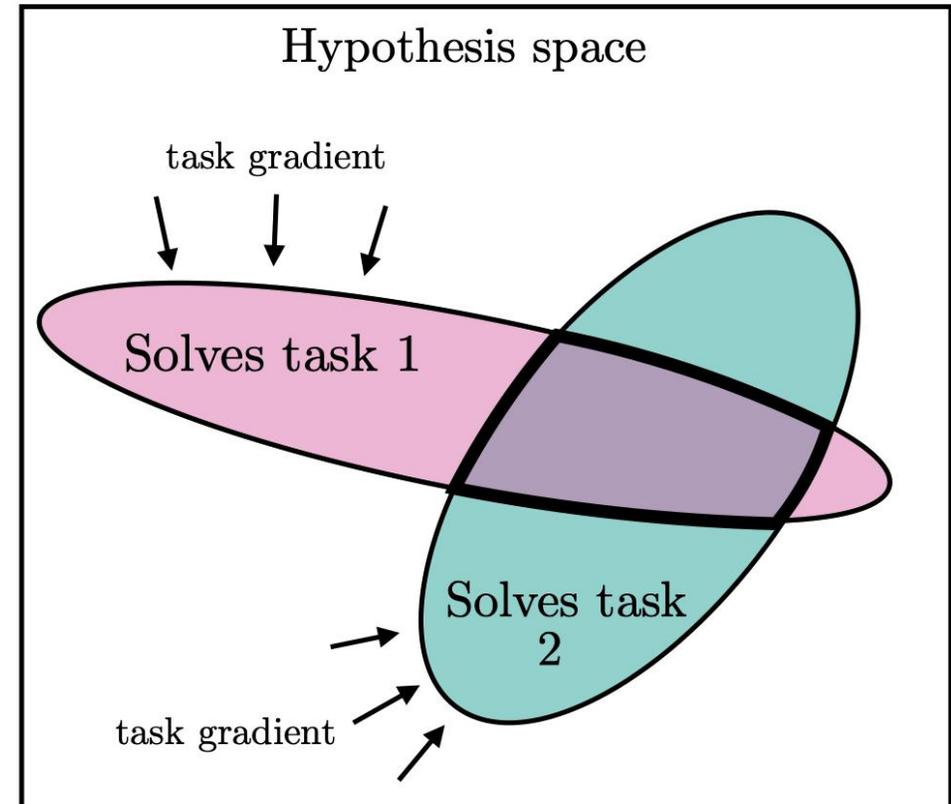
## The Platonic Representation Hypothesis

Minyoung Huh<sup>\*1</sup> Brian Cheung<sup>\*1</sup> Tongzhou Wang<sup>\*1</sup> Phillip Isola<sup>\*1</sup>

### ■ The Platonic Representation Hypothesis (ICML 2024)

#### The Multitask Scaling Hypothesis

There are fewer representations that are competent for  $N$  tasks than there are for  $M < N$  tasks. As we train more general models that solve more tasks at once, we should expect fewer possible solutions.



# Background

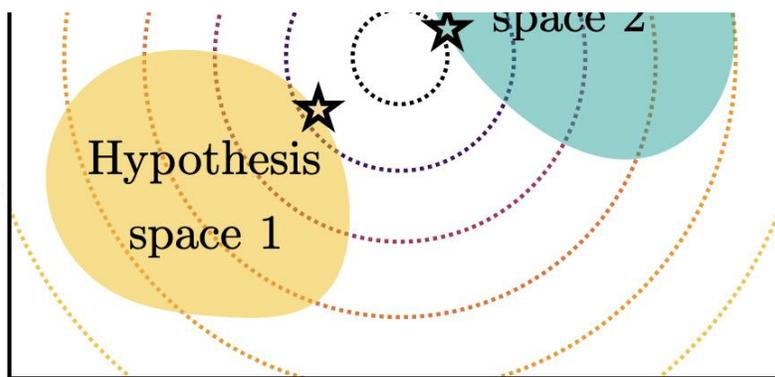
## The Platonic Representation Hypothesis

Minyoung Huh<sup>\*1</sup> Brian Cheung<sup>\*1</sup> Tongzhou Wang<sup>\*1</sup> Phillip Isola<sup>\*1</sup>

