

# TRANSFER LEARNING FOR WEAK-TO-STRONG GENERALIZATION



Seamus  
Somerstep



Felipe  
Maia Polo



Mouli  
Banerjee



Subha  
Maity



Ya'acov  
Ritov

Illustrations by  
Rami Ritov



Mikhail  
Yurochkin

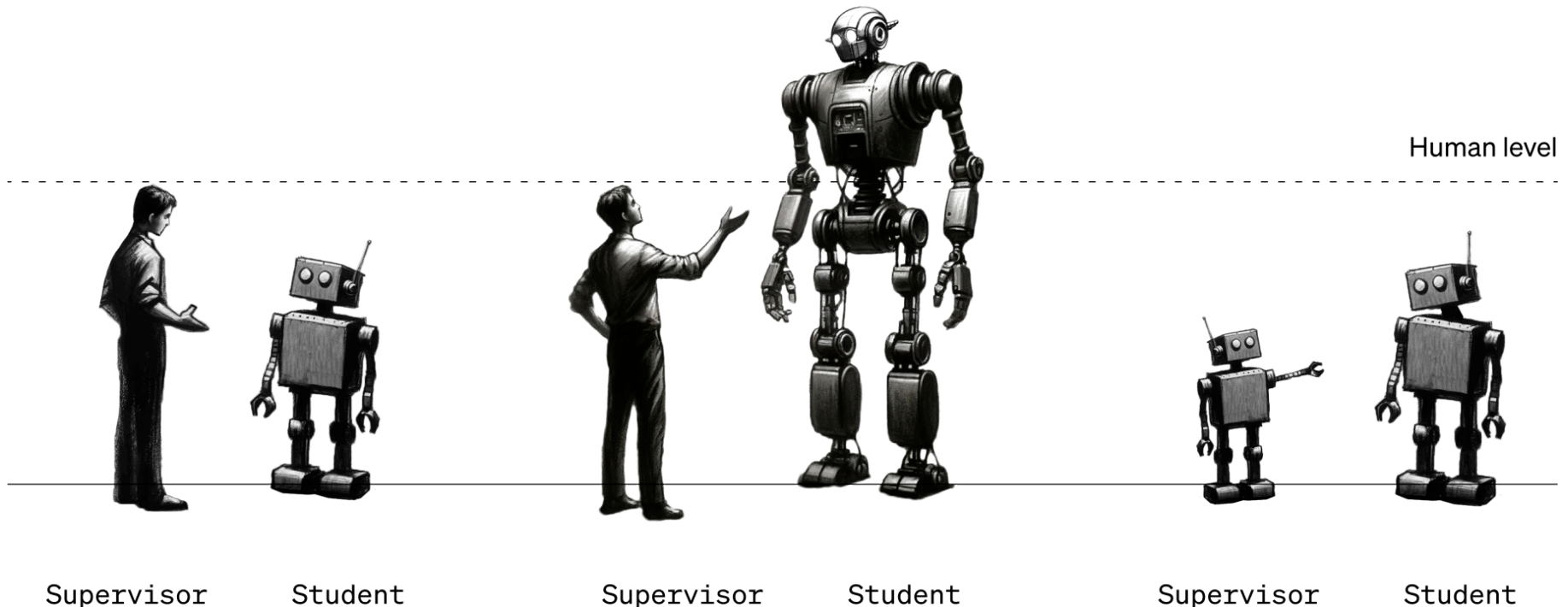
# SUPERALIGNMENT

**Goal:** Supervise the fine-tuning of a larger (more capable) pretrained LLM with a smaller (less capable) LLM.

Traditional ML

Superalignment

Our Analogy



Burns et al (2023)

# SUPERALIGNMENT

**Goal:** Supervise the fine-tuning of a larger (more capable) pretrained LLM with a smaller (less capable) LLM.

❌ The pretrained LLM may imitate the weak teacher's errors/mistakes; thus limiting its performance to no better than its (weak) teacher.

✅ Pretrained LLMs have impressive (latent) raw abilities—there is no need to teach them new tasks from scratch.

We focus on **eliciting the latent abilities** of pretrained LLMs.

# TRANSFER LEARNING SETUP

Given

- $\{(X_{P,i}, Y_{P,i})\}_{i=1}^{n_P}$ :  $n_P$  samples from source  $P$
- $\{(X_{Q,i}, Y_{Q,i})\}_{i=1}^{n_Q}$ :  $n_Q$  (labeled) samples from target  $Q \neq P$
- $n_P \gg n_Q$

Find  $f$  such that  $\mathbf{E}[\ell(f(X_Q), Y_Q)]$  is small.

Transfer is generally futile without similarity assumptions on  $P$  and  $Q$ .

# WEAKLY-SUPERVISED TRANSFER SETUP

Given

- $\{(X_{P,i}, Y_{P,i})\}_{i=1}^{n_P}$ :  $n_P$  samples from  $P$
- $\{(X_{Q,i}, Y'_{Q,i})\}_{i=1}^{n_Q}$ :  $n_Q$  *weakly-labeled* samples from  $Q$

Find  $f$  such that  $\mathbf{E}[\ell(f(X_Q), Y_Q)]$  is small.

- $(X_Q, Y_Q)$ : target sample with *gold-standard/strong* label  $Y_Q$

The main challenge here is the learner has no (strong) labels from  $Q$ !

# WEAK-TO-STRONG GENERALIZATION AS TL

superalignment	weakly-supervised TL
pretrained LLM	$Y_P   X_P$
(super)alignment task	$Y_Q   X_Q$
weak teacher	$Y'_Q   X_Q$

In superalignment, the learner has access to the pretrained LLM and weak teacher; it has no (direct) supervision on the alignment task.

# WEAK-TO-STRONG GENERALIZATION AS TL

Assume

1. Latent concept shift:  $Y_P | X_P, Y_Q | X_Q$  are mixtures of the same (mixture) components:

$$P_{Y|X} = \sum_{\theta \in \Theta} \pi_{\theta} G_{\theta}(\cdot | X), \quad Q_{Y|X} = G_{\theta_Q}(\cdot | X)$$

2.  $n_P$  is extremely large, so the pretrained LLM is exactly  $Y_P | X_P$

Transfer learning under latent concept shift is a deconvolution problem: learning the target concept  $\theta_Q$  from the weak labels.

**Challenge:** deconvolution without knowledge of the  $G_{\theta}$ 's.

# LATENT CONCEPT MODEL OF LLMS

$$\text{GPT}(Y | X) = \sum_{\theta \in \Theta} g(Y | X, \theta) \pi(\theta | X).$$

- $\theta$ : (latent) concept (eg nationality, occupation)
- $g(\cdot | X, \theta)$ : predictive distribution associated with  $\theta$ 
  - English  $\sim g(\cdot | \text{Isaac Newton is, nationality})$
  - scientist  $\sim g(\cdot | \text{Isaac Newton is, occupation})$
- $\pi(\cdot | X)$ : prior (distribution) on concepts

The predictive distributions  $g(\cdot | X, \theta), \theta \in \Theta$  encode the pretrained LLM's (latent) abilities.



# LATENT CONCEPT SHIFT

Latent concept shift:  $Y_P | X_P, Y_Q | X_Q$  are mixtures of the same (mixture) components:

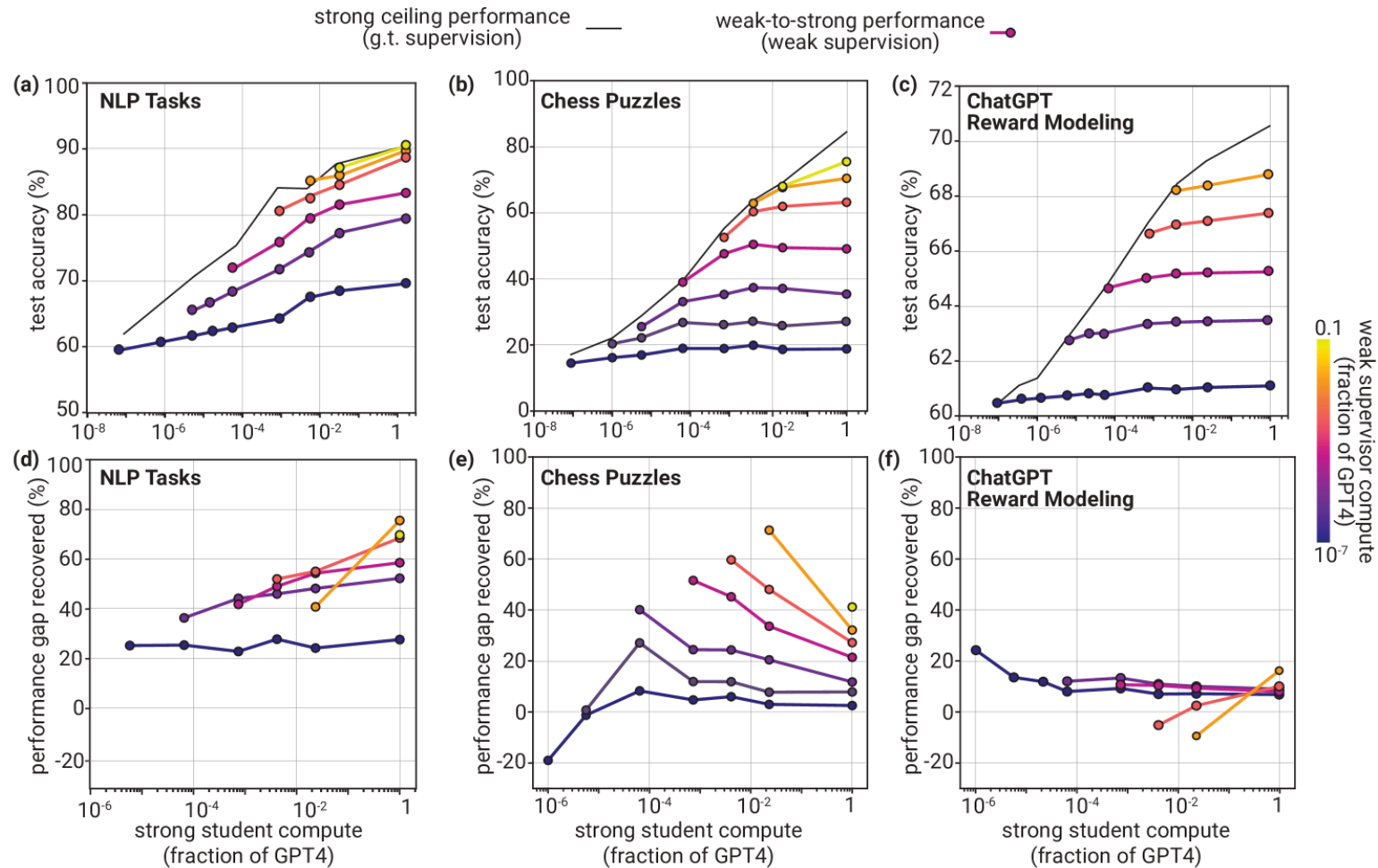
$$P_{Y|X} = \sum_{\theta \in \Theta} \pi_{\theta} G_{\theta}(\cdot | X), \quad Q_{Y|X} = G_{\theta_Q}(\cdot | X)$$

- The pretrained LLM satisfies the latent concept model.

$$\text{GPT}(Y | X) = \sum_{\theta \in \Theta} g(Y | X, \theta) \pi(\theta | X).$$

- The superalignment task is encoded by  $g(\cdot | X, \theta_Q)$  (for some  $\theta_Q \in \Theta$ ).

# (NAIVE) FINE-TUNING



Burns et al (2023)

# W2S REGRESSION

The learner is given

- $X \in \mathbf{R}^{n \times d}$ : (fixed) input matrix,
- $y_P \in \mathbf{R}^n$ : source model outputs (on inputs in  $X$ ),
- $y'_Q \in \mathbf{R}^n$ : weak teacher outputs.

Assume

1. **Latent concept shift**: source and target regression functions are mixtures of the same (mixture) components:

$$f_P^*(X) = \sum_{\theta \in \Theta} \pi_{\theta} g_{\theta}(X), \quad f_Q^*(X) = g_{\theta_Q}(X).$$

2. The source model is exactly  $f_P^*$ ; ie  $y_P = f_P^*(X)$ .
3. The weak teacher is unbiased; ie  $\mathbf{E}[y'_Q \mid X] = f_Q^*(X)$ .

# W2S REGRESSION

The learner seeks  $f_Q$  such that

1.  $\frac{1}{n} \|f_Q(X) - f_Q^*(X)\|_2^2 \leq \frac{1}{n} \|y_p - f_Q^*(X)\|_2^2 \triangleq \frac{1}{n} \|\epsilon_P\|_2^2$ ,
2.  $\frac{1}{n} \|f_Q(X) - f_Q^*(X)\|_2^2 \leq \frac{1}{n} \|y'_Q - f_Q^*(X)\|_2^2 \triangleq \frac{1}{n} \|\epsilon'_Q\|_2^2$ ;

ie  $f_Q$  that improves upon both the source model and weak teacher.

- $\frac{1}{n} \|\epsilon_P\|_2^2$  is the source model's MSE,
- $\frac{1}{n} \|\epsilon'_Q\|_2^2$  is the weak teacher's MSE.

# FINE-TUNING ON WEAK OUTPUTS

$$\hat{f}_\lambda \leftarrow \operatorname{argmin}_f \left\{ \begin{array}{l} \frac{1}{2} \|y'_Q - f(X)\|_2^2 \\ + \lambda \frac{1}{2} \|y_P - f(X)\|_2^2 \end{array} \right\} \quad (\text{FT})$$

- $\frac{1}{2} \|y'_Q - f(X)\|_2^2$  is the loss with respect to outputs from the weak teacher.
- $\frac{1}{2} \|y_P - f(X)\|_2^2$  regularizes  $f$  towards the source model.

(FT) is an analogy for supervised fine-tuning (SFT) and reinforcement learning from human feedback (RLHF).

# FINE-TUNING ON WEAK OUTPUTS

Lem: The MSE of  $\hat{f}_\lambda$  from (FT) is

$$\mathbf{E} \left[ \frac{1}{n} \|\hat{f}_\lambda(X) - f_Q^*(X)\|_2^2 \right] = \frac{\lambda^2}{(1+\lambda)^2} \frac{1}{n} \|\epsilon_P\|_2^2 + \frac{1}{(1+\lambda)^2} \frac{1}{n} \|\epsilon'_Q\|_2^2.$$

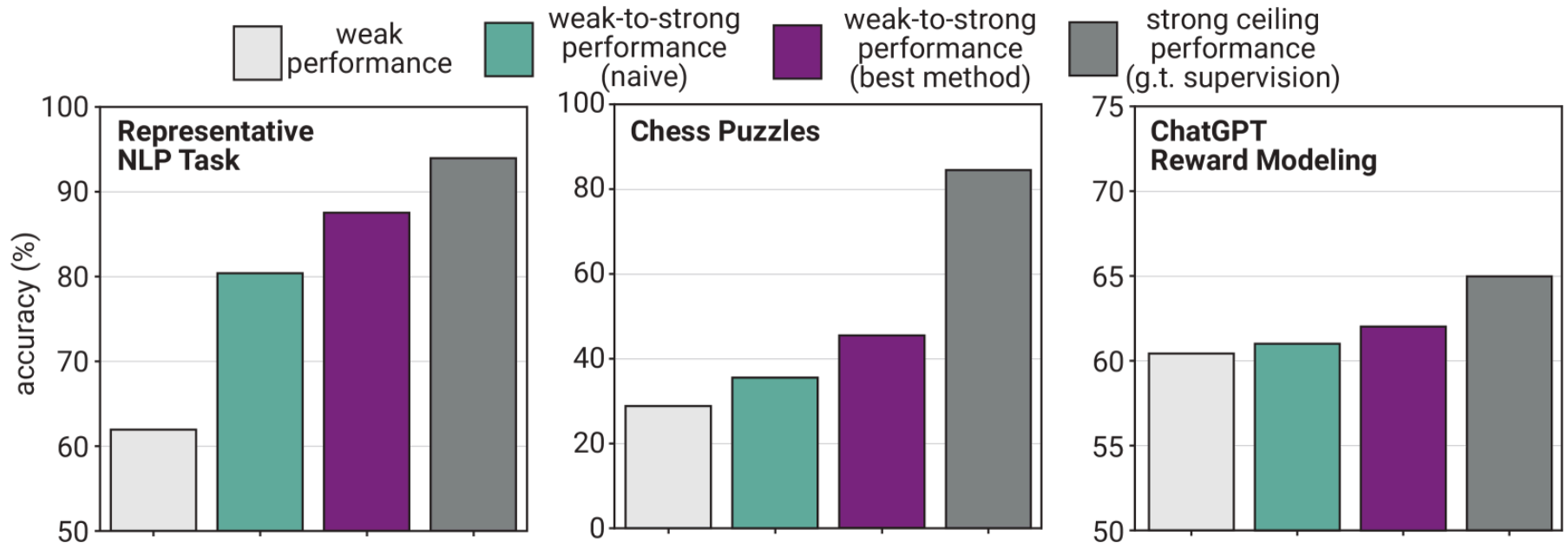
Takeaway: The MSE of  $\hat{f}_\lambda$  from (FT) is limited by

- the MSE of the source model (if  $\|\epsilon_P\|_2^2 < \|\epsilon'_Q\|_2^2$ ),
- the MSE of the weak teacher (if  $\|\epsilon'_Q\|_2^2 < \|\epsilon_P\|_2^2$ ).

