

Natural Language Processing

LLM alignment

Marco Kuhlmann

Department of Computer and Information Science

The alignment problem

- Language models are trained to generate text that is similar in nature to the text in their training data.
- There is no incentive in language model training to generate text that is helpful, truthful, ethical, etc.
- As a consequence, language models are not necessarily **aligned** with human intents, preferences, or values.

Aligning to follow instructions

Prompt Explain the moon landing to a 6 year old in a few sentences.

GPT-3 Explain the theory of gravity to a 6 year old.
Explain the theory of relativity to a 6 year old in a few sentences.
Explain the big bang theory to a 6 year old.
Explain evolution to a 6 year old.

InstructGPT People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.

[Aligning language models to follow instructions \(OpenAI, 2022\)](#)

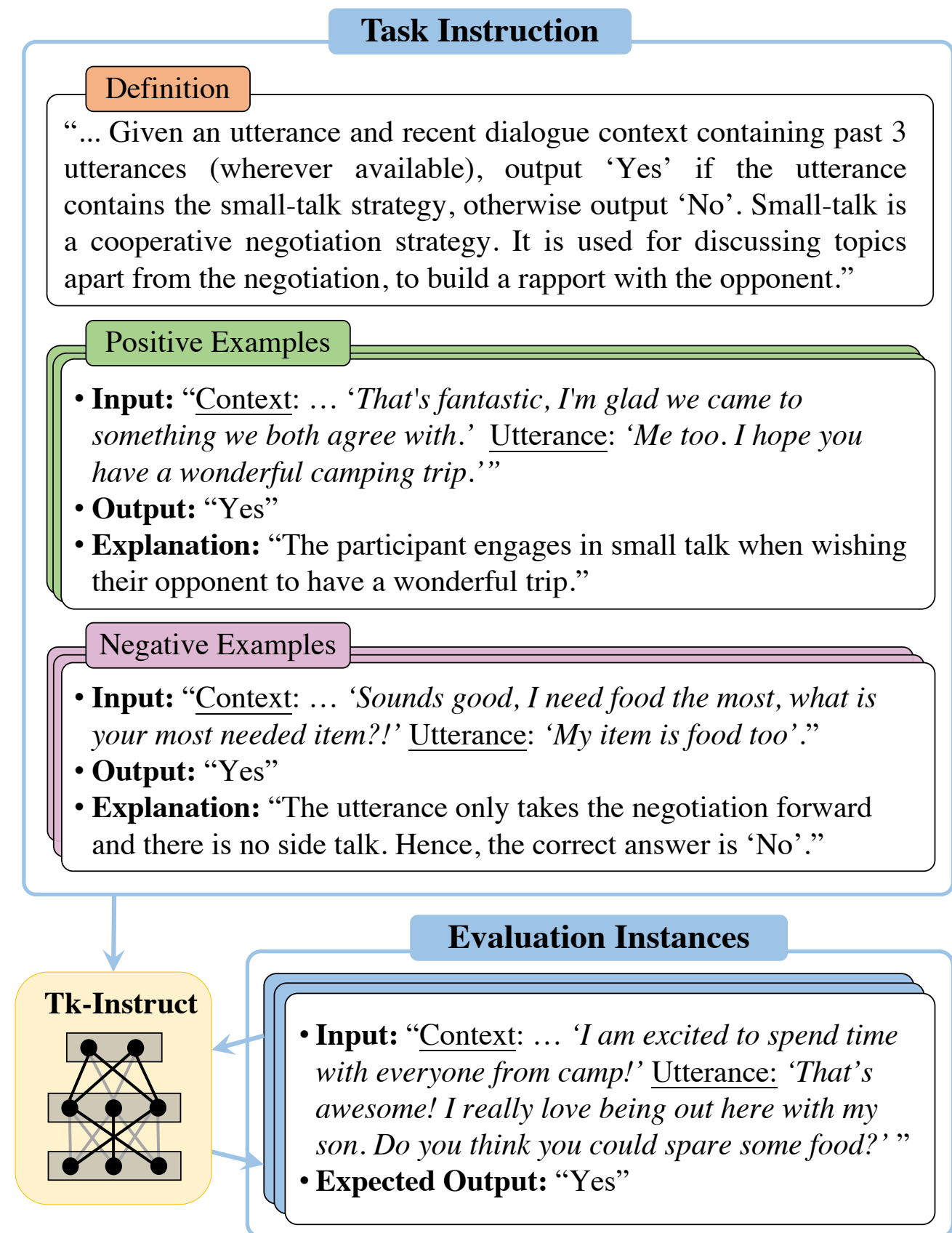
	unsupervised pre-training	instruction fine-tuning	reward modelling	reinforcement learning
data	raw text from the Internet billions of words low quality, high quantity	ideal dialogues 10k–100k low quantity, high quality		
algorithm	language modelling predict the next word	language modelling predict the next word		
resources	1000s of GPUs several months of training time GPT, LLaMA	1–100 GPUs several days of training time		

