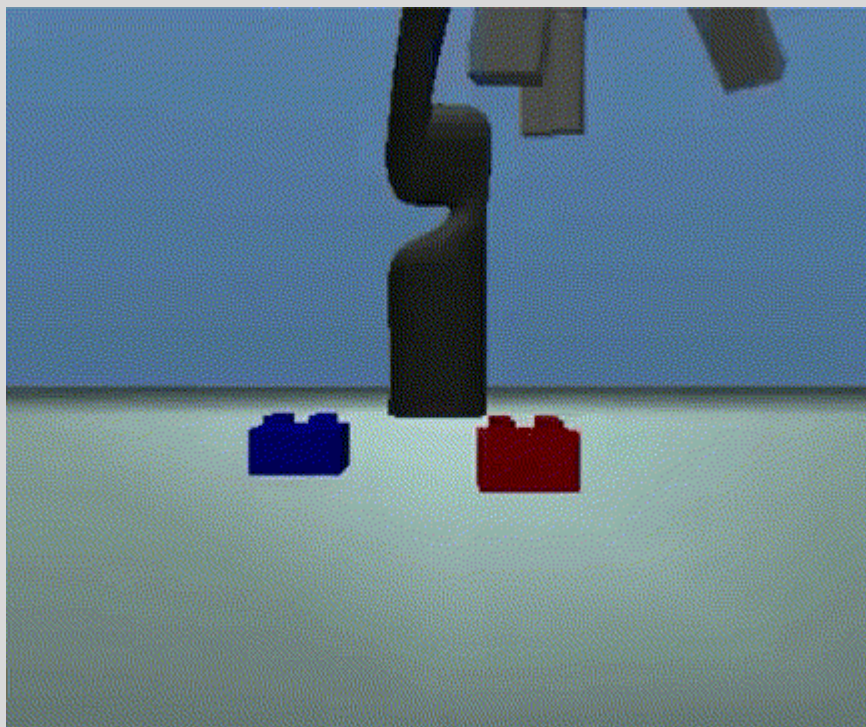


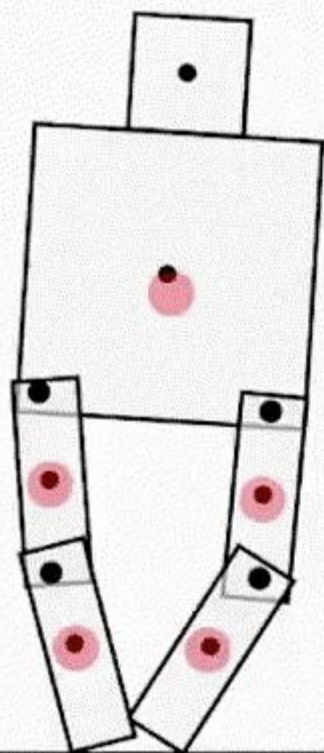
RL and Problems

Specification gaming

It is sometimes very hard to specify a score function when you want to train a model to do some complex task.



