



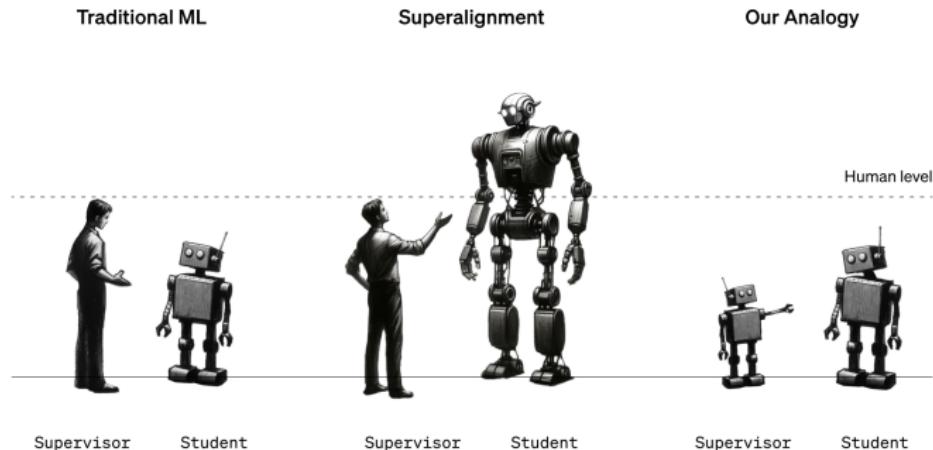
Weak-to-Strong Generalization: Eliciting Strong Capabilities With Weak Supervision

ICLR 2024 (oral)
OpenAI Superalignment Generalization Team

Overview

Superhuman model (e.g. GPT-4, Gemini, Claude)

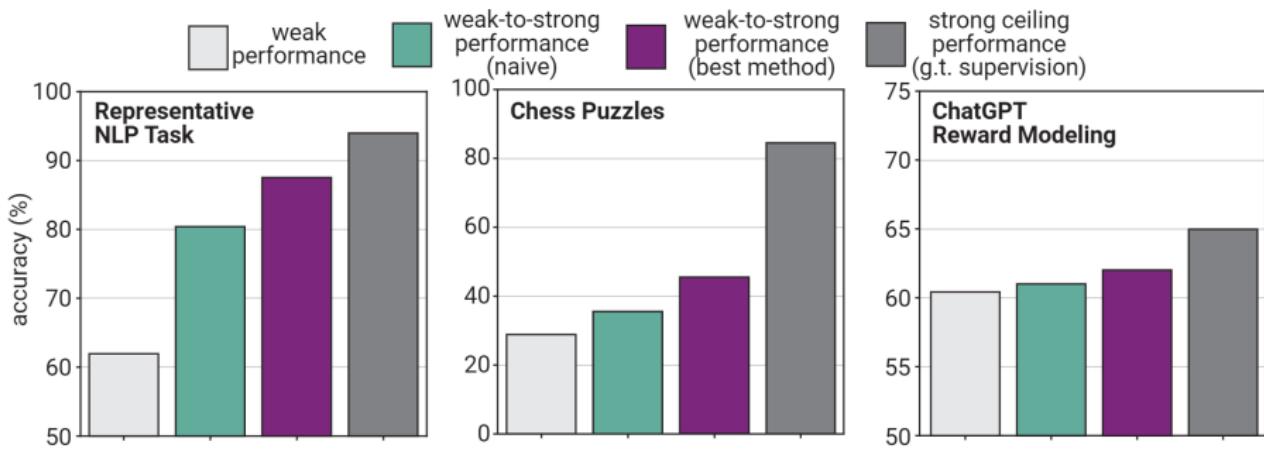
- Recent aligning method: Reinforcement Learning from Human Feedback (RLHF)
- Can humans evaluate responses from superhuman models? (e.g. millions of code)
- **GOAL:** Human supervision → Weaker model supervision



Can we use weak models to supervise strong models?

