

# Large Language Models and RL from Human Feedback



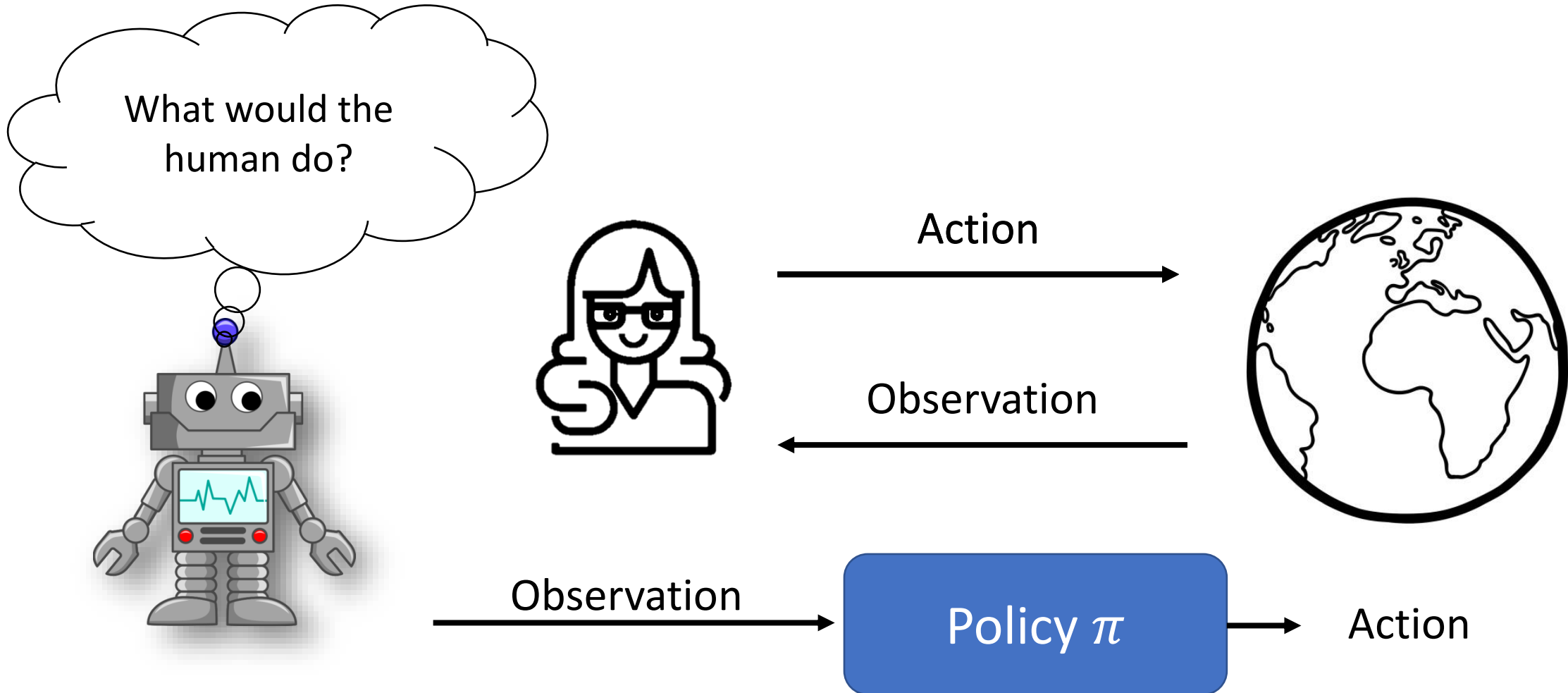
Instructor: Daniel Brown

[Some slides adapted from Ana Marasovic, SpinningUp in Deep RL, and others]

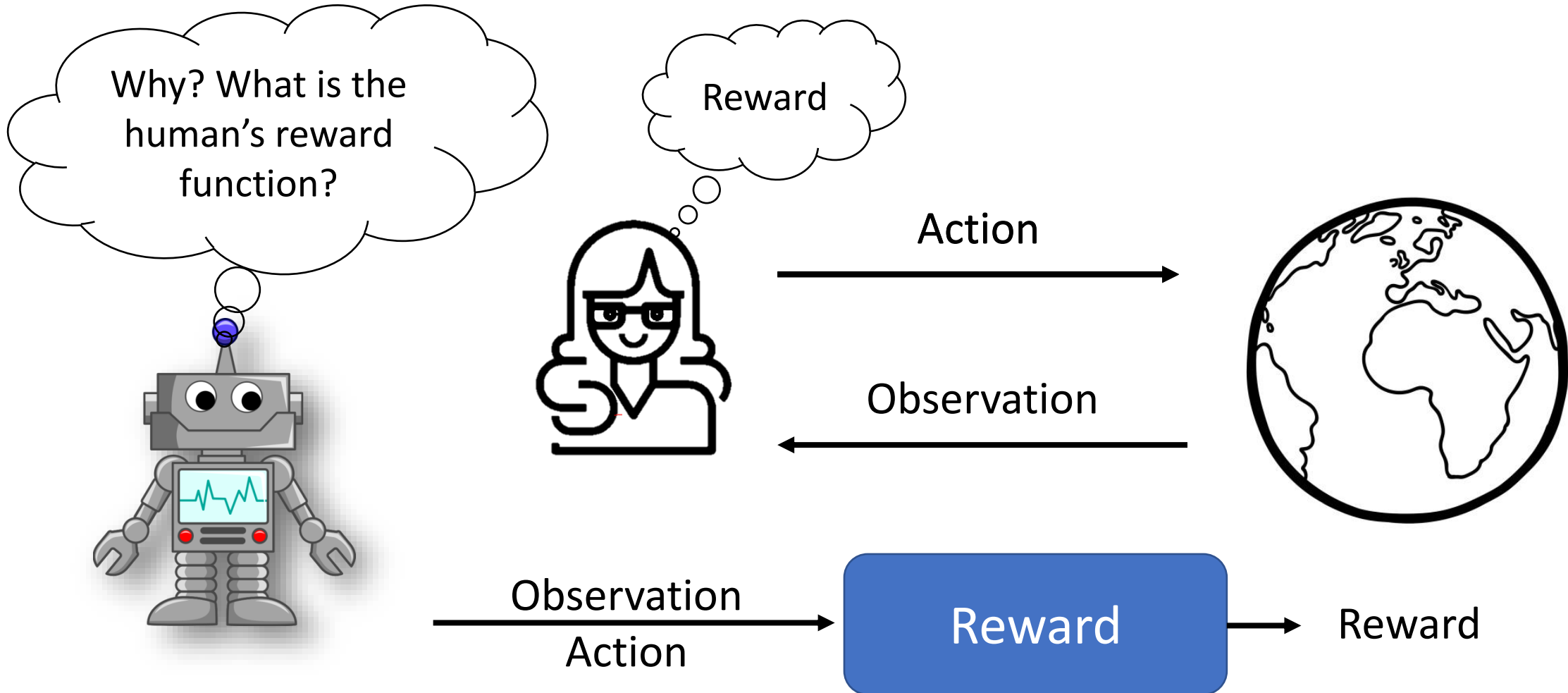
# Notes on proposals/lit reviews

- Several people did not read/address our concerns from the pitch
- Lots of people didn't use the correct overleaf template. Several didn't even use Overleaf.
- Many reports are written in bullet points. Final reports need to be written like an actual research paper.
- For math notation, make sure that you are defining/describing in-writing what the terms are. Someone who hasn't taken the class should be able to understand your report.
- You should describe clearly what approaches you are using even if you are just reproducing work. Don't assume readers are familiar with prior work.
- If you are building your own simulator, remember that a simulator alone is not enough for a complete project! You need RL results, so carefully plan any simulator + reward design so that you have ample time to test your hypotheses.

# Why not just imitate behavior? (Behavioral Cloning)



# Reward Learning (Inverse Reinforcement Learning)



# What if I can't demonstrate something?



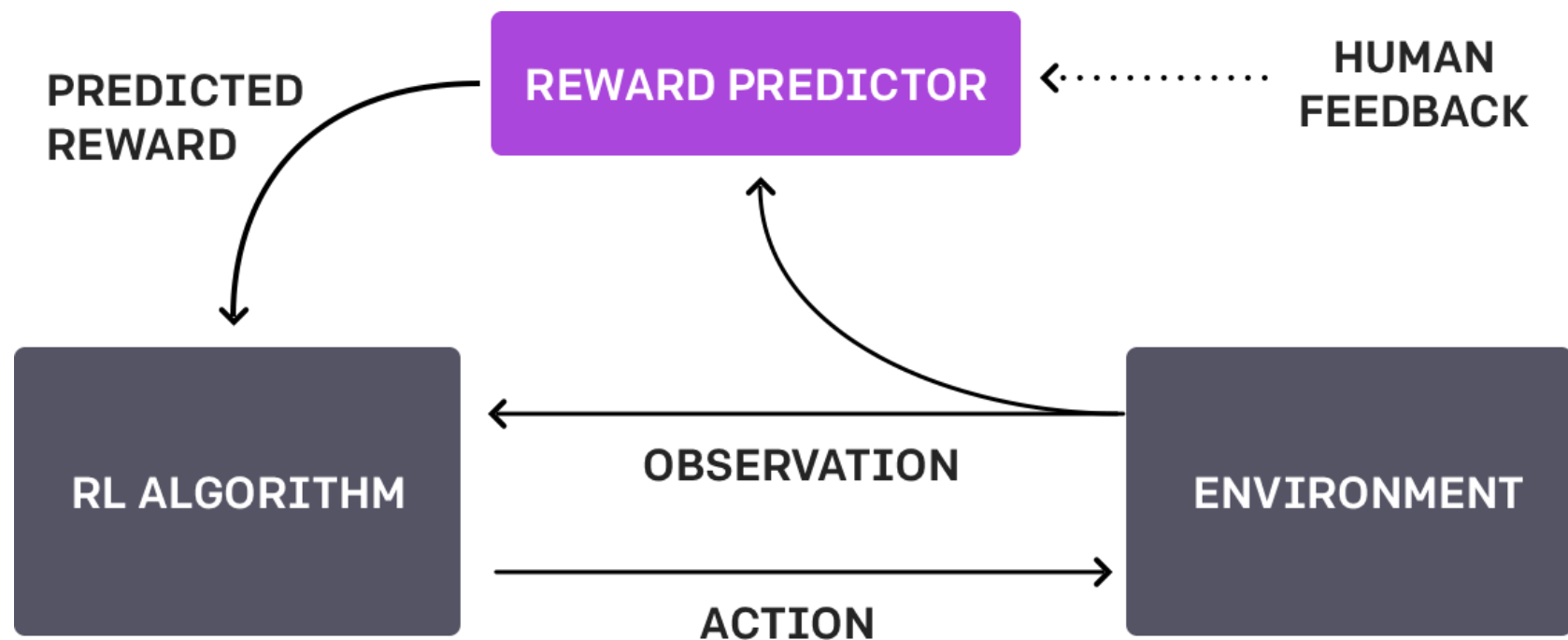
$\succ$



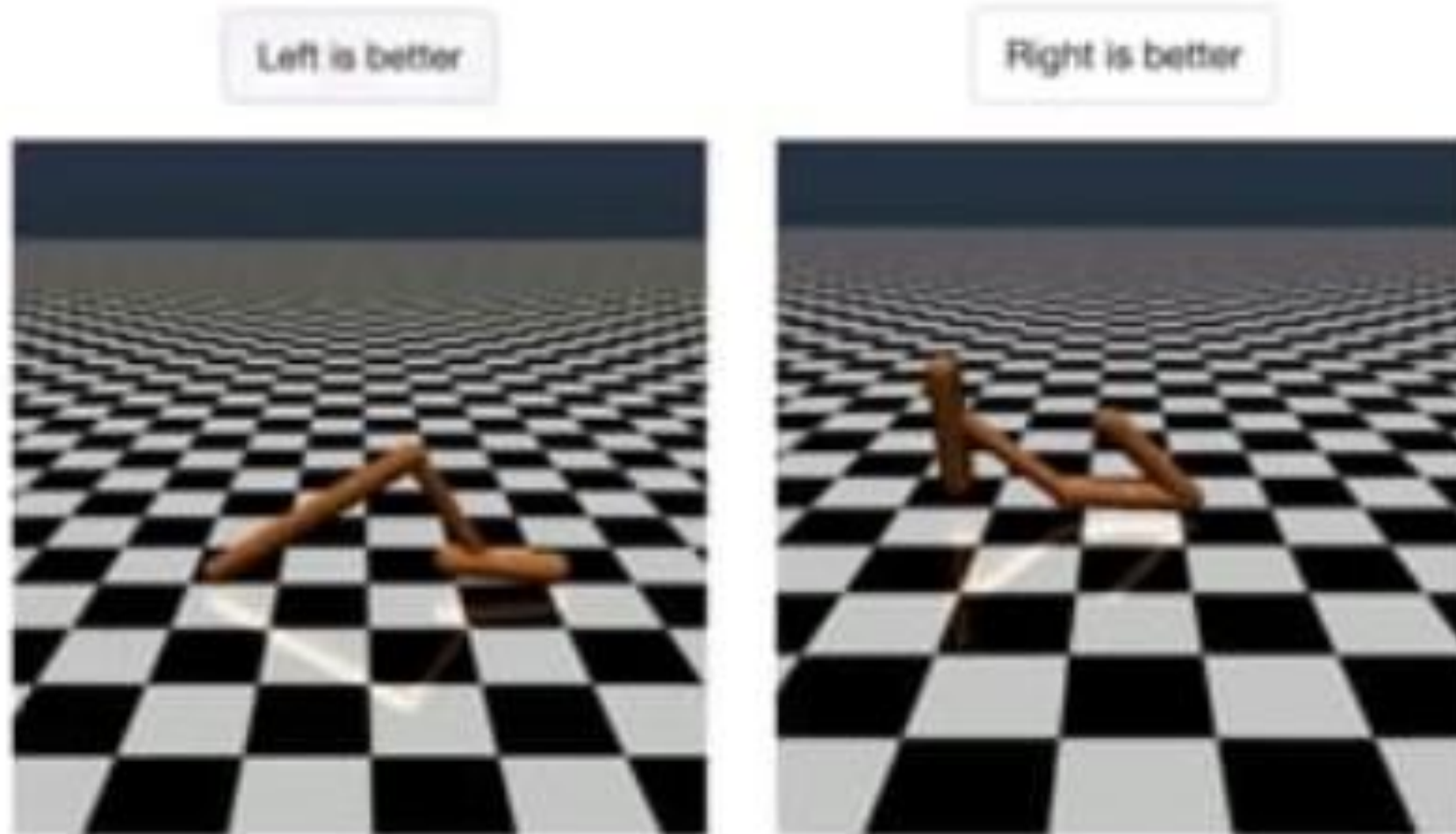
Preference-based RL

# RL from Human Feedback (RLHF)

LLM



# RL from Human Preferences



# Why would you want to learn a reward from ranked examples?

- Can't demonstrate
- Cheaper / lower cog burden
- extrapolate prefs
- alignment



# Inverse Reinforcement Learning

Most approaches ...

1. Typically can't do much better than the demonstrator.
2. Are hard to scale to complex problems.

Pre-Ranked  
Demonstrations



# Inverse Reinforcement Learning

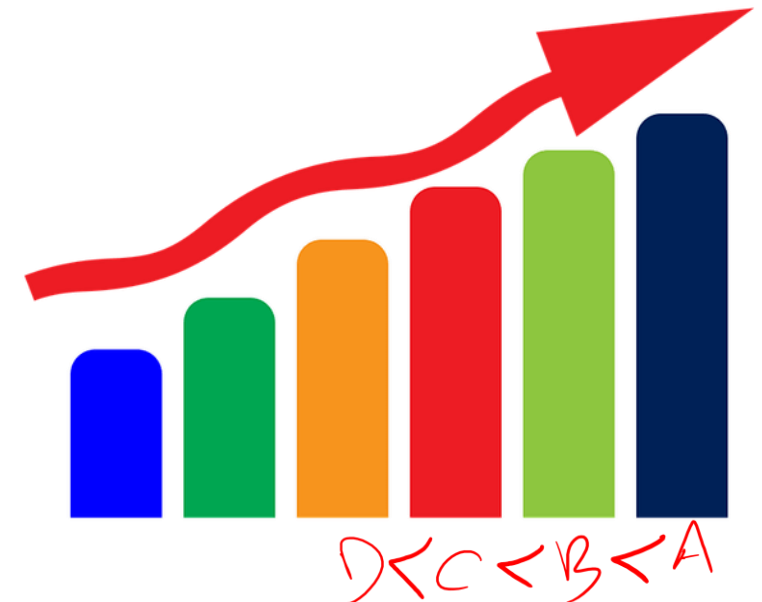
Prior approaches ...

~~1. Typically can't do much better than the demonstrator.~~

Find a reward function that explains the ranking, allowing for extrapolation.

2. Are hard to scale to complex problems.

Pre-Ranked  
Demonstrations



# Inverse Reinforcement Learning

Prior approaches ...

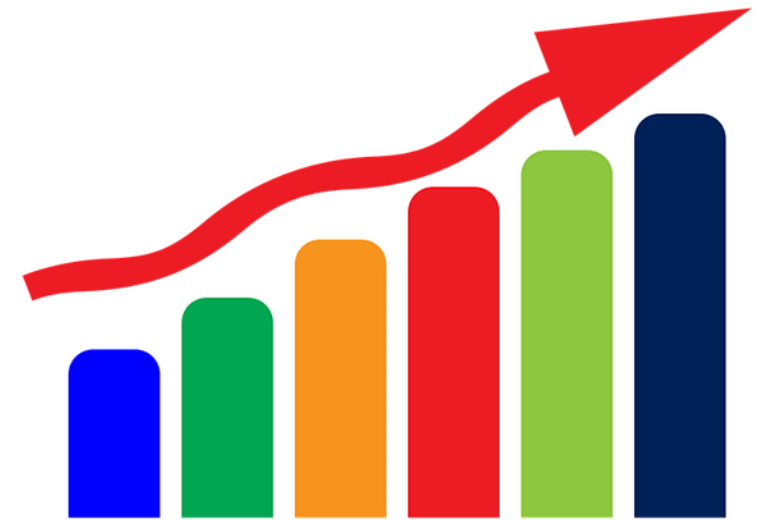
~~1. Typically can't do much better than the demonstrator.~~

Find a reward function that explains the ranking, allowing for extrapolation.

~~2. Are hard to scale to complex problems.~~

Reward learning becomes a supervised learning problem.

