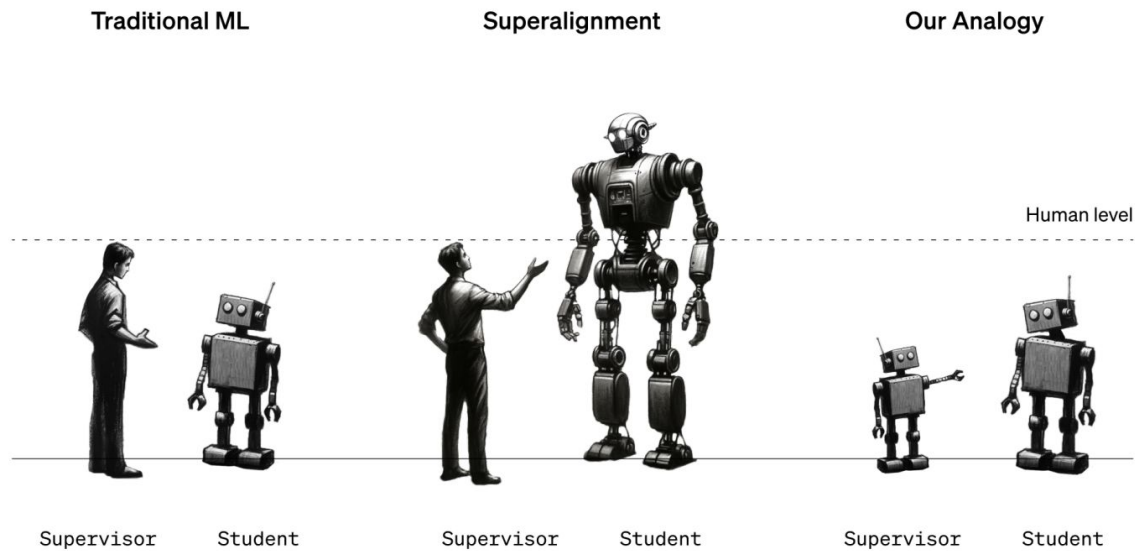


Weak-to-Strong Generalization: Eliciting Strong Capabilities With Weak Supervision

Введение



Что заметили

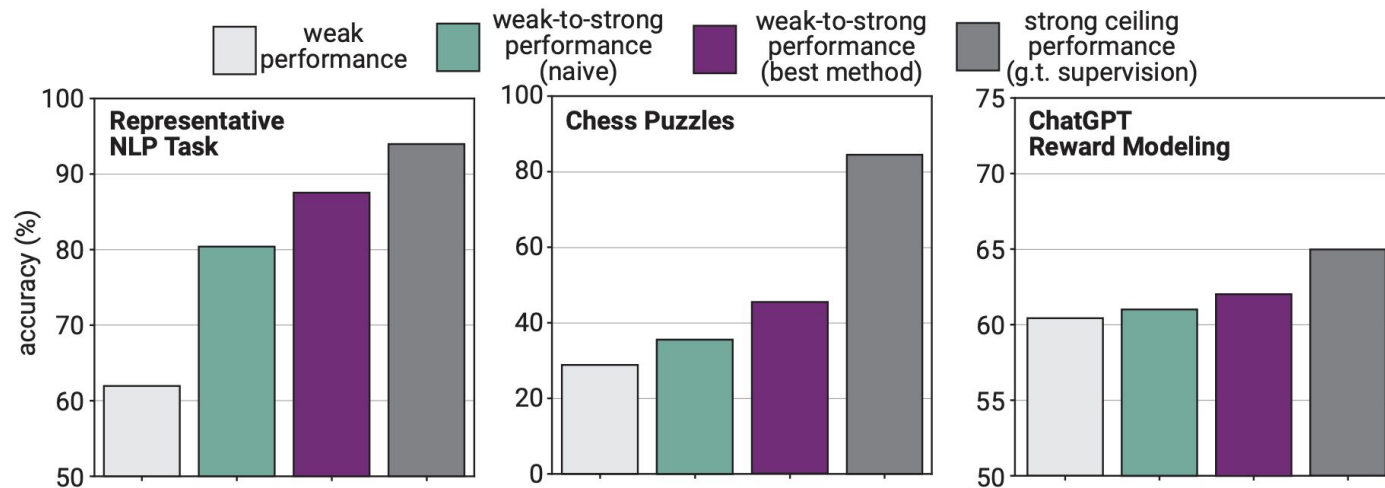
1. Сильные предварительно обученные модели, делают обобщения за пределами своих слабых наблюдателей
2. С ходу не работает
3. Методы улучшения позволяют значительно увеличивать качество