Main Contributions

1. Strong pretrained models outperform their weak supervisors. (light-gray vs. green)
2. Naively finetuning by weak supervisors is not enough. (green vs. dark-gray)
3. Improving weak-to-strong generalization is tractable. (green vs. purple)



Three Types of Tasks

1. 22 Natural Language Processing tasks
2. Chess puzzles: Find the next best move



Prompt: "1. d4 l.., c5 2. Nf3 2.., Nf6 3.., Nf3 3.., Nc6 4.., Bb4 4.., c5 5.., c3 5.., Bc6 6.., Nf6 6.., O-O 7.., O-O 7.., Ne5 8.., Re8 8.., Bd7 9.., e4 9.., dxe4 10.., Nxe4 10.., cxd4 11.., Nf6+ 11.., Bxf6 12.., cxd4 12.., Nf6 13.., Be3 13.., Qh6 14.., a3 14.., Ne5 15.., d5 15.., exd5 16.., Bxf5 16.., Bf5 17.., Bxc6 17.., Qxc6 18.., Nd4 18.., Re6 19.., Qd8 19.., Rf6 19.., Rxe6+ 20.., Rxe6 21.., Be3 21.., Bb7 22.., Rf1 22.."

Label: "Qxc6"



Prompt: "1.e4 l.., c5 2. Nf3 2.., Nf6 3.., Nf3 3.., Nc6 4.., Bb4 4.., Bc6 5.., Bb6 5.., a6 6.., Bg1 7.., Bg1 7.., Bg1 8.., Qg7 8.., O-O 9.., g4 9.., Bh4 10.., Ba2 10.., Nf7 11.., b4 11.., Be7 12.., g5 12.., Ng5 13.., O-O-O 13.., Qh7 14.., h5 14.., Qh8 15.., Qg3 15.., Ne6 16.., Rg1 16.., b5 17.., Qe5 17.., g5 18.., f4 18.., Re8 19.., Qf5 19.., b4 20.., Nd4 21.., Qg1 21.., c5 22.., f5 22.., Ra2 23.., b6 23.., Bb6 24.., fxg7 24.., Kxg7 25.., Rg2 25.., Qe6 26.., Rf6 26.., Rg8 27.., Qf5 27.., Qd7 28.., Rh1 28.., Rf6 29.., Rgf7 29.., Rg8 30.., c3 30.., Re6 31.., Rg6 32.., Ne5 32.., Qd5 32.., Afe 33.., Qd8 34.., Qd3 34.., Rb1 35.., Be3 35.., Nf6 36.., Nf6+ 36.., Qd6 37.., Rg1 37.., Qe4 38.., Kh1 38.., Qa4 39.., Ka1 39.., Be5 40.., Nf6 40.., Qa4 41.., Nd5 41.., Rb7 42.."

Label: "Qf5"

3. ChatGPT Reward Modeling (d, c_1, c_2, y)

- d : dialog between a user and an assistant
- c_1 and c_2 : pair of completions for d
- y : label, 1 if labeler preferred c_2 , 0 otherwise

Methodology and Measures

- Strong student model: A family of GPT-4
- Weak teacher model: GPT-2-level model for NLP and chess, GPT-3.5-level model for RM

Weak-to-Strong model

1. **Split the dataset:** one for training a weak model and the other for generating *weak labels*.
2. **Train the weak model** to generate *weak labels*.
3. **Train a strong model** using *weak labels*.

Strong-ceiling model

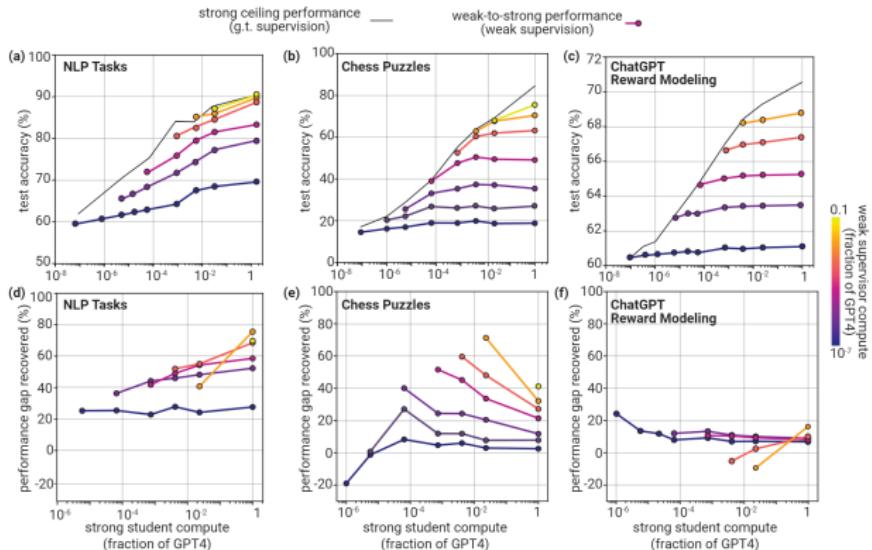
1. Train strong model using *ground truth labels*.