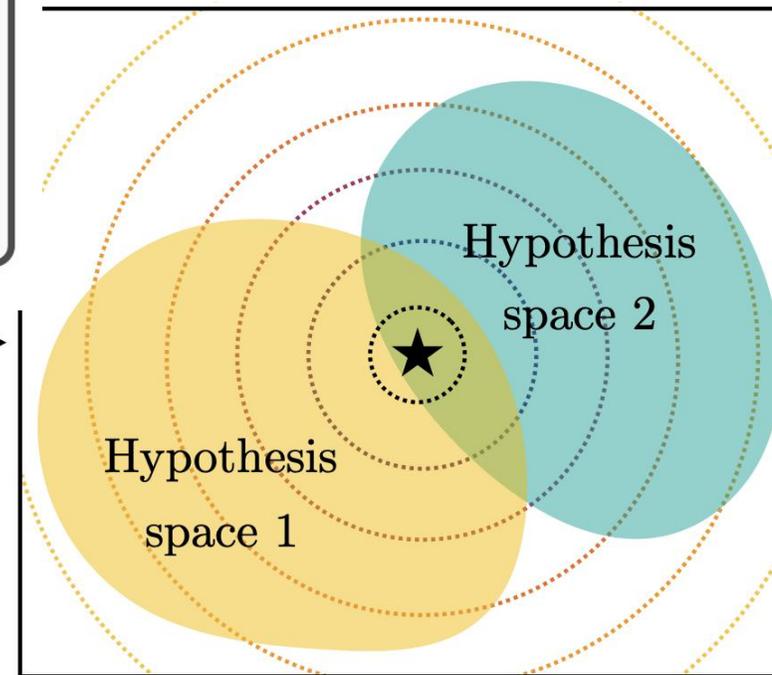
### ■ The Platonic Representation Hypothesis (ICML 2024)

#### The Capacity Hypothesis

Bigger models are more likely to converge to a shared representation than smaller models.



scale up  
architectures →



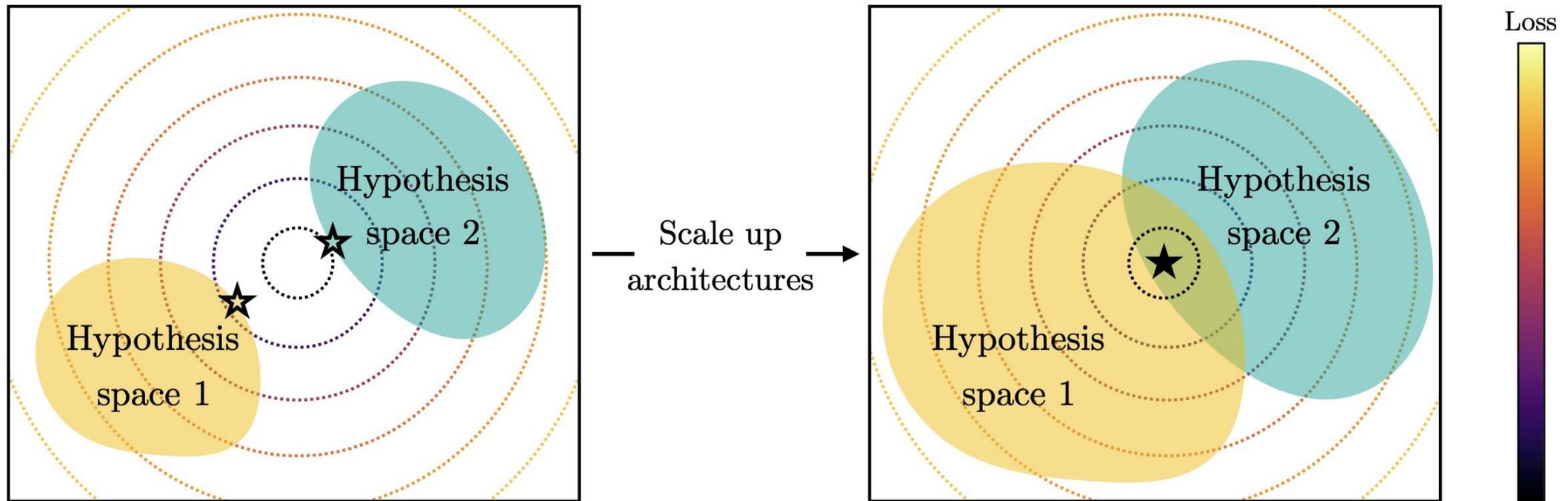
Loss

# Background

## The Platonic Representation Hypothesis

Minyoung Huh<sup>\*1</sup> Brian Cheung<sup>\*1</sup> Tongzhou Wang<sup>\*1</sup> Phillip Isola<sup>\*1</sup>

