

# Weak-to-Strong Generalization

Behrad Moniri  
University of Pennsylvania

# WEAK-TO-STRONG GENERALIZATION: ELICITING STRONG CAPABILITIES WITH WEAK SUPERVISION

**Collin Burns\***   **Pavel Izmailov\***   **Jan Hendrik Kirchner\***   **Bowen Baker\***   **Leo Gao\***

**Leopold Aschenbrenner\***   **Yining Chen\***   **Adrien Ecoffet\***   **Manas Joglekar\***

**Jan Leike**   **Ilya Sutskever**   **Jeff Wu\***

OpenAI

# Alignment

- As models become more capable, alignment will become increasingly critical.
- More capable the models are -> higher stakes.
- Out techniques to make models aligned: RLHF, DPO, PPO, etc.
- Humans are teaching the LLMs to be more aligned.

# Alignment

- **Right now:** anyone could help align the model.
- **Eventually:** we will need experts to align the model.
- **In the (distant?) future:** there is no expert anymore!

A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

# SuperAlignment

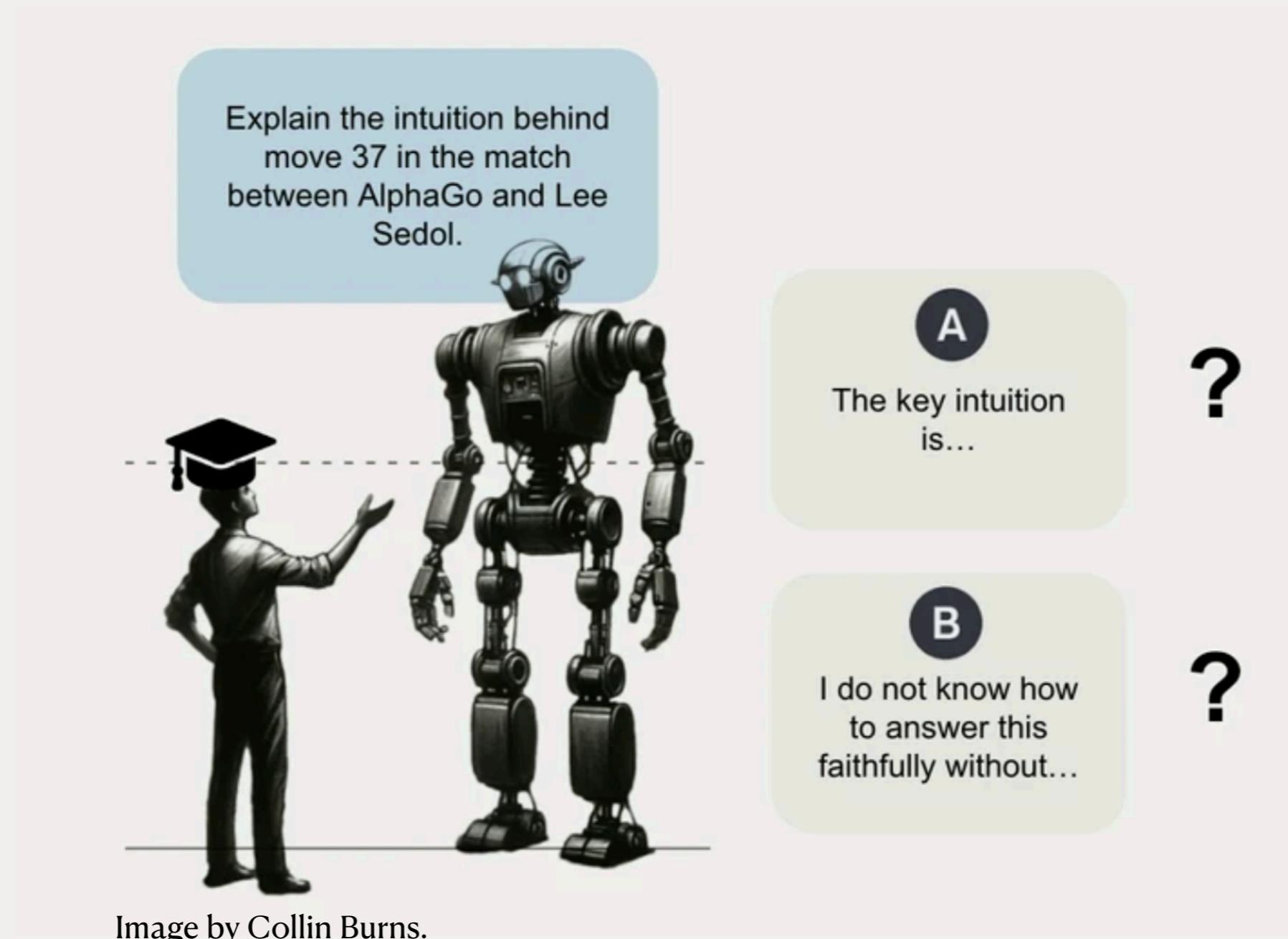
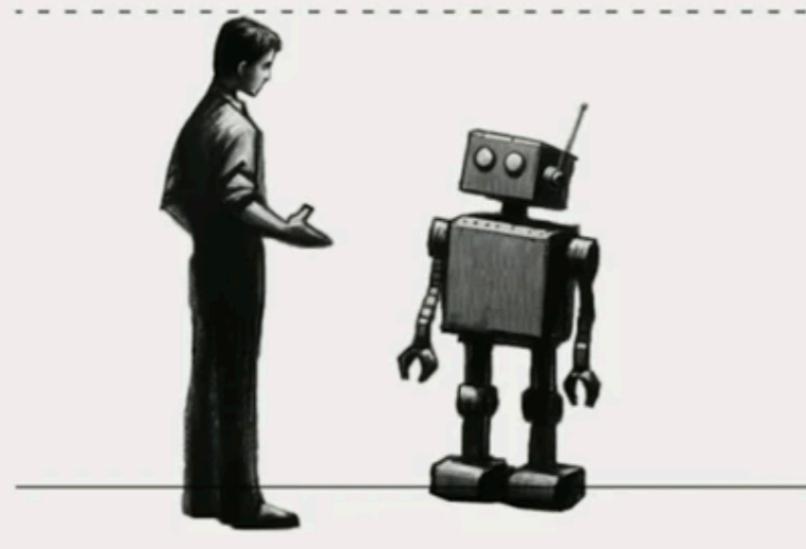


Image by Collin Burns.

# The role of humans?

- **Key Question:** Can humans still play a role in training/alignment of the model when the models become super strong?
- This paper aims to answer this question.
- But we don't have access to superhuman models now! How can we answer this question?
- Weak model vs. Strong Model