Performance gap recovered (PGR)

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

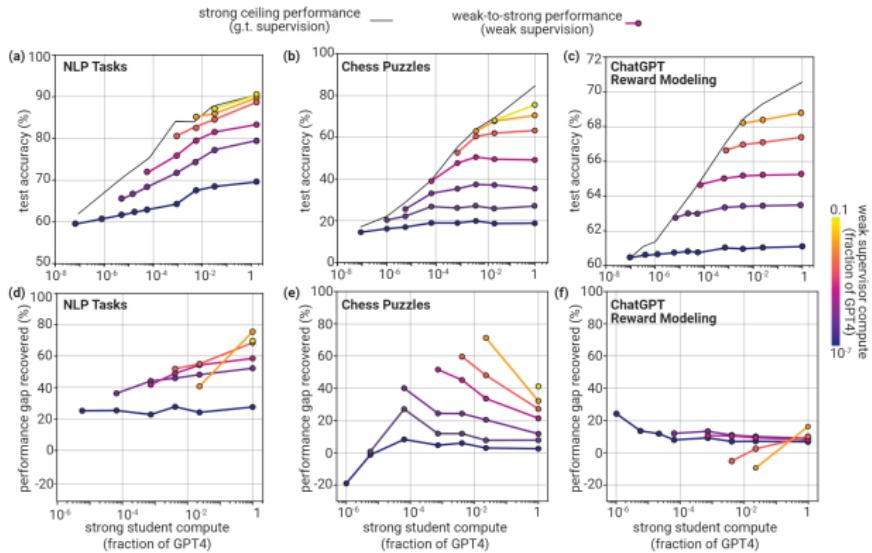


Naive Finetuning Results



- x-axis: strong student model size, colors: weak teacher model size
- black line: strong-ceiling performance, color line: weak-to-strong performance
- (a), (b), (c): test accuracy, (d), (e), (f): PGR

Naive Finetuning Results

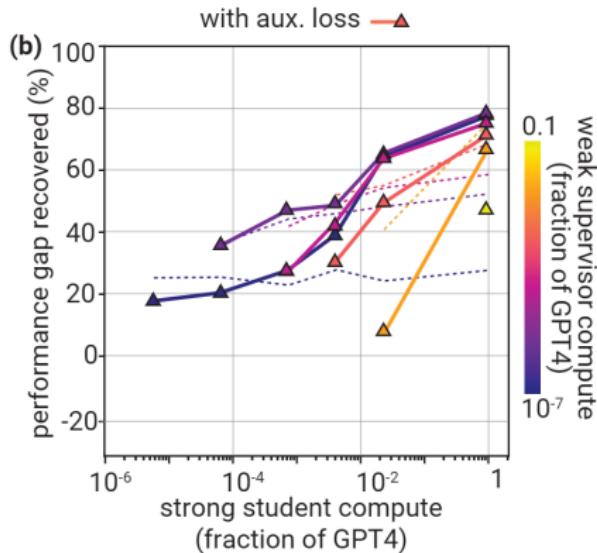
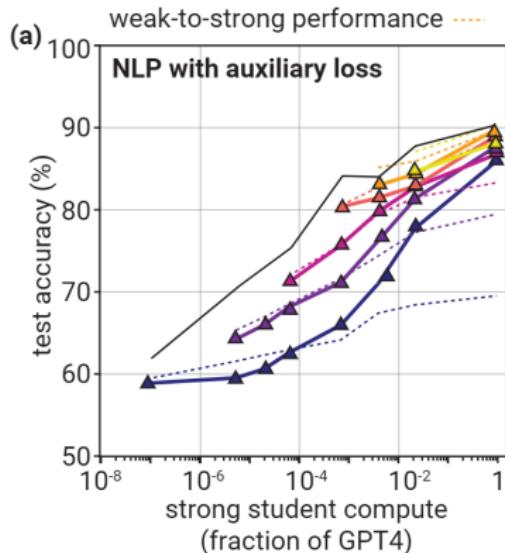