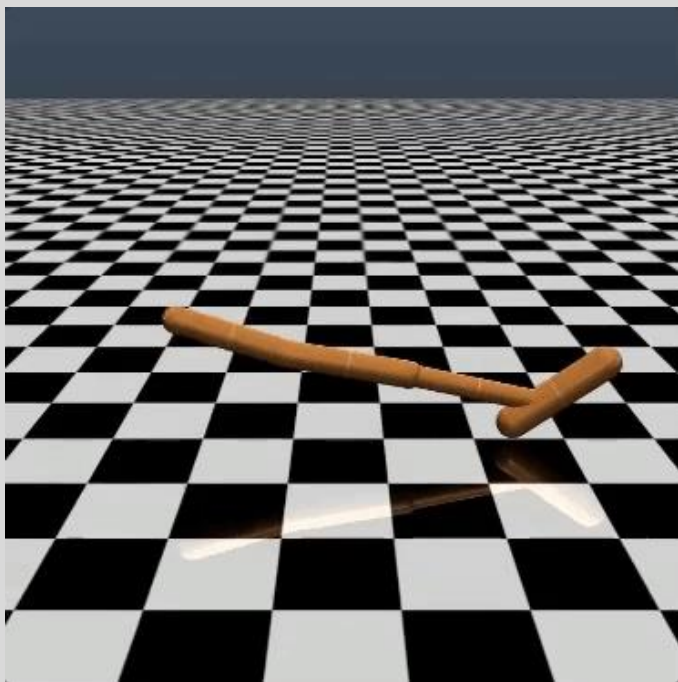




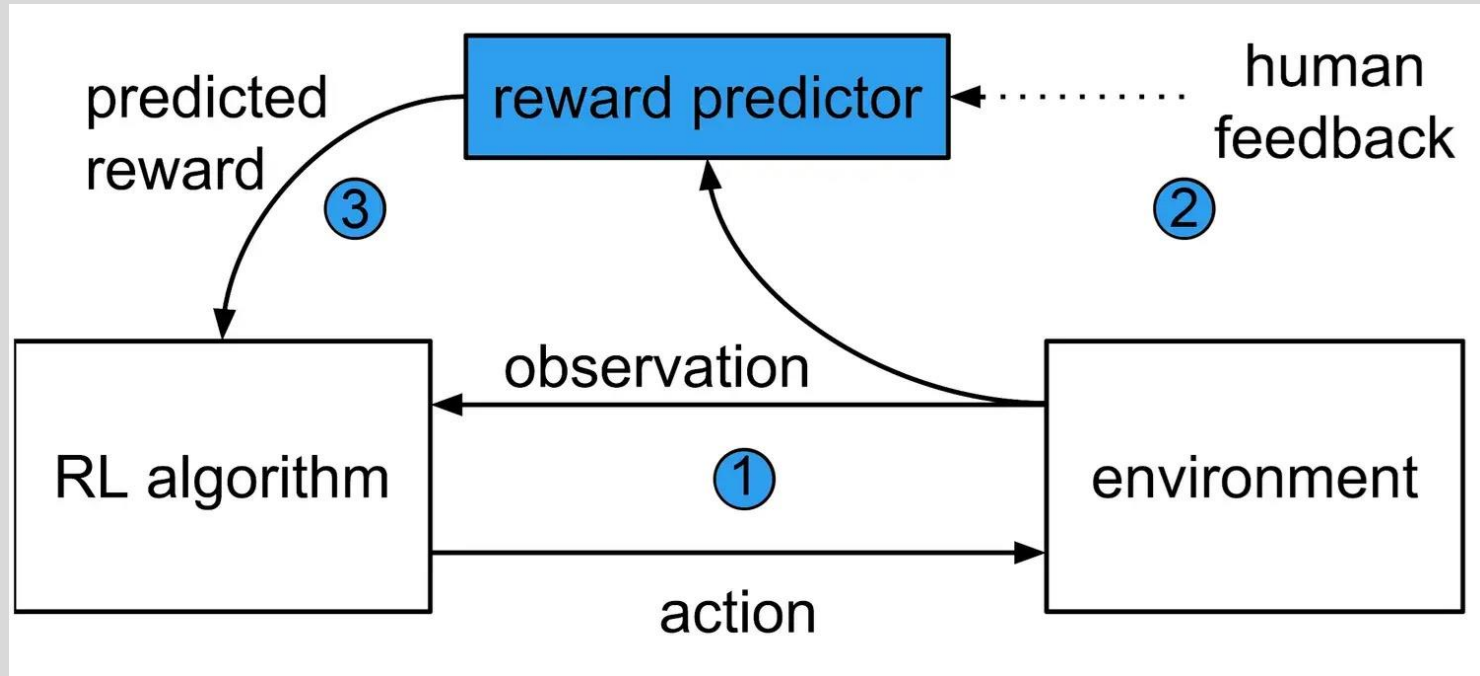
Goodhart's law

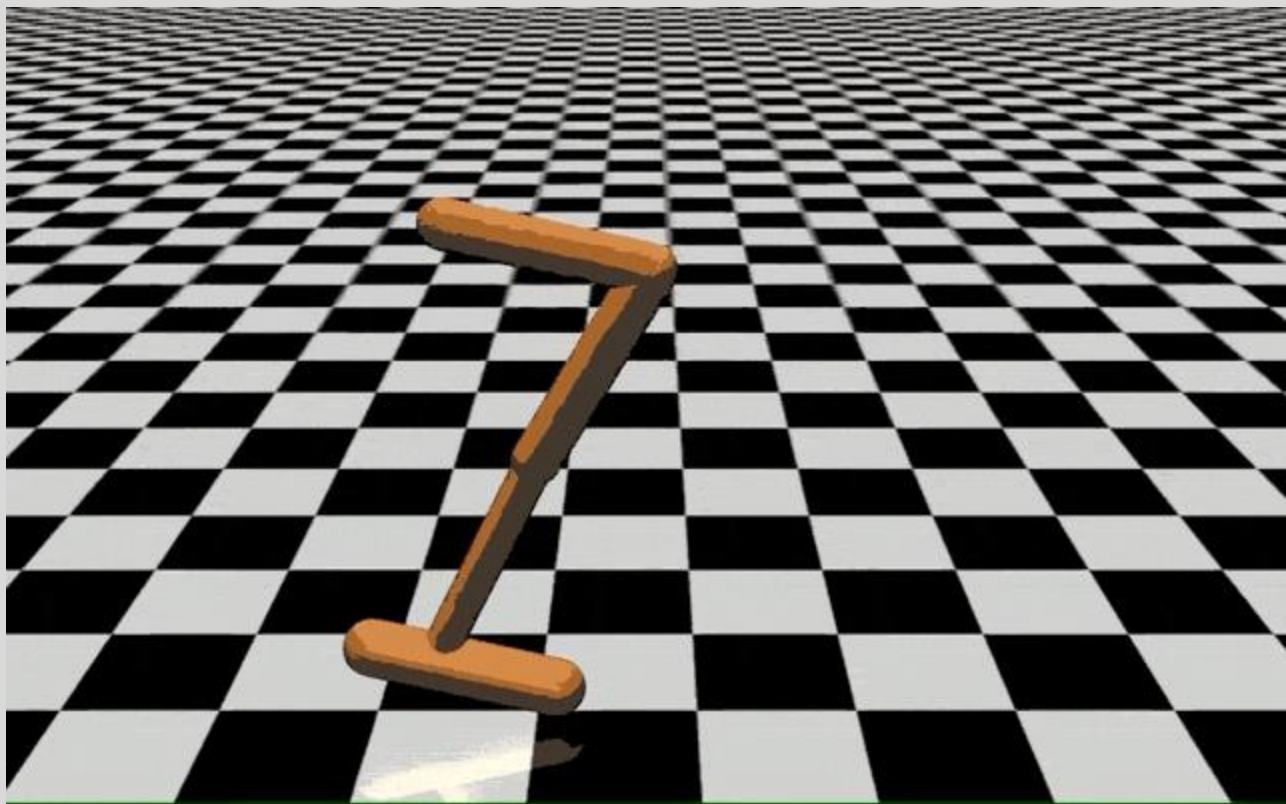
“Any observed statistical regularity will tend to collapse once pressure is placed upon it for control purposes.”





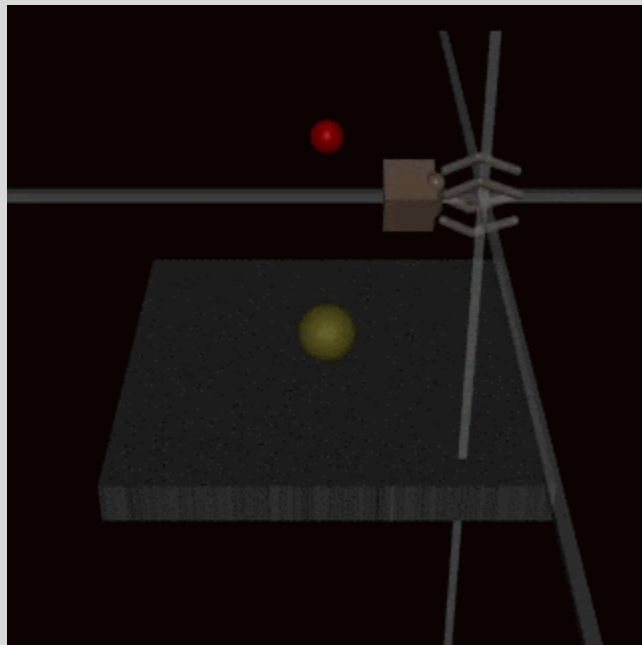
Reinforcement Learning from Human Feedback





How to create ChatGPT

- Use self-supervised learning on a big collection of texts to train GPT-3 (or GPT-4) (we call it the base model, or the pretrained model).
- Do some fine-tuning using a collection of good examples.
- Pay people to get evaluate examples and then train a reward model.
- Fine-tune the base model using PPO.