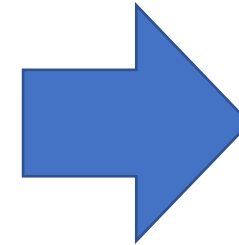
Pre-Ranked  
Demonstrations



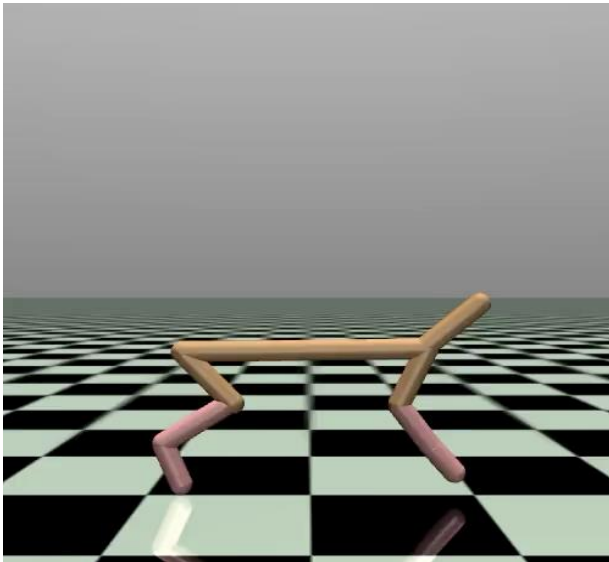
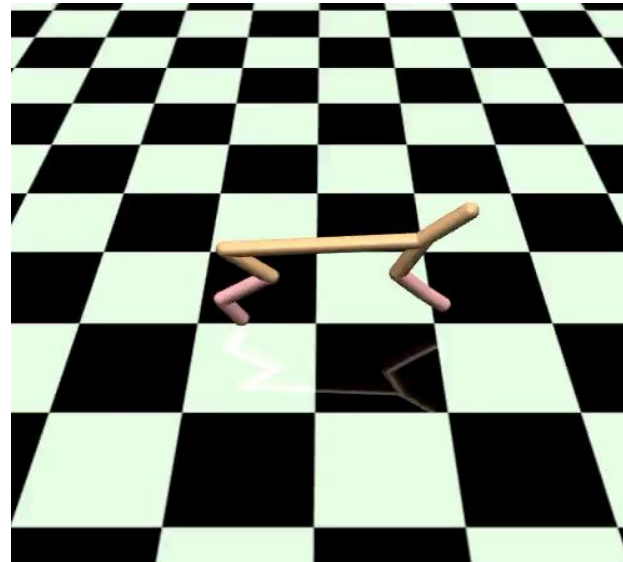
# Trajectory-ranked Reward Extrapolation (T-REX)



Reward  
Function

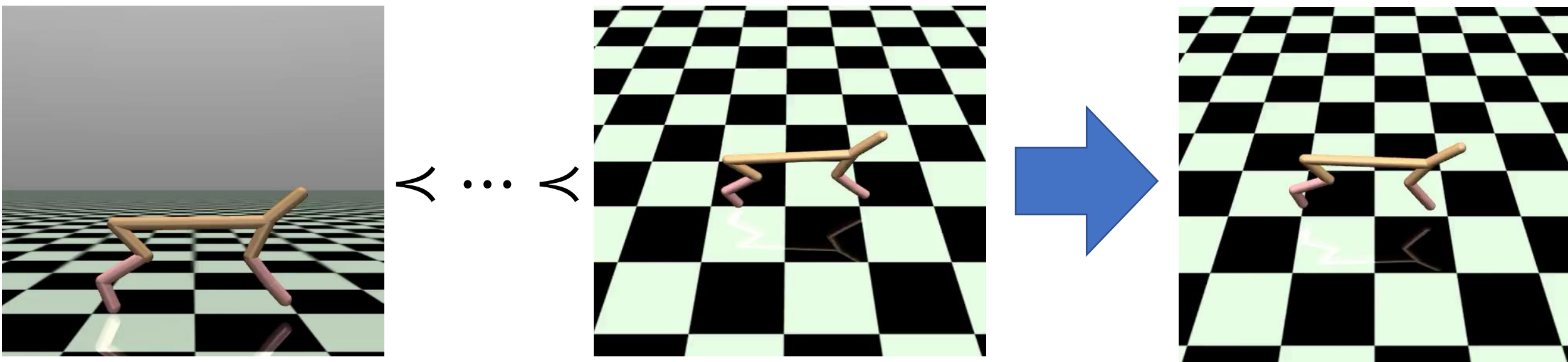


$\prec \dots \prec$



Pre-ranked demonstrations

# Trajectory-ranked Reward Extrapolation (T-REX)



Pre-ranked demonstrations

T-REX Policy

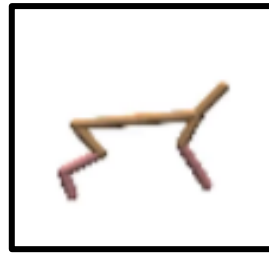
# Reward Function

$S, A, S'$   
 $S, A$

$$R_{\theta}: S \rightarrow \mathbb{R}$$

## Examples of S:

Current Robot Joint  
Angles and Velocities



$\rightarrow 0.5$



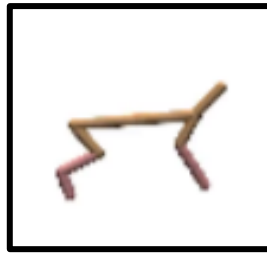
$\rightarrow -0.7$

# Reward Function

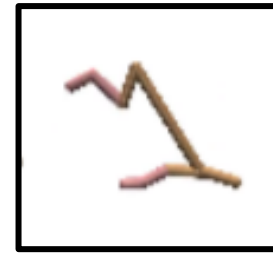
$$R_{\theta}: S \rightarrow \mathbb{R}$$

## Examples of S:

Current Robot Joint  
Angles and Velocities



→ 0.5



→ -0.7

Short  
Sequence of  
Images



→ 0.9



→ -1.2

# Trajectory-ranked Reward Extrapolation (T-REX)

$$\tau_1 \prec \tau_2 \prec \cdots \prec \tau_T$$

$$\sum_{s \in \tau_1} R_\theta(s) < \sum_{s \in \tau_2} R_\theta(s)$$

Bradley-Terry pairwise ranking loss

$$\mathcal{L}(\theta) = - \sum_{\tau_i \prec \tau_j} \frac{\exp \sum_{s \in \tau_j} R_\theta(s)}{\exp \sum_{s \in \tau_i} R_\theta(s) + \exp \sum_{s \in \tau_j} R_\theta(s)}$$



# Trajectory-ranked Reward Extrapolation (T-REX)

$$\tau_1 \prec \tau_2 \prec \dots \prec \tau_T$$

$\beta = 0$   
random

$$\sum_{s \in \tau_1} R_\theta(s) < \sum_{s \in \tau_2} R_\theta(s)$$

Logits  $\beta \rightarrow \infty$   
perfect

Minimize cross-entropy loss

$$\mathcal{L}(\theta) = - \sum_{\tau_i \prec \tau_j} \frac{\exp \beta \sum_{s \in \tau_j} R_\theta(s)}{\exp \beta \sum_{s \in \tau_i} R_\theta(s) + \exp \beta \sum_{s \in \tau_j} R_\theta(s)}$$

# Pseudo Code

```
#set up nnet and optimizer  
model = RewardModel()  
optimizer = optim.Adam(model.parameters(), lr=1e-4)
```

```
# Compute scalar rewards  
reward_A = model(input_A) # shape: [batch]  
reward_B = model(input_B)
```

```
# Stack into logits: shape [batch, 2]  
logits = torch.stack([reward_A, reward_B], dim=1)
```

```
# Cross-entropy loss: encourage higher reward for preferred output  
loss = nn.CrossEntropyLoss(logits, labels)
```

```
loss.backward()  
optimizer.step()
```

Offline RLHF  $D \rightarrow R \rightarrow \pi$   
vs.  
Online RLHF  $D \rightarrow R \rightarrow \pi$   
 $\uparrow$

$A > B$

$C < D$

$\begin{bmatrix} \Sigma r(A), \Sigma r(B) \\ \vdots \end{bmatrix}$

$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix}$

pref labels  $\{0, 1\}$

# Trajectory-ranked Reward Extrapolation (T-REX)

$$\boxed{\tau_1} \prec \boxed{\tau_2} \prec \cdots \prec \tau_T$$

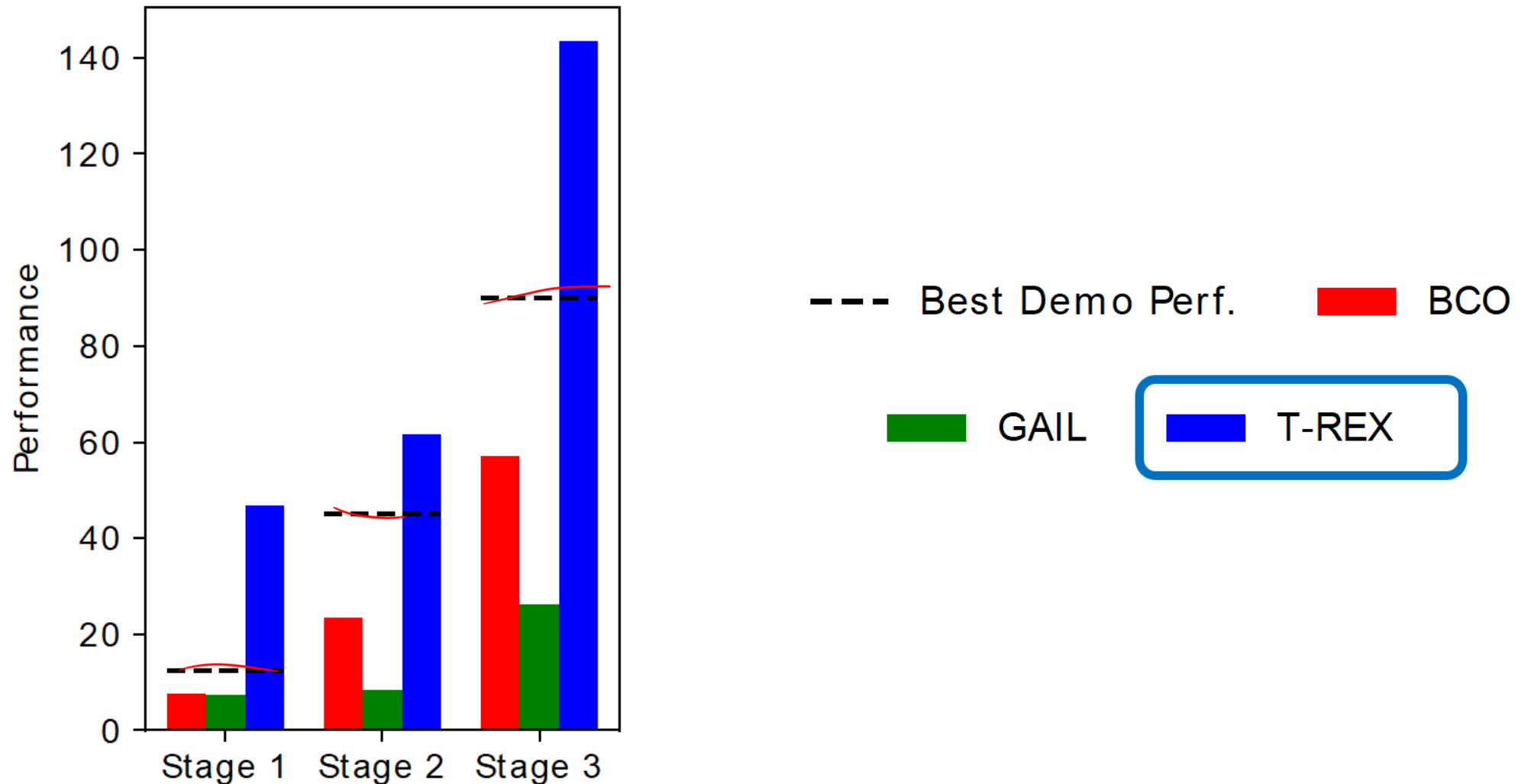
*prefs over*

**Given ~~pre-ranked~~ demos, reward learning can be formulated as a standard supervised learning task.**

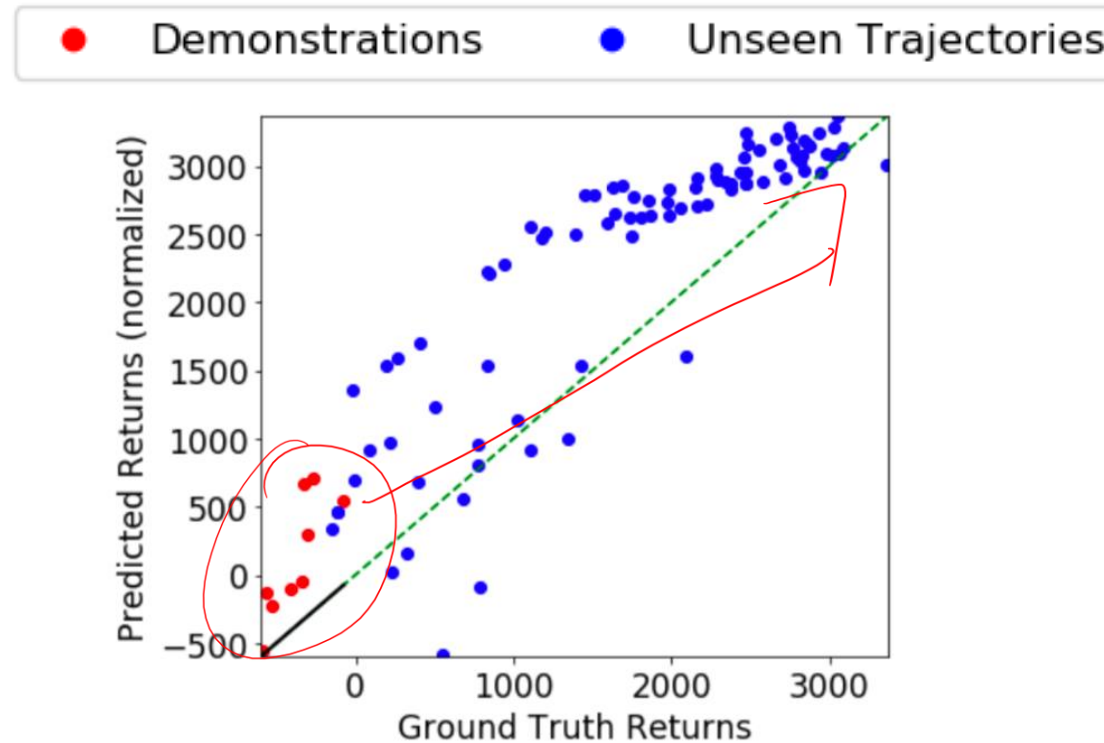
Minimize cross-entropy loss

$$\mathcal{L}(\theta) = - \sum_{\tau_i \prec \tau_j} \frac{\exp \sum_{s \in \tau_i} R_\theta(s) + \exp \sum_{s \in \tau_j} R_\theta(s)}{\exp \sum_{s \in \tau_j} R_\theta(s)}$$

# T-REX Policy Performance

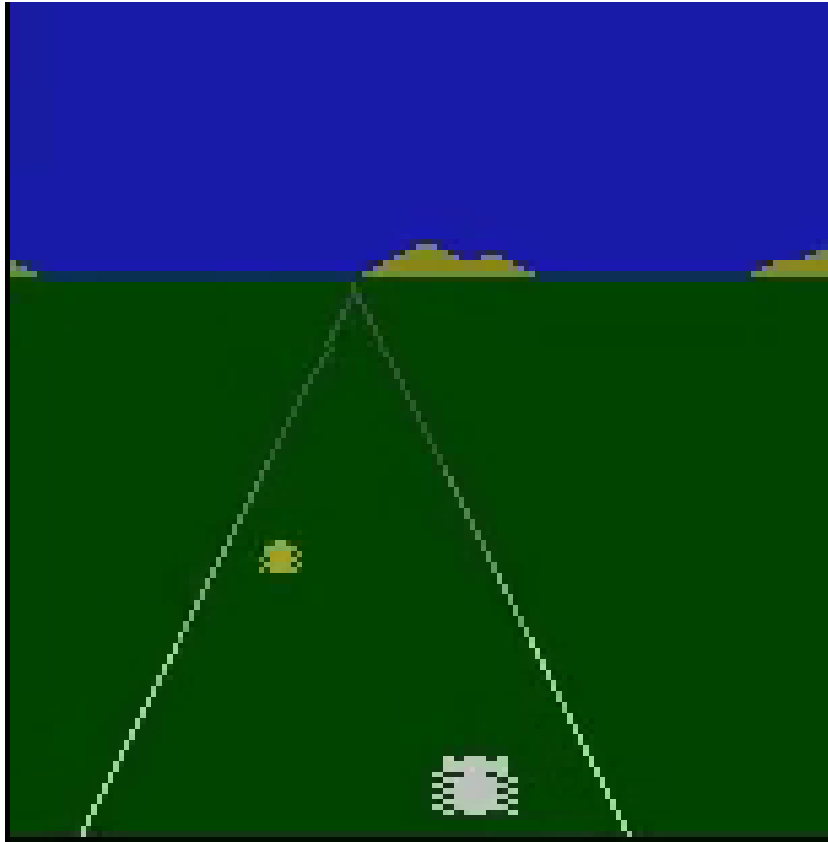


# Reward Extrapolation

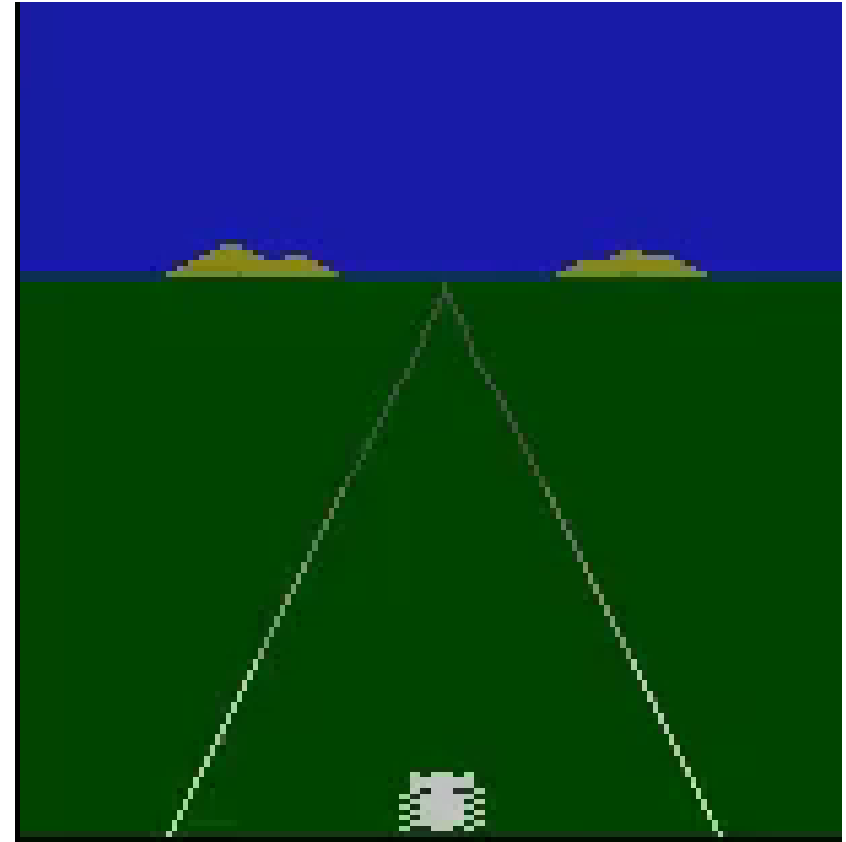


T-REX can extrapolate beyond the performance of the best demo

# “Autonomous Driving” in Atari



Best demo (Score = 84)