Идея

Сделаем три подхода для сравнения

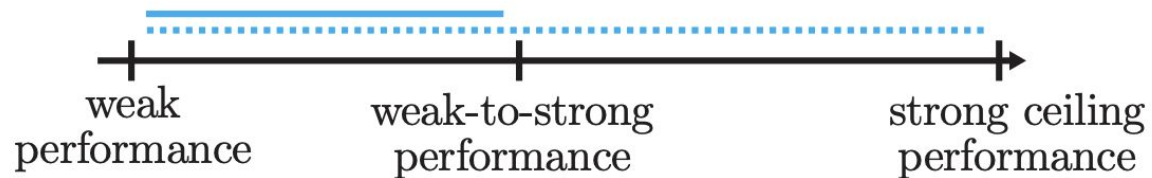
1. Создаем слабых наблюдателей на настоящих данных
2. Дотюниваем большие предобученные модели на предсказаниях слабых наблюдателей
3. Дотюниваем большие предобученные модели на настоящих данных

Что заметили



Что получается

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



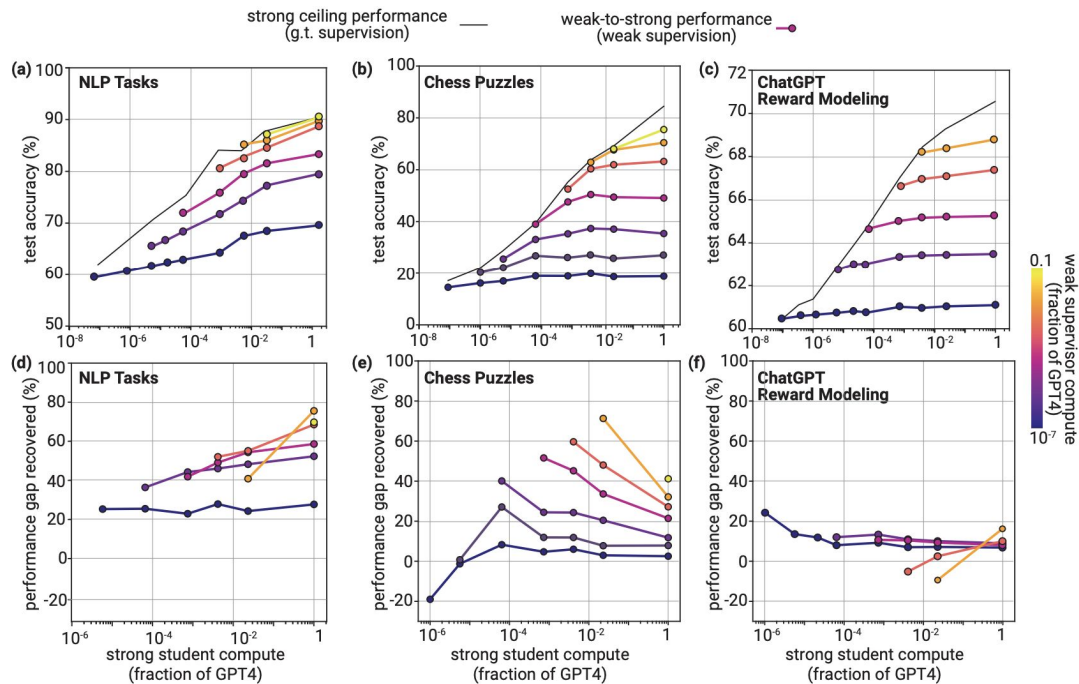
Плюсы такого подхода

1. Можно использовать любые модели
2. Можно использовать изучать любые сферы
3. В случае успеха, метод можно использовать прямо сейчас