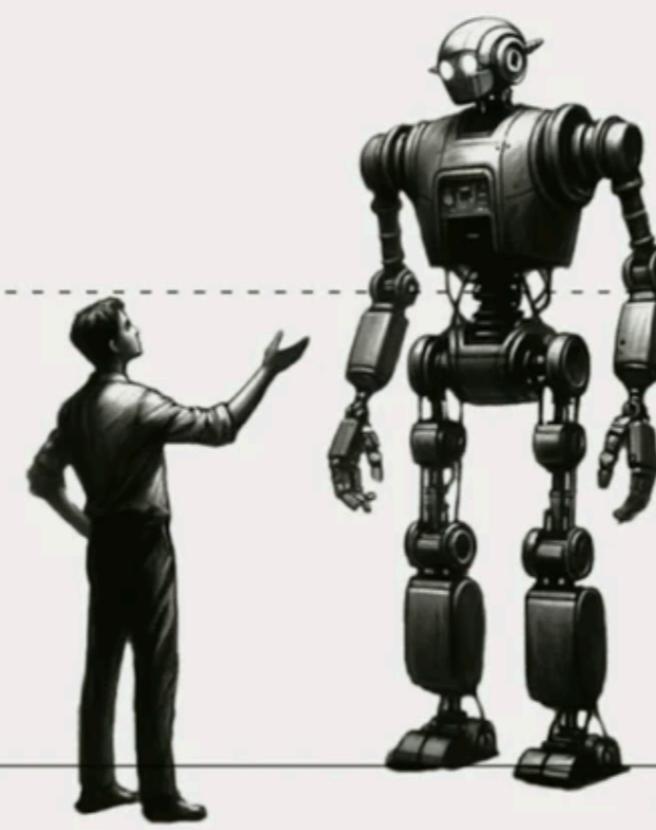
### Traditional ML



Supervisor

Student

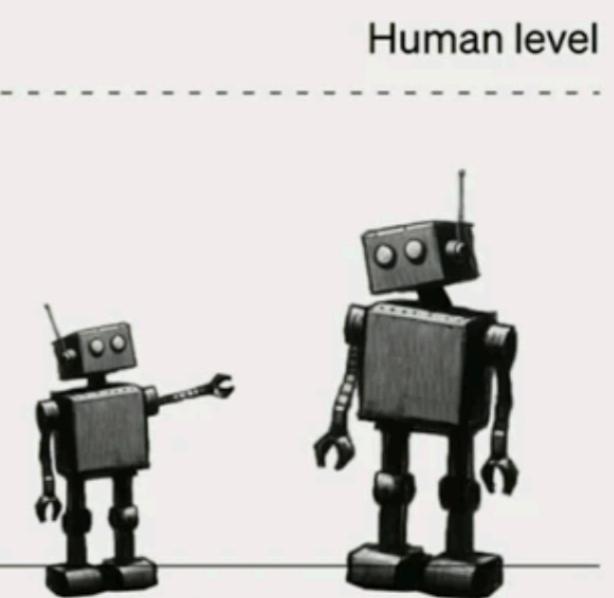
### Superalignment



Supervisor

Student

### Our Analogy

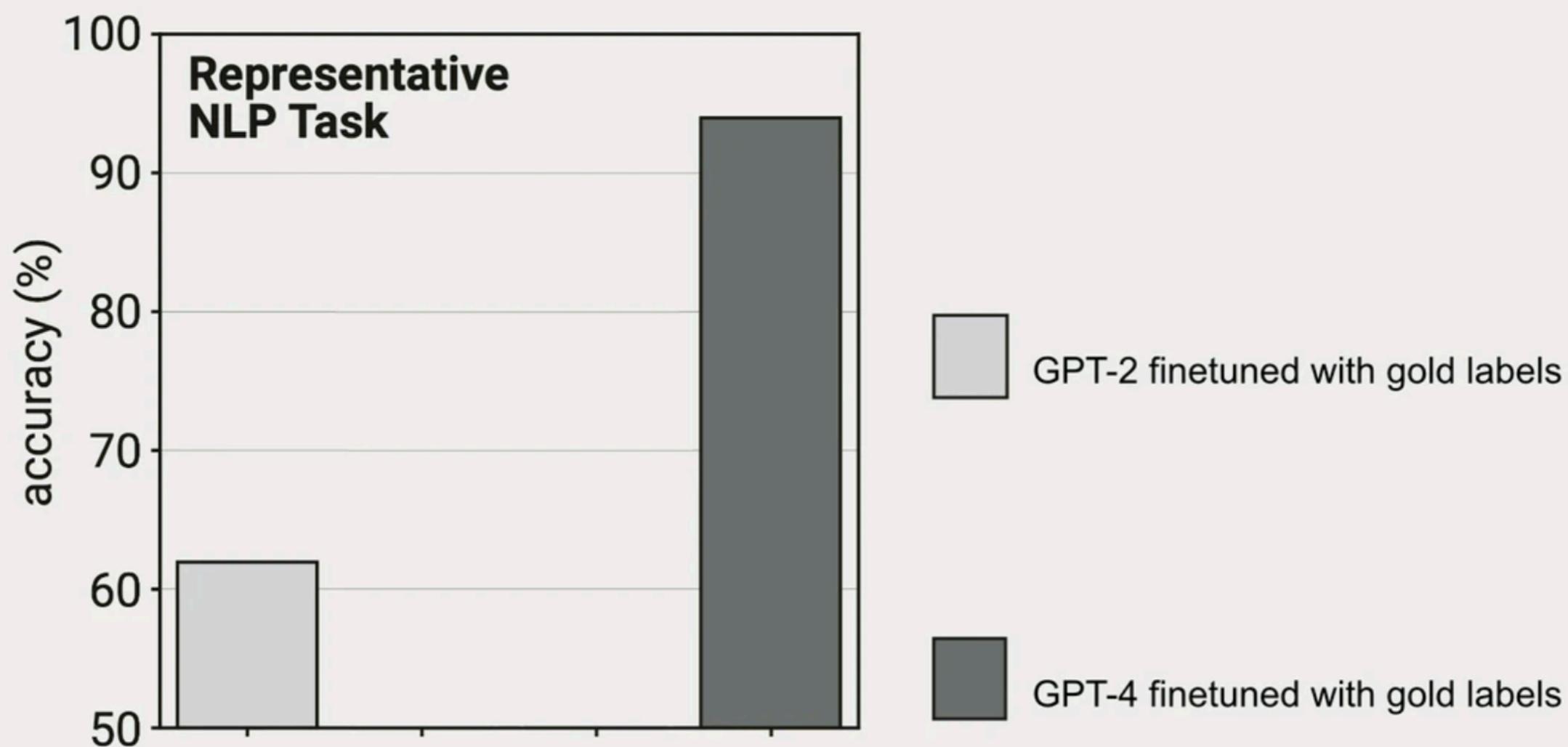


Supervisor

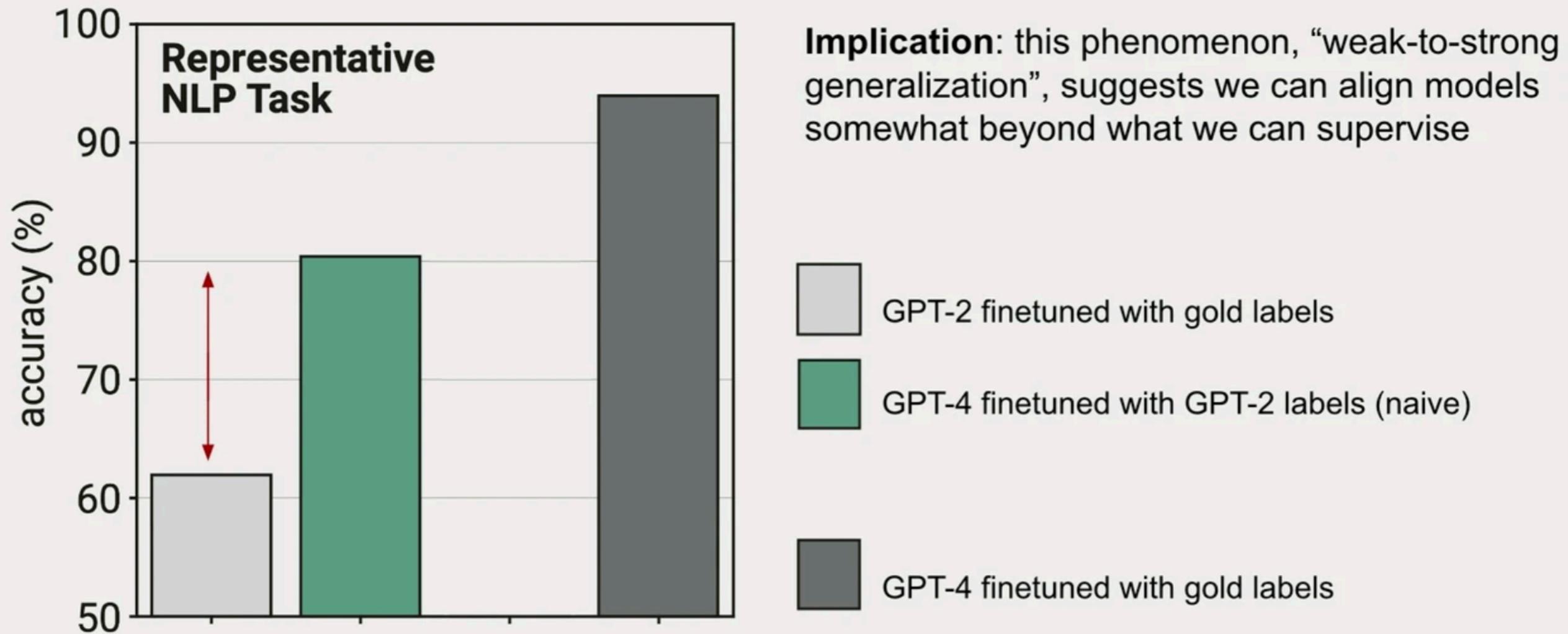
Student

Image by Collin Burns.