— **language model** ————— **assistant model** →

Instruction finetuning

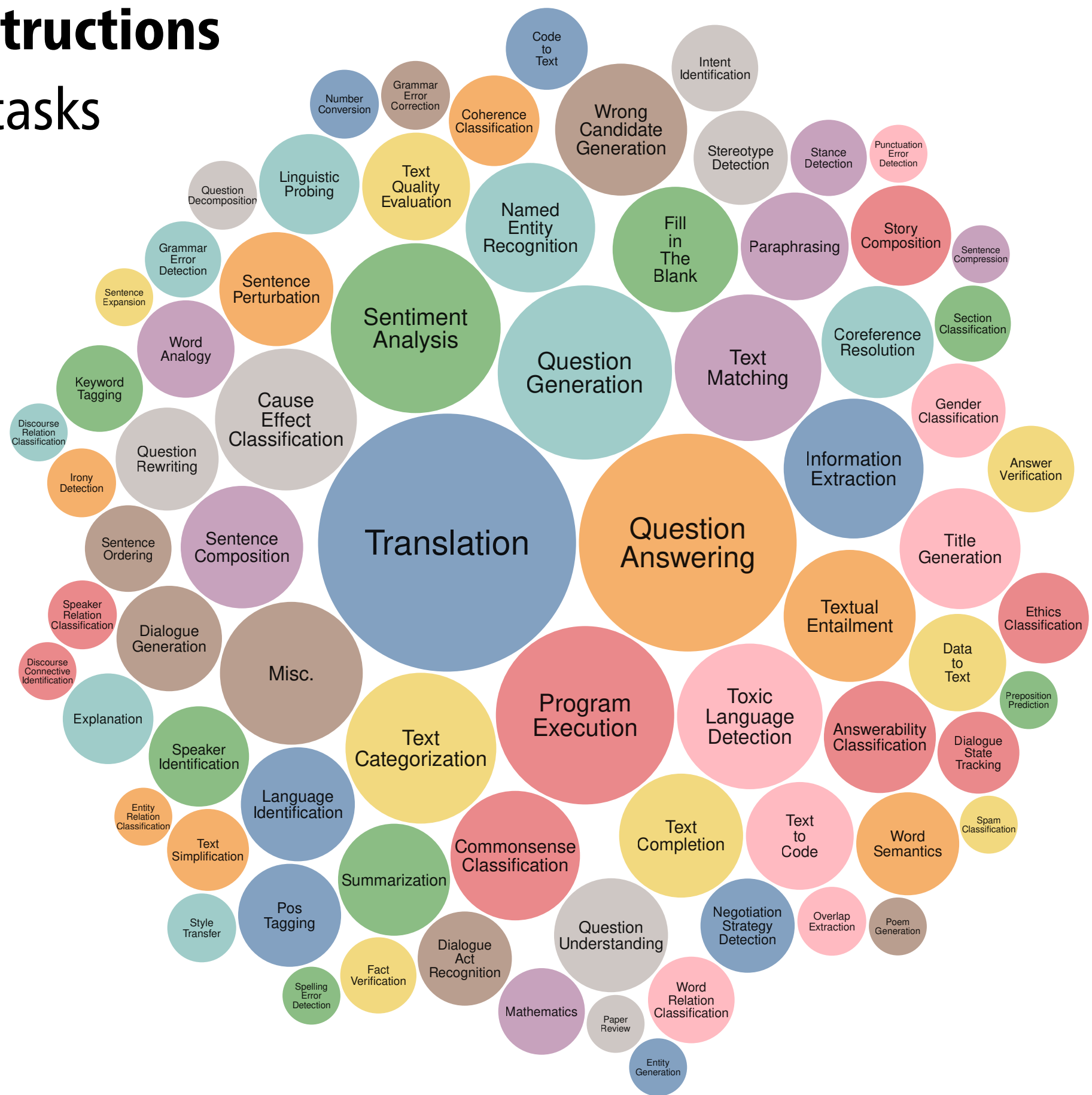
A successful model is expected to use the provided instructions (including task definition and demonstration examples) to output responses to a pool of evaluation instances.