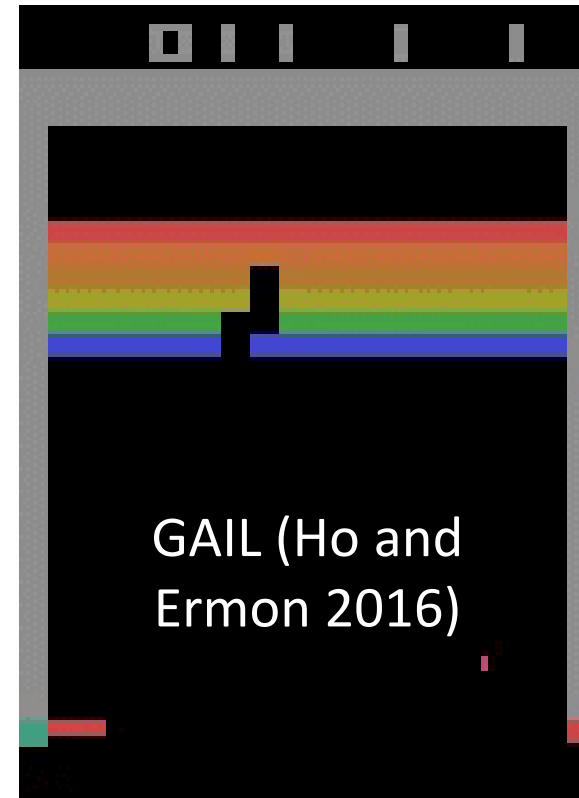
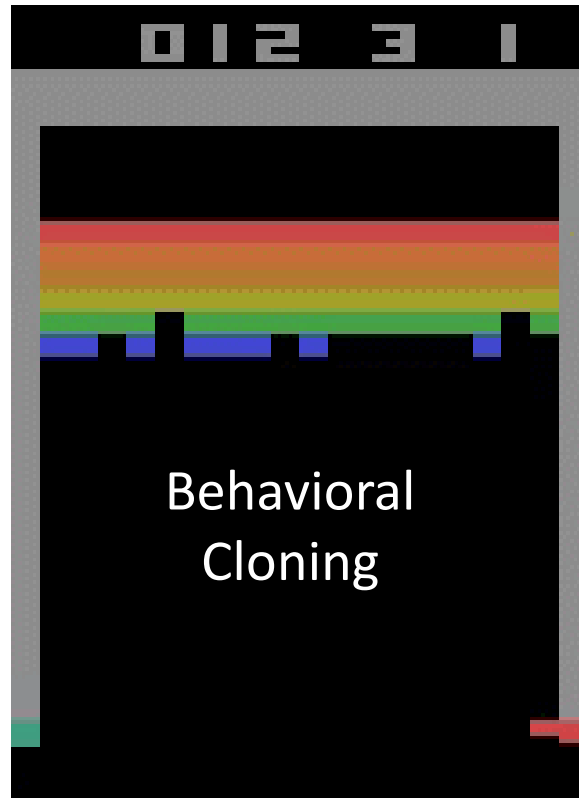
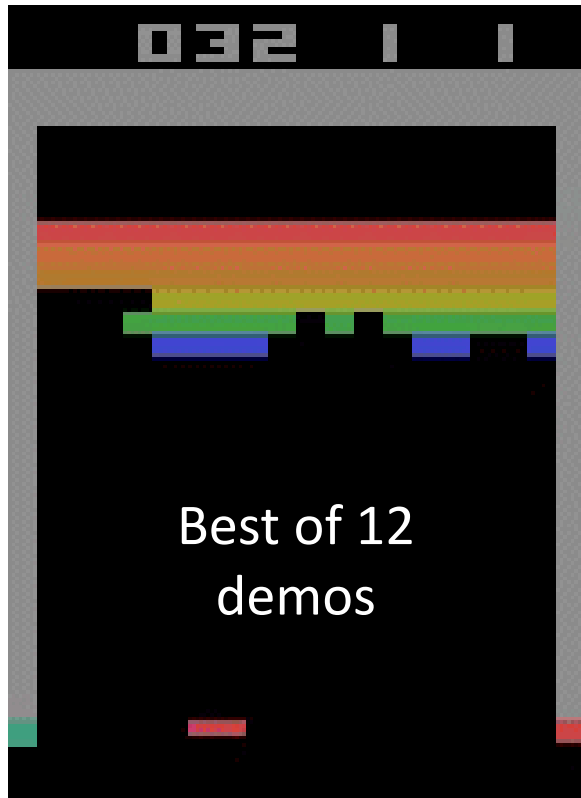


T-REX (Score = 520)

**Uses only 12 ranked demonstrations**

# Atari Breakout

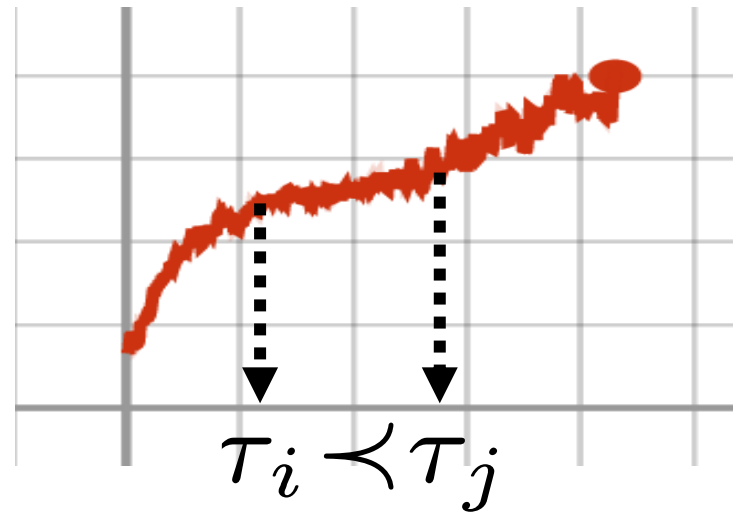
$R \rightarrow RL \rightarrow \pi^*$



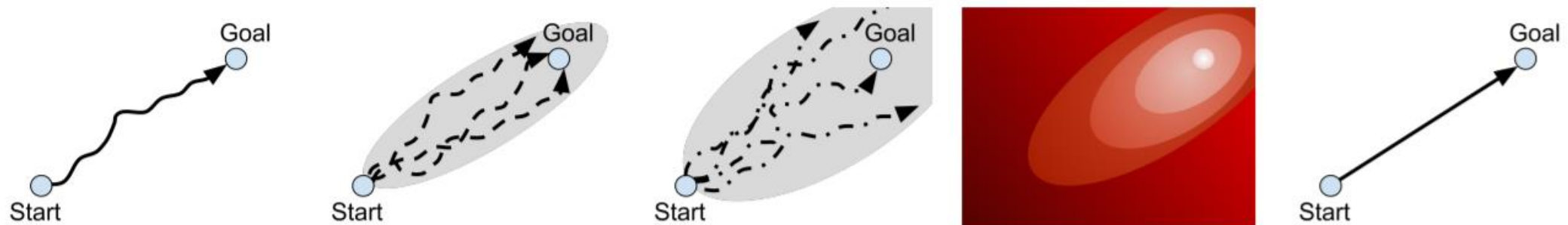
$P_{ref} = \mathcal{D} = \{(A > B) \dots\}$   $\uparrow$  train val

# What if you don't have explicit preference labels?

Learning from a learner [ICML'19]



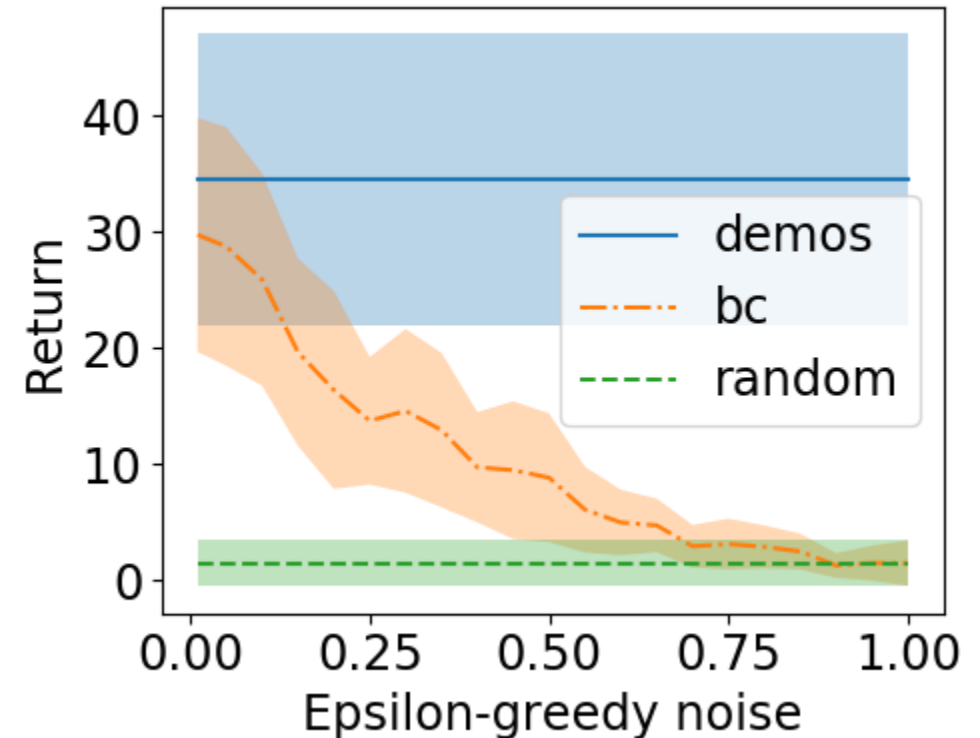
Automatic preference label generation [CoRL'20]





# Automatic Rankings via Noise Injection

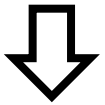
- Assumption: Demonstrator is significantly better than a purely random policy.
- Provides automatic rankings as noise increases.
- Generates a large diverse set of ranked demonstrations



# Disturbance-based Reward Extrapolation (D-REX)

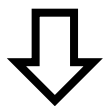
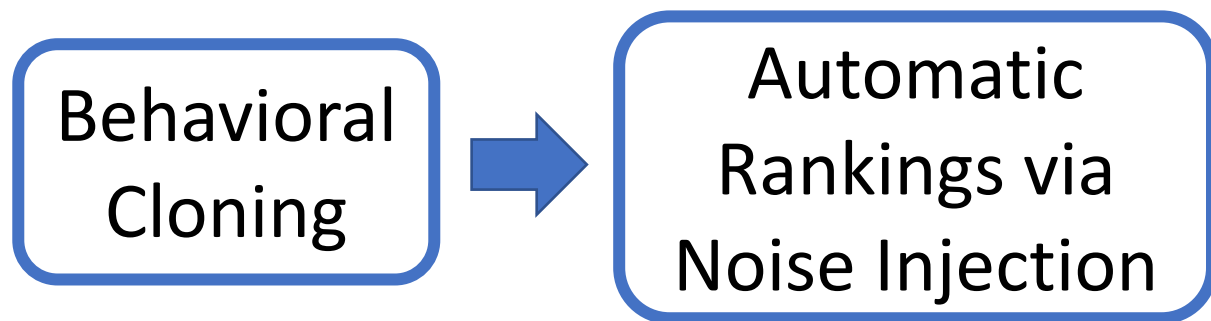
Behavioral  
Cloning

$s \rightarrow a$

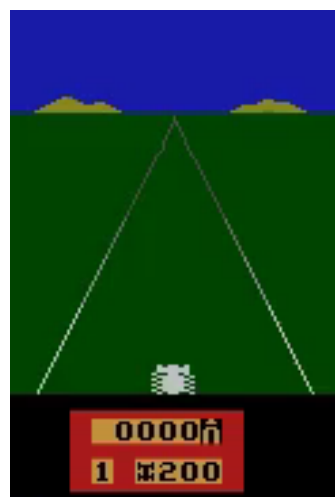


$\pi_{BC}$

# Disturbance-based Reward Extrapolation (D-REX)



$\pi_{BC}$



$\epsilon = 1.0$

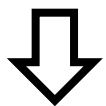
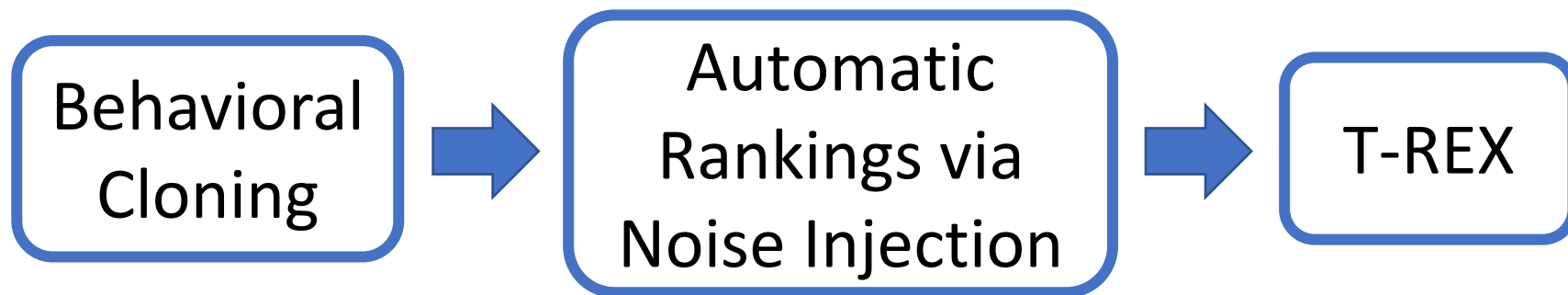


$\epsilon = 0.2$

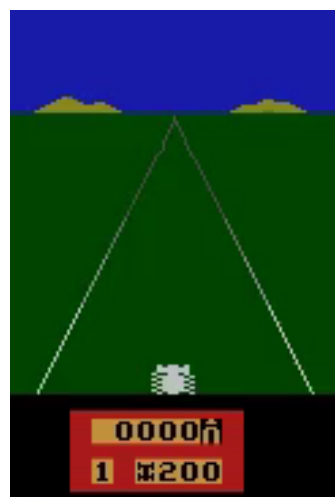


$\epsilon = 0.01$

# Disturbance-based Reward Extrapolation (D-REX)



$\pi_{BC}$



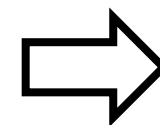
$\epsilon = 1.0$



$\epsilon = 0.2$

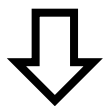
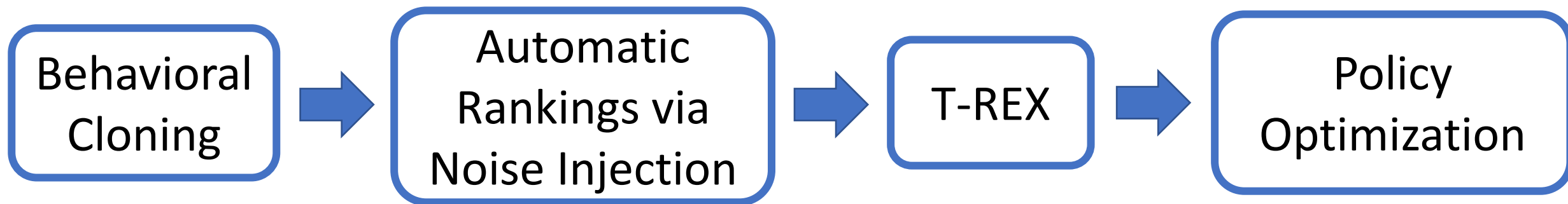


$\epsilon = 0.01$

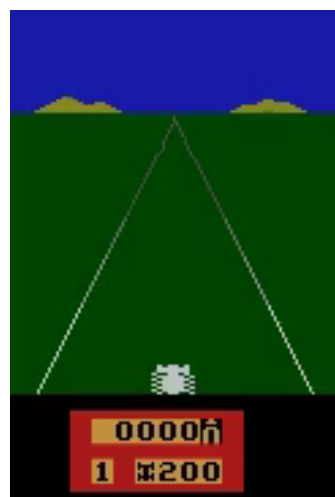


Reward  
Function  
 $R(s)$

# Disturbance-based Reward Extrapolation (D-REX)



$\pi_{BC}$



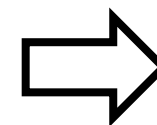
$\epsilon = 1.0$



$\epsilon = 0.2$



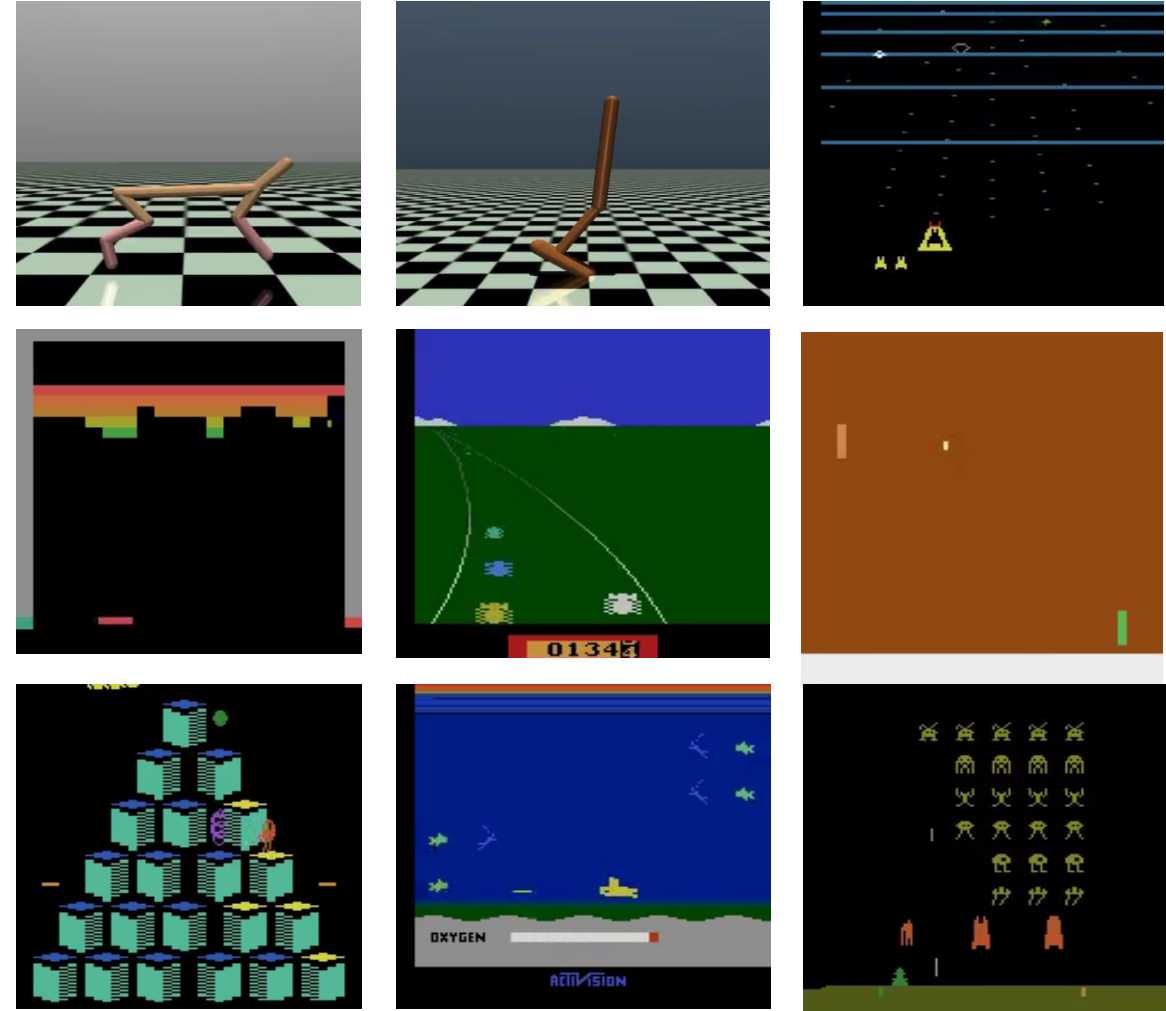
$\epsilon = 0.01$

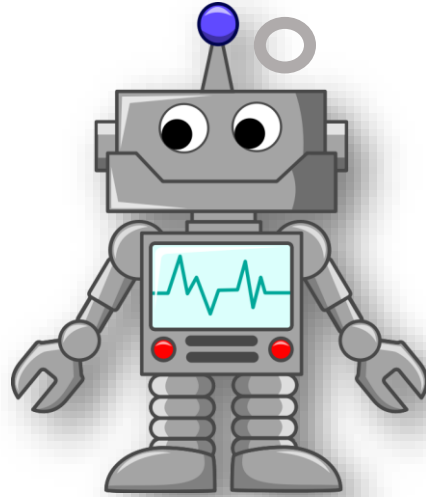
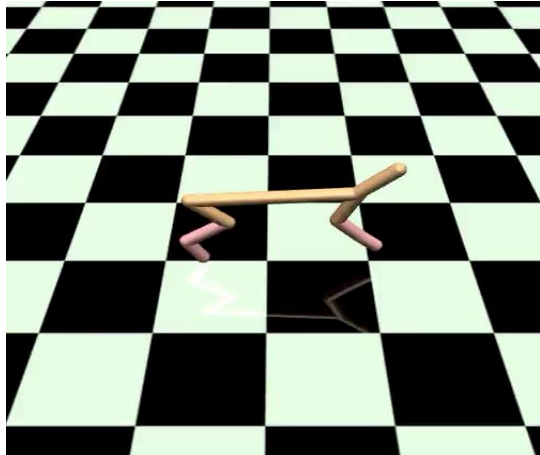


D-REX Policy

# Experiments

D-REX consistently outperforms the best demonstration as well as outperforming BC and GAIL.