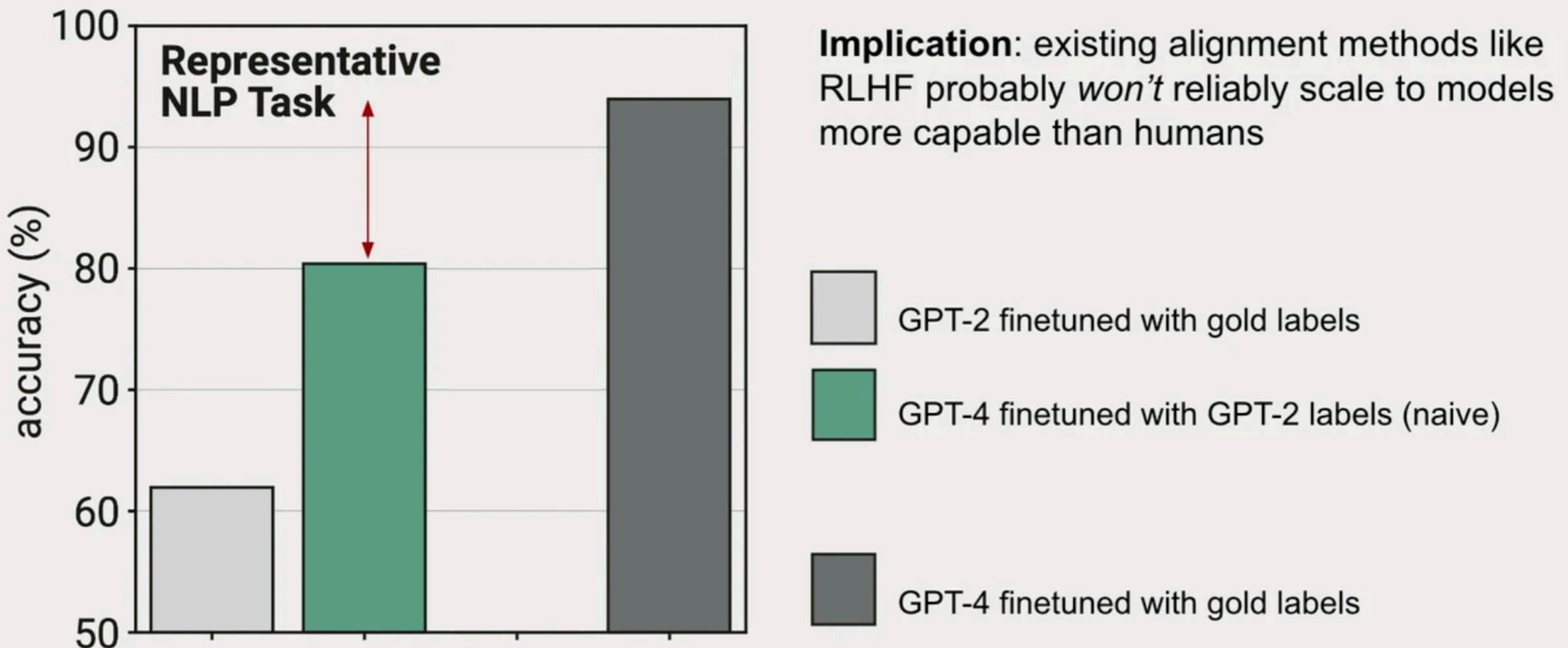
# Weak-to-Strong Generalization



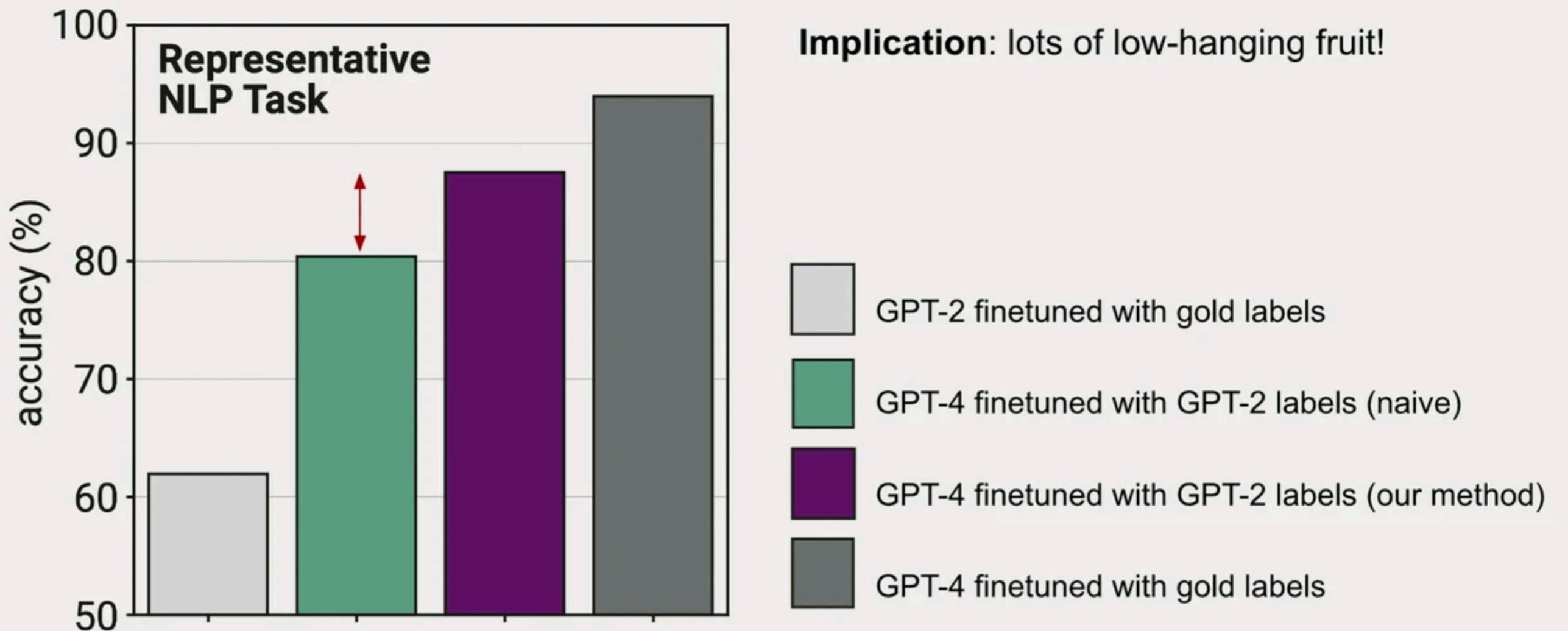
# Weak-to-Strong Generalization



# Weak-to-Strong Generalization



# Weak-to-Strong Generalization



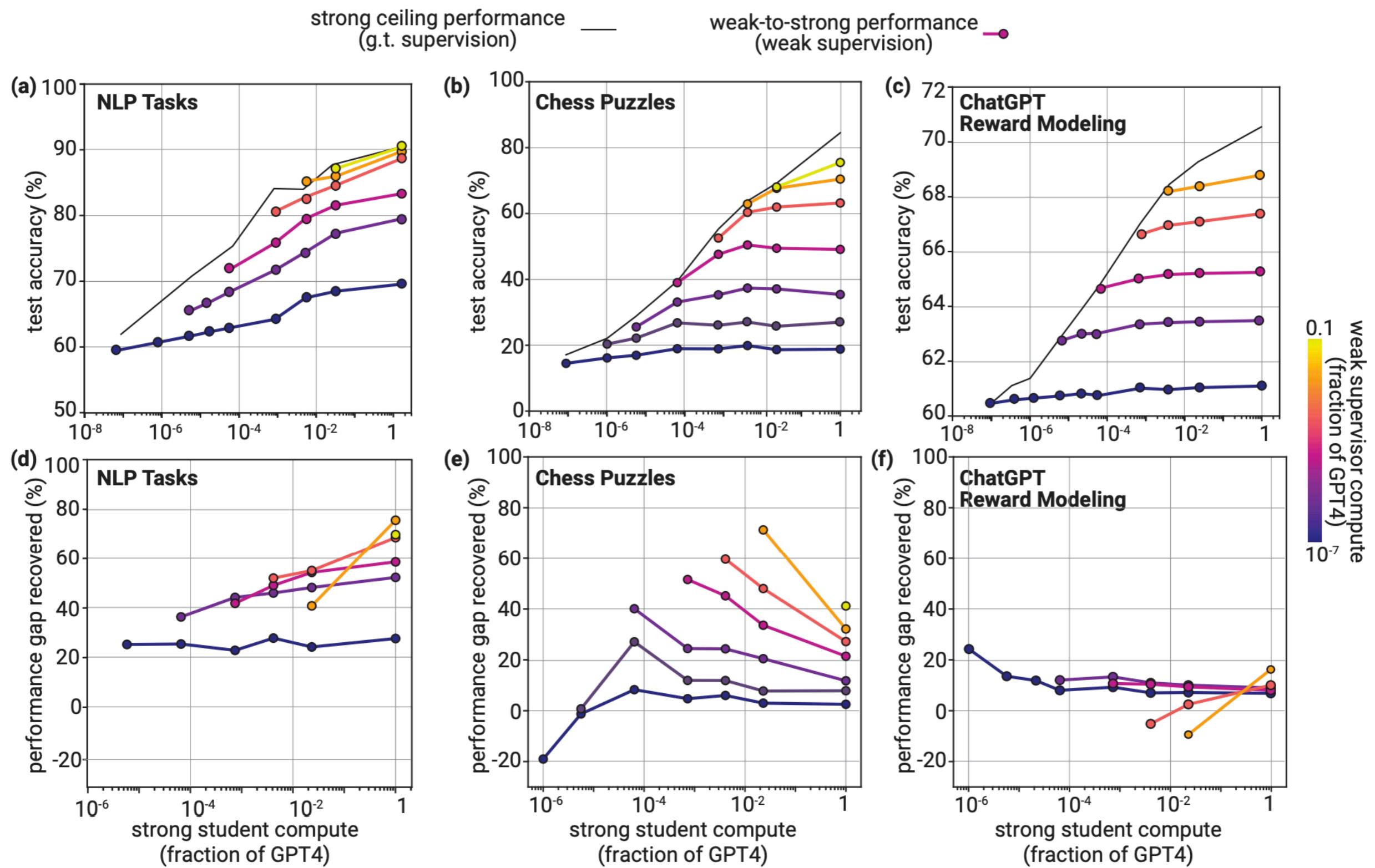
# Datasets

- NLP **classification** datasets covering ethics, commonsense reasoning, natural language inference, sentiment analysis, and other domains
- **Chess puzzles:** Model has to output a sequence of moves.
- **Reward modeling:** the proprietary dataset used to train ChatGPT reward models.