### ■ The Platonic Representation Hypothesis (ICML 2024)



# Background

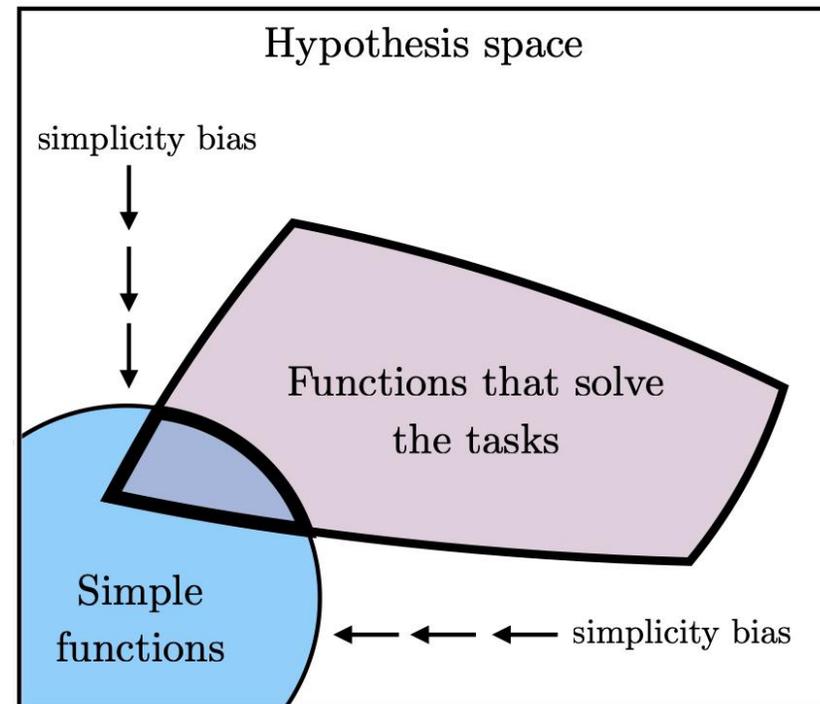
## The Platonic Representation Hypothesis

Minyoung Huh<sup>\*1</sup> Brian Cheung<sup>\*1</sup> Tongzhou Wang<sup>\*1</sup> Phillip Isola<sup>\*1</sup>

### ■ The Platonic Representation Hypothesis (ICML 2024)

#### The Simplicity Bias Hypothesis

Deep networks are biased toward finding simple fits to the data, and the bigger the model, the stronger the bias. Therefore, as models get bigger, we should expect convergence to a smaller solution space.