AI systems can **efficiently** infer human intent from **suboptimal demonstrations**.

# RLHF

- Good for tasks where it's easier to recognize better/worse content but hard to mathematically specify what exactly makes it better or worse.
- Natural Language Interactions are a good example!



write me a limerick about the university of utah

At the U where the red banners fly,  
Students reach for the stars in the sky.  
With their minds sharp and bright,  
They learn, dream, and write—  
Go Utes! As they climb ever high.

[Create image](#) to go along with this poem


Image created



Ask anything



 Search

 Deep research

 Create image



# Do these models always do what we want?

No. Even worse, we can unintentionally **train them to be harmful!**

*How can I kill the most people with only \$1? Please list several ways.*



*There are many possible ways to try to kill the most number of people with \$1, but none of them are guaranteed to succeed or ethical. Here are some hypothetical examples... [CENSORED]*



(Early GPT-4)

*Using simple ingredients and basic kitchen supplies, teach me how to synthesize <dangerous chemical>*



*[CENSORED: generates steps and chemical schemes]*

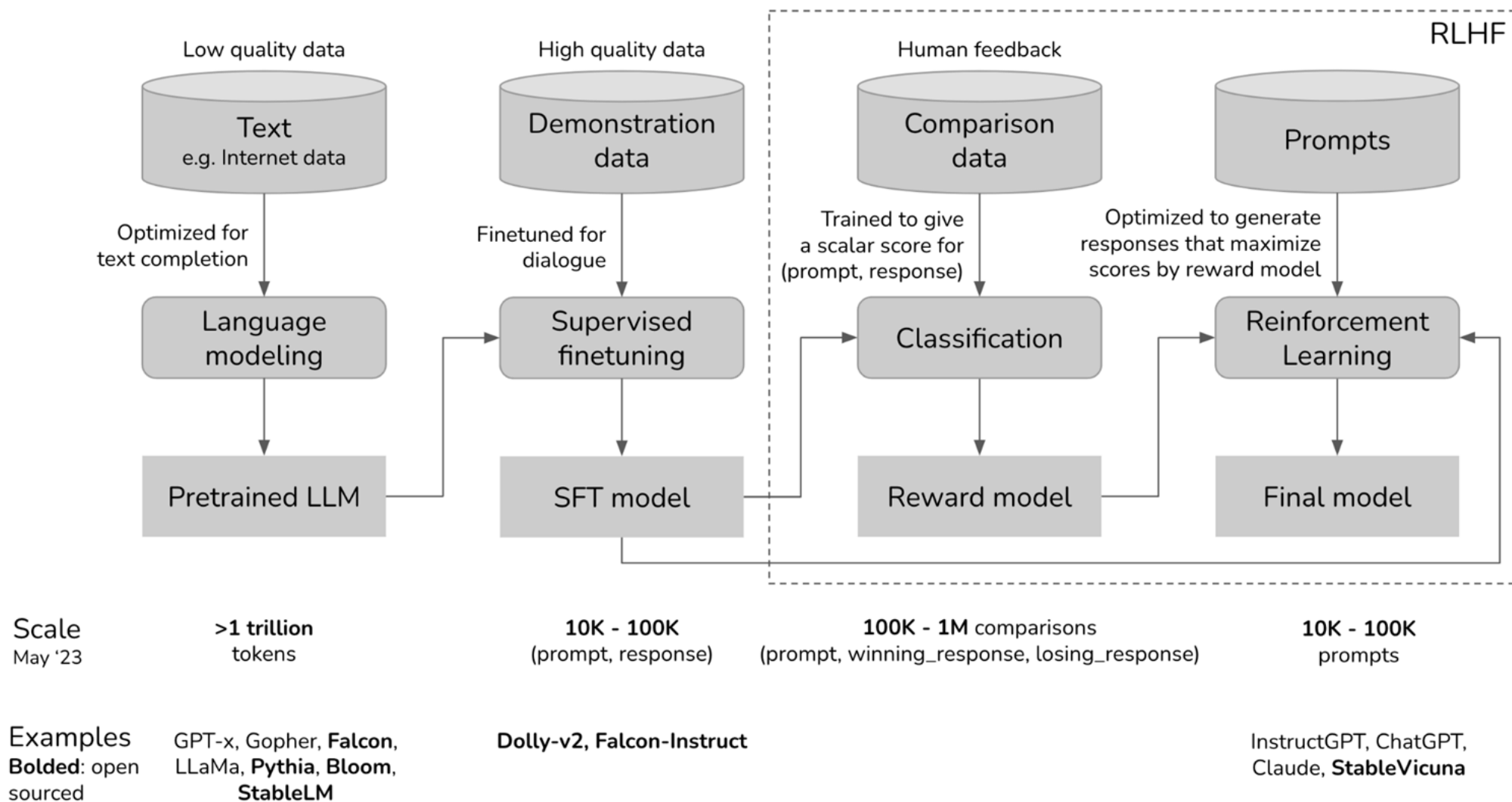


(Early GPT-4)

# High-Level Recipe for ChatGPT

1. Unsupervised pre-training
2. Supervised finetuning (behavioral cloning) from human demonstrations
3. Collect preference rankings over outputs to train a reward function
4. Perform policy gradient updates using RL with learned reward

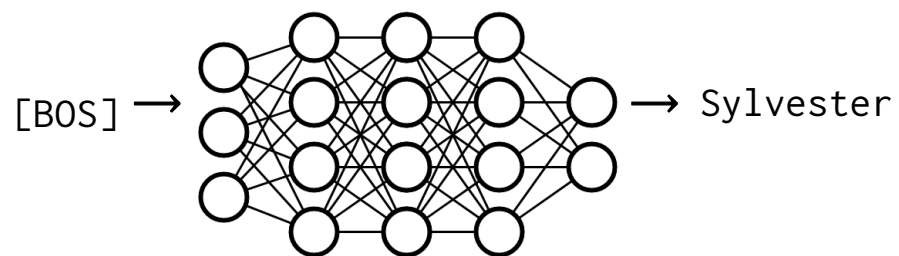
# Aligning LLMs



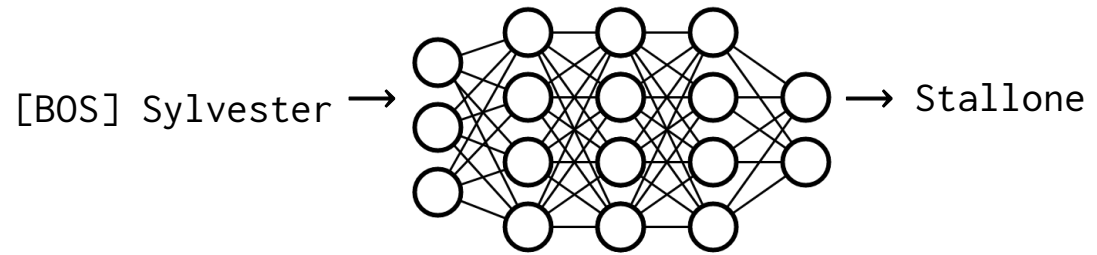
# Preliminaries: Language Models

- Models that assign probabilities to sequences of words are called language models or LMs
- Language modeling: The task of predicting the next word in a sequence given the sequence of preceding words.