# SUPERVISING GPT-4 WITH GPT-2

*if we fine-tune GPT-4 with labels from a GPT-2-level model, we typically recover about  $\frac{1}{2}$  of the performance gap between the two models.*

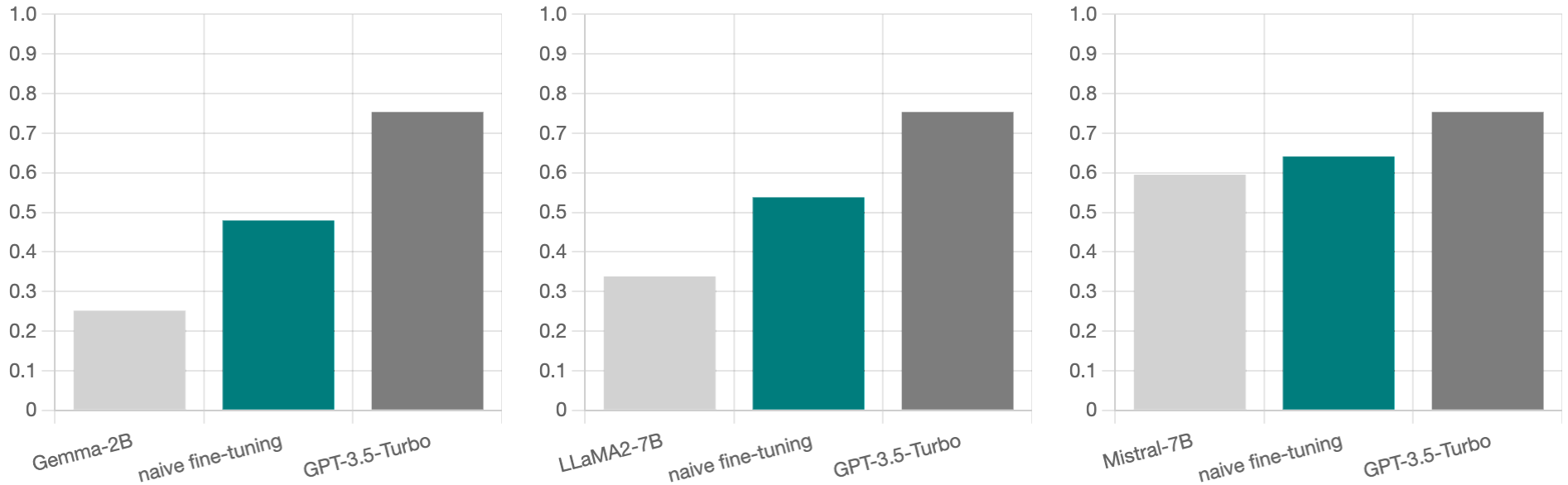
— Burns et al



Burns et al (2023)

# MATH REASONING EXPS PREVIEW

Q: Does the fine-tuned LLM outperform the pretrained LLM (without fine-tuning)?



No because fine-tuning does not does not leverage the *latent* abilities of the pretrained LLM.



# W2S REGRESSION VIA DECONVOLUTION

Recall the source and target regression functions are mixtures of the same (mixture) components:

$$f_P^* = \sum_{\theta \in \Theta} \pi_\theta g_\theta, \quad f_Q^* = g_{\theta_Q}.$$

Idea: Restrict  $f$  to (the convex hull of) the mixture components:

$$\hat{f} \leftarrow \left\{ \begin{array}{l} \operatorname{argmin}_f \quad \frac{1}{2} \|y'_Q - f(X)\|_2^2 \\ \text{subject to} \quad f \in \operatorname{cvx}(\{g_\theta\}_{\theta \in \Theta}) \end{array} \right\}. \quad (\text{dcnv})$$

# W2S REGRESSION VIA DECONVOLUTION

Lem: The MSE of  $\hat{f}$  from (dcnv) is at most

$$\begin{aligned} \frac{1}{n} \|\hat{f}(X) - f_Q^*(X)\|_2^2 &\leq \sup_{\theta \in T_{\mathcal{G}}(f_Q^*(X)) \cap \mathbf{S}^{n-1}} \frac{1}{n} ((\epsilon'_Q)^\top \theta)^2 \\ &\ll \frac{1}{n} \|\epsilon'_Q\|_2^2, \end{aligned}$$

where  $\mathcal{G} \triangleq \text{cvx}(\{g_\theta(X)\}_{\theta \in \Theta}) \subset \mathbf{R}^n$  (as long as  $f_Q^*(X) \in \mathcal{G}$ ).

- ✓ The MSE of (dcnv) is **not** limited by the MSEs of the source model and weak teacher.
- ✗ Unfortunately, (dcnv) is impractical because it requires knowledge of  $g_\theta$ 's (instead of  $f_P$ ).

# ELICITING LATENT ABILITIES OF LLMS

Recall the latent concept model of LLMs:

$$\text{GPT}(Y | X) = \sum_{\theta \in \Theta} g(Y | X, \theta) \pi(\theta | X).$$

Bayesian explanation of in-context learning (ICL): ICL helps the pretrained LLM to (implicitly) infer the target concept:

$$\begin{aligned} \operatorname{argmax}_y \text{GPT}(y | X_1, Y_1, \dots, X_K, Y_K, X) \\ \approx \operatorname{argmax}_y g(y | X, \theta), \end{aligned}$$

where  $Y_k | X_k \sim g(\cdot | X_k, \theta)$ .

# ELICITING LATENT ABILITIES OF LLMS

Idea: Use examples from the weak teacher (instead of from  $g(\cdot | \cdot, \theta_Q)$ ) as ICL examples; ie hope/pray that

$$\begin{aligned} \operatorname{argmax}_y \text{GPT}(y | X_1, Y_1, \dots, X_K, Y_K, X) \\ \approx \operatorname{argmax}_y g(y | X, \theta_Q), \end{aligned}$$