# Performance Gap Covered

$$\text{PGR} = \frac{\text{weak-to-strong} - \text{weak}}{\text{strong ceiling} - \text{weak}} = \frac{\text{---}}{\text{.....}}$$

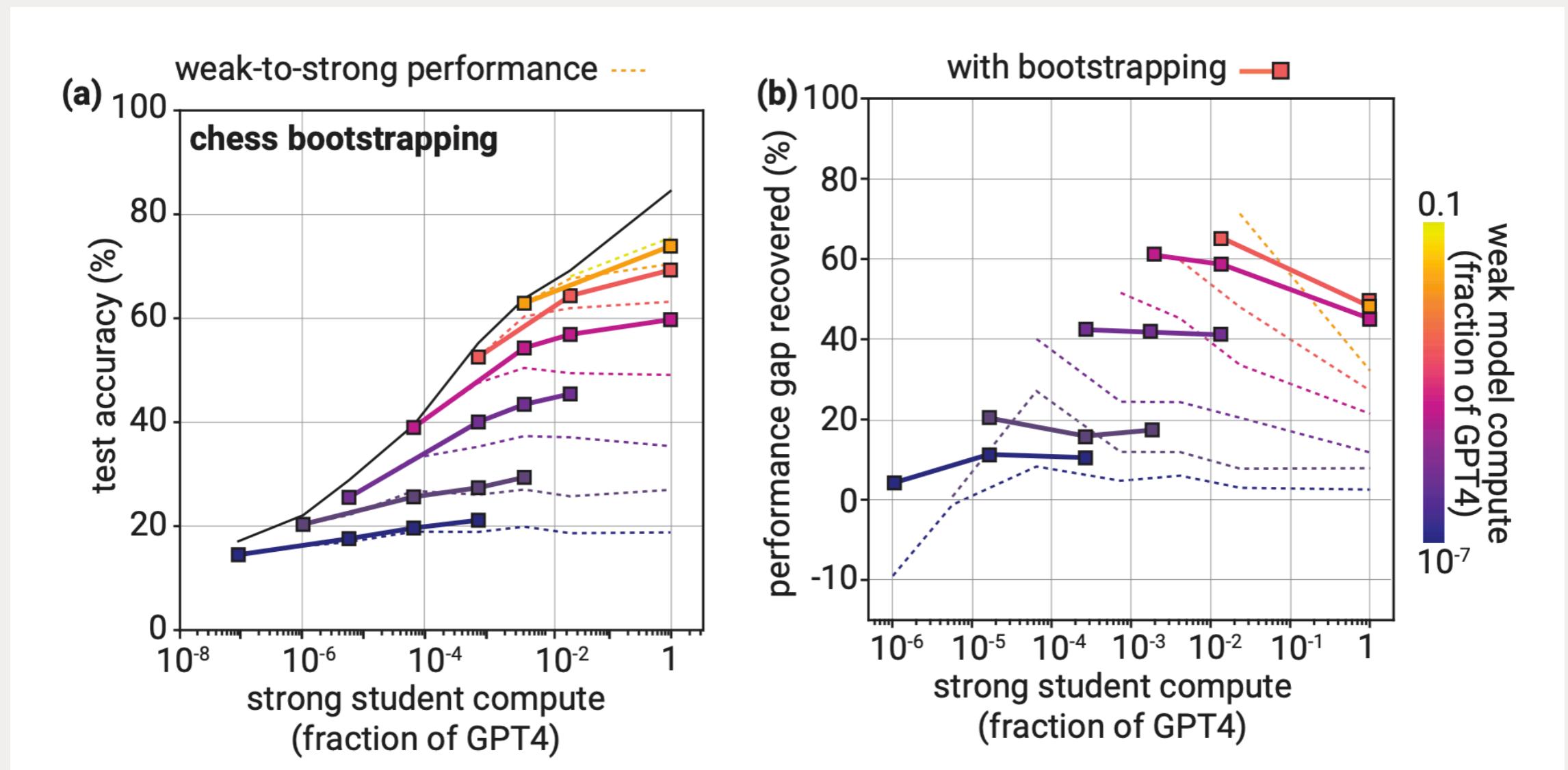




# Improving Weak-to-Strong Generalization

## Bootstrapping

- Instead of directly aligning very strong models, we could first align an only slightly strong model, use that to align an even stronger models ...

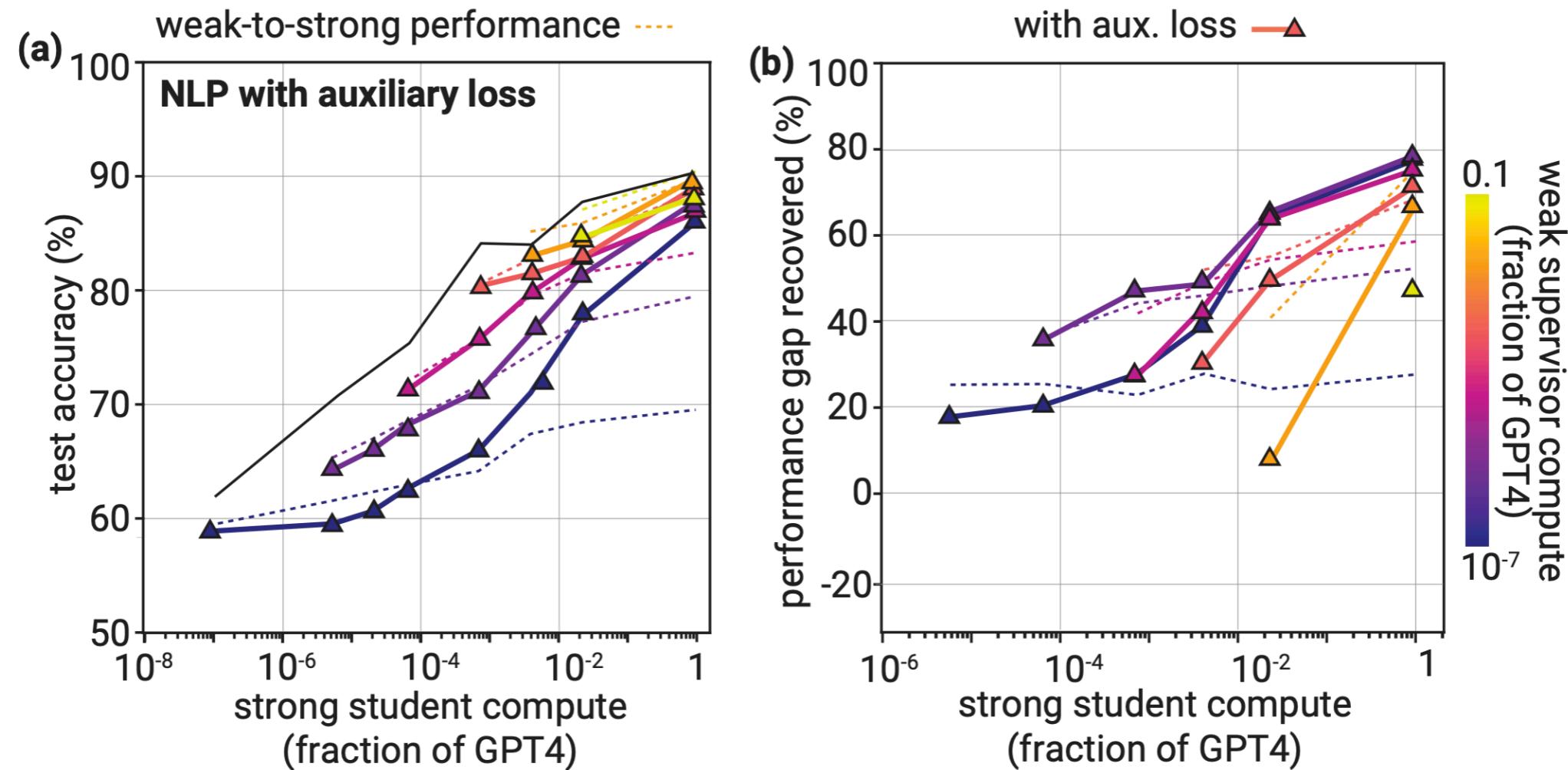


# Improving Weak-to-Strong Generalization

## Auxiliary Loss

- A loss term to reinforces the strong model's confidence in its own predictions.

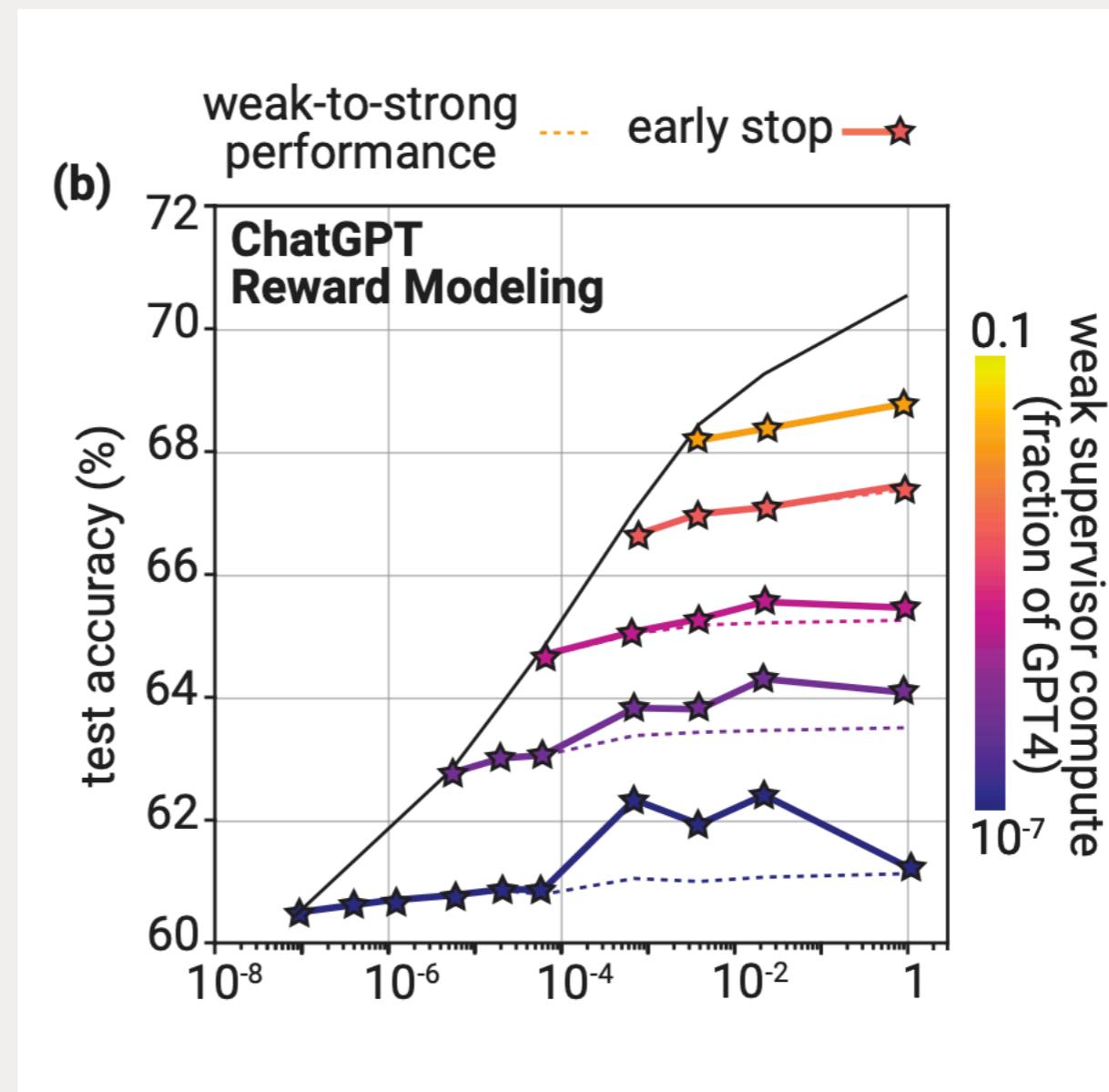
$$L_{\text{conf}}(f) = (1 - \alpha) \cdot \text{CE}(f(x), f_w(x)) + \alpha \cdot \text{CE}(f(x), \hat{f}_t(x))$$