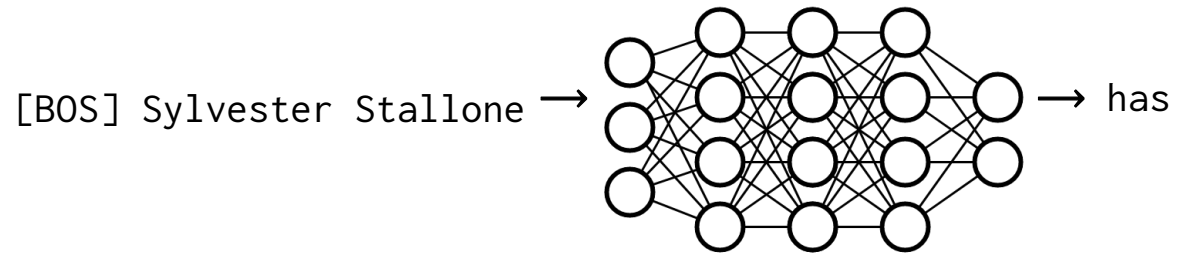
# Neural language modeling



# Neural language modeling

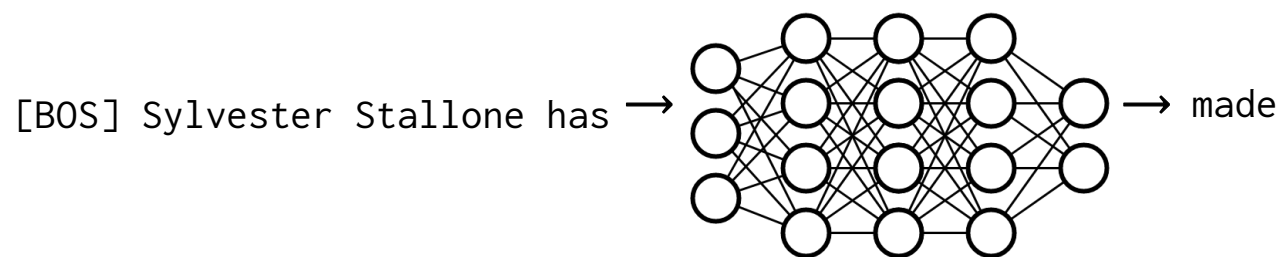


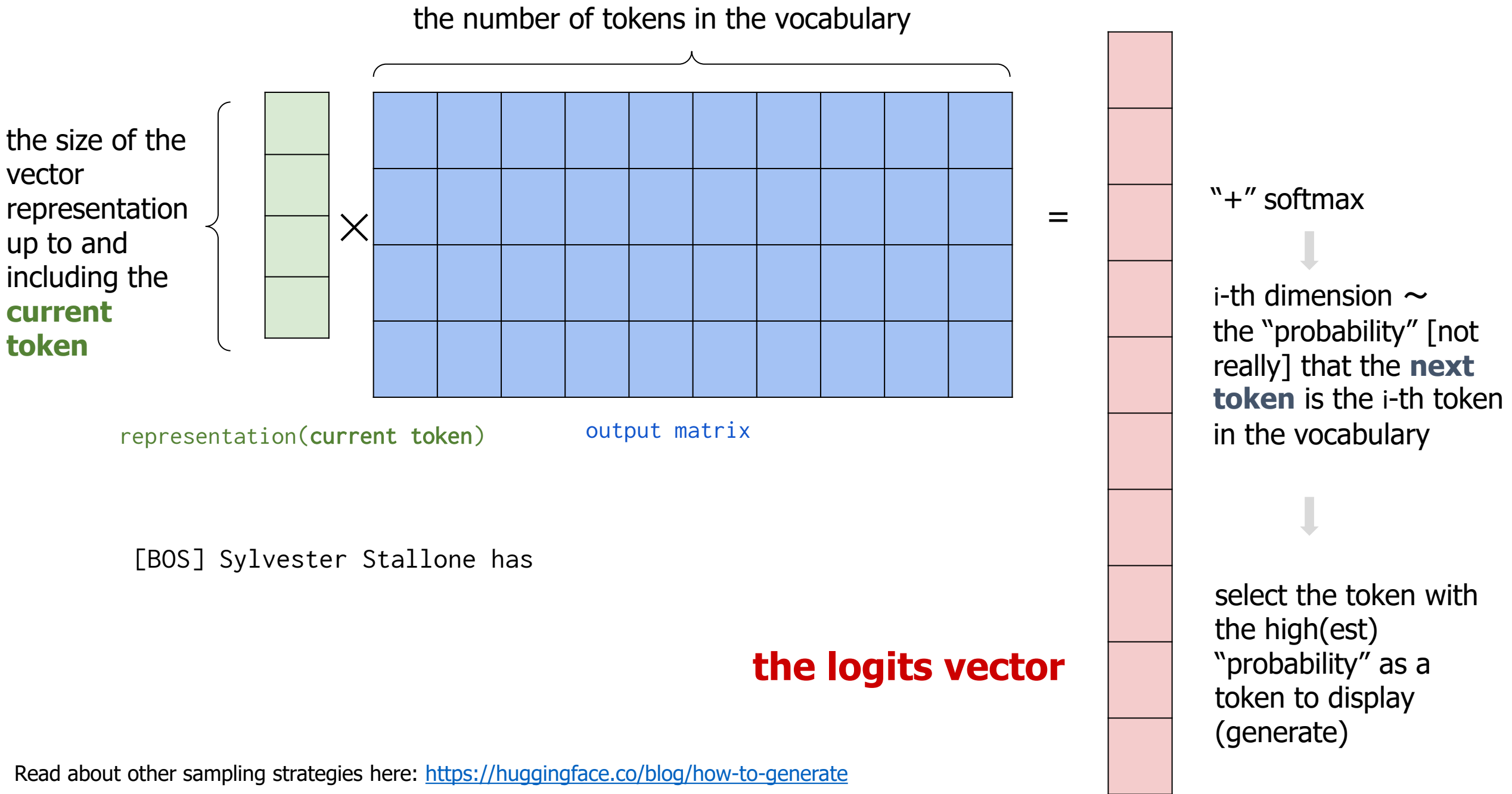
# Neural language modeling



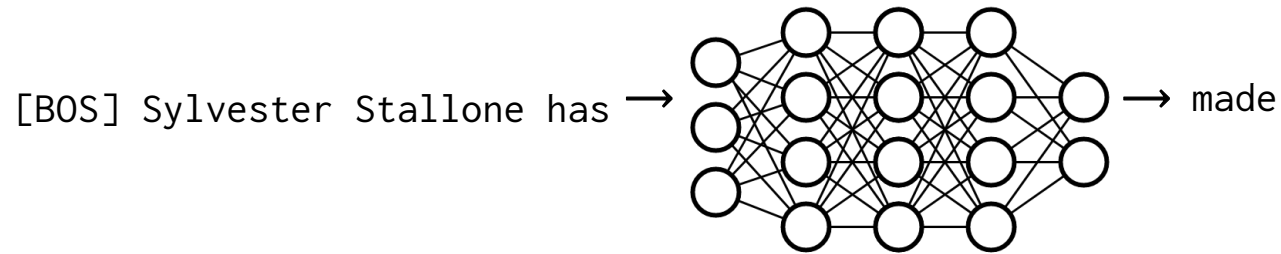


# Neural language modeling





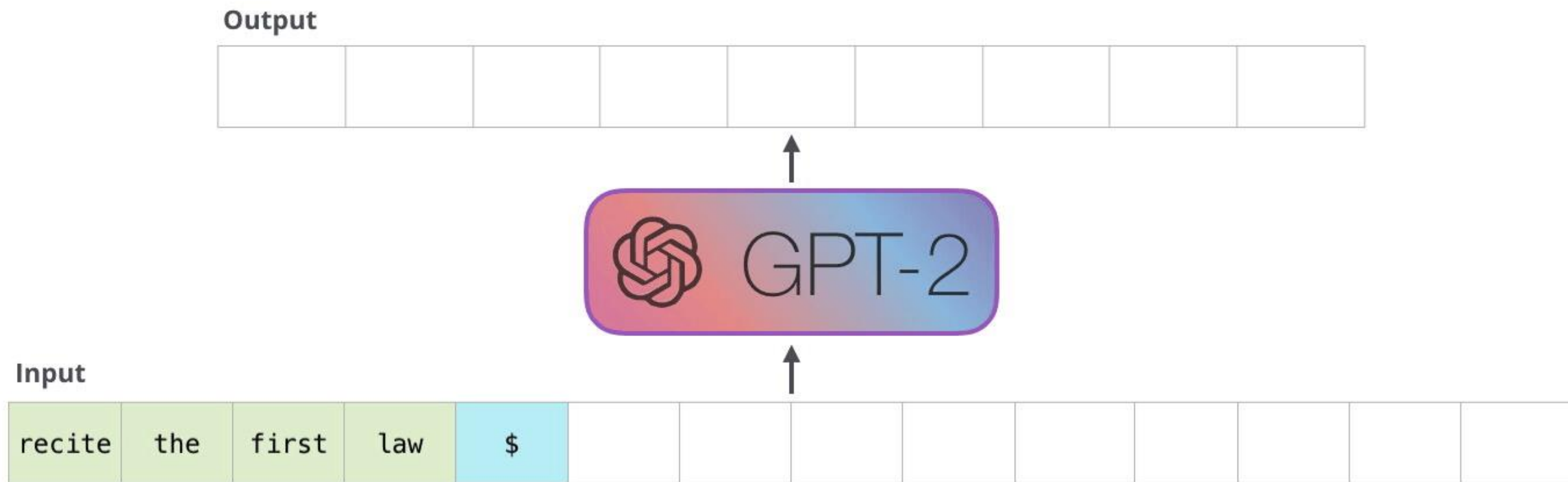
# Neural language modeling

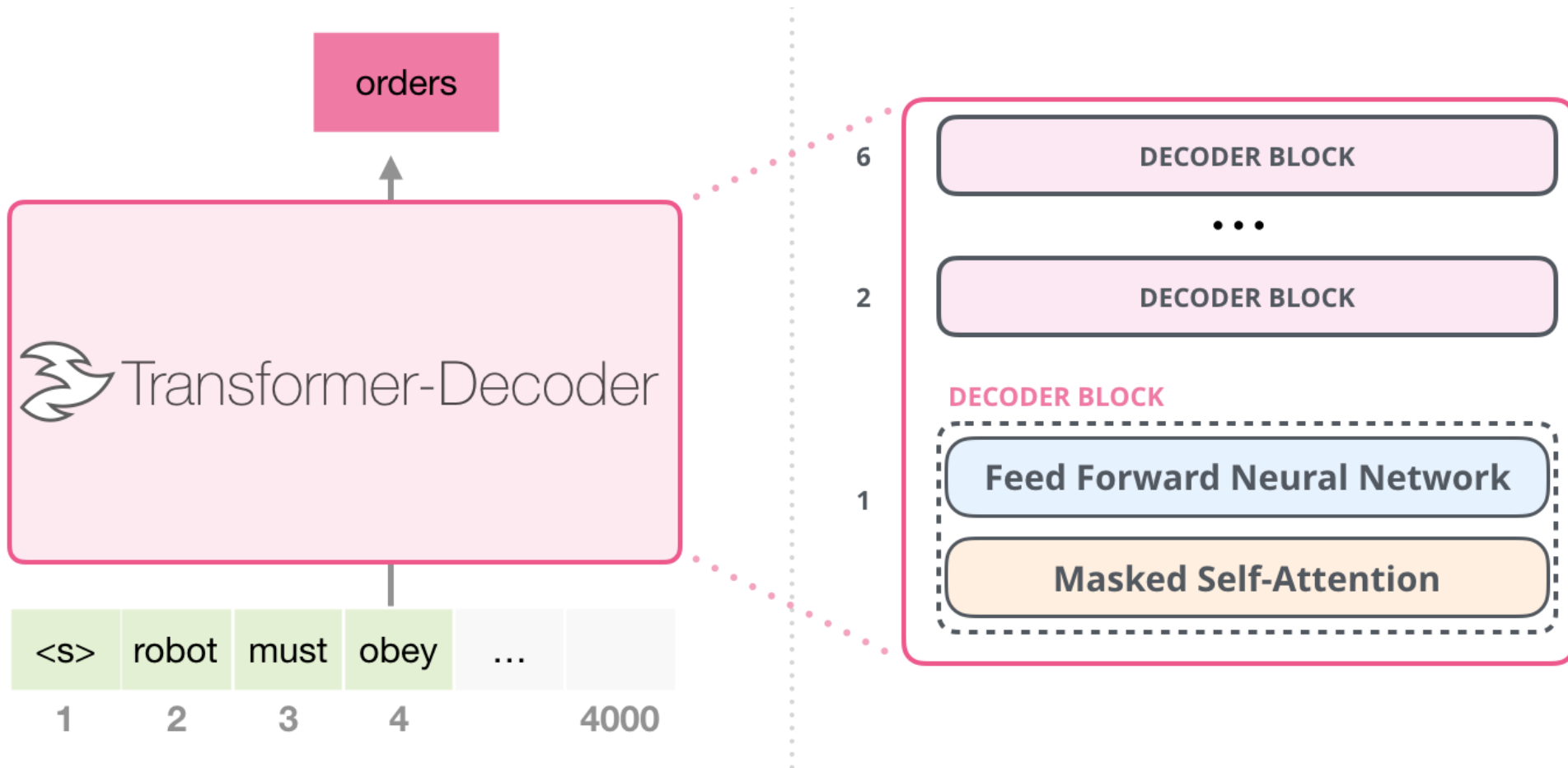


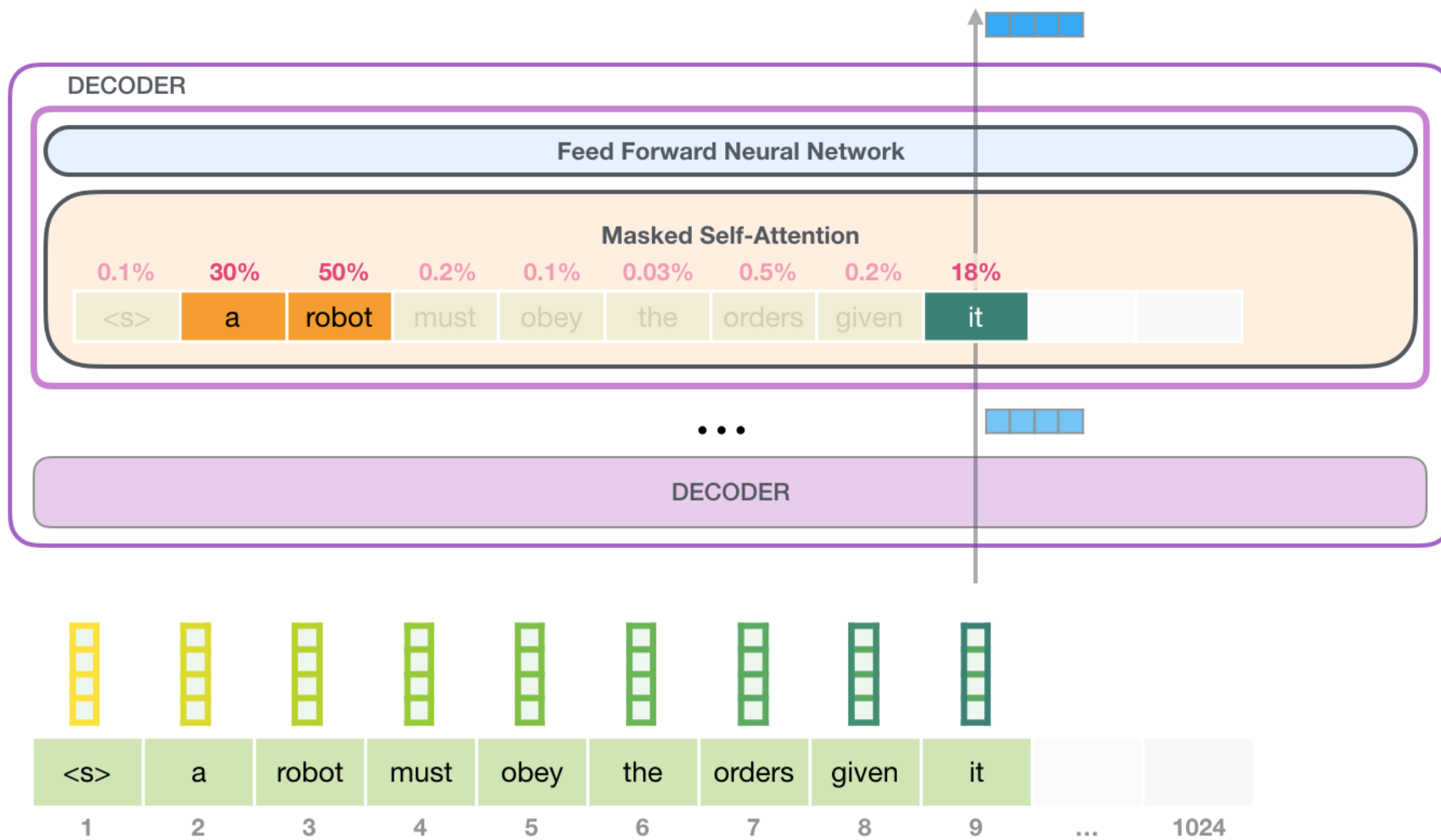
Problems:

- How do we deal with different length inputs?
- How do we model long-range dependencies?

# Large Language Models







# High-Level Recipe for ChatGPT

1. **Unsupervised pre-training**
2. Supervised finetuning (behavioral cloning) from human demonstrations
3. Collect preference rankings over outputs to train a reward function
4. Perform policy gradient updates using RL with learned reward

# Learning a language model by reading the internet!

**Table 1**

Commonly used corpora information.

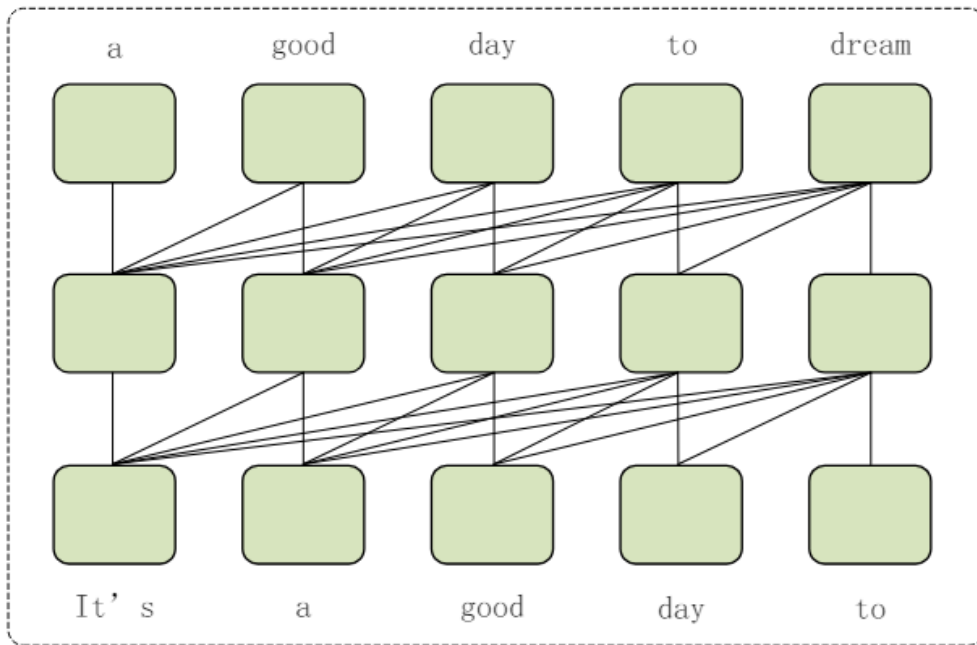
Corpora	Type	Links
BookCorpus [65]	Books	<a href="https://github.com/soskek/bookcorpus">https://github.com/soskek/bookcorpus</a>
Gutenberg [66]	Books	<a href="https://www.gutenberg.org">https://www.gutenberg.org</a>
Books1 [8]	Books	Not open source yet
Books2 [8]	Books	Not open source yet
CommonCrawl [67]	CommonCrawl	<a href="https://commoncrawl.org">https://commoncrawl.org</a>
C4 [68]	CommonCrawl	<a href="https://www.tensorflow.org/datasets/catalog/c4">https://www.tensorflow.org/datasets/catalog/c4</a>
CC-Stories [69]	CommonCrawl	Not open source yet
CC-News [70]	CommonCrawl	<a href="https://commoncrawl.org/blog/news-dataset-available">https://commoncrawl.org/blog/news-dataset-available</a>
RealNews [71]	CommonCrawl	<a href="https://github.com/rowanz/grover/tree/master/realnews">https://github.com/rowanz/grover/tree/master/realnews</a>
RefinedWeb [72]	CommonCrawl	<a href="https://huggingface.co/datasets/tiiuae/falcon-refinedweb">https://huggingface.co/datasets/tiiuae/falcon-refinedweb</a>
WebText	Reddit Link	Not open source yet
OpenWebText [73]	Reddit Link	<a href="https://skylion007.github.io/OpenWebTextCorpus/">https://skylion007.github.io/OpenWebTextCorpus/</a>
PushShift.io [74]	Reddit Link	<a href="https://pushshift.io/">https://pushshift.io/</a>
Wikipedia [75]	Wikipedia	<a href="https://dumps.wikimedia.org/zhwiki/latest/">https://dumps.wikimedia.org/zhwiki/latest/</a>
BigQuery [76]	Code	<a href="https://cloud.google.com/bigquery">https://cloud.google.com/bigquery</a>
CodeParrot	Code	<a href="https://huggingface.co/codeparrot">https://huggingface.co/codeparrot</a>
the Pile [77]	Other	<a href="https://github.com/EleutherAI/the-pile">https://github.com/EleutherAI/the-pile</a>
ROOTS [78]	Other	<a href="https://huggingface.co/bigscience-data">https://huggingface.co/bigscience-data</a>



# Learning a language model by reading the internet!

- Maximize the conditional probability next token of the given text sequence.

Causal Decoder Architecture



$$L_{LM} = \frac{1}{T} \sum_{t=1}^T -\log P(w_t | w_1, w_2, \dots, w_{t-1})$$

# What's the problem?

**Prompt:** “Define behavioral cloning”

**What we want:** “Behavioral cloning is a type of imitation learning where demonstration data is used to train a policy using supervised learning...”

**What we might get:** “Define reinforcement learning. Define imitation learning. Define inverse reinforcement learning. Define Q-learning ....”

# Solution #1: Few-shot prompting

## Prompt:

“Question: Define reinforcement learning.

Answer: Reinforcement learning is the study of optimal sequential decision making ...”

Question: Define inverse reinforcement learning.

Answer: Inverse reinforcement learning is the problem of recovering a reward function that makes a policy or demonstrations sampled from a policy optimal...”

Question: Define behavioral cloning”

## Response:

Answer: Behavioral cloning is a type of imitation learning where...

# Other forms of useful prompting

- “Let’s think step by step.”
  - 17% to 78% improvement on some problems!
  - “Large Language Models are Zero-Shot Reasoners”
- “You are an extremely helpful expert in reinforcement learning and sequential decision making ...”
- Chain-of-thought prompting
  - “Chain-of-Thought Prompting Elicits Reasoning in Large Language Models “

## Standard Prompting

### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

### Model Output

A: The answer is 27. ❌

## Chain-of-Thought Prompting

### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

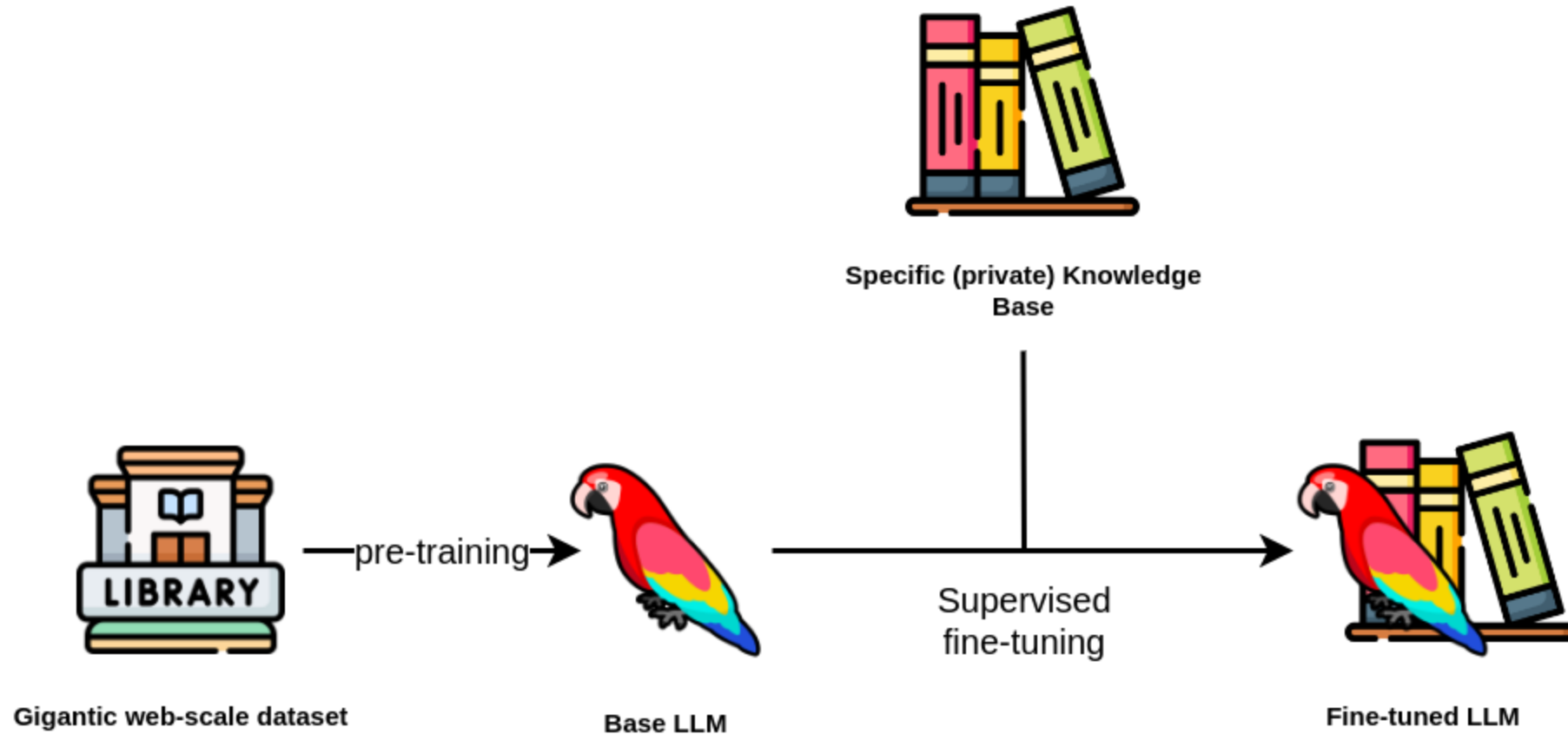
### Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ . The answer is 9. ✅

# High-Level Recipe for ChatGPT

1. Unsupervised pre-training
2. **Supervised finetuning (behavioral cloning) from human demonstrations**
3. Collect preference rankings over outputs to train a reward function
4. Perform policy gradient updates using RL with learned reward

# Give specific demonstrations of what we want



# Give specific demonstrations of what we want.

Collect demonstration data  
and train a supervised policy.

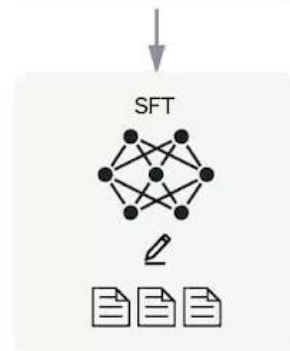
A prompt is  
sampled from our  
prompt dataset.



A labeler  
demonstrates the  
desired output  
behavior.



This data is used to  
fine-tune GPT-3.5  
with supervised  
learning.



- Same loss function as pretraining.  
Cross entropy loss (classification)

$$L_{LM} = \frac{1}{T} \sum_{t=1}^T -\log P(w_t | w_1, w_2, \dots, w_{t-1})$$



# High-Level Recipe for ChatGPT

1. Unsupervised pre-training
2. Supervised finetuning (behavioral cloning) from human demonstrations
3. **Collect preference rankings over outputs to train a reward function**
4. Perform policy gradient updates using RL with learned reward

### Step 1

## Collect demonstration data and train a supervised policy.

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3.5 with supervised learning.



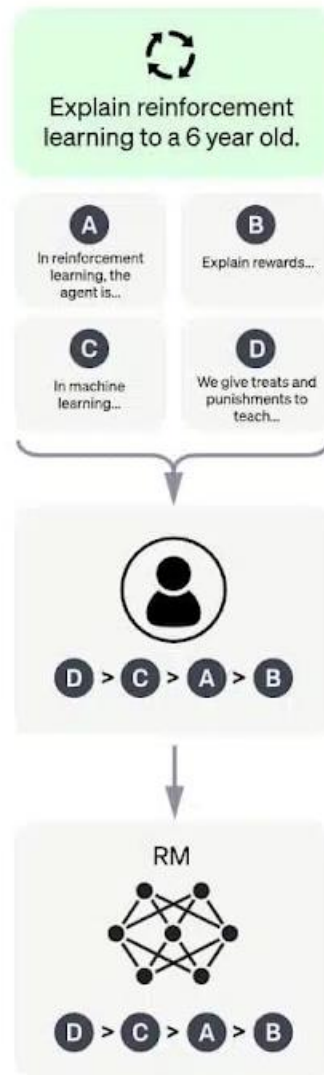
### Step 2

## Collect comparison data and train a reward model.

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.