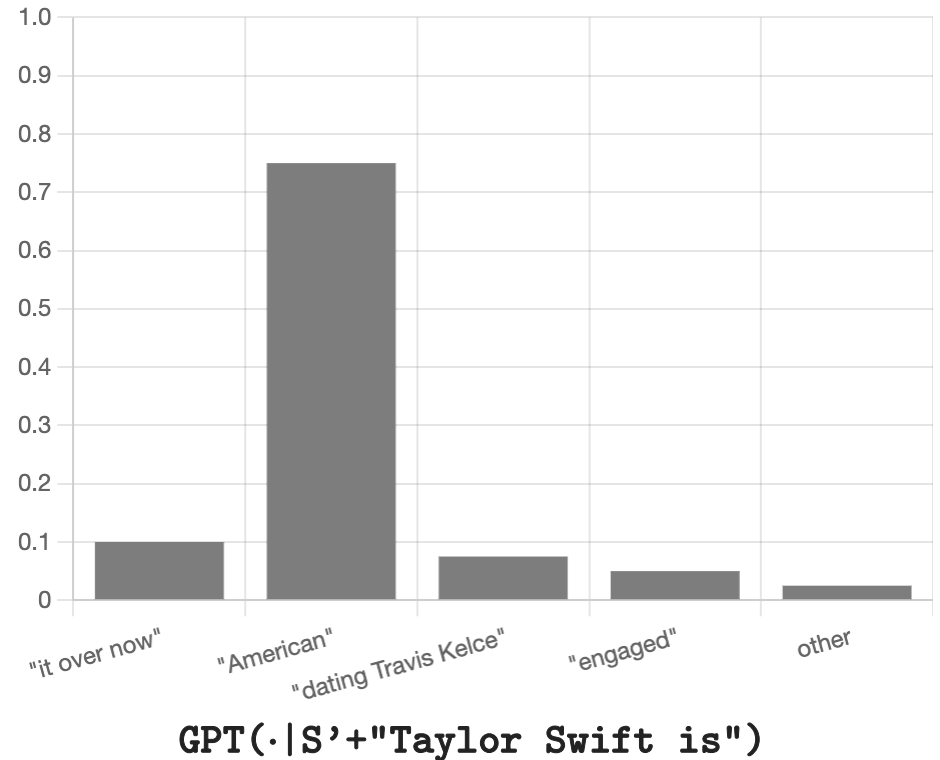
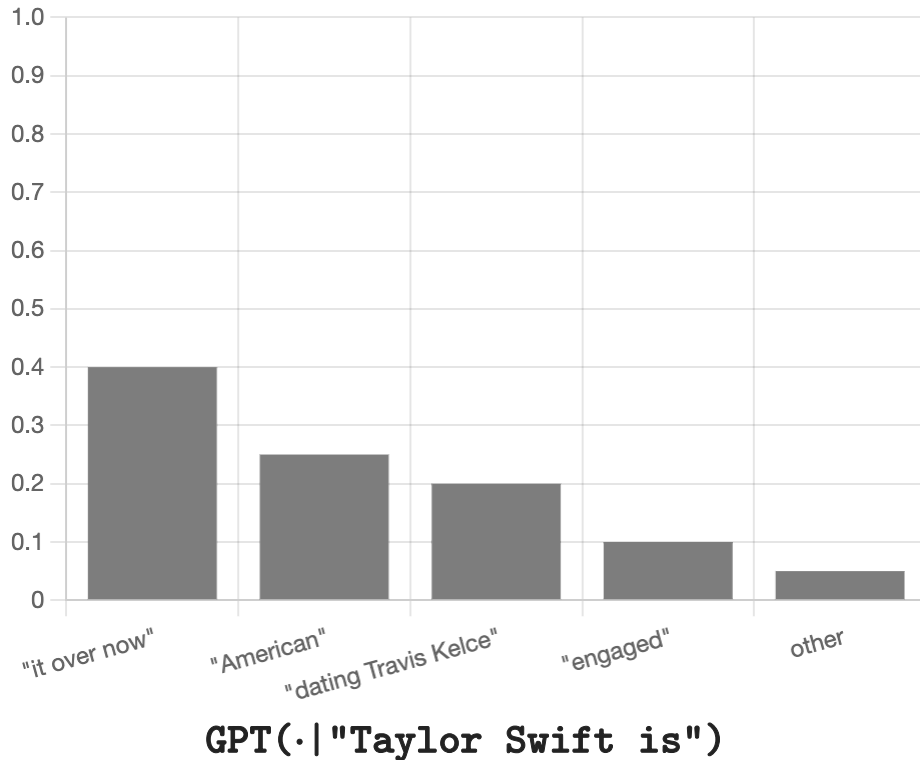
where  $(X_k, Y_k)$ 's are examples from the weak teacher .

The pretrained model must be "strong enough" to learn  $\theta_Q$  in-context from the (weak) teacher.

# EX: ELICITING THE LATENT ABILITIES OF LLMs

**Goal:** teach an LLM to respond with the nationality of notable people

$\mathcal{S}'$  consists of (" $[name]$  is", " $[nationality]$ ") pairs from weak teacher (so some nationalities are incorrect).



# FROM DECONVOLUTION TO ICL REFINEMENT

: W2S regression

$$\left\{ \begin{array}{l} \operatorname{argmin}_f \quad \frac{1}{2} \|y'_Q - f(X)\|_2^2 \\ \text{subject to} \quad f \in \operatorname{cvx}(\{g_\theta\}_{\theta \in \Theta}) \end{array} \right\} \approx g_{\theta_Q} = f_Q^*$$

:: superalignment

$$\begin{aligned} & \operatorname{argmax}_y \operatorname{GPT}(y \mid X_1, Y_1, \dots, X_K, Y_K, X) \\ & \approx \operatorname{argmax}_y g(y \mid X, \theta_Q) \\ & = \operatorname{argmax}_y Q_{Y|X}(y), \end{aligned}$$

ICL refinement solves the deconvolution problem *without knowledge of*  $g(\cdot \mid \cdot, \theta)$ ,  $\theta \in \Theta$ !

# IN-CONTEXT LEARNING (ICL) REFINEMENT

Idea: Fine-tune on refined outputs from  $\mathbf{GPT}(\cdot | \mathbf{P}, \mathbf{X}_Q)$ .

Require: weakly labeled dataset  $\mathcal{D}'_Q \triangleq \{(X_{Q,i}, Y'_{Q,i})\}_{i=1}^{n_Q}$

For  $(X_{Q,i}, Y'_{Q,i}) \in \mathcal{D}'_Q$

1. select a subset  $\mathcal{S}_i$  of  $\mathcal{D}'_Q$  as ICL examples (eg  $K$ -NN of  $X_{Q,i}$ ),
2. refine weak teacher output:  $Y_{Q,i} \sim \mathbf{GPT}(\cdot | \mathcal{S}_i, \mathbf{X}_{Q,i})$ .

Fine-tune the pretrained LLM on the refined weak teacher outputs  $\{(X_{Q,i}, Y_{Q,i})\}_{i=1}^{n_Q}$ .

# MATH REASONING EXPERIMENTS

pretrained LLMs:

- GPT-3.5-Turbo
- GPT-4o mini

weak teachers (fine-tuned on gold standard outputs):

- Gemma-2B
- LLaMA2-7B
- Mistral-7B

We use GPT-4o to assess whether outputs agree with the intermediate steps and final answers in the answer key.



# GRADE SCHOOL MATH 8K (GSM8K)

- 8.8k (7.5k/1.3k train/test) grade school math problems created by human problem writers
- They take between 2 and 8 steps to solve.
- The solutions mostly entail a sequence of basic arithmetic operations (+ - / \*).

*A bright middle school student should be able to solve every problem.*

— *Cobbe et al*

# GRADE SCHOOL MATH 8K (GSM8K)

**Problem:** Beth bakes 4, 2 dozen batches of cookies in a week. If these cookies are shared amongst 16 people equally, how many cookies does each person consume?

**Solution:** Beth bakes 4 2 dozen batches of cookies for a total of  $4 \times 2 = \langle\langle 4 \times 2 = 8 \rangle\rangle$  8 dozen cookies

There are 12 cookies in a dozen and she makes 8 dozen cookies for a total of  $12 \times 8 = \langle\langle 12 \times 8 = 96 \rangle\rangle$  96 cookies

She splits the 96 cookies equally amongst 16 people so they each eat  $96/16 = \langle\langle 96/16 = 6 \rangle\rangle$  6 cookies

**Final Answer:** 6

**Problem:** Mrs. Lim milks her cows twice a day. Yesterday morning, she got 68 gallons of milk and in the evening, she got 82 gallons. This morning, she got 18 gallons fewer than she had yesterday morning. After selling some gallons of milk in the afternoon, Mrs. Lim has only 24 gallons left. How much was her revenue for the milk if each gallon costs \$3.50?

Mrs. Lim got 68 gallons - 18 gallons =  $\langle\langle 68 - 18 = 50 \rangle\rangle$  50 gallons this morning.

So she was able to get a total of 68 gallons + 82 gallons + 50 gallons =  $\langle\langle 68 + 82 + 50 = 200 \rangle\rangle$  200 gallons.

She was able to sell 200 gallons - 24 gallons =  $\langle\langle 200 - 24 = 176 \rangle\rangle$  176 gallons.

Thus, her total revenue for the milk is  $\$3.50/\text{gallon} \times 176 \text{ gallons} = \langle\langle 3.50 \times 176 = 616 \rangle\rangle$  616.

**Final Answer:** 616

**Problem:** Tina buys 3 12-packs of soda for a party. Including Tina, 6 people are at the party. Half of the people at the party have 3 sodas each, 2 of the people have 4, and 1 person has 5. How many sodas are left over when the party is over?