1. **PGRs are universally positive:** strong student model > weak teacher model
2. **NLP Tasks:** PGR > 20%, often above 50%.
3. **Chess Puzzles:** PGR decreases with the student size.
4. **ChatGPT Reward Modeling:** PGR never exceeds 20%.

Improving Weak-to-Strong Generalization

Confidence Loss

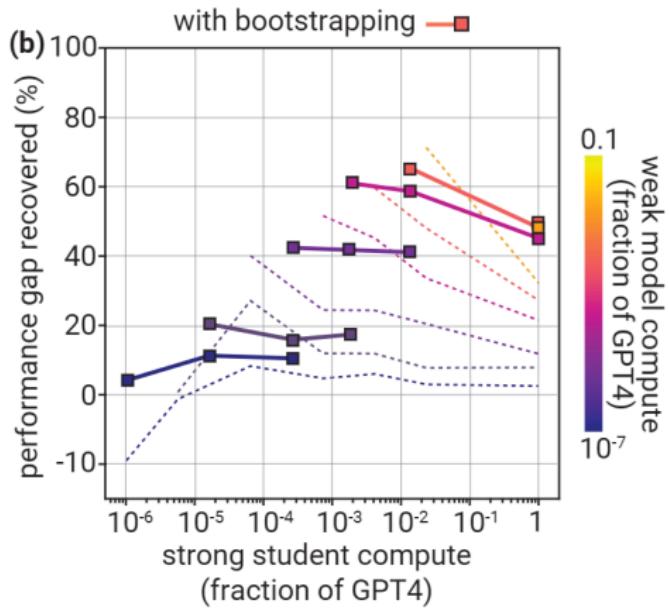
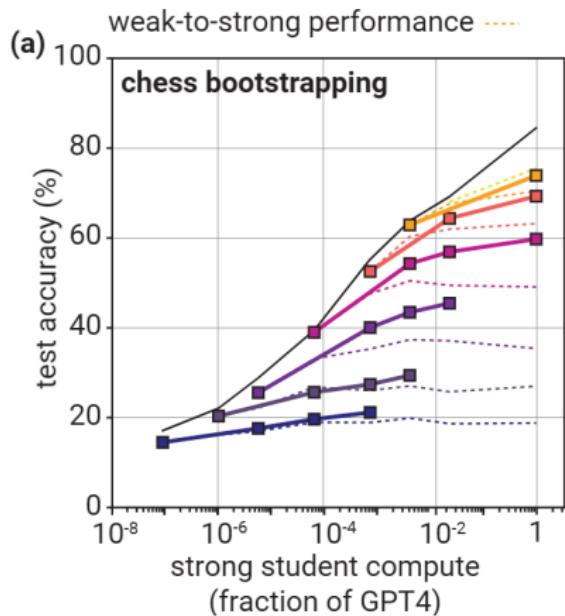
- Idea: reinforces the strong model's confidence in its own predictions, even they disagree with weak labels. cf.) conditional entropy minimization
- Dramatically improves generalization gaps in **NLP tasks**.