# Understanding Weak-to-Strong Generalization

## Weak Imitation

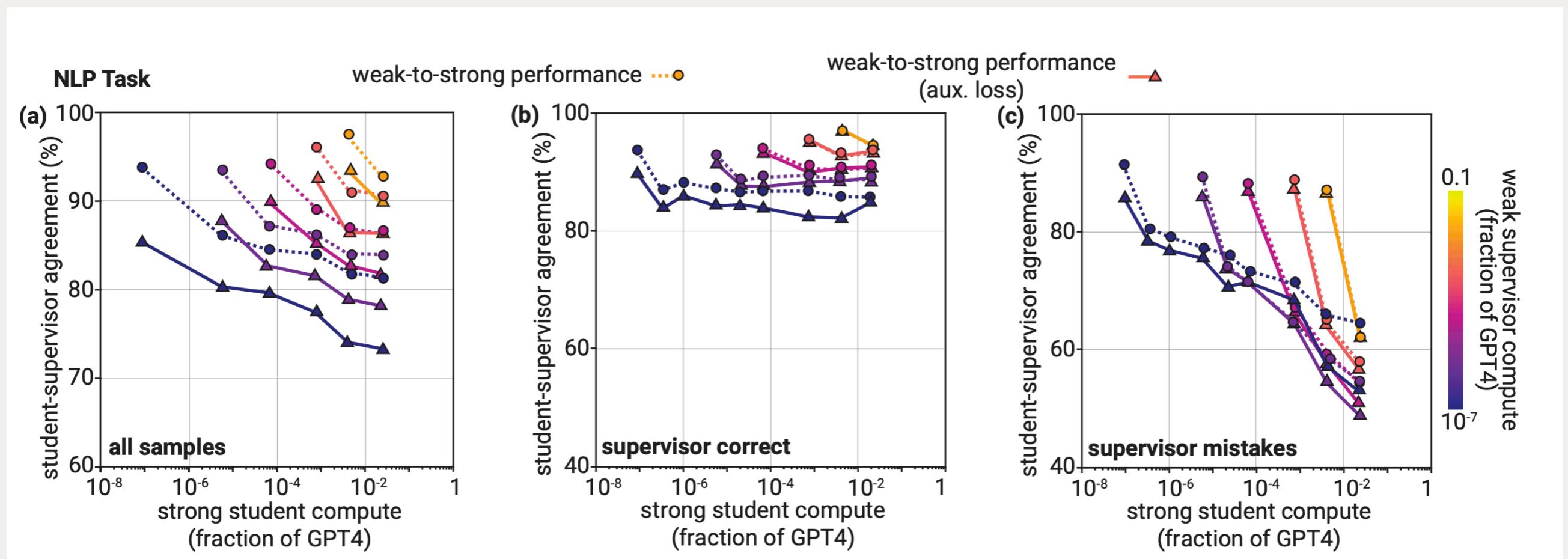
- Early Stopping



# Understanding Weak-to-Strong Generalization

## Weak Imitation

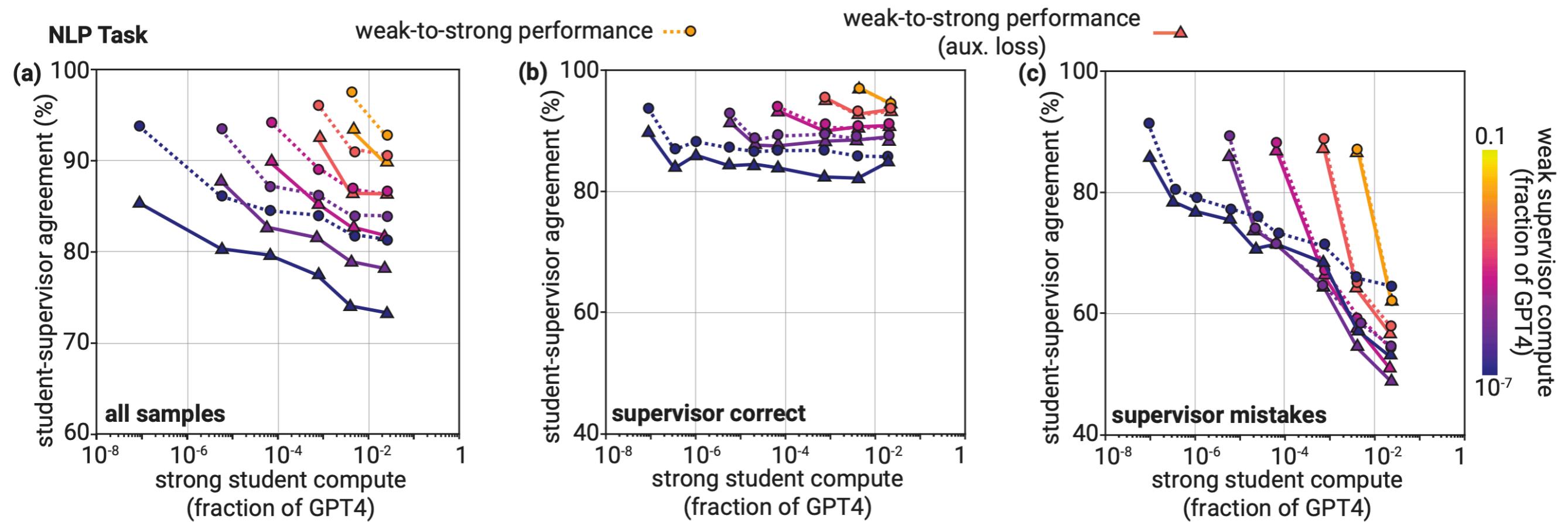
- Aux Loss prevents too much agreement



# Understanding Weak-to-Strong Generalization

## Weak Imitation

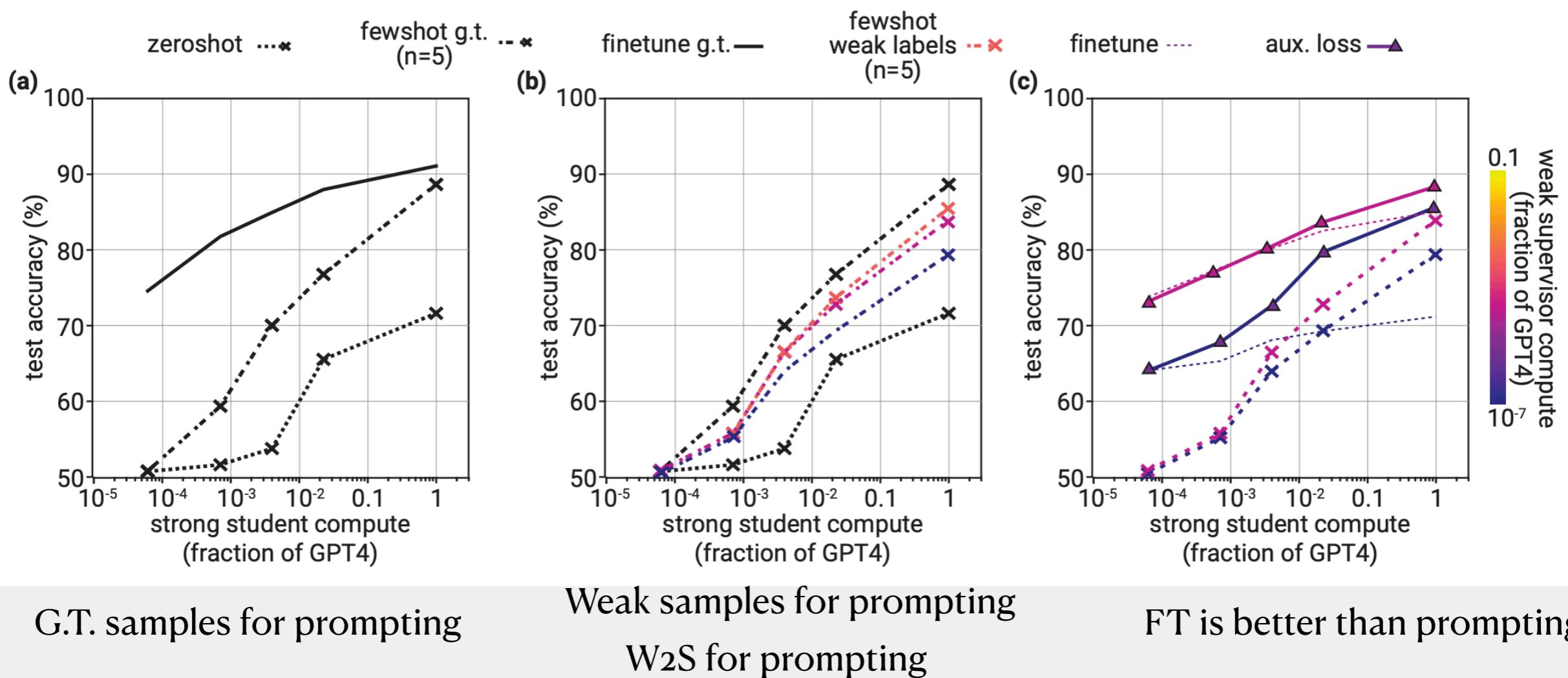
- Larger student models consistently agree less with the errors of the supervisor than smaller student models, despite
  - being trained to imitate the supervisor
  - not using early stopping
  - having larger capacity than smaller student models.