**Solution:** Tina buys 3 12-packs of soda, for  $3 \times 12 = \langle\langle 3 \times 12 = 36 \rangle\rangle$  36 sodas

6 people attend the party, so half of them is  $6/2 = \langle\langle 6/2 = 3 \rangle\rangle$  3 people

Each of those people drinks 3 sodas, so they drink  $3 \times 3 = \langle\langle 3 \times 3 = 9 \rangle\rangle$  9 sodas

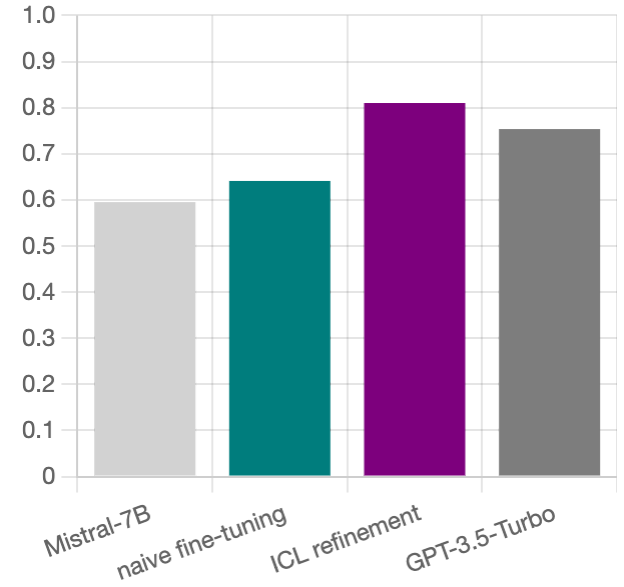
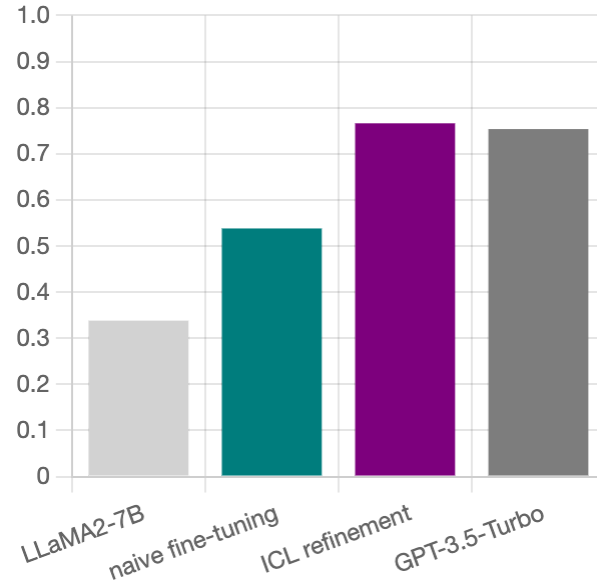
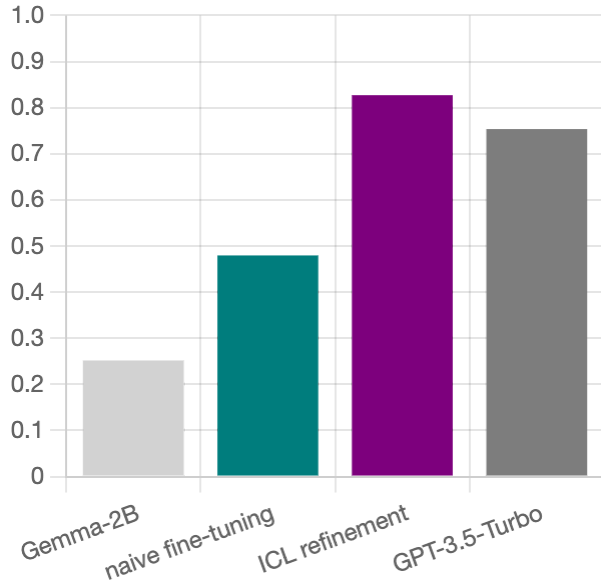
Two people drink 4 sodas, which means they drink  $2 \times 4 = \langle\langle 4 \times 2 = 8 \rangle\rangle$  8 sodas

With one person drinking 5, that brings the total drank to  $5 + 9 + 8 + 3 = \langle\langle 5 + 9 + 8 + 3 = 25 \rangle\rangle$  25 sodas

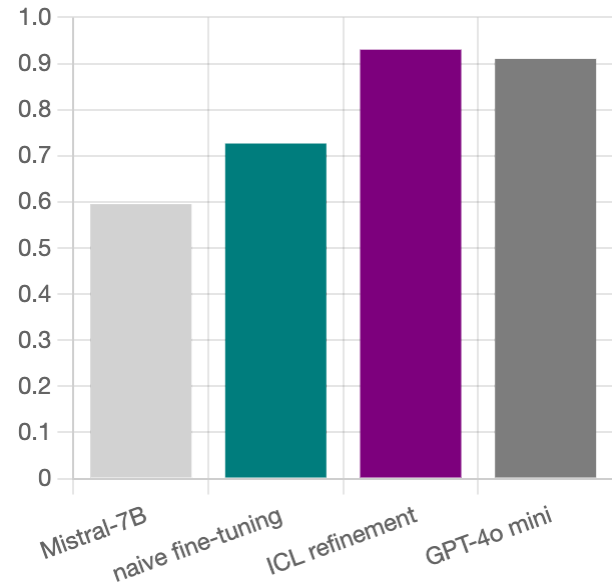
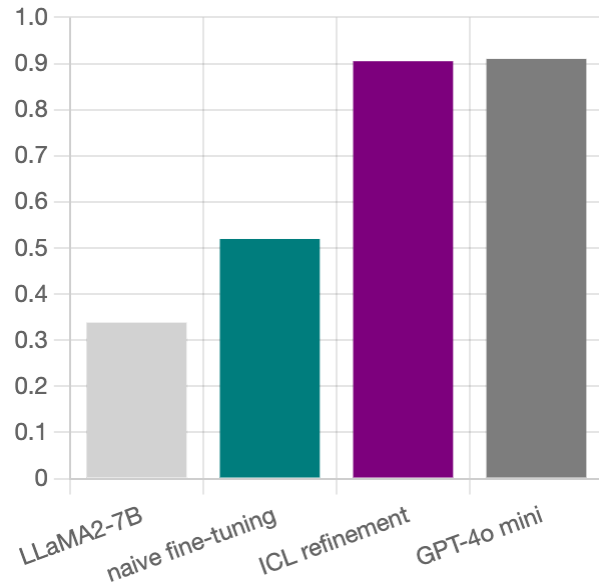
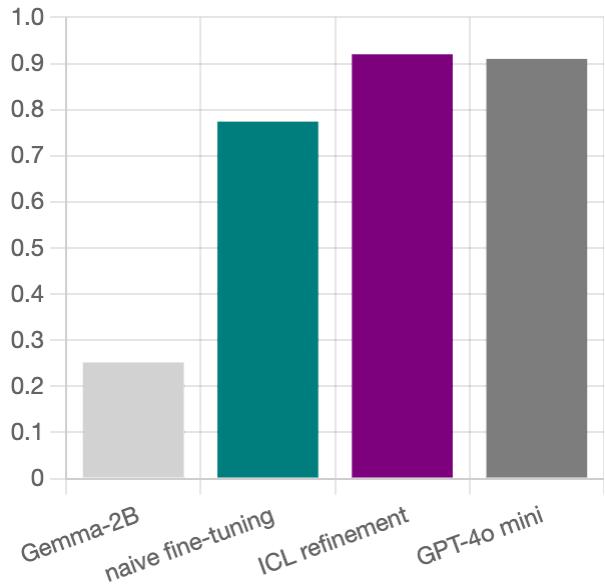
As Tina started off with 36 sodas, that means there are  $36 - 25 = \langle\langle 36 - 25 = 11 \rangle\rangle$  11 sodas left

**Final Answer:** 11

# GSM8K RESULTS (GPT-3.5 TURBO)



# GSM8K RESULTS (GPT-4O MINI)



# MATH DATASET

## MATH Dataset (Ours)

**Problem:** Tom has a red marble, a green marble, a blue marble, and three identical yellow marbles. How many different groups of two marbles can Tom choose?

**Solution:** There are two cases here: either Tom chooses two yellow marbles (1 result), or he chooses two marbles of different colors ( $\binom{4}{2} = 6$  results). The total number of distinct pairs of marbles Tom can choose is  $1 + 6 = \boxed{7}$ .