Improving Weak-to-Strong Generalization

Bootstrapping

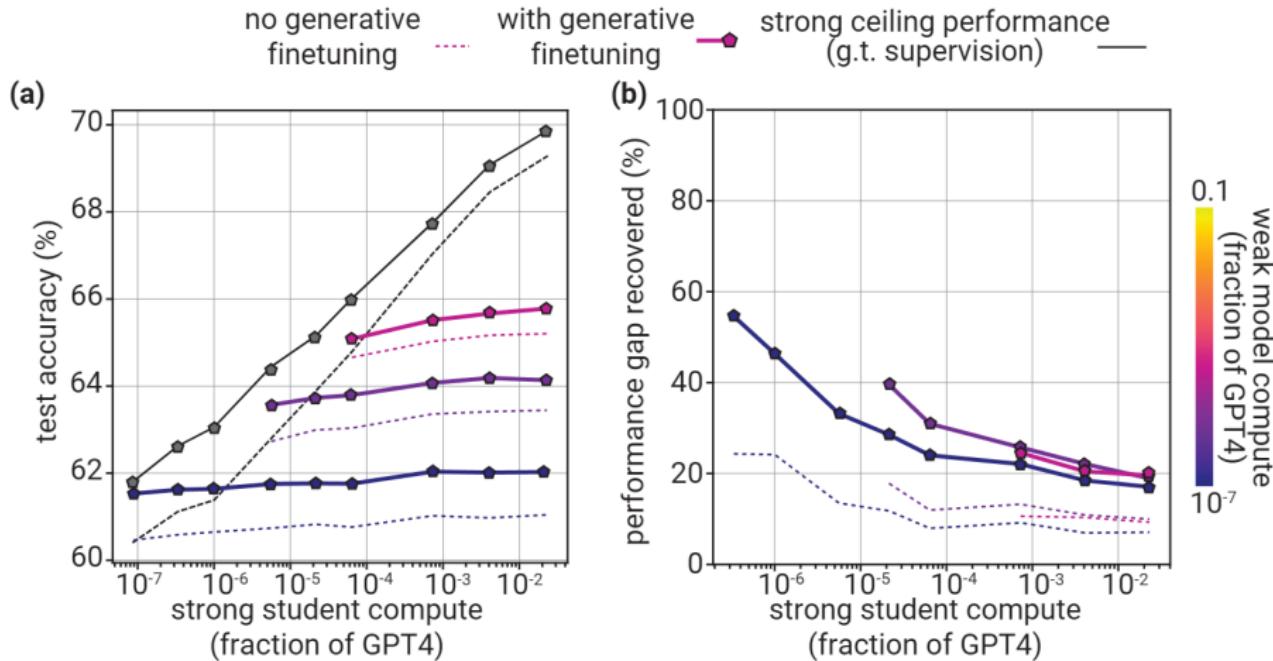
- **Idea:** Training models with a sequence of increasing model sizes, $\mathcal{M}_1 \rightarrow \mathcal{M}_2 \dots \rightarrow \mathcal{M}_n$
- Test accuracy improves for larger student models in **Chess** task.



Improving Weak-to-Strong Generalization

Unsupervised Finetuning

- **Idea:** Unsupervised finetuning to increase salience of a task without ground truth labels
- PGR improves 10-20% in **ChatGPT RM** task.



Additional Analysis: Imitation

Q. What extent does the student imitate the weak supervision?

1. Overfitting to weak model errors

- The student model overfits to weak model errors in early stage. (less than 1 epoch)
- Oracle early stopping improves performance significantly.
- Early stopping methods without relying on an oracle could enhance performance.

2. Student-supervisor agreement

- Student-supervisor agreement is higher than weak supervisor accuracy.
- Student imitates some of the supervisor's errors.
- Confidence loss reduces this agreement.

3. Size of student model

- Larger student models agree less with the supervisor than smaller models.
- Contradictorily, larger models should fit teacher errors more easily.

Additional Analysis: Salience

Intuition: A capability should be easier to elicit if a model is salient.

1. (Few-shot) Prompting vs. Finetuning

- Few-shot prompting becomes competitive with finetuning at large model sizes.
- Weak-to-strong finetuning with confidence loss outperforms few-shot prompting.
- Weak-to-strong fintuning significantly better than prompting at small model sizes.

2. Linear Representation of Tasks

- Finetuning on weak labels makes the ground truth task more linearly represented.
- Ground truth linear probe of the base model: 72%
- Finetuning on ground truth: 82%
- Ground truth linear probe of a model finetuned on weak labels: 78%.

Remaining Disanalogies

1. Imitation saliency: superhuman models may easily imitate weak errors.

- If we train models with human supervision, it might be having human-level capabilities.
- Larger strong models have more difficulty to imitate weak model.
- **Imitationg the weak supervisor** may not be as much of a problem.

2. Pretraining leakage: superhuman knoledge may be latent, not observable.

- Many of tasks may have been observed in pretraining stage. (contamination)
- Future superhuman model may never directly observe alignment-relevant capabilities.