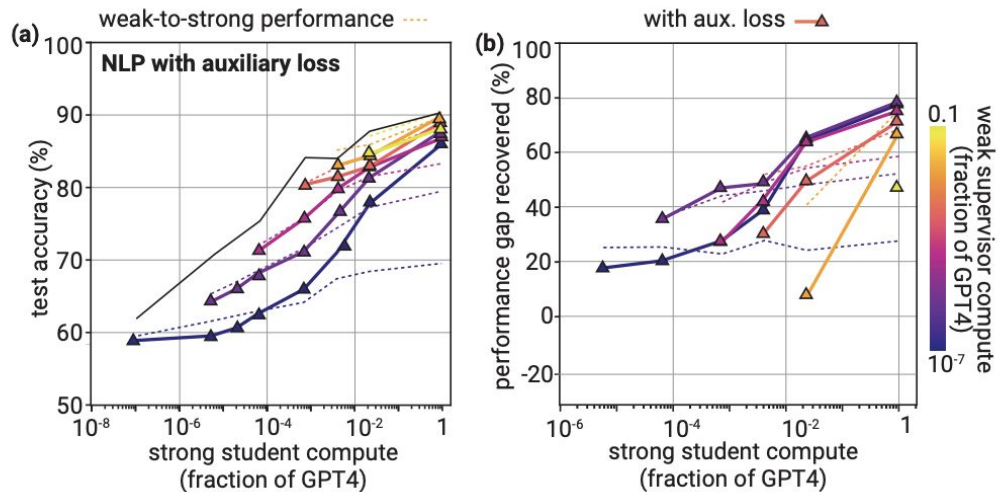
Минусы такого подхода

1. Слабые модели могут делать ошибки, отличающиеся от ошибок людей
2. В тренировочных данных все равно есть данные, которые были размечены людьми