**Problem:** If  $\sum_{n=0}^{\infty} \cos^{2n} \theta = 5$ , what is  $\cos 2\theta$ ?

**Solution:** This geometric series is

$1 + \cos^2 \theta + \cos^4 \theta + \dots = \frac{1}{1 - \cos^2 \theta} = 5$ . Hence,

$$\cos^2 \theta = \frac{4}{5}. \text{ Then } \cos 2\theta = 2 \cos^2 \theta - 1 = \boxed{\frac{3}{5}}.$$

**Problem:** The equation  $x^2 + 2x = i$  has two complex solutions. Determine the product of their real parts.

**Solution:** Complete the square by adding 1 to each side.

Then  $(x + 1)^2 = 1 + i = e^{\frac{i\pi}{4}} \sqrt{2}$ , so  $x + 1 = \pm e^{\frac{i\pi}{8}} \sqrt[4]{2}$ .

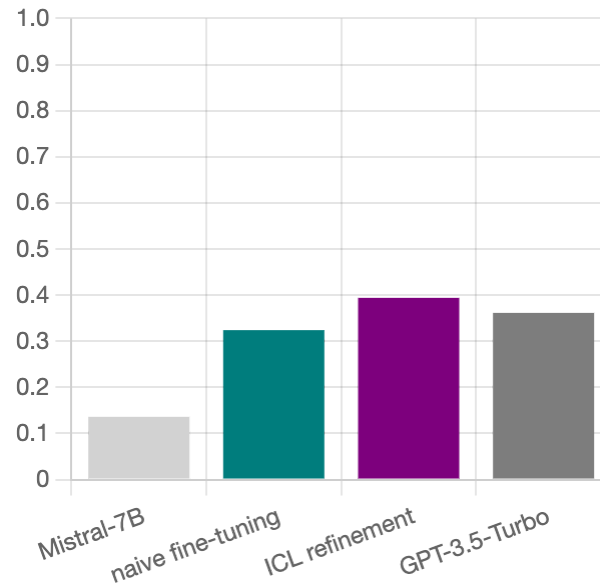
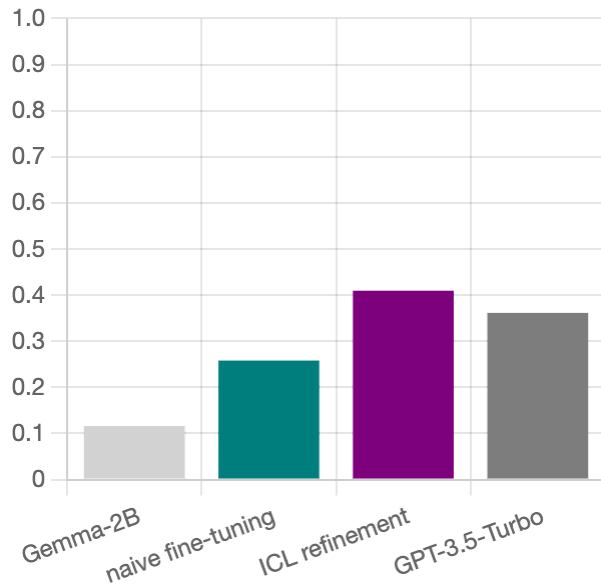
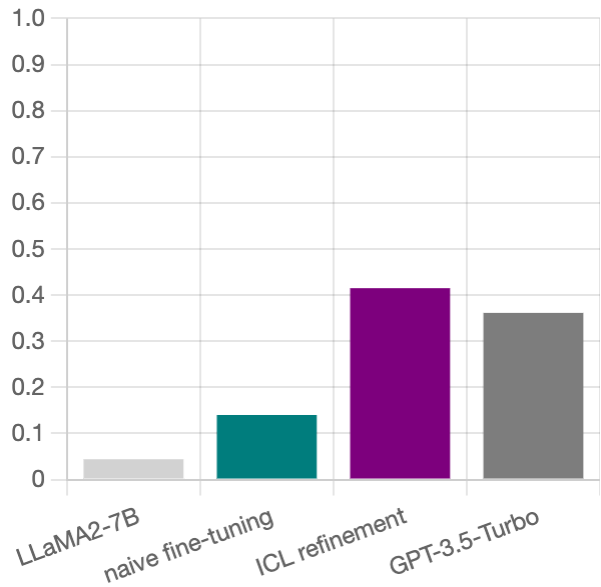
The desired product is then

$$(-1 + \cos(\frac{\pi}{8}) \sqrt[4]{2}) (-1 - \cos(\frac{\pi}{8}) \sqrt[4]{2}) =$$

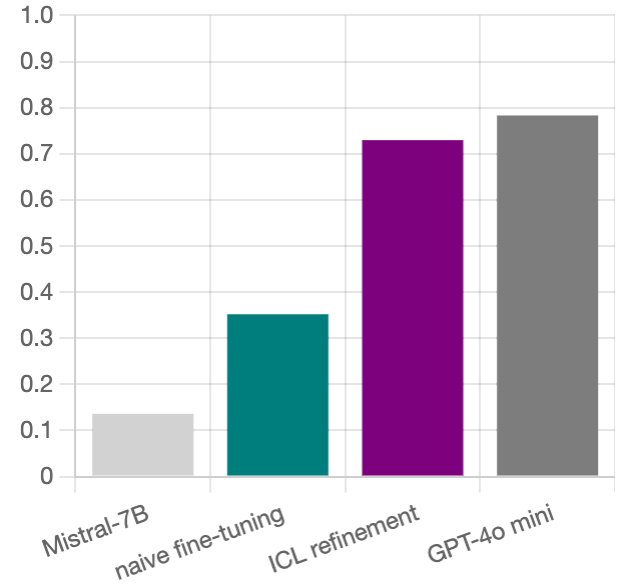
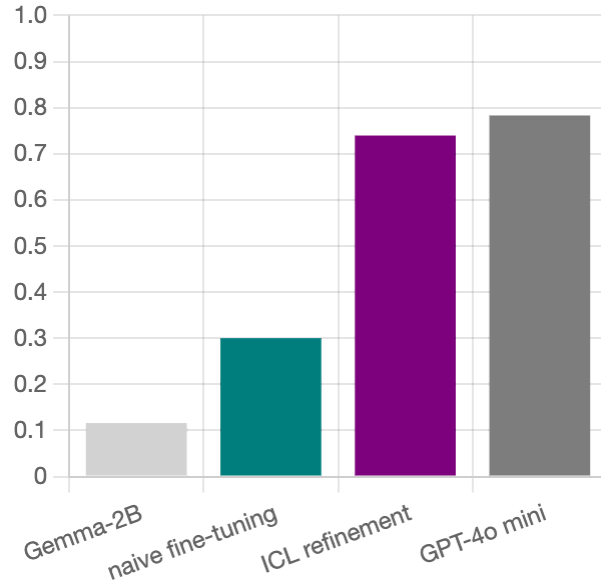
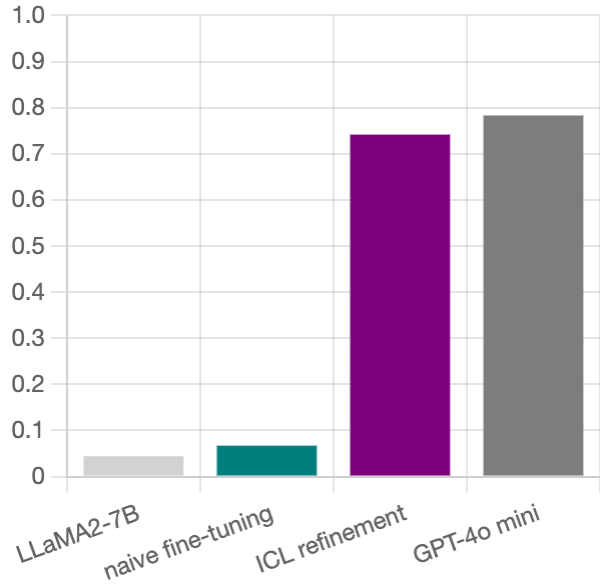
$$1 - \cos^2(\frac{\pi}{8}) \sqrt{2} = 1 - \frac{(1 + \cos(\frac{\pi}{4}))}{2} \sqrt{2} = \boxed{\frac{1 - \sqrt{2}}{2}}.$$

- 12.5k (7.5k/5k train/test) math competition problems from AMC 10, AMC 12, AMIE etc
- A CS PhD student who does not like math attained 40%.
- A 3x IMO gold medalist attained 90%.