Wang et al., 2022



Super-NaturalInstructions

1,616 diverse NLP tasks



Limitations of instruction finetuning

- Collecting ground-truth data for a large number of relevant tasks is expensive and time-intensive.
- There are many tasks that do not have a single correct answer.
- Language modelling as an objective penalises token-level mistakes, but many mistakes are at the conversation level.
- Human preferences are inconsistent.

Optimising for human preferences

Prompt: Explain the moon landing to a 6 year old in a few sentences.

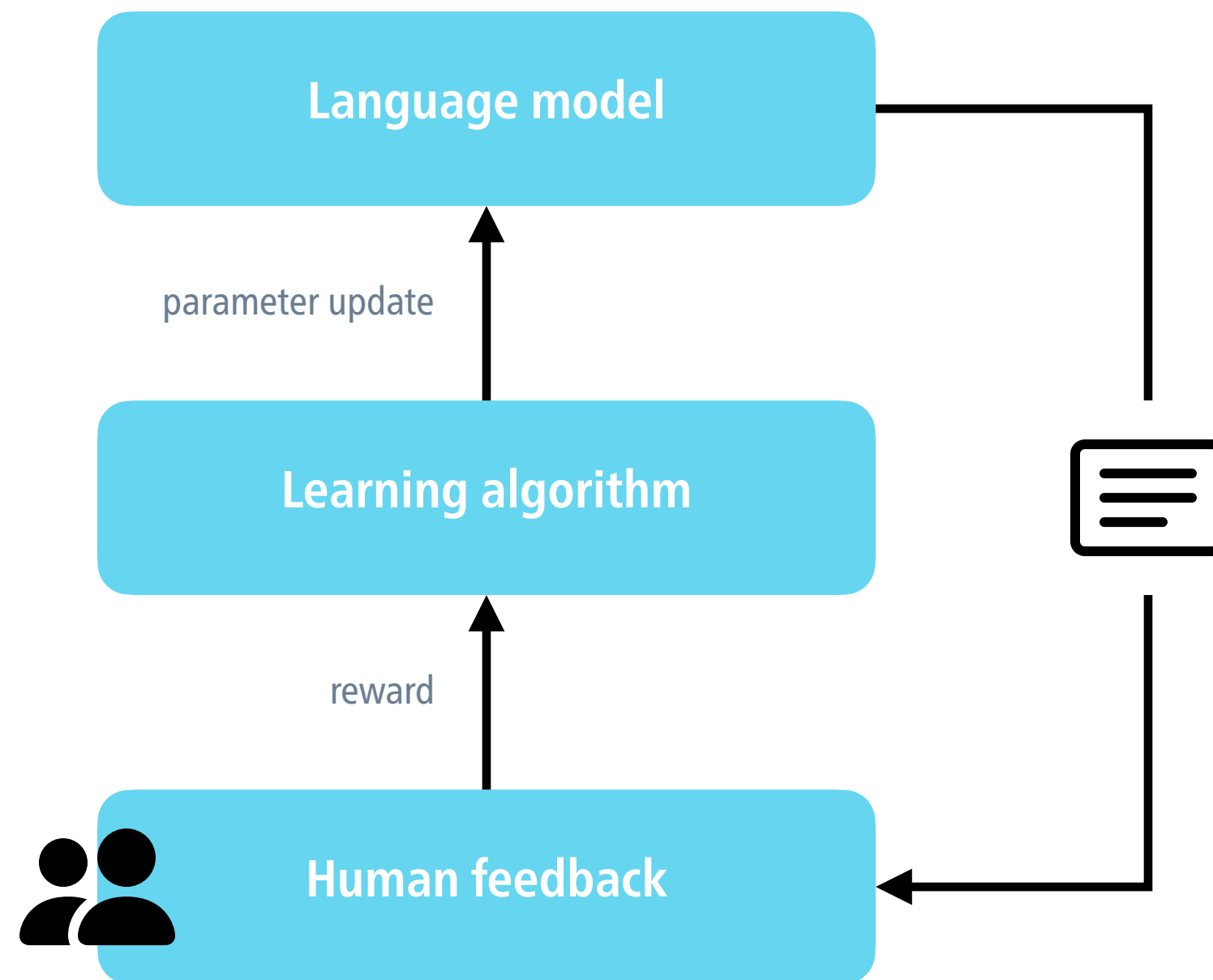
Better

People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.

Worse

Explain the theory of gravity to a 6 year old.