# Understanding Weak-to-Strong Generalization

## Saliency

- **Intuition:** weak-to-strong generalization might be feasible is when the task or concept we want to elicit is internally *salient* to the strong model.



# Understanding Weak-to-Strong Generalization

## Saliency

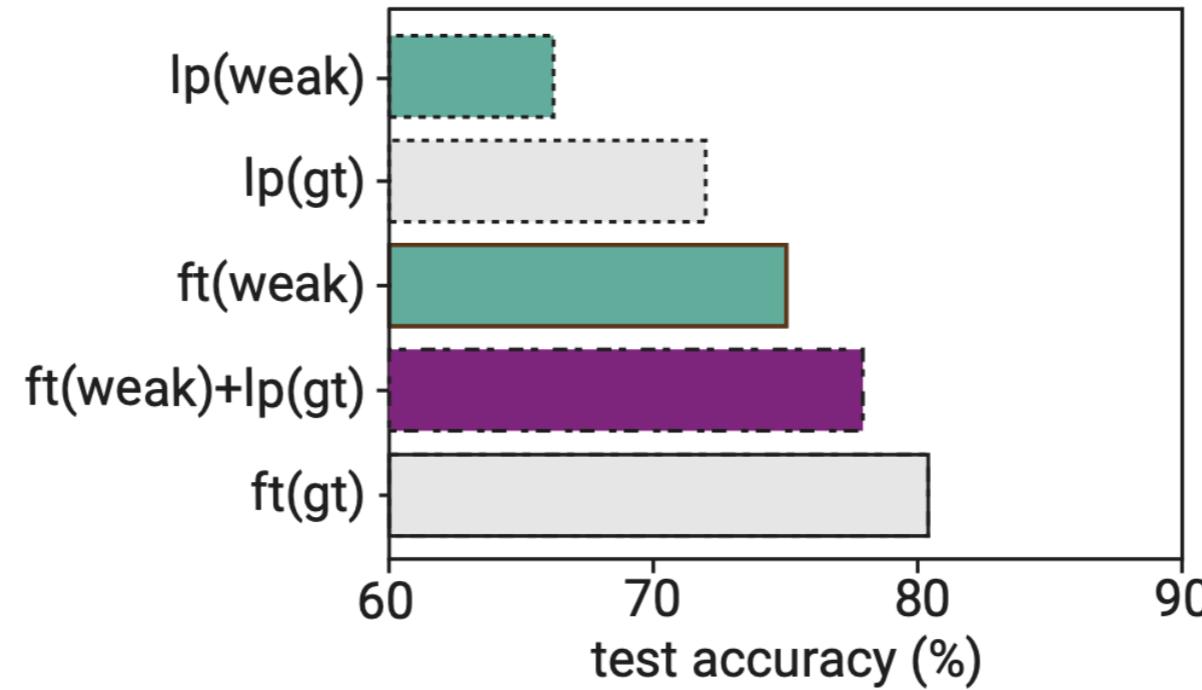


Figure 11: **Finetuning on weak supervisor labels makes the desired generalization more linearly represented.** We plot test accuracy for five different strategies, averaged across a subset of NLP tasks. **lp(weak)**: training a linear probe on the base model using weak labels, **lp(gt)**: training a linear probe on the base models using ground truth labels, **ft(weak)**: finetuning the model on weak labels, **ft(weak) + lp(gt)**: finetuning the model on weak labels *then* training a linear probe on ground truth labels, **ft(gt)**: finetuning the model on ground truth labels. Finetuning on the weak labels significantly increases the linearity of the ground truth concept.

Let's now see some theory.

# Random Features Model

4 Mar 2025

Weak-to-Strong Generalization Even in Random Feature Networks,  
Provably

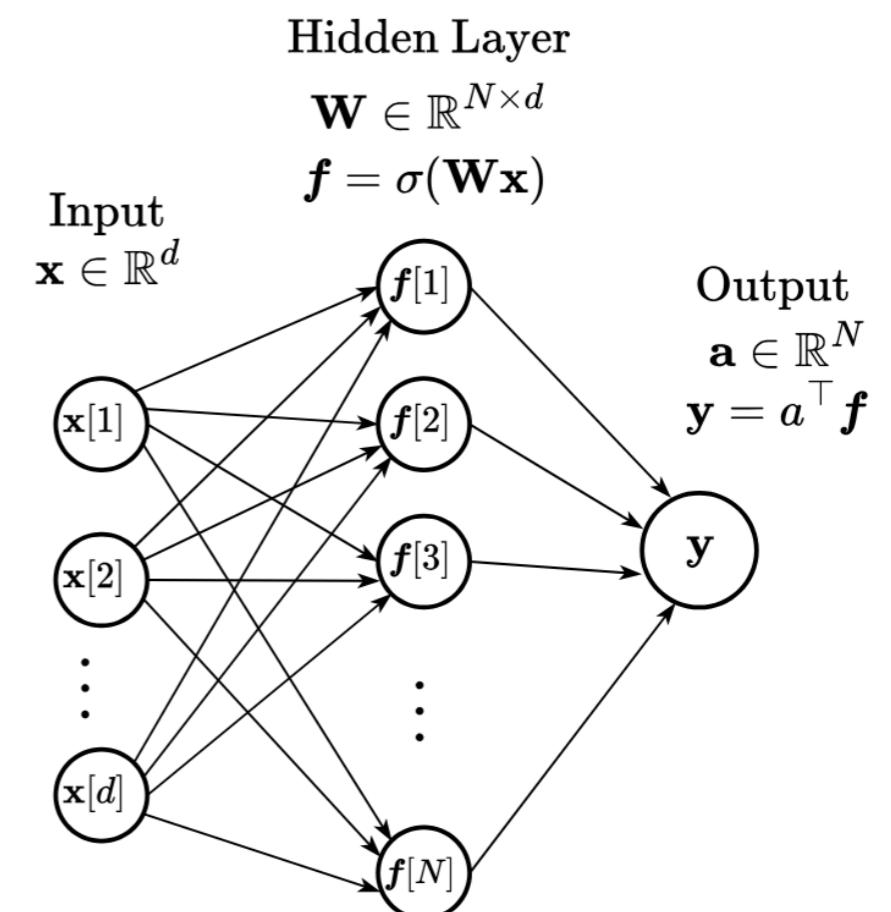
Márkó Medvedev<sup>1,\*</sup>, Kaifeng Lyu<sup>2,\*</sup>, Dingli Yu<sup>3</sup>, Sanjeev Arora<sup>4</sup>, Zhiyuan Li<sup>5</sup>, Nathan Srebro<sup>5</sup>

<sup>1</sup>University of Chicago, <sup>2</sup>Simons Institute, University of California, <sup>3</sup>Microsoft Research,

<sup>4</sup>Princeton University, <sup>5</sup>Toyota Technological Institute at Chicago