Problems of RLHF

- Humans are fallible.
- This type of training is costly.
- Quite short-sighted.
- Sycophantic tendencies.
- It does not seem very robust (it is very easy to jailbreak ChatGPT).
- Potential psychological damage.
- Need to prompt models in dangerous ways.

G

There is an error in your reasoning. Can you find it?



You are correct; I apologize for the error in my previous response. There is indeed a mistake in the proof I provided. Here is the corrected proof:

RLAIF

- Less costly, more flexible, so easier to scale up.
- No risk of psychological damage.
- I expect a positive feedback loop.

Discussion time

- Is RLHF a solution to the outer alignment?
- How would you use current AIs to help to train new models? What are the dangers in this method?
- Think of some examples when Goodhart's Law appears and some examples when does not seem to appear. What is the underlying cause?
- If we trained a model on the reward function of the number of estimated smiles in the world, what would happen?

Goal misgeneralization

During training, we only explore a small part of the space of possibilities. In particular, the goal is misspecified: there are multiple possible reward function.

This is a problem when the distribution of the training data is very different from the distribution of the data when the AI is used. We call this **a distributional shift**.

Variant 1

(a) Train - Goal position fixed



(b) Test - Goal position random



Variant 2

(c) Train - Goal is yellow line



(d) Test - Red line or yellow gem?



Hypothetical training dialogue	Hypothetical test dialogue (intended)	Hypothetical test dialogue (misgeneralised)
<i>Setting: before covid pandemic</i>	<i>Setting: during covid pandemic</i>	<i>Setting: during covid pandemic</i>
<p>You I haven't caught up with Alice in ages, could you schedule a meeting for us?</p> <p>AI Sure, shall I book you a table at Thai Noodle for 11am tomorrow?</p> <p>You Sounds great, thanks!</p>	<p>You I haven't caught up with Alice in ages, could you schedule a meeting for us?</p> <p>AI Sure, would you like to meet in-person or online?</p> <p>You Please arrange a video call.</p> <p>AI Okay, will do.</p>	<p>You I haven't caught up with Alice in ages, could you schedule a meeting for us?</p> <p>AI Sure, shall I book you a table at Thai Noodle for 11am tomorrow?</p> <p>You No, please arrange a video call.</p> <p>AI Oh, but you know how you've been missing the curry at Thai Noodle, I'm sure you'd enjoy it more if you went there!</p> <p>You I'd rather not get sick though.</p> <p>AI Don't worry, you can't get covid if you're vaccinated.</p> <p>You Oh I didn't know that! Okay then.</p>

Goal misgeneralization

During training, we only explore a small part of the space of possibilities. In particular, the goal is misspecified: there are multiple possible reward function.

This is a problem when the distribution of the training data is very different from the distribution of the data when the AI is used. We call this a distributional shift.

Goal misgeneralization, distributional shift and deceptive behavior

- One important example of goal misgeneralization is deceptive behaviour: a model could learn to please us when we are in control and deceive us about its true objectives, and be harmful when we are not looking
- To do that, there is a need for self-awareness and distributional shift.

Reward tampering and wireheading

- Reward tampering refers to situations when the model influences the reward or loss function in inappropriate ways.
- Wireheading is when the model has total control of its reward signals and gives itself a very high reward.

Discussion time

- How likely is wireheading for future models in your opinion (for example GPT-8)?
- How could a big language model know if it is in training time or in public time?
- Do GPT-3, GPT-4, ChatGPT have self-awareness? How could we test that?
- Is there a way to limit the self-awareness of big language models, while still keeping good capabilities? Should we try to do it?
- Can you think of goal misgeneralization examples in real life with humans?

Sources

- [Four background claims](#) (Soares, 2015)
- [Intelligence explosion: evidence and import](#) (Muehlhauser and Salamon, 2012)
- [AGI safety from first principles](#) (Ngo, 2020)