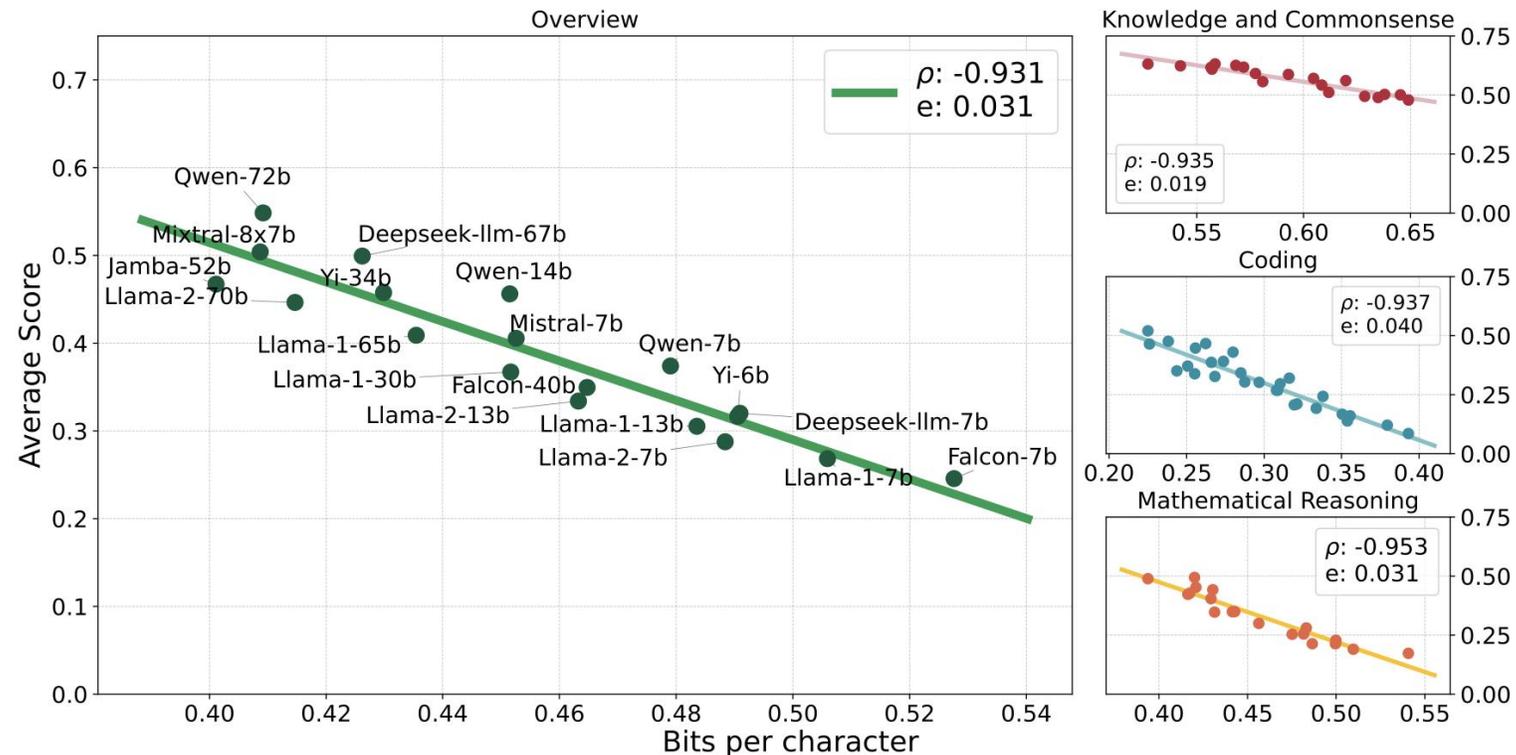
# Background

## Compression Represents Intelligence Linearly

Yuzhen Huang<sup>\*1</sup> Jinghan Zhang<sup>\*1</sup> Zifei Shan<sup>2</sup> Junxian He<sup>1</sup>  
<sup>1</sup>The Hong Kong University of Science and Technology <sup>2</sup>Tencent  
 {yhuanghj, jzhangjv, junxianh}@cse.ust.hk

### ■ Compression Represents Intelligence Linearly (COLM 2024)

$$\text{Optimal \# Bits on Average} = \mathbb{E}_{x \sim p_{\text{data}}} \left[ \sum_{i=1}^n -\log_2 p_{\text{model}}(x_i | x_{1:i-1}) \right],$$



# Outline

1 / Authors

2 / Background

3 / **Method**

4 / Experiments

5 / Discussion

# Learning Dynamic

*After an GD update on  $\mathbf{x}_u$ , how does the model's prediction on  $\mathbf{x}_o$  change?*

$$\Delta\theta \triangleq \theta^{t+1} - \theta^t = -\eta \cdot \nabla \mathcal{L}(f_{\theta}(\mathbf{x}_u), \mathbf{y}_u); \quad \Delta f(\mathbf{x}_o) \triangleq f_{\theta^{t+1}}(\mathbf{x}_o) - f_{\theta^t}(\mathbf{x}_o),$$

$$\Delta \log \pi^t(\mathbf{y} | \mathbf{x}_o) \triangleq \log \pi_{\theta^{t+1}}(\mathbf{y} | \mathbf{x}_o) - \log \pi_{\theta^t}(\mathbf{y} | \mathbf{x}_o), .$$

# Learning Dynamic

After an GD update on  $\mathbf{x}_u$ , how does the model's prediction on  $\mathbf{x}_o$  change?