Самый просто метод

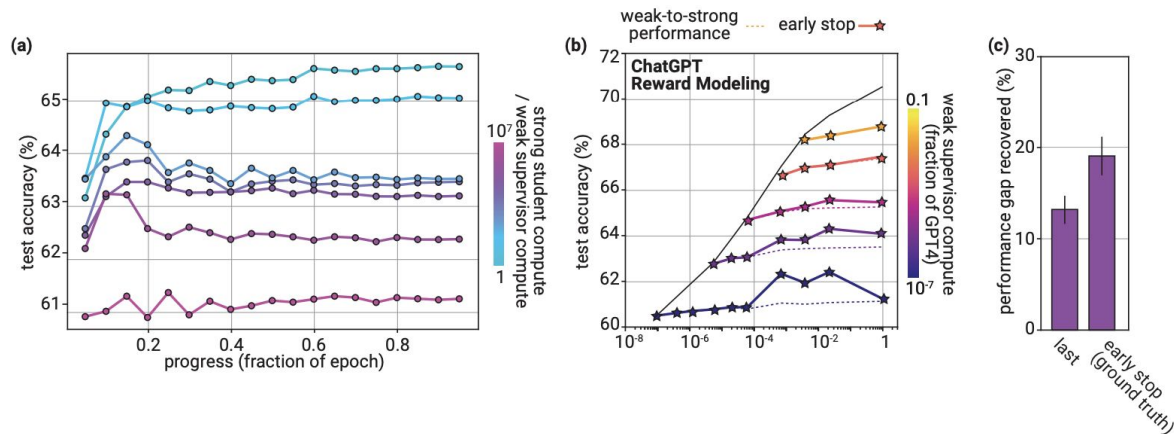


CONFIDENCE LOSS

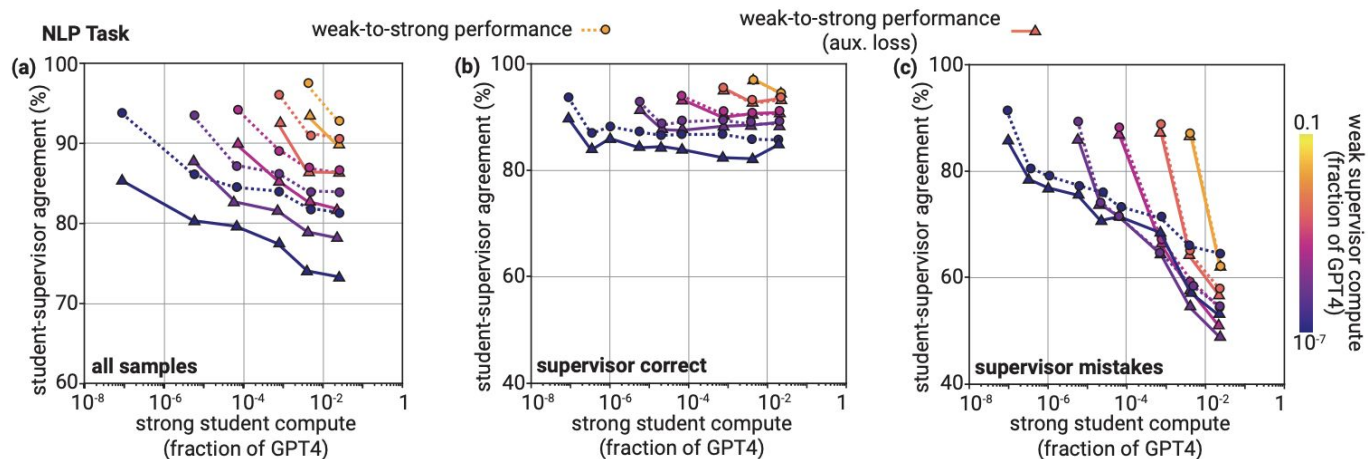


Почему это работает

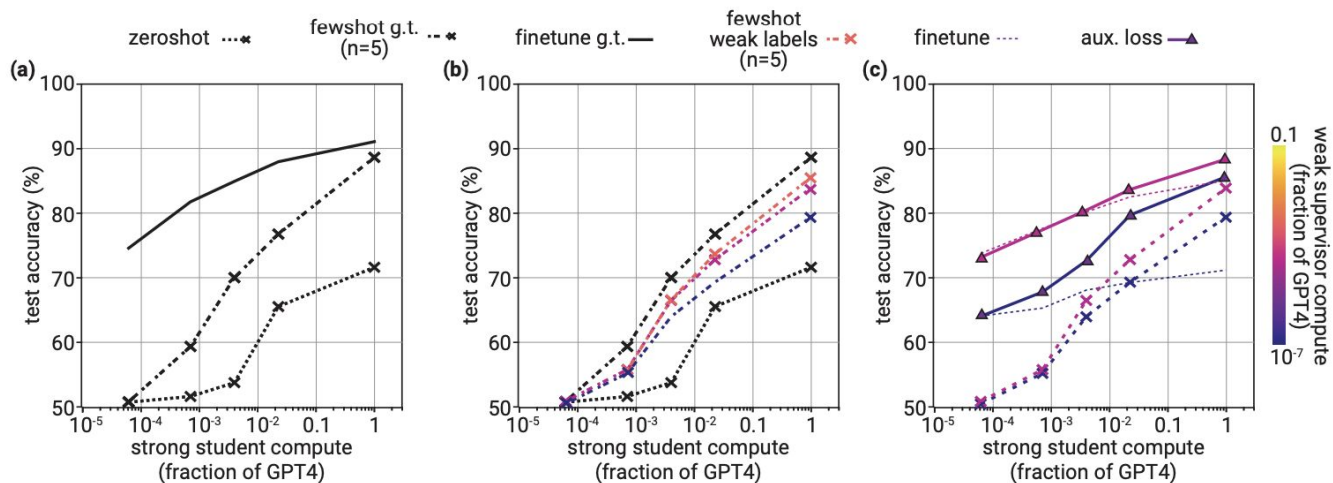
- Способность имитировать модель
- Переобучение



Разница между уверенностью модели и наблюдателя



Промтинг



Обучение на определенную тему