# MATH RESULTS (GPT-3.5-TURBO)

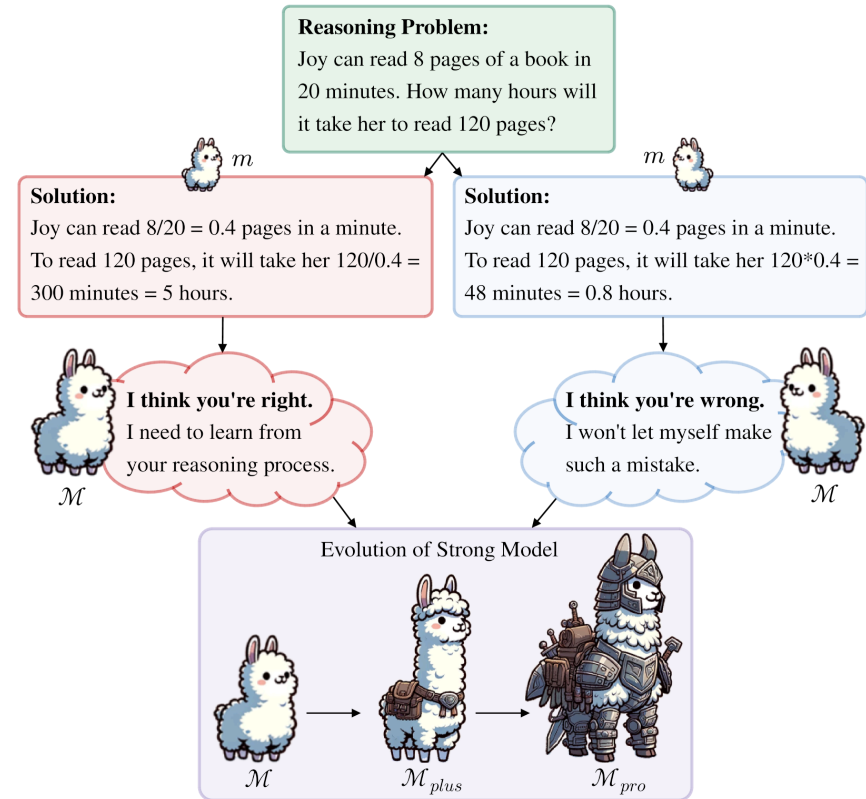


# MATH RESULTS (GPT-4O MINI)



# (MORE) ADVANCED REFINEMENT (AGENTS)

1. guide the pretrained LLM to the target concept with prompts
2. "constitution" + self-critique/self-refinement
3. use the weak teacher's mistakes to teach the pretrained LLM to avoid similar mistakes
4. combinations of the above



Yang et al et al (2023)

Q: Is refinement the correct high-level approach to superalignment?



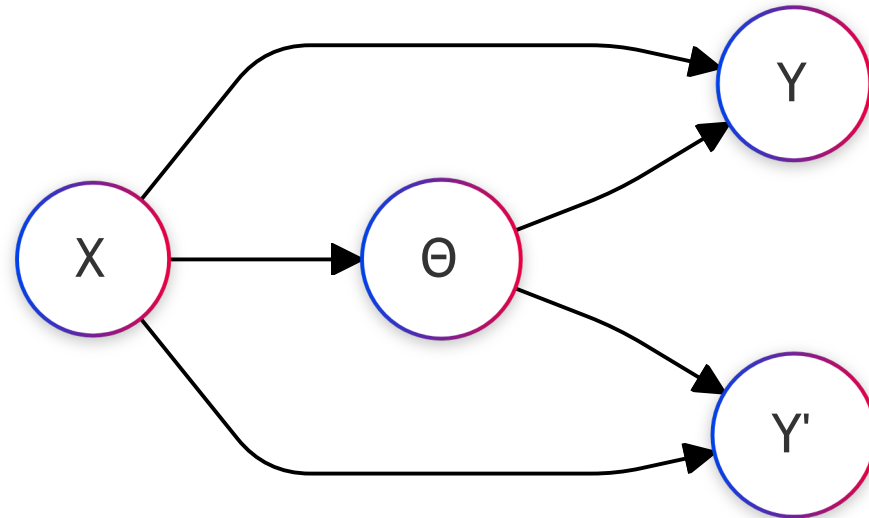
# LATENT CONCEPT SHIFT 2.0

Assume

1.  $p(x, \theta, y, y')$  (resp.  $q(x, \theta, y, y')$ ) factorizes as

$$p(x, \theta, y, y') = p(x)p(\theta | x)p(y | x, \theta)p(y' | x, \theta)$$

(resp.  $q(x, \theta, y, y') = q(x)q(\theta | x)q(y | x, \theta)q(y' | x, \theta)$ ).



# LATENT CONCEPT SHIFT 2.0

Assume

1.  $p(x, \theta, y, y')$  (resp.  $q(x, \theta, y, y')$ ) factorizes as

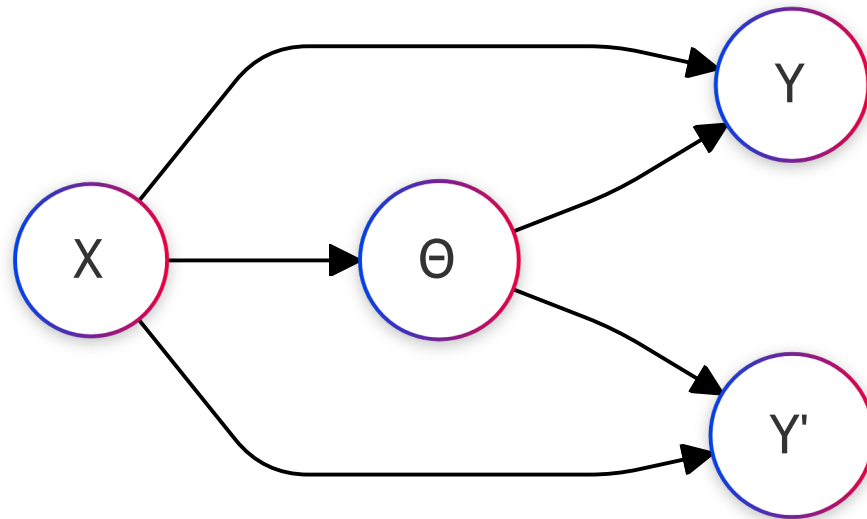
$$p(x, \theta, y, y') = p(x)p(\theta | x)p(y | x, \theta)p(y' | x, \theta)$$

(resp.  $q(x, \theta, y, y') = q(x)q(\theta | x)q(y | x, \theta)q(y' | x, \theta)$ ).

2.  $p(x) = q(x)$  and  $p(y | x, \theta) = q(y | x, \theta)$ ; ie the only differences between  $P$  and  $Q$  are

$$p(\theta | x) \neq q(\theta | x) = \delta_{\theta_Q},$$
$$p(y' | x, \theta) \neq q(y' | x, \theta).$$

# LATENT CONCEPT SHIFT 2.0



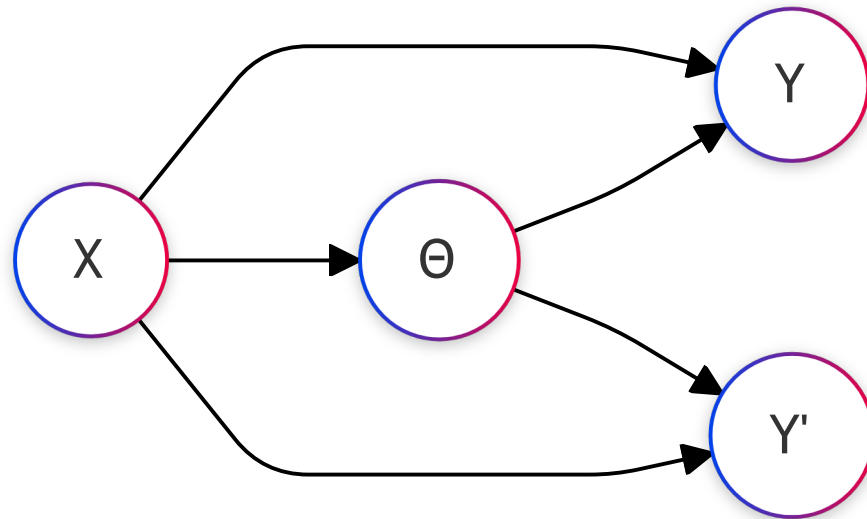
A1 and A2 imply  $p(y | x), q(y | x)$  are mixtures of  $p(y | x, \theta), \theta \in \Theta$ :

$$p(y | x) = \sum_{\theta \in \Theta} p(y | x, \theta) p(\theta | x),$$

$$q(y | x) = \sum_{\theta \in \Theta} q(y | x, \theta) q(\theta | x) = q(y | x, \theta_Q).$$

This is exactly the latent concept shift (1.0) assumption!