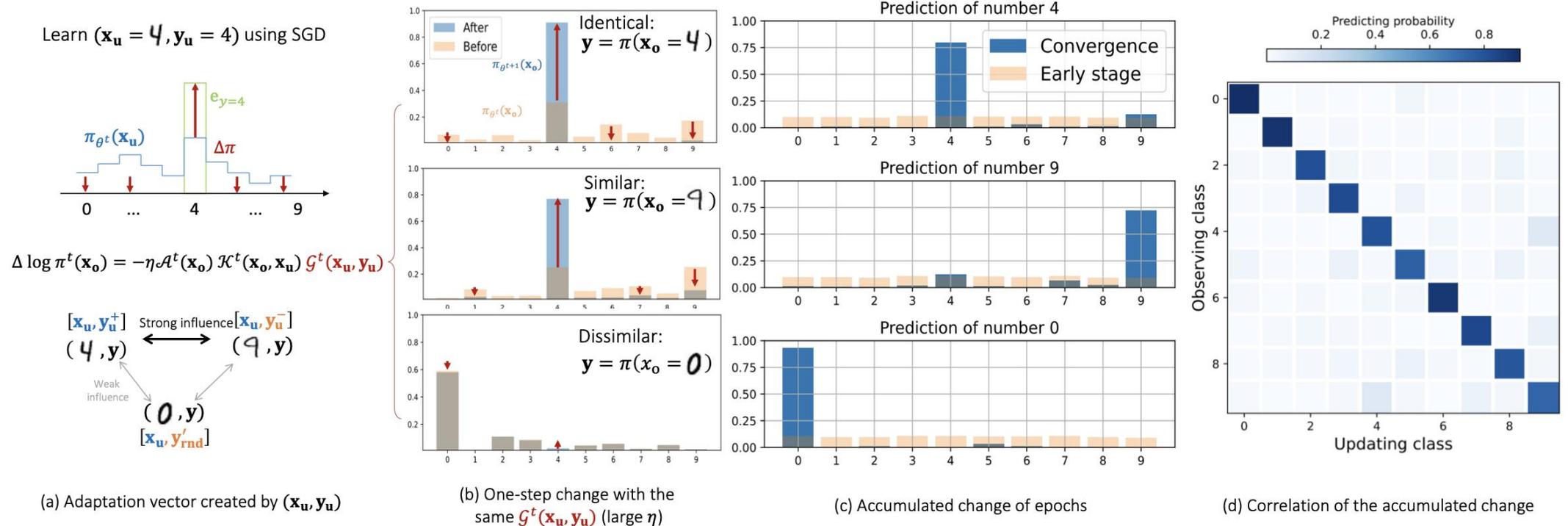
**Proposition 1.** Let  $\pi = \text{Softmax}(\mathbf{z})$  and  $\mathbf{z} = h_\theta(\mathbf{x})$ . The one-step learning dynamics decompose as

$$\underbrace{\Delta \log \pi^t(\mathbf{y} \mid \mathbf{x}_o)}_{V \times 1} = -\eta \underbrace{\mathcal{A}^t(\mathbf{x}_o)}_{V \times V} \underbrace{\mathcal{K}^t(\mathbf{x}_o, \mathbf{x}_u)}_{V \times V} \underbrace{\mathcal{G}^t(\mathbf{x}_u, \mathbf{y}_u)}_{V \times 1} + \mathcal{O}(\eta^2 \|\nabla_\theta \mathbf{z}(\mathbf{x}_u)\|_{\text{op}}^2), \quad (3)$$

where  $\mathcal{A}^t(\mathbf{x}_o) = \nabla_{\mathbf{z}} \log \pi_{\theta^t}(\mathbf{x}_o) = I - \mathbf{1} \pi_{\theta^t}^\top(\mathbf{x}_o)$ ,  $\mathcal{K}^t(\mathbf{x}_o, \mathbf{x}_u) = (\nabla_\theta \mathbf{z}(\mathbf{x}_o)|_{\theta^t})(\nabla_\theta \mathbf{z}(\mathbf{x}_u)|_{\theta^t})^\top$  is the empirical neural tangent kernel of the logit network  $\mathbf{z}$ , and  $\mathcal{G}^t(\mathbf{x}_u, \mathbf{y}_u) = \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{x}_u, \mathbf{y}_u)|_{\mathbf{z}^t}$ .

# Learning Dynamic

After an GD update on  $\mathbf{x}_u$ , how does the model's prediction on  $\mathbf{x}_o$  change?