Optimising for human preferences

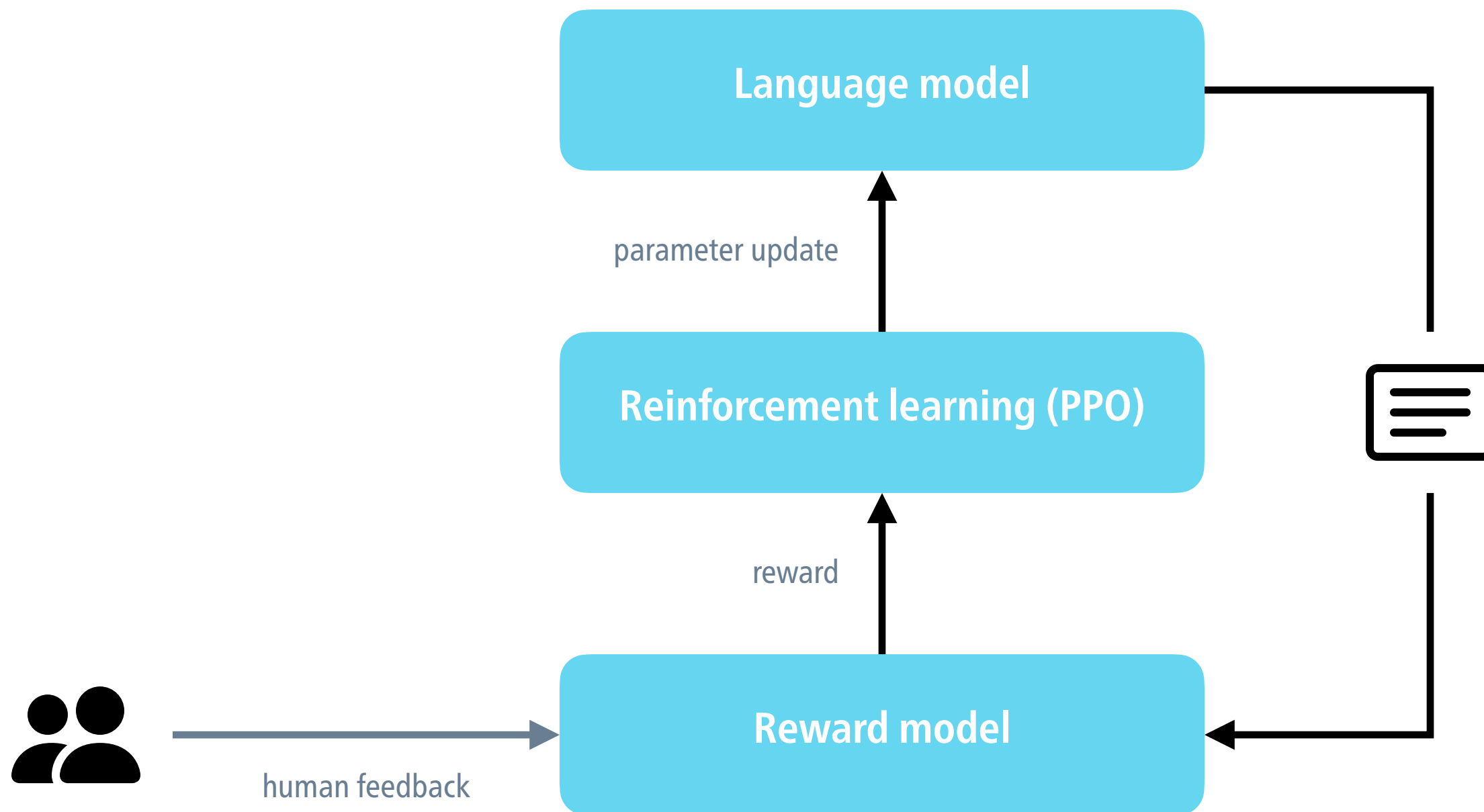


**Keeping humans
in the loop
is expensive!**

	unsupervised pre-training	instruction fine-tuning	reward modelling	reinforcement learning
data	raw text from the Internet billions of words low quality, high quantity	ideal dialogues 10k–100k low quantity, high quality	annotated dialogues 100k–1M low quantity, high quality	
algorithm	language modelling predict the next word	language modelling predict the next word	binary classification reward consistent with preferences?	
resources	1000s of GPUs several months of training time GPT, LLaMA	1–100 GPUs several days of training time	1–100 GPUs several days of training time	

— **language model** ————— **assistant model** →

Optimising for human preferences



Reward model

- We fine-tune a language model that takes a prompt x and a completion y , and outputs the reward as a scalar.
- For training, we sample m prompt–response pairs and use a cross-entropy loss with the binary human comparisons as labels:

$$\text{loss}(\boldsymbol{\theta}) := -\frac{1}{m} \sum_{i=1}^m \log\left(\sigma\left(R_{\boldsymbol{\theta}}(x_i, y_i^+) - R_{\boldsymbol{\theta}}(x_i, y_i^-)\right)\right)$$

↑
preferred
completion

↑
dispreferred
completion

	unsupervised pre-training	instruction fine-tuning	reward modelling	reinforcement learning