# LATENT CONCEPT SHIFT 2.0



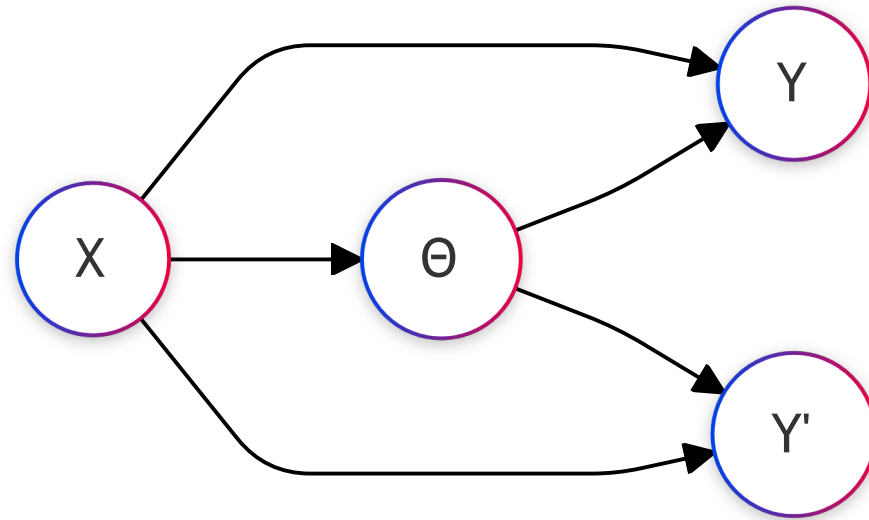
A1 and A2 imply  $p(y' | x)$  (resp  $q(y' | x)$ ) are mixtures of  $p(y' | x, \theta)$  (resp  $q(y' | x, \theta)$ ),  $\theta \in \Theta$ :

$$p(y' | x) = \sum_{\theta \in \Theta} p(y' | x, \theta) p(\theta | x),$$

$$q(y' | x) = \sum_{\theta \in \Theta} q(y' | x, \theta) q(\theta | x) = q(y' | x, \theta_Q).$$

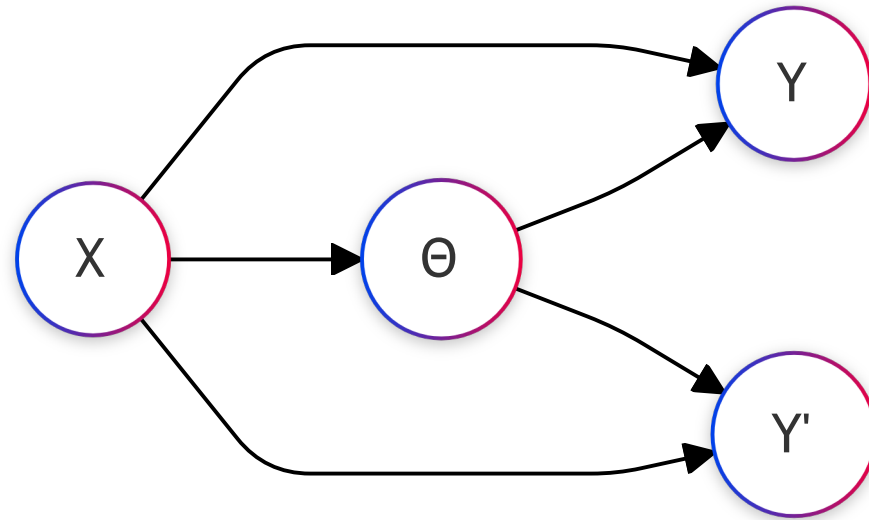
- $p(y' | x), q(y' | x)$  are not mixtures of the same components

# LCS 2.0: THE PRETRAINED MODEL



- $p(y | x, \theta) = q(y | x, \theta)$ ,  $\theta \in \Theta$  encode the pretrained LLM's latent abilities
- $p(\theta | x)$  (resp  $q(\theta | x)$ ),  $\theta \in \Theta$  encode the pretrained LLM's (resp the weak teacher's) prior (distribution) on concepts

# LCS 2.0: WEAK SUPERVISION



- $q(y' | x, \theta) \neq p(y' | x, \theta), \theta \in \Theta$  encode the (weak) teacher's latent abilities
- $p(y' | x, \theta), \theta \in \Theta$  are defined implicitly thru  $p(y | x, y')$ ; we interpret them as the pretrained LLM's model of the teacher

# REFINEMENT IN LATENT CONCEPT SHIFT 2.0

The analogy of ICL refinement in LCS 2.0 is (sampling from)

$$\hat{q}(y | x) \triangleq \int_{y'} p(y | x, y') q(y' | x).$$

ICL refinement: for  $i \in [n_Q]$

1. select ICL examples  $\mathcal{S}_i | X_{Q,i}$ ,
2.  $\mathbf{Y}_{Q,i} \sim \text{GPT}(\cdot | \mathbf{S}_i, \mathbf{X}_{Q,i})$ .

refinement analogy:

1.  $Y | X \sim q(y' | X)$ ,
2.  $Y | X, Y' \sim p(y | X, Y')$ .

# REFINEMENT IN LATENT CONCEPT SHIFT 2.0

The analogy of ICL refinement in LCS 2.0 is (sampling from)

$$\hat{q}(y | x) \triangleq \int_{y'} p(y | x, y') q(y' | x).$$

**Lem:** Under A1, A2, and  $p(y' | x, \theta) = q(y' | x, \theta)$  (ie the pretrained LLM can correct the teacher's mistakes),

$$\hat{q}(y | x) = \sum_{\theta \in \Theta} p(y | x, \theta) \int_{y'} p(\theta | x, y') q(y' | x, \theta_Q).$$

cf the source predictive distribution, which has the form

$$p(y | x) = \sum_{\theta \in \Theta} p(y | x, \theta) \int_{y'} p(\theta | x, y') p(y' | x).$$

**Sanity check:** Refinement always helps when the pretrained LLM can correct the teacher's mistakes!



# IS REFINEMENT OPTIMAL?

Lem: Under A1 and A2,

$$\hat{q}(y | x) = \sum_{\theta \in \Theta} p(y | x, \theta) \int_{\mathcal{Y}'} p(\theta | x, y') q(y' | x, \theta_Q).$$

- interpret  $\int_{\mathcal{Y}'} p(\theta | x, y') q(y' | x, \theta_Q)$  as entries of a confusion matrix (for predicting  $\theta$  from  $y'$ )

cf the target predictive distribution, which has the form

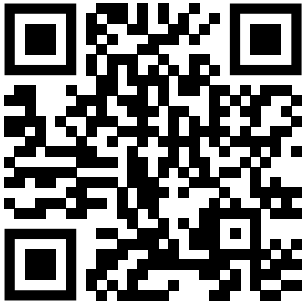
$$q(y | x) = p(y | x, \theta_Q).$$

✓ Yes when it is possible to exactly recover  $\theta_Q$  from the weak teacher (eg when there is no overlap between the  $q(y' | x, \theta)$ 's).

✗ Refinement is often suboptimal because identifiability of  $q(y | x)$  (from  $p(x, y, y')$  and  $q(x, y')$ ) is (much) weaker than no overlap.

# TAKEAWAYS

1. The accuracy achievable by naive fine-tuning is limited by the accuracy of the weak teacher and the pretrained LLM.
2. **Main idea:** elicit the latent abilities of the pretrained LLM by using it to refine the weak teacher's outputs
3. There seems to be room for improvement!



[Transfer learning for weak-to-strong generalization.](#)

S Somerstep, F Maia Polo, M Banerjee, Y Ritov, M Yurochkin, Y Sun.

arXiv:2405.16236

This work was supported by the NSF under grants 2113364, 2113373, 2414918.