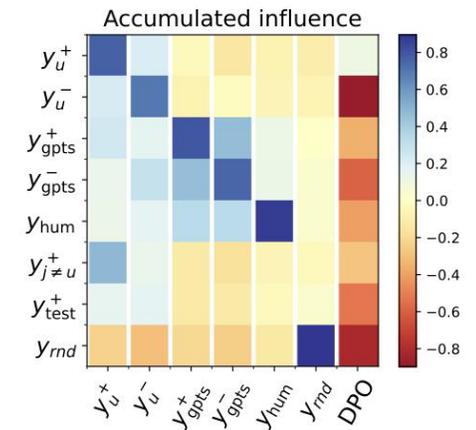
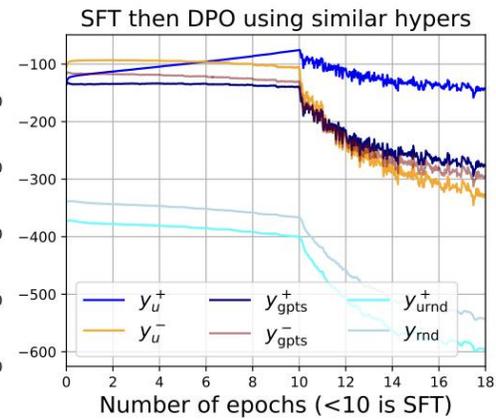
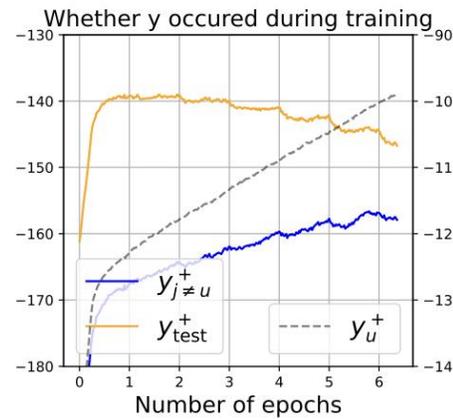
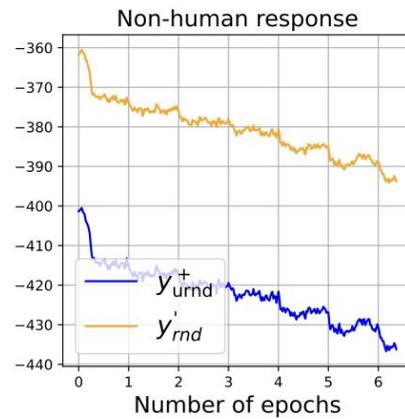
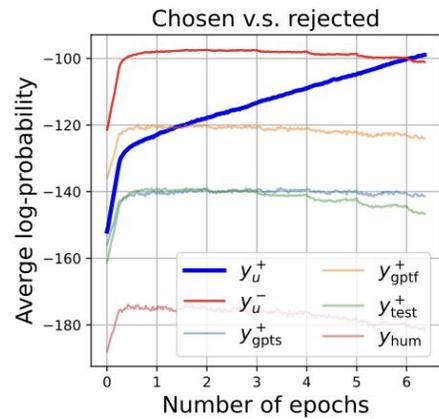
# Learning Dynamic

$$\mathcal{L}_{\text{SFT}}(\mathbf{x}_u, \mathbf{y}_u^+) \triangleq - \sum_{l=1}^L \log \pi(y = y_l^+ | \mathbf{y}_{<l}^+, \mathbf{x}_u) = - \sum_{l=1}^L \mathbf{e}_{y_l^+} \cdot \log \pi(\mathbf{y} | \mathbf{x}_u, \mathbf{y}_{<l}^+).$$

$$\underbrace{[\Delta \log \pi^t(\mathbf{y} | \boldsymbol{\chi}_o)]_m}_{V \times M} = - \sum_{l=1}^L \eta \underbrace{[\mathcal{A}^t(\boldsymbol{\chi}_o)]_m}_{V \times V \times M} \underbrace{[\mathcal{K}^t(\boldsymbol{\chi}_o, \boldsymbol{\chi}_u)]_l}_{V \times V \times L} \underbrace{[\mathcal{G}^t(\boldsymbol{\chi}_u)]_l}_{V \times L} + \mathcal{O}(\eta^2),$$

# Experiments

## ■ LEARNING DYNAMICS OF SFT



# Learning Dynamic

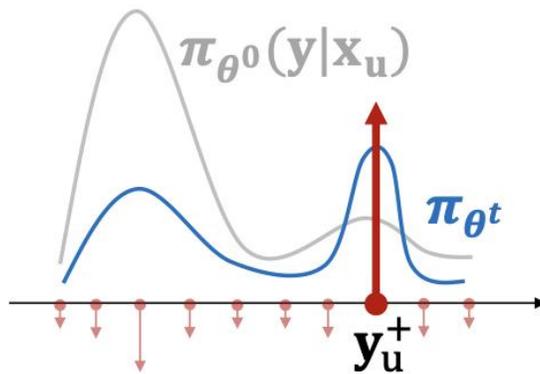
$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(\mathbf{x}_u, \mathbf{y}_u^+, \mathbf{y}_u^-) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta^t}(\mathbf{y}_u^+ | \boldsymbol{\chi}_u^+)}{\pi_{\text{ref}}(\mathbf{y}_u^+ | \boldsymbol{\chi}_u^+)} - \beta \log \frac{\pi_{\theta^t}(\mathbf{y}_u^- | \boldsymbol{\chi}_u^-)}{\pi_{\text{ref}}(\mathbf{y}_u^- | \boldsymbol{\chi}_u^-)} \right) \right],$$