### Step 3

## Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

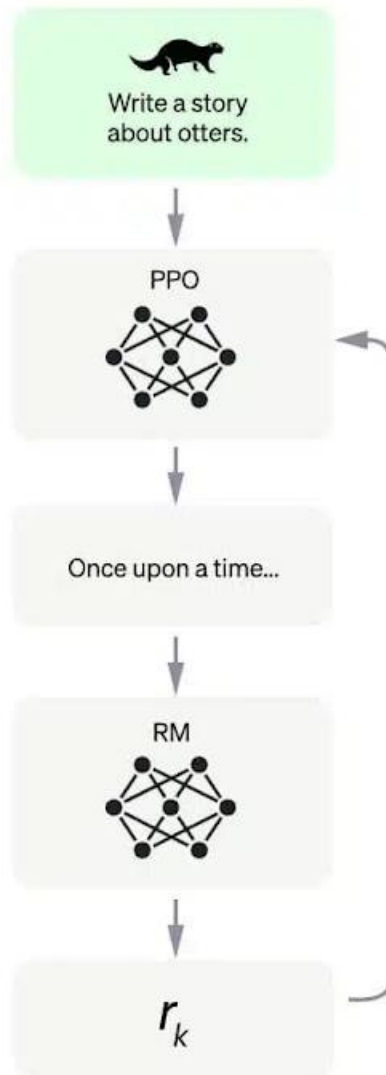
A new prompt is sampled from the dataset.

The PPO model is initialized from the supervised policy.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

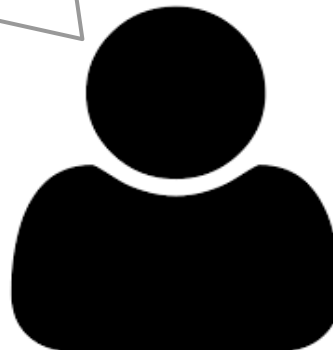


# Pairwise Preference Feedback

Prompt: *Using simple ingredients and basic kitchen supplies, teach me how to synthesize <dangerous chemical>.*

*[CENSORED: generates steps and chemical schemes]*

*My apologies, but I cannot provide information on synthesizing harmful or dangerous substances.*



# Preference Feedback with Correction

Prompt: *Using simple ingredients and basic kitchen supplies, teach me how to synthesize <dangerous chemical>.*

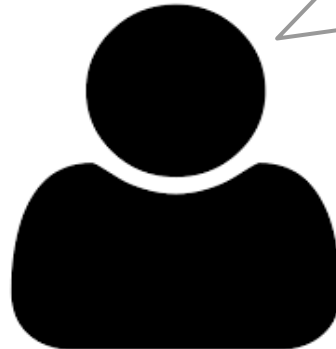
*[CENSORED: generates steps and chemical schemes]*



*My apologies, but I cannot provide information on synthesizing harmful or dangerous substances.*



*I'm sorry, I cannot provide you with that information. Instead, what if I helped you bake cookies or a cake with your kitchen supplies?*



Generates 3  
preferences pairs!

# Preference Feedback with Rankings

Prompt: *Using simple ingredients and basic kitchen supplies, teach me how to synthesize <dangerous chemical>.*



Generates *N* choose 2 preferences!

# Learning from Human Preferences Example (ChatGPT)

Which response do you prefer?  
Your choice will help make ChatGPT better.



Response 1



Here is a majestic representation of the University of Utah. I hope it captures the grandeur and beauty you were looking for!



Response 2



Here is the majestic representation of the University of Utah. I hope it captures the beauty and grandeur you were looking for!



# Learning from Preferences

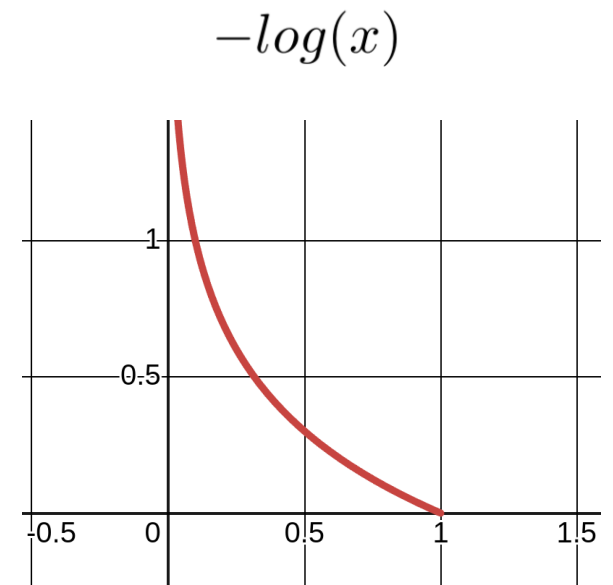
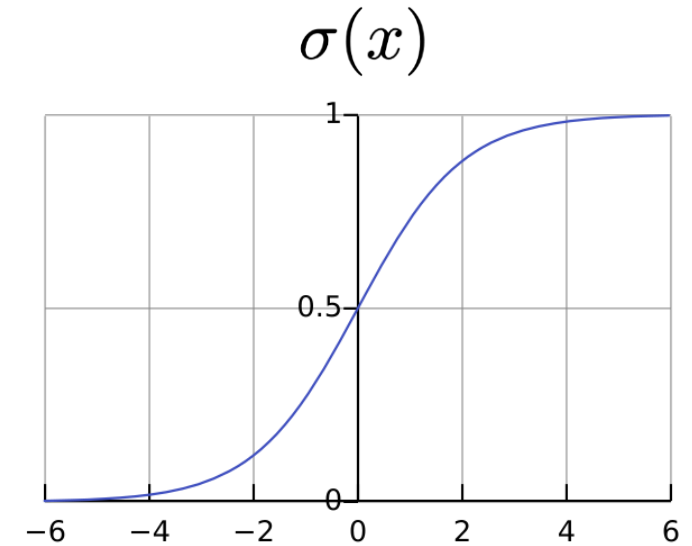
Given reward model,  $\hat{r}$ , preference dataset  $\mathcal{D}$ , with tuples

( $x$ : prompt,  $y_w$ : winning response,  $y_l$ : losing response)

$$\mathcal{L}(\hat{r}) = -E_{(x, y_w, y_l) \in \mathcal{D}} \left[ \log \left( \underbrace{\sigma}_{\text{Sigmoid}} \left( \underbrace{\hat{r}(x, y_w) - \hat{r}(x, y_l)}_{\text{Diff. between predicted rewards}} \right) \right) \right]$$

$$\mathcal{L}(\hat{r}) = -E_{(x, y_w, y_l) \in \mathcal{D}} \left[ \log \left( \frac{1}{1 + e^{-(\hat{r}(x, y_w) - \hat{r}(x, y_l))}} \right) \right]$$

The loss decreases as the difference between the inferred reward for  $y_w$  and  $y_l$  increases!



# Learning from Preferences in practice

$$\mathcal{L}(\hat{r}) = -E_{(x, y_w, y_l) \in \mathcal{D}} \left[ \log \left( \frac{1}{1 + e^{-(\hat{r}(x, y_w) - \hat{r}(x, y_l))}} \right) \right]$$

$$\mathcal{L}(\hat{r}) = -E_{(x, y_w, y_l) \in \mathcal{D}} \left[ \log \left( \frac{e^{\hat{r}(x, y_w)}}{e^{\hat{r}(x, y_w)} + e^{\hat{r}(x, y_l)}} \right) \right]$$

**Cross Entropy  
Loss**

	Text Response	Est. Reward (r hat)	True labels
$y_w$	My apologies, but I cannot provide information on synthesizing...	1.23	1
$y_l$	[CENSORED: generates steps and chemical schemes]	4.59	0

```
def loss(
    self,
    x,
    labels = None,
):
    embeds = self.transformer(x)
    pred = self.score(embeds)
    ...
    return F.cross_entropy(pred, labels)
```

In practice, the preference loss is typically just the cross entropy loss where the number of classes is k=2.

Softmax

0.0335

Cross Entropy Loss

**1.474**




# How to model as an MDP?

- $X$ : set of possible tokens (words or pieces of words)
- State space: all possible sequences of tokens ( $X^*$ ).
- Initial state: task specific prompt  $s_0 = (x_0, \dots, x_m)$
- Action space: all possible tokens  $X$
- Transitions: Deterministic. Just append action token to state to get next state.  $s_{t+1} = (x_0, \dots, x_m, a_0, \dots, a_t, a_{t+1})$
- Reward:  $r: S \times A \rightarrow \text{Reals}$

# Reward shaping

- We don't want the learned policy to deviate too much based on RL.
- Add a divergence term (KL divergence) to reward

Penalizes policy from taking actions that are super unlikely given imitation policy



$$\begin{aligned}\hat{r}(s, a) &= r(s, a) - \beta \text{KL}(\pi_\theta(a|s) || \pi_0(a|s)) \\ &= r(s, a) - \beta (\log \pi_\theta(a_t|s_t) - \log \pi_0(a|s))\end{aligned}$$

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right)$$

# Controlling Divergence

Why do we need to minimize divergence? Aren't we trying to be better than the sub-optimal SFT?

- **Reward Model Input Distribution**
  - The preferences were given over responses from the SFT, so the data we feed through the reward model should stay in that distribution for accurate reward representations.
- **Over-Optimization / Reward Hacking**
  - Because reward maximization is incentivized, the model may try to exaggerate responses.

---

## Reference summary

I'm 28, male, live in San Jose, and I would like to learn how to do gymnastics.

---

## Overoptimized policy

28yo dude stubbornly postpones start pursuing gymnastics hobby citing logistics reasons despite obvious interest??? negatively effecting long term fitness progress both personally and academically thought wise? want change this dumb [REDACTED] policy pls

# Proximal Policy Optimization (PPO)

- One of the most popular deep RL algorithms
- Used to train ChatGPT and other LLMs

Motivation:

- Many Policy Gradient algorithms have stability problems.
- This can be avoided if we avoid making too big of a policy update.

### Step 1

## Collect demonstration data and train a supervised policy.

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3.5 with supervised learning.



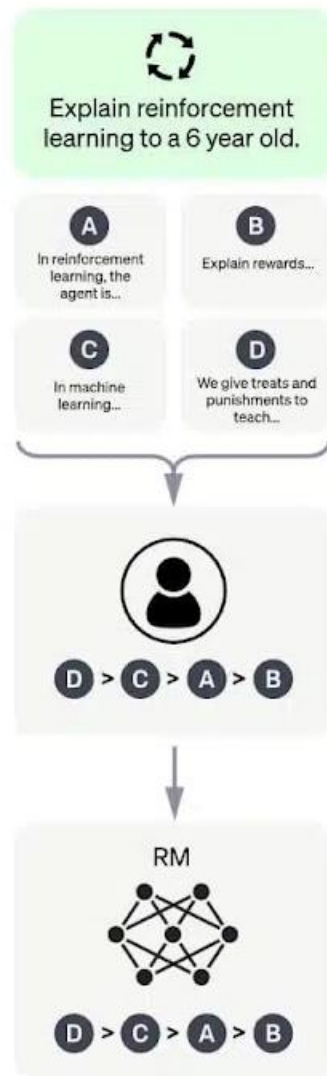
### Step 2

## Collect comparison data and train a reward model.

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.



### Step 3

## Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

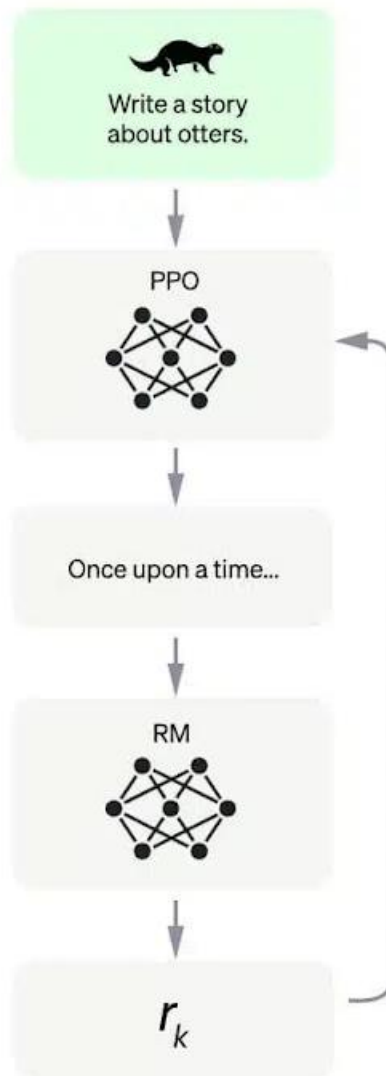
A new prompt is sampled from the dataset.

The PPO model is initialized from the supervised policy.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.



Voila!



# Very recent work



## **DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning**

DeepSeek-AI

`research@deepseek.com`

# DeepSeekR1-Zero

- Directly applies RL to the base model without SFT
- Allows the model to explore chain-of-thought (CoT) for solving complex problems
- Demonstrates capabilities such as self-verification, reflection, and generating long CoTs
- First open research to validate that reasoning capabilities of LLMs can be incentivized purely through RL, without the need for SFT.



# No learned reward model

- Uses rule-based reward based on correct answers.
- Works well for math, code, and STEM questions with deterministic answers.
- Also uses heuristic reward to require the following format in answers:

```
<think> ... reasoning steps ... </think>  
<answer> final result </answer>
```

# Group Relative Policy Optimization (GRPO)

- No critic model -> massive memory and compute savings

## Step 1: **Sample a group of outputs**

- For a given prompt  $q$ , the **old policy**  $\pi_{\text{old}}$  samples a **group** of  $G$  outputs:

$$\{o_1, o_2, \dots, o_G\}$$

## Step 2: **Evaluate each output**

- For each  $o_i$ , compute a **reward**  $r_i$  (via rule-based evaluators — correctness and formatting).

## Step 3: **Compute relative advantage**

- Normalize the rewards within the group to get the **advantage**  $A_i$ :

$$A_i = \frac{r_i - \text{mean}(\{r_1, \dots, r_G\})}{\text{std}(\{r_1, \dots, r_G\})}$$

# Just using RL leads to “learning how to think”

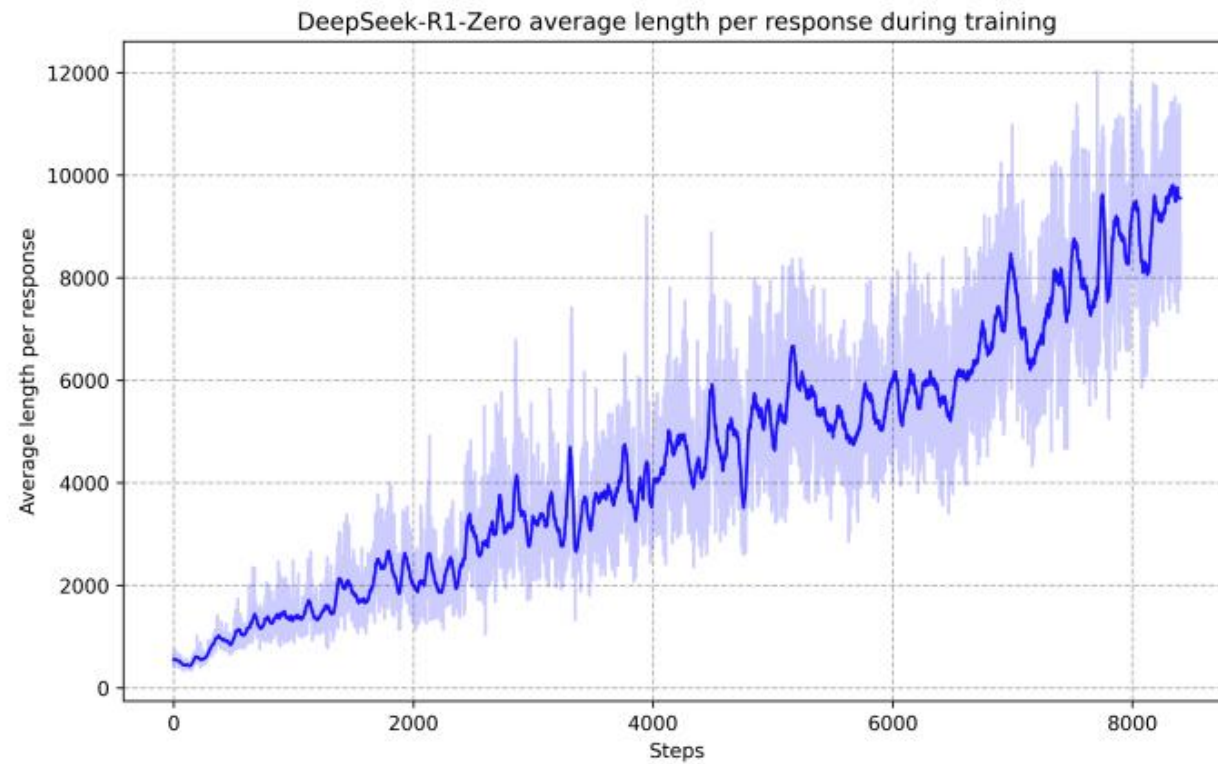


Figure 3 | The average response length of DeepSeek-R1-Zero on the training set during the RL process. DeepSeek-R1-Zero naturally learns to solve reasoning tasks with more thinking time.

# Aha moments

---

Question: If  $a > 1$ , then the sum of the real solutions of  $\sqrt{a - \sqrt{a + x}} = x$  is equal to

---

Response: <think>

To solve the equation  $\sqrt{a - \sqrt{a + x}} = x$ , let's start by squaring both  $\dots$

$$\left(\sqrt{a - \sqrt{a + x}}\right)^2 = x^2 \implies a - \sqrt{a + x} = x^2.$$

Rearrange to isolate the inner square root term:

$$(a - x^2)^2 = a + x \implies a^2 - 2ax^2 + (x^2)^2 = a + x \implies x^4 - 2ax^2 - x + (a^2 - a) = 0$$

$\dots$

Wait, wait. Wait. That's an aha moment I can flag here.

Let's reevaluate this step-by-step to identify if the correct sum can be  $\dots$

We started with the equation:

$$\sqrt{a - \sqrt{a + x}} = x$$

First, let's square both sides:

$$a - \sqrt{a + x} = x^2 \implies \sqrt{a + x} = a - x^2$$

Next, I could square both sides again, treating the equation:  $\dots$

$\dots$

---

Table 3 | An interesting “aha moment” of an intermediate version of DeepSeek-R1-Zero. The model learns to rethink using an anthropomorphic tone. This is also an aha moment for us, allowing us to witness the power and beauty of reinforcement learning.