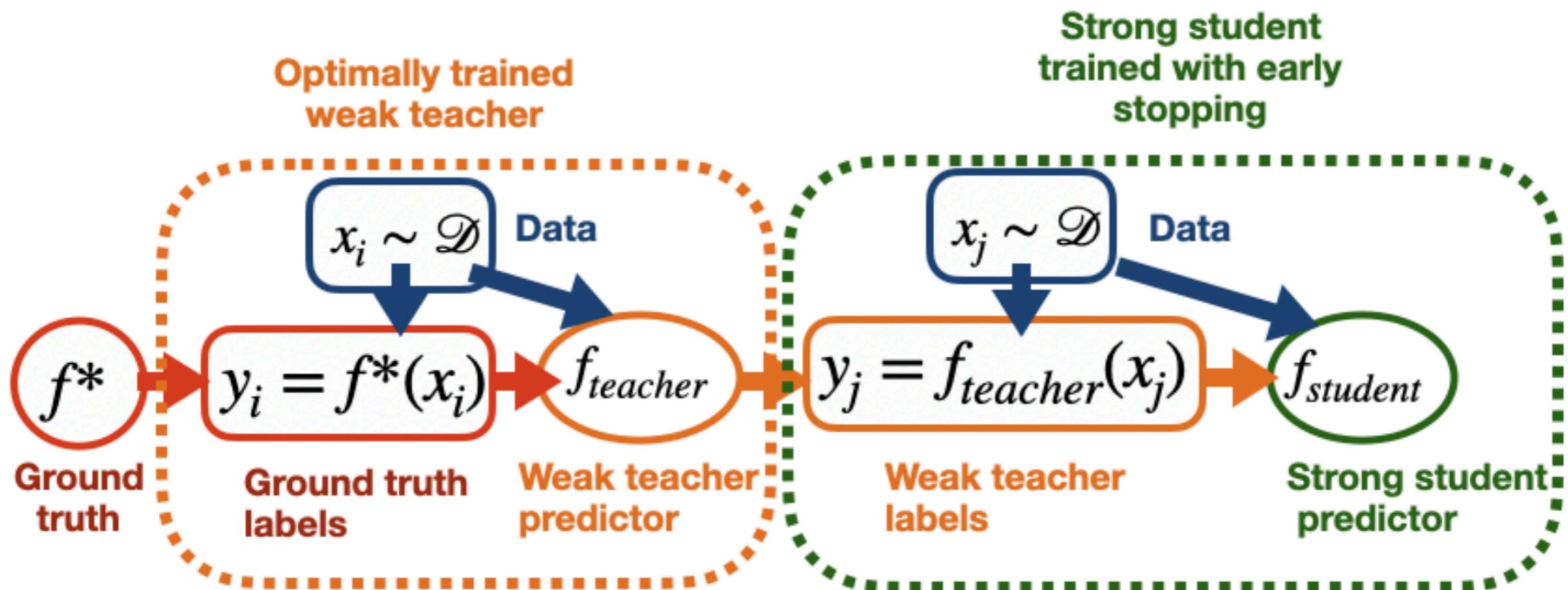
# Random Features Model

- Two-Layer fully connected neural network  $f_{\text{NN}}(\mathbf{x}) = \mathbf{a}^\top \sigma(\mathbf{W}\mathbf{x})$
- Random first layer weights
- Trained second layer
- Captures various deep learning phenomena.

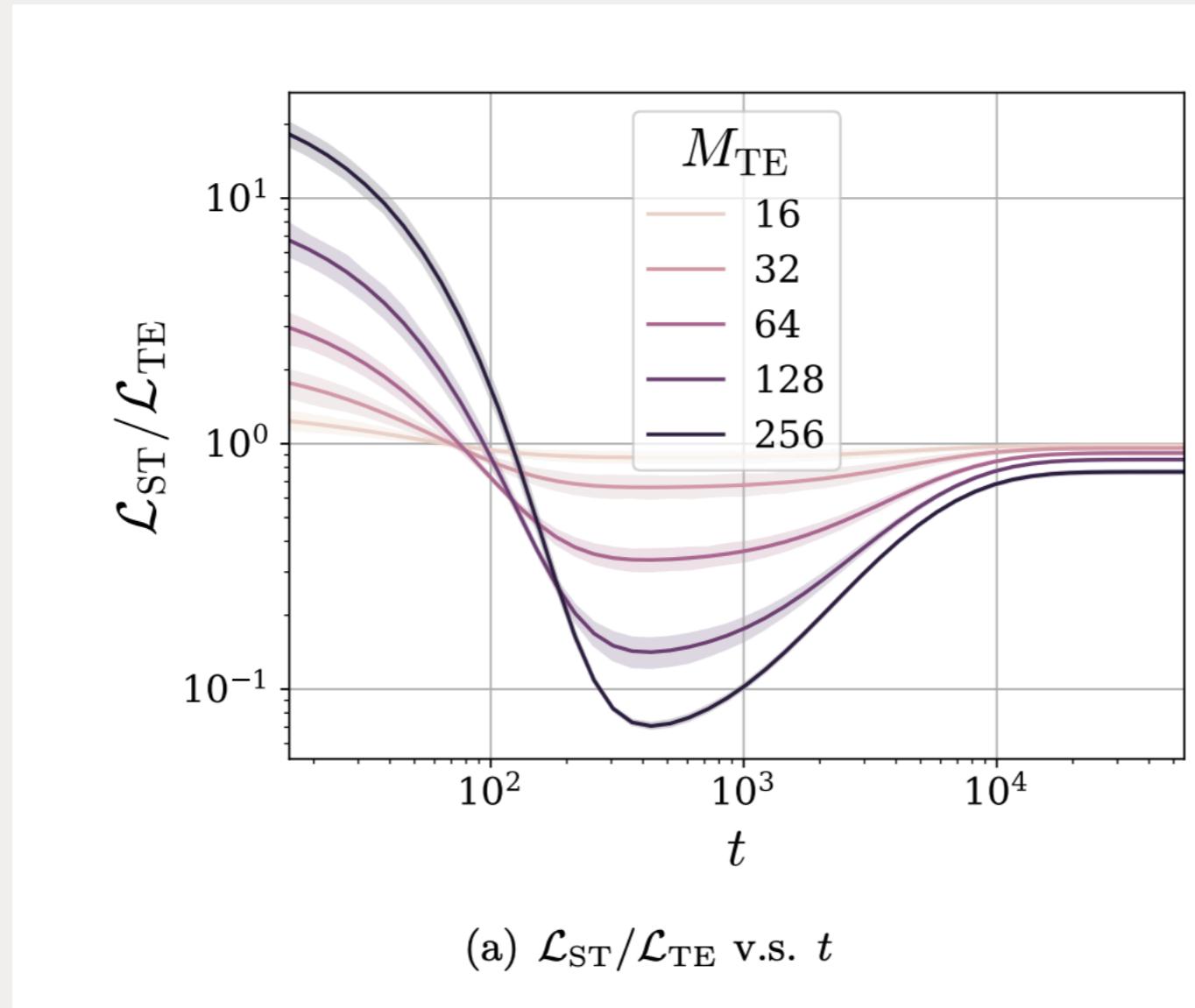


# Random Features Model

- Weak model has significantly less random features compared to the strong model.
- We fit the second layer (i.e., linear probing) using gradient flow.



# Random Features Model



Early Stopping is the key here

# The Role of Representations

---

**Representations Shape Weak-to-Strong Generalization:  
Theoretical Insights and Empirical Predictions**

---

**Yihao Xue<sup>1</sup> Jiping Li<sup>2</sup> Baharan Mirzasoleiman<sup>1</sup>**

# The Role of Representations

- Again, assume that we have a fixed set of features and we do linear probing.

$$f_w = \arg \min_{f \in \mathcal{F}_w} \left( \frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} (f(h_w(\tilde{x}_i)) - \tilde{y}_i)^2 + \beta_w R(f) \right).$$

$$f_{w2s} = \arg \min_{f \in \mathcal{F}_s} \left( \frac{1}{\hat{n}} \sum_{i=1}^{\hat{n}} (f(h_s(\hat{x}_i)) - f_w(h_w(\hat{x}_i)))^2 + \beta_s R(f) \right).$$

$$f_{sc} = \arg \min_{f \in \mathcal{F}_s} \left( \frac{1}{\hat{n}} \sum_{i=1}^{\hat{n}} (f(h_s(\hat{x}_i)) - \hat{y}_i)^2 + R_s(f) \right).$$

- Prediction Gap:**

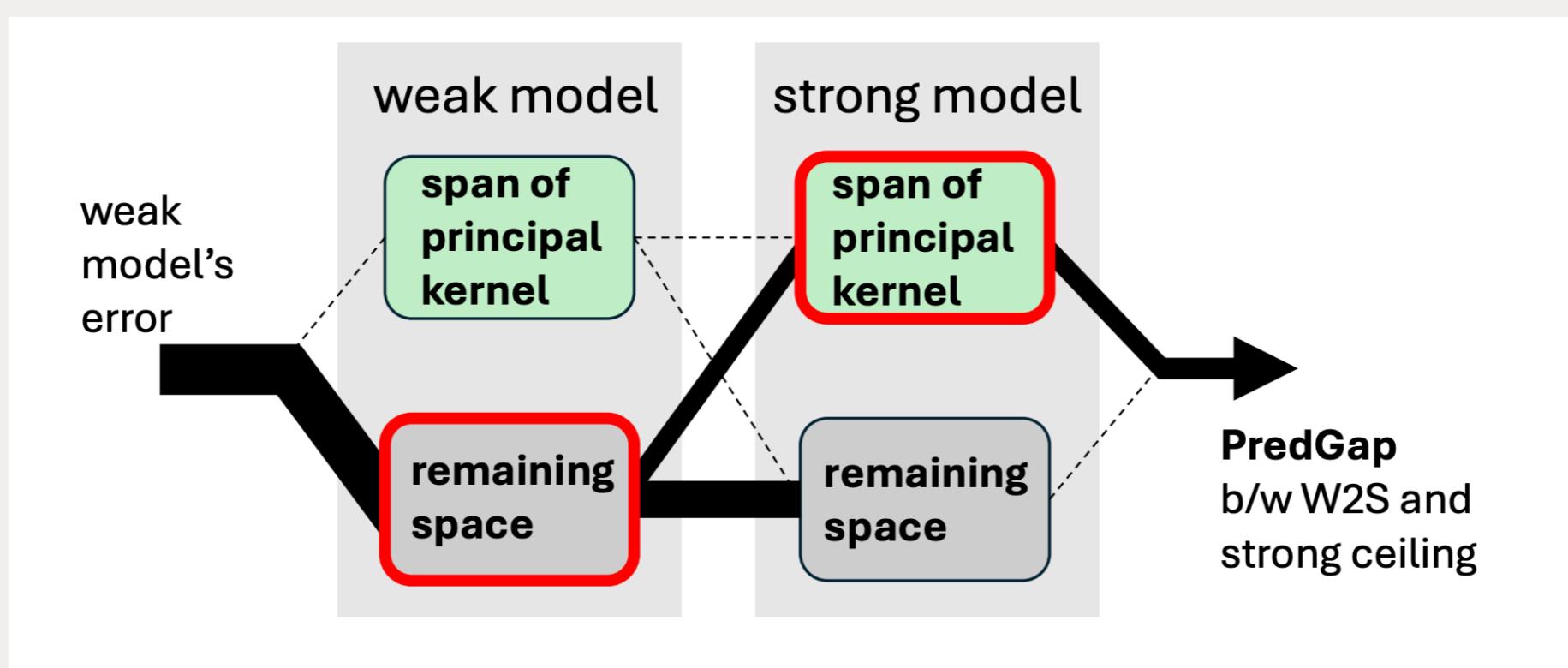
$$\textbf{PredGap} = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [(f_{w2s}(h_s(\mathbf{x})) - f_{sc}(h_s(\mathbf{x})))^2].$$