$$[\Delta \log \pi^t(\mathbf{y} | \boldsymbol{\chi}_o)]_m = - \sum_{l=1}^L \eta [\mathcal{A}^t(\boldsymbol{\chi}_o)]_m \left( [\mathcal{K}^t(\boldsymbol{\chi}_o, \boldsymbol{\chi}_u^+)]_l [\mathcal{G}_{\text{DPO}^+}^t]_l - [\mathcal{K}^t(\boldsymbol{\chi}_o, \boldsymbol{\chi}_u^-)]_l [\mathcal{G}_{\text{DPO}^-}^t]_l \right) + \mathcal{O}(\eta^2)$$

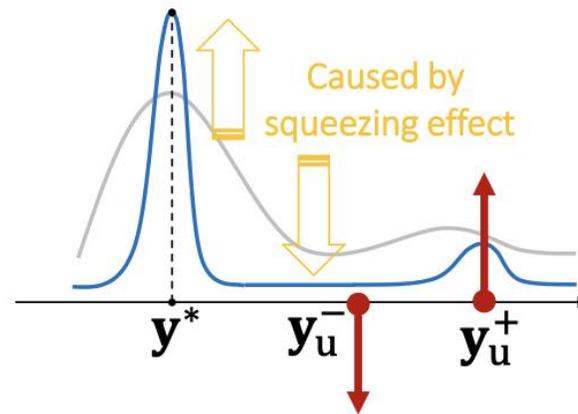
$$\mathcal{G}_{\text{DPO}^+}^t = \beta(1-a) (\pi_{\theta^t}(\mathbf{y} | \boldsymbol{\chi}_u^+) - \mathbf{y}_u^+); \quad \mathcal{G}_{\text{DPO}^-}^t = \beta(1-a) (\pi_{\theta^t}(\mathbf{y} | \boldsymbol{\chi}_u^-) - \mathbf{y}_u^-),$$

# Learning Dynamic

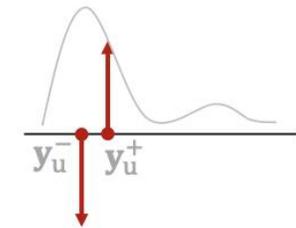
- SFT



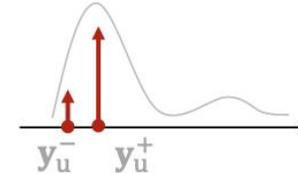
- Off-policy DPO, IPO



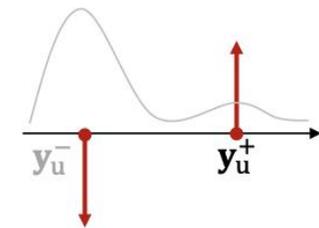
- On-policy DPO, IPO



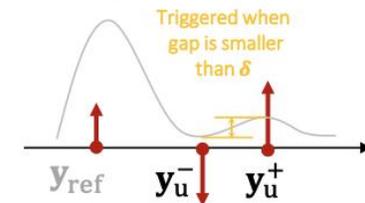
- SPPO



- SPIN

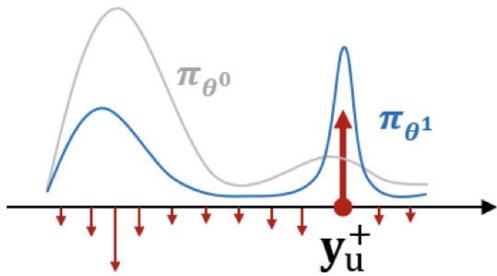


- SLIC

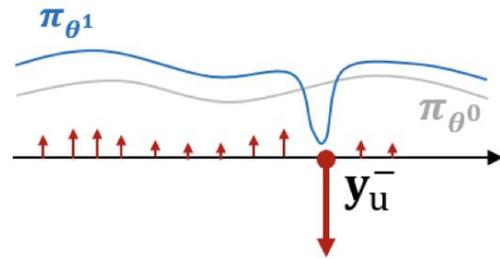


# Learning Dynamic

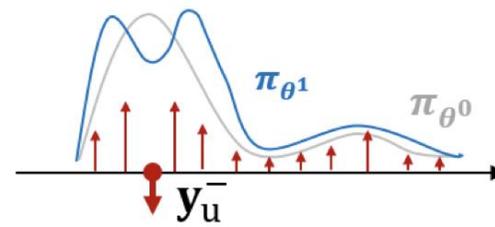
- A: Positive gradient ( $\pi_{\theta^0} - e_{\hat{y}}$ )