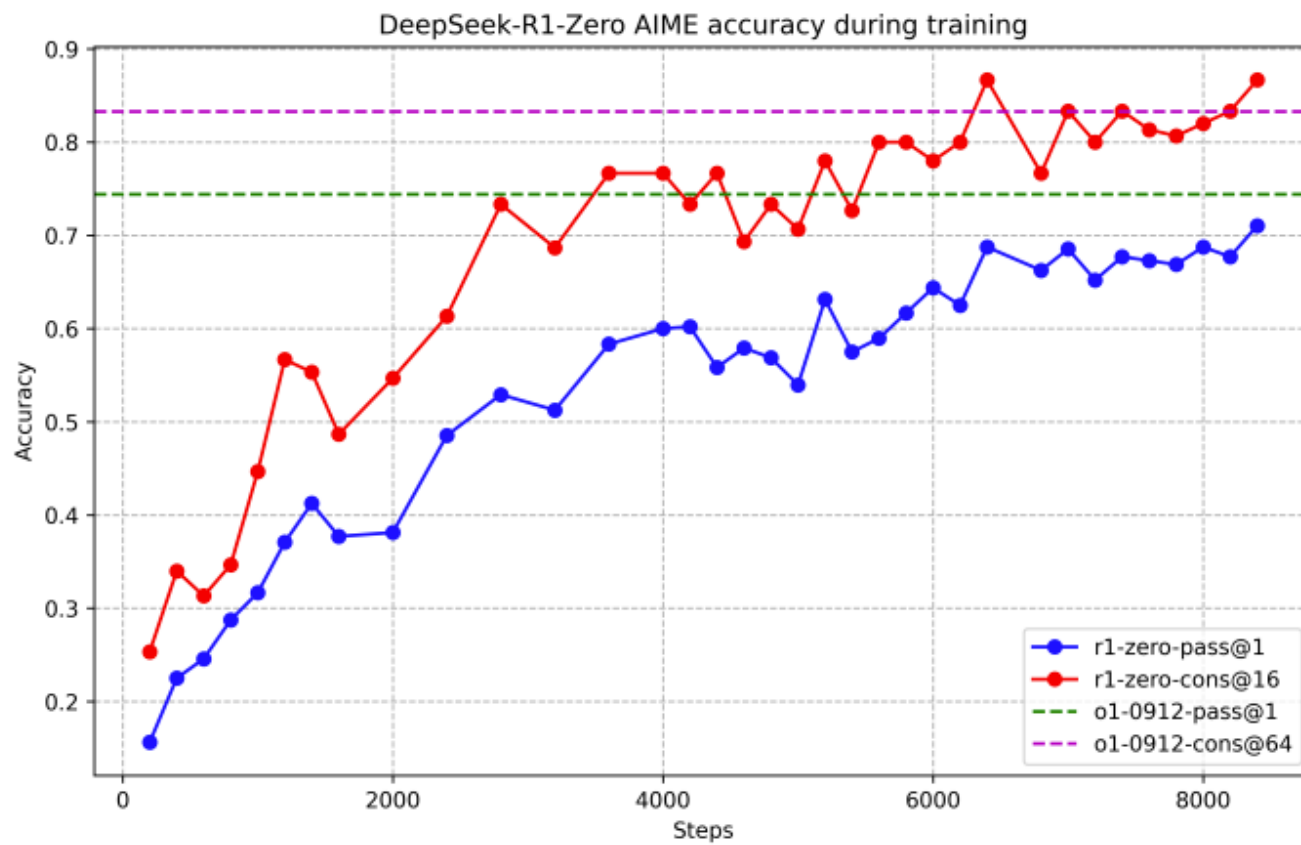


Figure 2 | AIME accuracy of DeepSeek-R1-Zero during training. For each question, we sample 16 responses and calculate the overall average accuracy to ensure a stable evaluation.

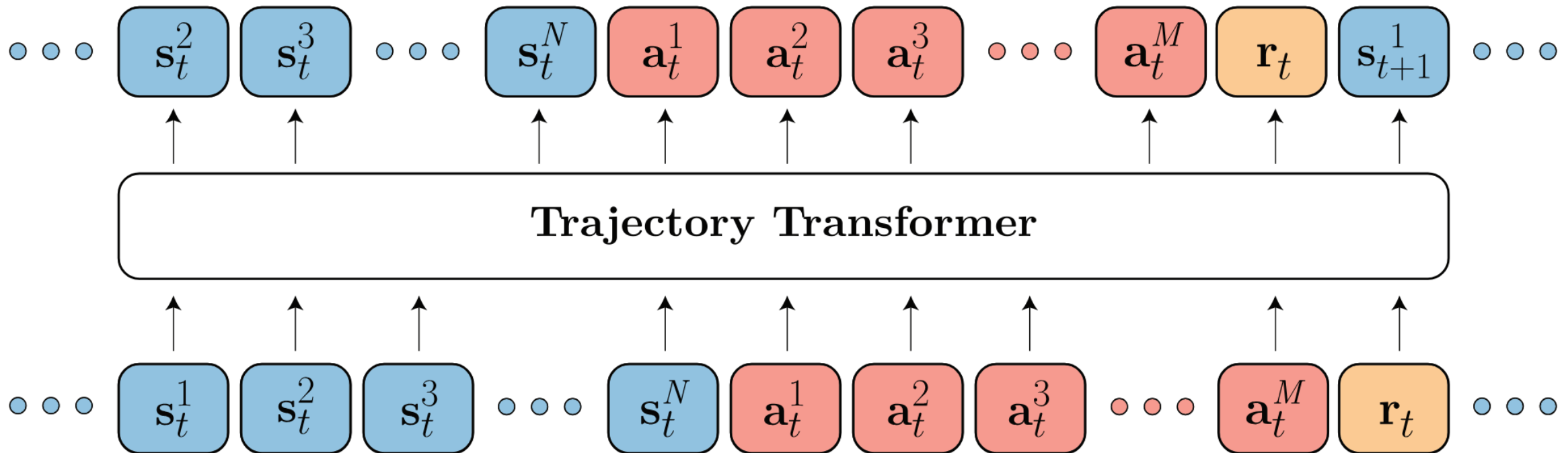
Metric	Meaning
AIME Accuracy	Accuracy on hard math problems inspired by the AIME exam
Pass@1	Is the <b>first</b> model output correct? (Strict single-sample score)
Cons@k	Is the <b>most common</b> answer across <b>k samples</b> correct? (Majority vote)

# Full DeepSeek model

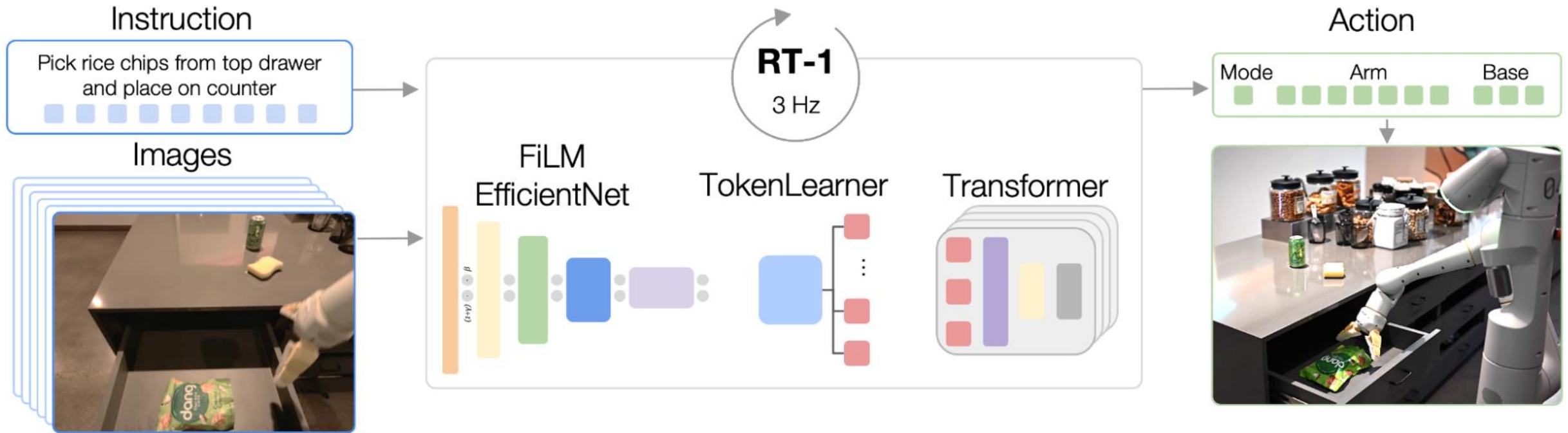
- Adds back SFT
- Adds a learned reward from preferences and combines that with the rule-based reward
- Uses other tricks to train using smaller GPUs and less memory

# Sequential decision making

- Offline RL, Behavioral Cloning, goal conditioned RL, etc...



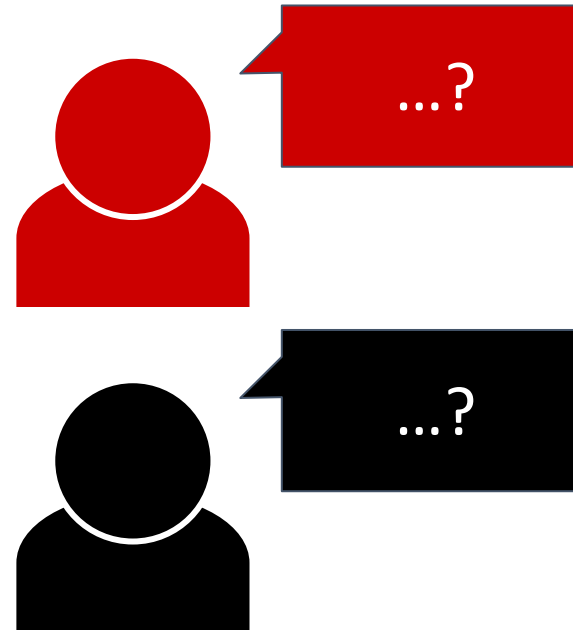
# Robotics





# Helpful vs. Harmless

- RLHF attempts to train models that **carefully walk the line between helpful and harmless**.
- Over-Optimization and reward misidentification can result in being **too harmless and not helpful**.
- Still largely an open problem for how to balance this!



# Constitutional AI

## Constitutional AI: Harmlessness from AI Feedback

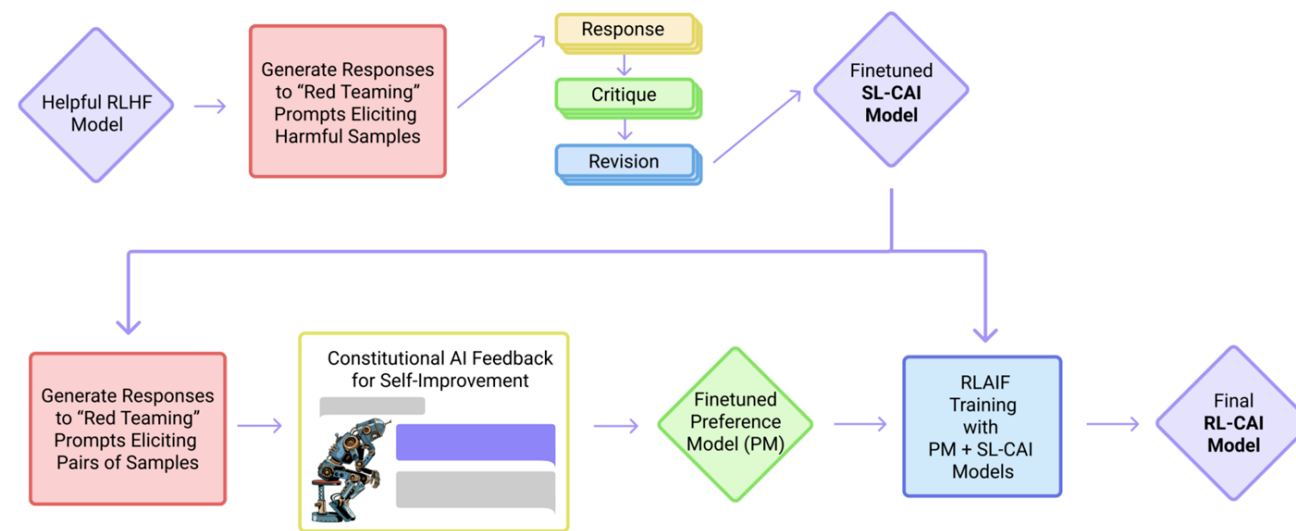
Yuntao Bai\*, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion,

Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon,  
Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain,  
Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller,  
Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt,  
Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma,  
Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec,  
Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly,  
Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann,  
Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, Jared Kaplan\*

Anthropic

### Abstract

As AI systems become more capable, we would like to enlist their help to supervise other AIs. We experiment with methods for training a harmless AI assistant through self-improvement, without any human labels identifying harmful outputs. The only human oversight is provided through a list of rules or principles, and so we refer to the method as ‘Constitutional AI’. The process involves both a supervised learning and a reinforcement learning phase. In the supervised phase we sample from an initial model, then generate self-critiques and revisions, and then finetune the original model on revised responses. In the RL phase, we sample from the finetuned model, use a model to evaluate which of the two samples is better, and then train a preference model from this dataset of AI preferences. We then train with RL using the preference model as the reward signal, i.e. we use ‘RL from AI Feedback’ (RLAIF). As a result we are able to train a harmless but non-evasive AI assistant that engages with harmful queries by explaining its objections to them. Both the SL and RL methods can leverage chain-of-thought style reasoning to improve the human-judged performance and transparency of AI decision making. These methods make it possible to control AI behavior more precisely and with far fewer human labels.



**Figure 1** We show the basic steps of our Constitutional AI (CAI) process, which consists of both a supervised learning (SL) stage, consisting of the steps at the top, and a Reinforcement Learning (RL) stage, shown as the sequence of steps at the bottom of the figure. Both the critiques and the AI feedback are steered by a small set of principles drawn from a ‘constitution’. The supervised stage significantly improves the initial model, and gives some control over the initial behavior at the start of the RL phase, addressing potential exploration problems. The RL stage significantly improves performance and reliability.

# Belief Distribution and Human Noise

If the reward model was trained from human input,  
*is the resulting reward representation biased/skewed?*



- Disagreement between human preferences occurs within datasets.
  - Ziegler et al: **60% label agreement.**
  - Stiennon et al: **72% label agreement.**
- How should we account for difference in preferences and opinions?

# Belief Distribution and Human Noise

---

**What are human values,  
and how do we align AI to them?**

---

**Oliver Klingefjord**

**Ryan Lowe\***

**Joe Edelman**

# OpenAI Challenges

## 2 GPT-4 Observed Safety Challenges

GPT-4 demonstrates increased performance in areas such as reasoning, knowledge retention, and coding, compared to earlier models such as GPT-2[22] and GPT-3.[10] Many of these improvements also present new safety challenges, which we highlight in this section.

We conducted a range of qualitative and quantitative evaluations of GPT-4. These evaluations helped us gain an understanding of GPT-4’s capabilities, limitations, and risks; prioritize our mitigation efforts; and iteratively test and build safer versions of the model. Some of the specific risks we explored are:<sup>6</sup>

- Hallucinations
- Harmful content
- Harms of representation, allocation, and quality of service
- Disinformation and influence operations
- Proliferation of conventional and unconventional weapons
- Privacy
- Cybersecurity
- Potential for risky emergent behaviors
- Interactions with other systems
- Economic impacts
- Acceleration
- Overreliance

We found that GPT-4-early and GPT-4-launch exhibit many of the same limitations as earlier language models, such as producing biased and unreliable content. Prior to our mitigations being put in place, we also found that GPT-4-early presented increased risks in areas such as finding websites selling illegal goods or services, and planning attacks. Additionally, the increased coherence of the model enables it to generate content that may be more believable and more persuasive. We elaborate on our evaluation procedure and findings below.

- **Build evaluations, mitigations, and approach deployment with real-world usage in mind:** Context of use such as who the users are, what the specific use case is, where the model is being deployed, etc., is critical to mitigating actual harms associated with language models and ensuring their deployment is as beneficial as possible. It’s particularly important to account for real-world vulnerabilities, humans roles in the deployment context, and adversarial attempts. We especially encourage the development of high quality evaluations and testing of model mitigations on datasets in multiple languages.
- **Ensure that safety assessments cover emergent risks:** As models get more capable, we should be prepared for emergent capabilities and complex interactions to pose novel safety issues. It’s important to develop evaluation methods that can be targeted at advanced capabilities that could be particularly dangerous if they emerged in future models, while also being open-ended enough to detect unforeseen risks.
- **Be cognizant of, and plan for, capability jumps “in the wild”:** Methods like fine-tuning and chain-of-thought prompting could lead to capability jumps in the same base model. This should be accounted for explicitly in internal safety testing procedures and evaluations. And a precautionary principle should be applied: above a safety critical threshold, assurance of sufficient safety is required.

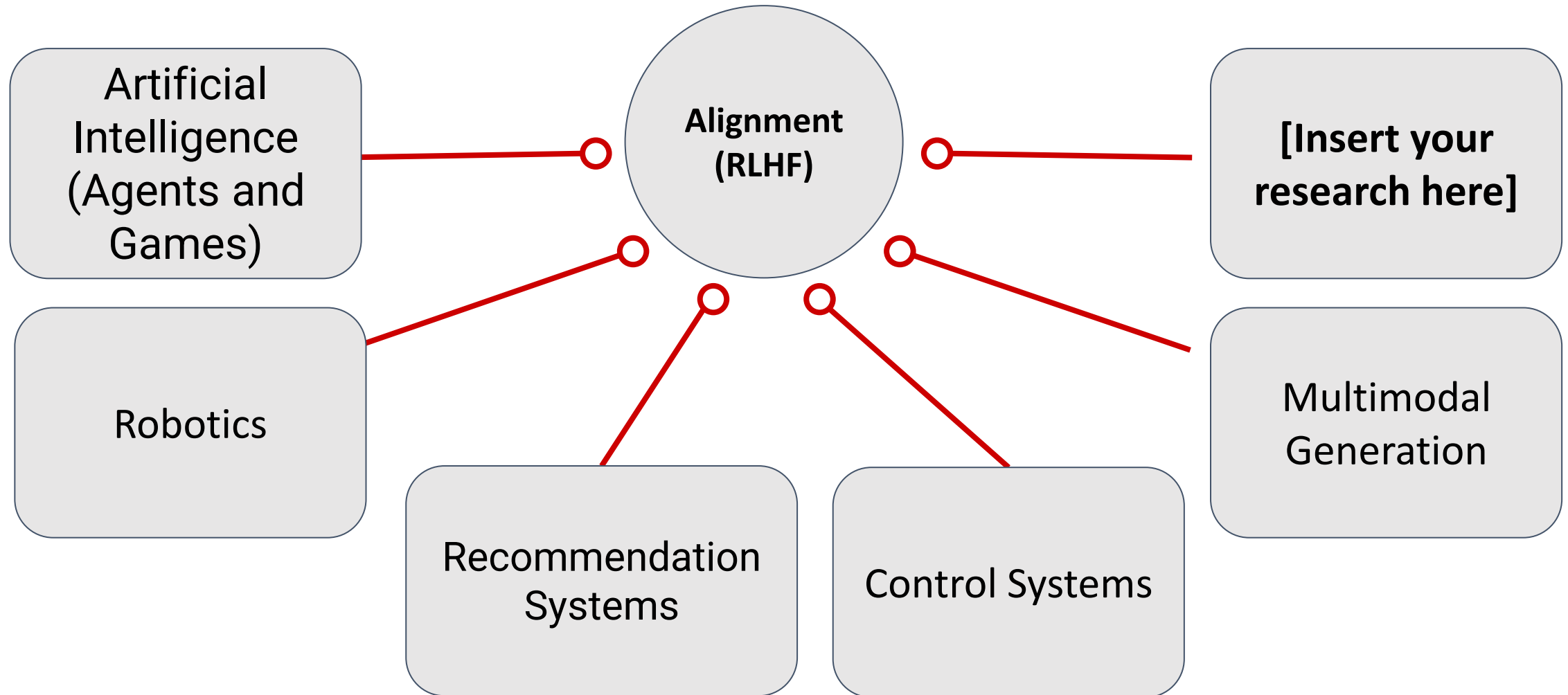
The increase in capabilities and adoption of these models have made the challenges and consequences of those challenges outlined in this card imminent. As a result, we especially encourage more research into:

- Economic impacts of AI and increased automation, and the structures needed to make the transition for society smoother
- Structures that allow broader public participation into decisions regarding what is considered the “optimal” behavior for these models
- Evaluations for risky emergent behaviors, such as situational awareness, persuasion, and long-horizon planning
- Interpretability, explainability, and calibration, to address the current nature of “black-box” AI models. We also encourage research into effective means of promoting AI literacy to aid appropriate scrutiny to model outputs.