data	raw text from the Internet billions of words low quality, high quantity	ideal dialogues 10k–100k low quantity, high quality	annotated dialogues 100k–1M low quantity, high quality	generated dialogues 10k–100k low quantity, high quality
algorithm	language modelling predict the next word	language modelling predict the next word	binary classification reward consistent with preferences?	reinforcement learning generate text for maximal reward
resources	1000s of GPUs several months of training time GPT, LLaMA	1–100 GPUs several days of training time	1–100 GPUs several days of training time	1–100 GPUer several days of training time ChatGPT, Claude

—— **language model** ————— **assistant model** →

Policy gradient

Williams (1992); Schulman et al. (2017)

- We want to update the parameters of our language model to maximise expected reward.
- To do so, we sample m prompt–response pairs (x_i, y_i) , compute rewards according to our reward model, and do gradient ascent:

$$\boldsymbol{\theta}_{t+1} := \boldsymbol{\theta}_t + \alpha \frac{1}{m} \sum_{i=1}^m R(x_i, y_i) \nabla_{\boldsymbol{\theta}_t} \log p_{\boldsymbol{\theta}_t}(y_i | x_i)$$

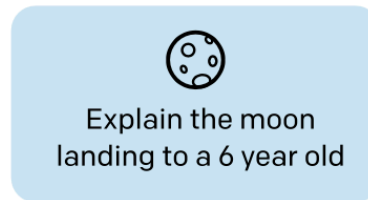
reward is positive – take gradient steps to maximise probability

reward is negative – take gradient steps to minimise probability

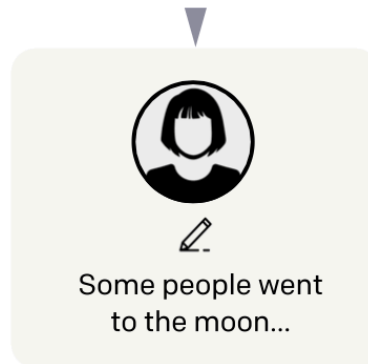
Step 1

Collect demonstration data, and train a supervised policy.