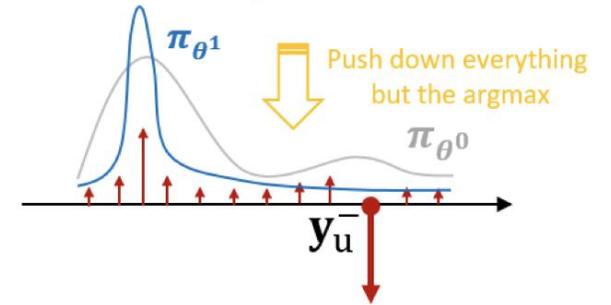
- B: Big negative gradient for flat  $\pi_{\theta^0}$



- C: Big negative gradient on the “peak” of a non-flat  $\pi_{\theta^0}$  (on-policy DPO)

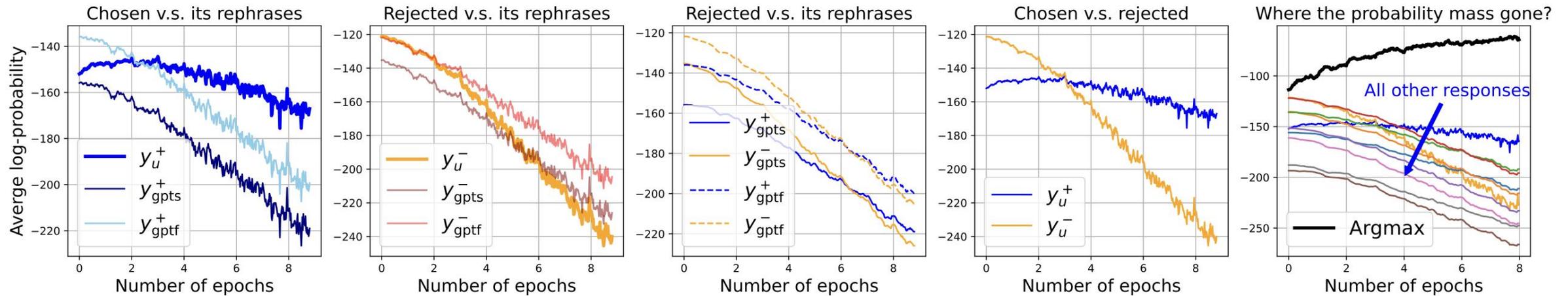


- D: Big negative gradient on the “valley” of a non-flat  $\pi_{\theta^0}$



# Experiments

## ■ LEARNING DYNAMICS OF OFF-POLICY DPO



# Outline

1 / Authors

2 / Background

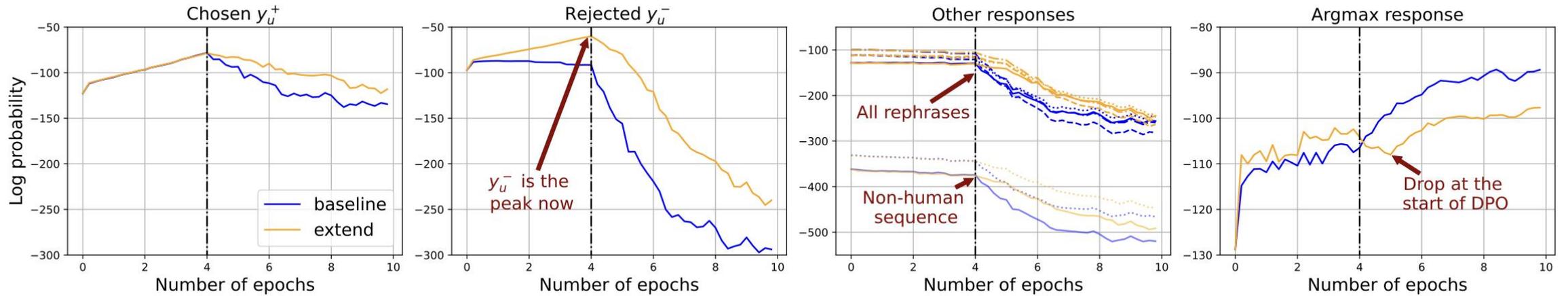
3 / Method

4 / **Experiments**

5 / Discussion

# Experiments

## ■ MITIGATING THE SQUEEZING EFFECT



# Outline

**1** / Authors

**2** / Background

**3** / Method

**4** / Experiments

**5** / Discussion

# Discussion

- SFT & RL
- Pos & Neg

**Thanks!**