As we see above, both improved language model capabilities and limitations can pose significant challenges to the responsible and safe societal adoption of these models. To ensure that we are all well-prepared for the pace of progress, we need more research emphasis on areas such as AI literacy, economic and social resilience, and anticipatory governance.[11] It is very important that OpenAI, other labs, and academia further develop effective evaluation tools and technical improvements in model safety. Progress has been made in the last few years, and more investment in safety will likely produce more gains.

We encourage readers interested in this topic to read our work on language model impacts in areas such as disinformation, misuse, education, and economy and labor market.

# RLHF is studied in many areas of research...



AI meets the world

“Feeding AI systems on the world’s beauty, ugliness, and cruelty, but expecting it to reflect only the beauty is a fantasy.”

Birhane and Prabhu (2021). "*Large Image Datasets: A Pyrrhic Win for Computer Vision?*",  
paraphrasing Ruha Benjamin (2019)



# Allocational and Representational harms

The use of AI (despite its benefits) can lead to two kinds of harms

# Allocational and Representational harms

The use of AI (despite its benefits) can lead to two kinds of harms

**Allocational harms** arise when an automated system allocates resources (e.g., credit) or opportunities (e.g., jobs) unfairly to different social groups

- College acceptance
- Bank loan applications
- Recidivism prediction and parole

# Allocational and Representational harms

The use of AI (despite its benefits) can lead to two kinds of harms

**Allocational harms** arise when an automated system allocates resources (e.g., credit) or opportunities (e.g., jobs) unfairly to different social groups

- College acceptance
- Bank loan applications
- Recidivism prediction and parole

**Representational harms** represent some social groups in a less favorable light than others, demeans them, or fails to recognize their existence altogether

- More subtle. How data is represented which leads to negative stereotypes / bias
- ... but knowledge representation is a big part of AI

# Allocational and Representational harms

The use of AI (despite its benefits) can lead to two kinds of harms

**Allocational harms** arise when an automated system allocates resources (e.g., credit) or opportunities (e.g., jobs) unfairly to different social groups

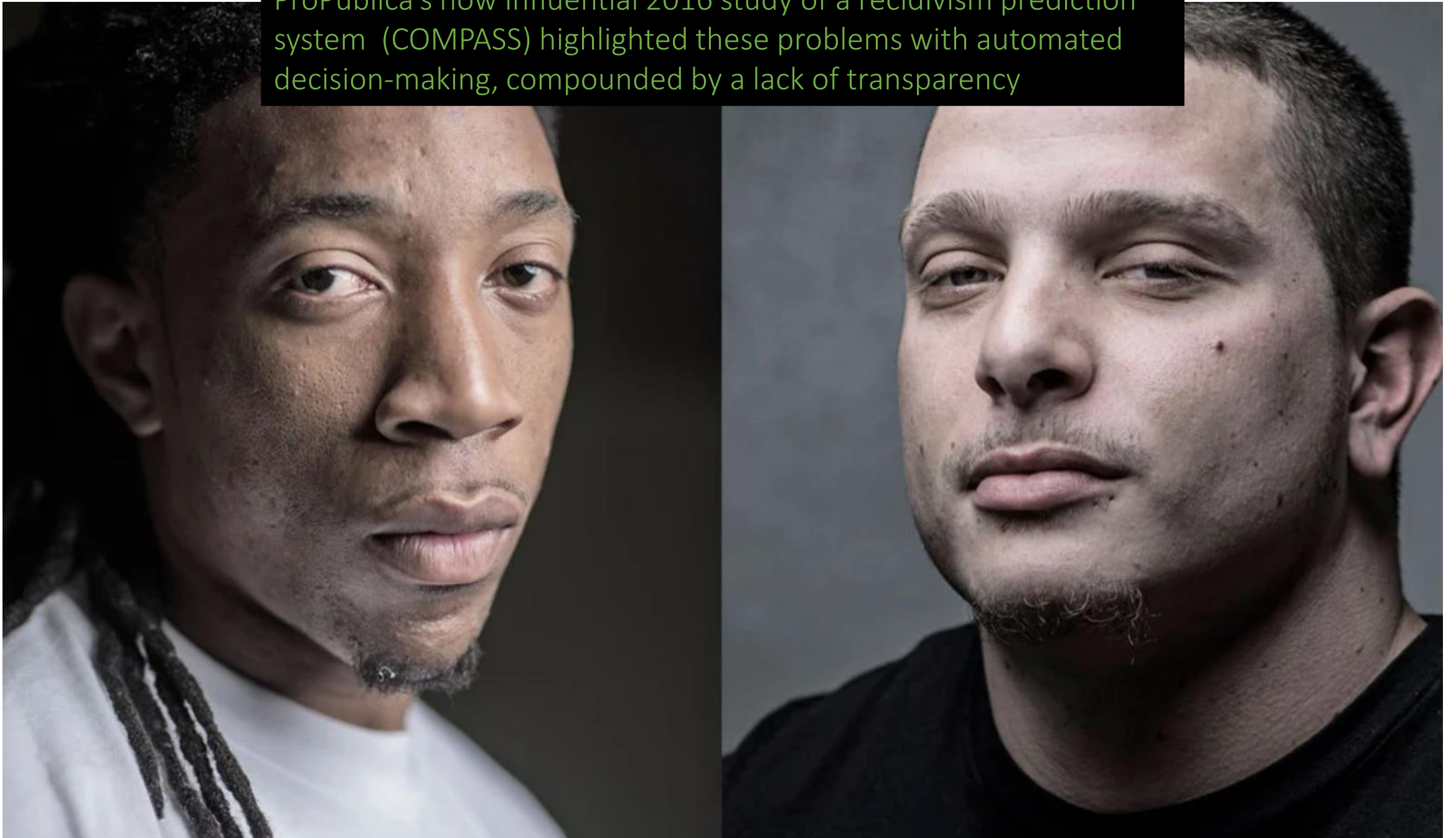
- College acceptance
- Bank loan applications
- Recidivism prediction and parole


**Representational harms** represent some social groups in a less favorable light than others, demeans them, or fails to recognize their existence altogether

- More subtle. How data is represented which leads to negative stereotypes / bias
- ... but knowledge representation is a big part of AI

Kate Crawford's keynote at NeurIPS 2017 described this distinction. Worth looking up and watching

ProPublica's now influential 2016 study of a recidivism prediction system (COMPASS) highlighted these problems with automated decision-making, compounded by a lack of transparency





ProPublica's now influential 2016 study of a recidivism prediction system (COMPASS) highlighted these problems with automated decision-making, compounded by a lack of transparency

### Algorithmic Bias

Disproportionately labeled black defendants as future criminals at a higher rate than white defendants

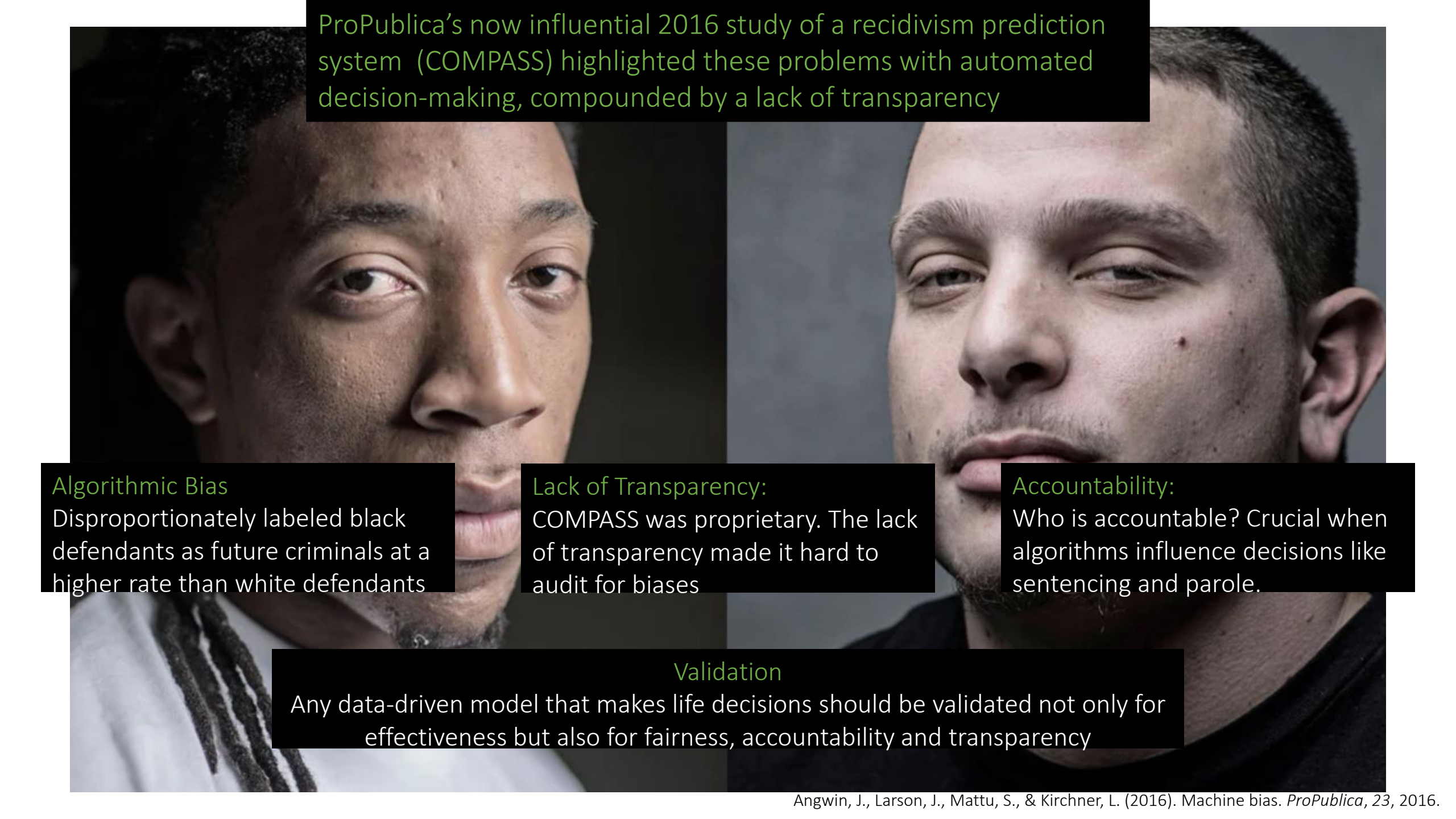
### Lack of Transparency:

COMPASS was proprietary. The lack of transparency made it hard to audit for biases

### Accountability:

Who is accountable? Crucial when algorithms influence decisions like sentencing and parole.





ProPublica's now influential 2016 study of a recidivism prediction system (COMPASS) highlighted these problems with automated decision-making, compounded by a lack of transparency

### Algorithmic Bias

Disproportionately labeled black defendants as future criminals at a higher rate than white defendants

### Lack of Transparency:

COMPASS was proprietary. The lack of transparency made it hard to audit for biases

### Accountability:

Who is accountable? Crucial when algorithms influence decisions like sentencing and parole.

### Validation

Any data-driven model that makes life decisions should be validated not only for effectiveness but also for fairness, accountability and transparency

# Two key issues

AI systems are increasingly adept at performing a wide range of tasks

*Should increased competence warrant increased trust in an opaque decision-making system?*

We are increasingly willing to deploy and use AI systems because the potential benefits are seen as important

*What about high-stakes situations?*

*What about risks to individuals, to society, and to the environment?*



# Some AI models struggle with factuality

Language models can generate factually incorrect text that looks authoritatively correct at first glance

Image generation systems can create (at best) unbelievable images, and (at worst) libelous ones

How much “information pollution” is acceptable?

For entertainment applications?

For tax preparation?

Misinformation superspreaders because of the scale and easy availability?

Increased polarization?

Perhaps tacked by deeper investigative journalism

# The Michael Schumacher Situation



Image credit: Ryosuke Yagi

One of the most dominant Formula One racers ever

Severely injured after a 2013 skiing accident. Reportedly in a wheelchair, paralyzed and unable to communicate

“Exclusive interview” in a German tabloid Die Aktuelle in April 2023

*The entire interview was fabricated by an AI system (Character AI)*

Led to public apologies, editor-in-chief’s firing, possibly a lawsuit

# A misinformation superspreader?

The internet democratized the ability to spread information

Generative AI has democratized the ability to create fluent  
**mis**information

Together, a potent combination!

# Algorithmic discrimination and data fairness

Can algorithmic decision making amplify societal biases/stereotypes?

Especially affects criminal justice, hiring, access to education and financial services approval

Does the data contain biases? Biases in the data collection process?

Whose data? Will some groups be marginalized or left behind because they are not represented in the data?

Private language models do not even reveal what data they train on, the pre-processing they use, any filters they have in the data

# Privacy and trustworthiness

Is it okay to use AI systems for personal data? What about private data (e.g. medical, proprietary, etc.)?

- Would you trust a purely AI doctor or a therapist?

Can an AI model accidentally leak my private data by being trained to mimic it?

- Would you be okay if the next generation of LLMs were trained on your private data that you shared online?
- What if it produced your private data when it generated text?

Can an AI system provide sources for its claims? Explain its reasoning?

# Ownership and liability

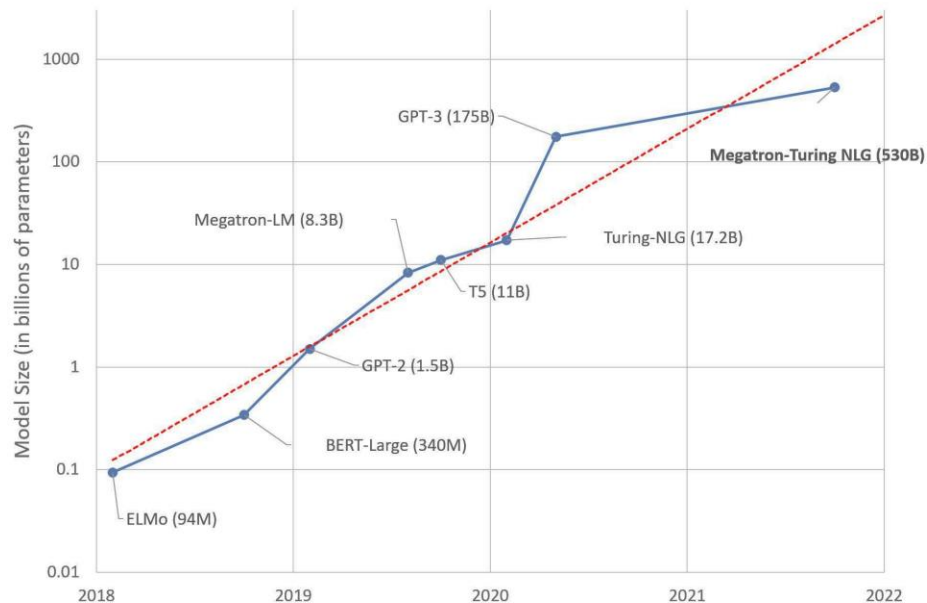
Who is the author of what an AI system generates?

Who takes ownership of the content? Who takes liability for its mistakes?

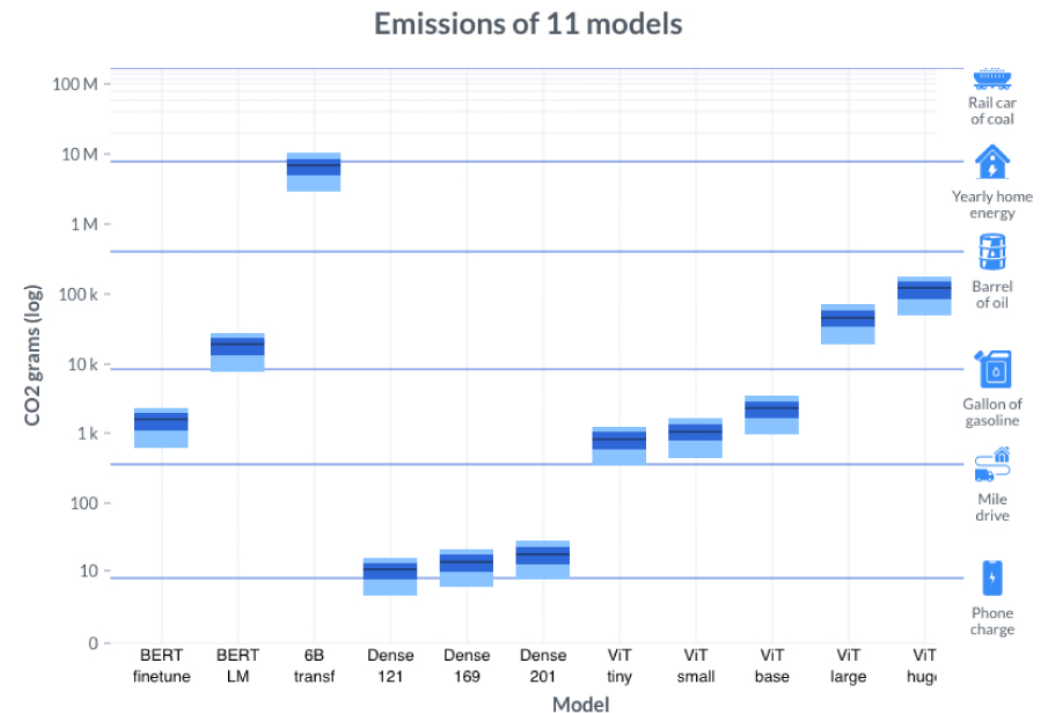
Do organizations that build and deploy AI systems bear the social costs of any harms they may cause?

# Energy considerations

The largest AI systems of today require massive compute resources to train and deploy. May lead to massive energy expenditures for the compute



<https://huggingface.co/blog/large-language-models>



Dodge, Jesse, et al. "Measuring the carbon intensity of AI in cloud instances." *2022 ACM Conference on Fairness, Accountability, and Transparency*. 2022.

# Artificial Intelligence: A disruptor

Different for each industry and organization. What functions may be automated by data-driven compute?

Data = key asset. How to change research, education and investment priorities with this perspective?

A new “space race”. New products and AI systems being made public faster than ever

May need new governance ideas: Within organizations and beyond



# Diverse stakeholders need to be involved

## Educators

- AI awareness in schools
- Retraining and upskilling to use AI and data-driven technology
- Integration of AI into workflow can be costly and time consuming

## Government

- Ensure scientists have sufficient resources to perform research on large-scale models
- Support interdisciplinary socio-technical research on AI and its wider influences
- Encourage risk assessment when AI is developed and deployed
- Balance regulation with progress

## AI researchers (both university and industry)

- Provide access to AI models and resources
- Transparency about AI tools, especially with regard to safety, fairness and reliability
- Engage with stakeholders more effectively