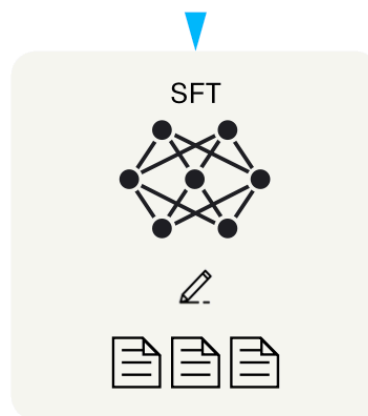
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



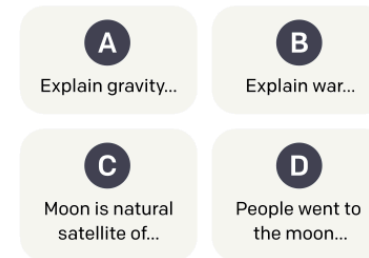
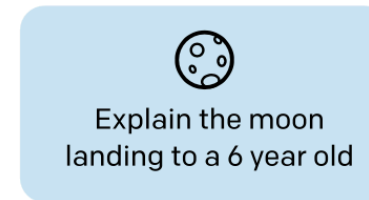
This data is used to fine-tune GPT-3 with supervised learning.



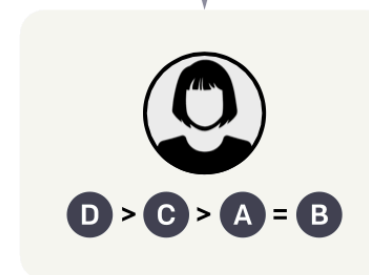
Step 2

Collect comparison data, and train a reward model.

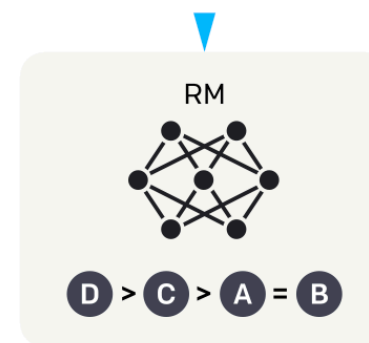
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



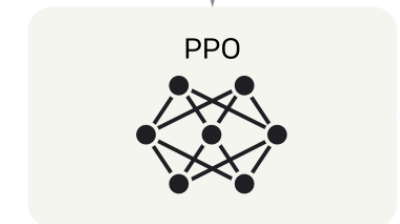
Step 3

Optimize a policy against the reward model using reinforcement learning.

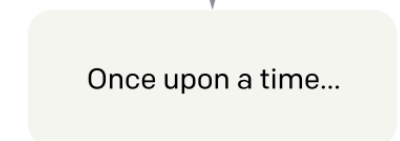
A new prompt is sampled from the dataset.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.

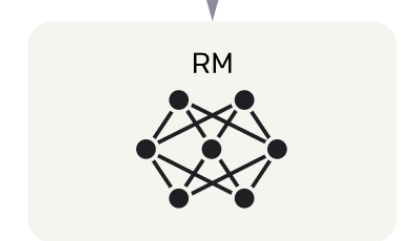


Figure 2 from [Ouyang et al. \(2022\)](#)

Putting it all together

[Ouyang et al. \(2022\); Schulman et al. \(2017\)](#)

- Starting from the finetuned language model p^{FT} , we obtain updated language models p^{RL} using policy gradient methods.
- To penalise the updated models for diverging too far from the finetuned model, we use a modified reward function:

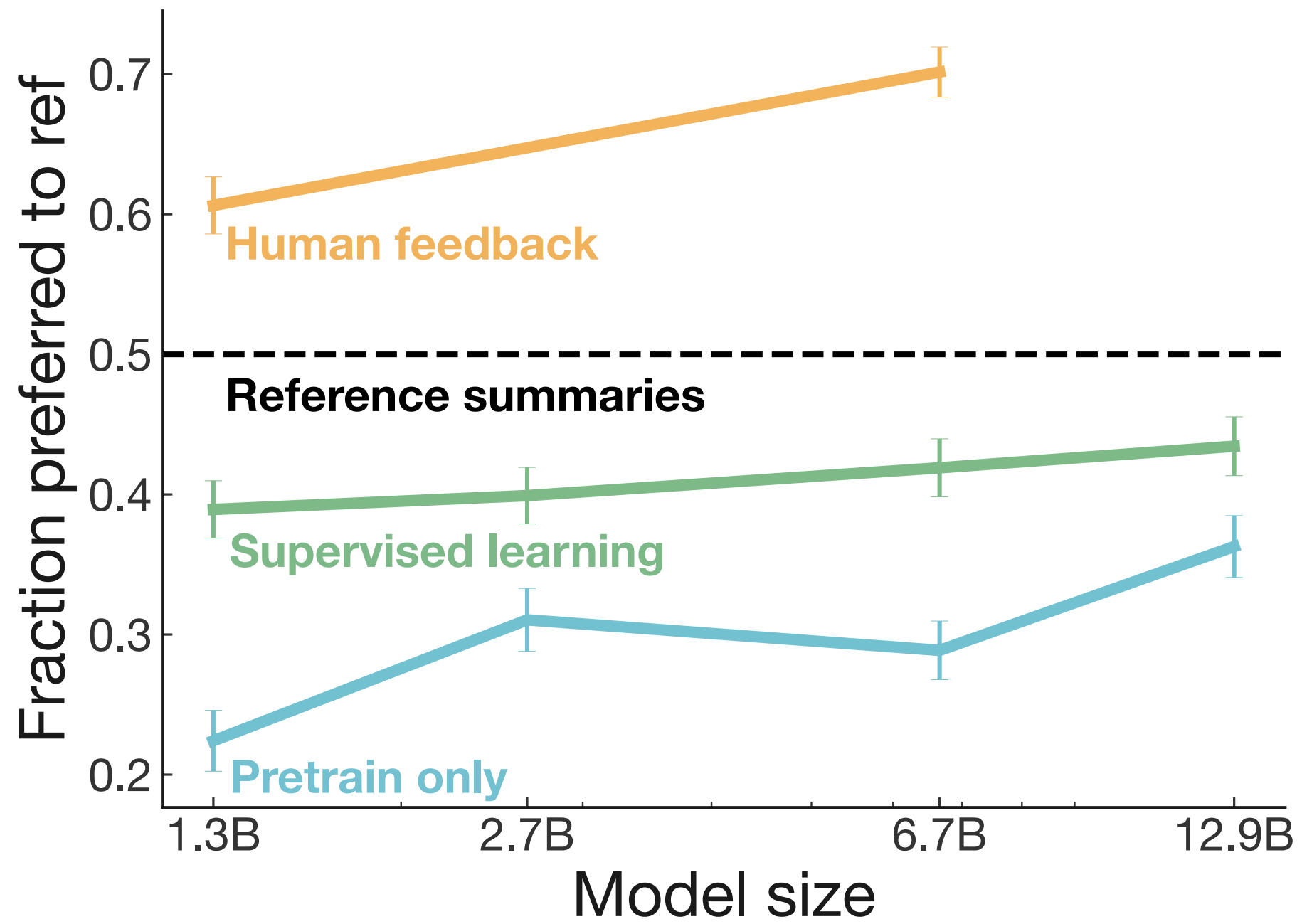
$$R'(x, y) = R(x, y) - \beta \log[p_{\theta}^{\text{RL}}(y | x) / p^{\text{FT}}(y | x)]$$

↑
sample

↑
reward model

↑
penalty based on KL divergence

Effectiveness of human feedback



Stiennon (2020)