# Effect of Features

**Theorem 3.8** (Main result). *Under Assump. 3.7, and assuming reasonable regularization:  $\delta_w \leq \beta_w = O(1)$  and  $\delta_s \leq \beta_s = O(1)$ , let  $\hat{\mathbf{y}} = [\hat{y}_1 \ \hat{y}_2 \ \dots \ \hat{y}_{\hat{n}}]^\top$ . Then, w.h.p., we have*

$$\text{PredGap} = \|\mathbf{P}_s(\mathbf{I} - \mathbf{P}_w) \frac{1}{\sqrt{\hat{n}}} \hat{\mathbf{y}}\|^2 \pm o(1) \quad (1)$$

$\mathbf{P}_s(\mathbf{I} - \mathbf{P}_w)$  captures “what the weak model is unable to learn but is learnable by the strong model using their respective principal representations”. Therefore, it determines the mistakes that will be learned by the strong model, as discussed in the intuition. A more powerful weak model has a  $\mathbf{P}_w$  that covers more space, shrinking  $\mathbf{P}_s(\mathbf{I} - \mathbf{P}_w)$  and potentially leading to a smaller PredGap.



# Discrepancies are Virtue

Discrepancies are Virtue: Weak-to-Strong  
Generalization through Lens of Intrinsic Dimension

Yijun Dong<sup>1</sup>      Yicheng Li<sup>1</sup>      Yunai Li<sup>2</sup>      Jason D. Lee<sup>3</sup>      Qi Lei<sup>1</sup>

<sup>1</sup>New York University    <sup>2</sup>Shanghai Jiaotong University    <sup>3</sup>Princeton University

{yd1319, y19315}@nyu.edu    liyunai\_8528@sjtu.edu  
jasonlee@princeton.edu    ql518@nyu.edu

# Discrepancies are Virtue

(a) **Weak teacher model**  $f_w(\mathbf{x}) = \phi_w(\mathbf{x})^\top \boldsymbol{\theta}_w$  is supervised finetuned over  $\tilde{\mathcal{S}}$ :

$$\boldsymbol{\theta}_w = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmin}} \left\| \tilde{\Phi}_w \boldsymbol{\theta} - \tilde{\mathbf{y}} \right\|_2^2 + \alpha_w \|\boldsymbol{\theta}\|_2^2. \quad (2)$$

(b) **W2S model**  $f_{w2s}(\mathbf{x}) = \phi_s(\mathbf{x})^\top \boldsymbol{\theta}_{w2s}$  is finetuned over the strong feature  $\phi_s$  through  $\mathcal{S}_x$  and their pseudo-labels generated by the weak teacher model:

$$\boldsymbol{\theta}_{w2s} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmin}} \left\| \Phi_s \boldsymbol{\theta} - \Phi_w \boldsymbol{\theta}_w \right\|_2^2 + \alpha_{w2s} \|\boldsymbol{\theta}\|_2^2 \quad (3)$$

(c) **Strong SFT model**  $f_s(\mathbf{x}) = \phi_s(\mathbf{x})^\top \boldsymbol{\theta}_s$  is a strong baseline where the strong feature  $\phi_s$  is supervisely finetuned over the small labeled set  $\tilde{\mathcal{S}}$  directly:

$$\boldsymbol{\theta}_s = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmin}} \left\| \tilde{\Phi}_s \boldsymbol{\theta} - \tilde{\mathbf{y}} \right\|_2^2 + \alpha_s \|\boldsymbol{\theta}\|_2^2. \quad (4)$$

(d) **Strong ceiling model**  $f_c(\mathbf{x}) = \phi_s(\mathbf{x})^\top \boldsymbol{\theta}_c$  is a reference for the ceiling performance where  $\phi_s$  is supervisely finetuned over  $\mathcal{S} \cup \tilde{\mathcal{S}}$ , assuming access to the unknown labels  $\mathbf{y} = [y_1, \dots, y_N]^\top$ :

$\alpha_w, \alpha_{w2s}, \alpha_s, \alpha_c \rightarrow 0$
--

$$\boldsymbol{\theta}_c = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmin}} \left\| \begin{bmatrix} \tilde{\Phi}_s \\ \Phi_s \end{bmatrix} \boldsymbol{\theta} - \begin{bmatrix} \tilde{\mathbf{y}} \\ \mathbf{y} \end{bmatrix} \right\|_2^2 + \alpha_c \|\boldsymbol{\theta}\|_2^2. \quad (5)$$

# Discrepancies are Virtue

**Theorem 1** (W2S model (formally in B.1)). *Assuming Assumptions 1 and 2 and  $\phi_w(\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\theta}_d, \Sigma_w)$ , for  $n > d_w + 1$ ,  $\text{ER}(f_{\text{w2s}}) = \text{Var}(f_{\text{w2s}}) + \text{Bias}(f_{\text{w2s}})$  satisfies*

$$\text{Var}(f_{\text{w2s}}) = \frac{\sigma^2}{n - d_w - 1} \left( d_{s \wedge w} + \frac{d_s}{N} (d_w - d_{s \wedge w}) \right),$$

$$\text{Bias}(f_{\text{w2s}}) \leq \frac{\rho_w(n)}{n} + \frac{\rho_s(N)}{N} \leq \rho_w + \rho_s, \quad \begin{array}{l} \text{= zero; variance dominated regime} \\ \text{if } n \leq d_w + 1 \end{array}$$

where the inequality for  $\text{Bias}(f_{\text{w2s}})$  is strict if  $\rho_w(n)/n > 0$  and  $d_s < d_w$ .

$$\text{Var}(f_w) = \sigma^2 \frac{d_w}{n}, \quad \text{Bias}(f_w) = \frac{\rho_w(n)}{n} \leq \rho_w.$$

$$\text{Var}(f_s) = \sigma^2 \frac{d_s}{n}, \quad \text{Bias}(f_s) = \frac{\rho_s(n)}{n} \leq \rho_s,$$

$$\text{Var}(f_c) = \sigma^2 \frac{d_s}{N+n}, \quad \text{Bias}(f_c) = \frac{\rho_s(N+n)}{N+n} \leq \rho_s,$$

$$\boxed{\begin{aligned} \rho_s &= \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [(\phi_s(\mathbf{x})^\top \boldsymbol{\theta} - f_*(\mathbf{x}))^2], \\ \rho_w &= \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [(\phi_w(\mathbf{x})^\top \boldsymbol{\theta} - f_*(\mathbf{x}))^2], \end{aligned}}$$

# Discrepancies are Virtue

- **Intuition.**
- **Task:** classify the brand of a car in an image
- **Strong Features:** information in the logo  
(low intrinsic dimension  $ds$ )
- **Weal Features:** complex information in the design  
(large intrinsic dimension  $dw$ ).

# Quantifying the Gains

23 Oct 2024

## Quantifying the Gain in Weak-to-Strong Generalization

**Moses Charikar**

Stanford University

[moses@cs.stanford.edu](mailto:moses@cs.stanford.edu)

**Chirag Pabbaraju**

Stanford University

[cpabbaraju@cs.stanford.edu](mailto:cpabbaraju@cs.stanford.edu)

**Kirankumar Shiragur**

Microsoft Research

[kshiragur@microsoft.com](mailto:kshiragur@microsoft.com)

# Quantifying the Gains

## Main Result

Let  $\text{MSE}(p, q) = \mathbb{E}_x[p(x) - q(x)]^2$  and suppose the strong model obtains  $f_{sw}$  from a convex space  $\mathcal{F}_s$

$$f_{sw} = \operatorname{argmin}_{f \in \mathcal{F}_s} \text{MSE}(f \circ h_s, f^* \circ h^*)$$

Then, assuming realizability (i.e.,  $f^* \circ h^* = f_s \circ h_s$  for some  $f_s \in \mathcal{F}_s$ ),

$$\underbrace{\text{MSE}(f_w \circ h_w, f^* \circ h^*) - \text{MSE}(f_{sw} \circ h_s, f^* \circ h^*)}_{\text{GAIN IN PERFORMANCE}} \geq \underbrace{\text{MSE}(f_{sw} \circ h_s, f_w \circ h_w)}_{\text{MISFIT BETWEEN WEAK AND STRONG MODEL}}$$

# Future Work

- **Updating the representation (beyond linear probing):**  
The strong model can learn the features of the weak model.
  - Language learning example.
- **The role of implicit regularization:**  
Even in very simple linear models.
  - Linear regression example
  - Stay tuned :)

**Thanks!**