

LLMs 2

Deep Learning for Computational Linguistics

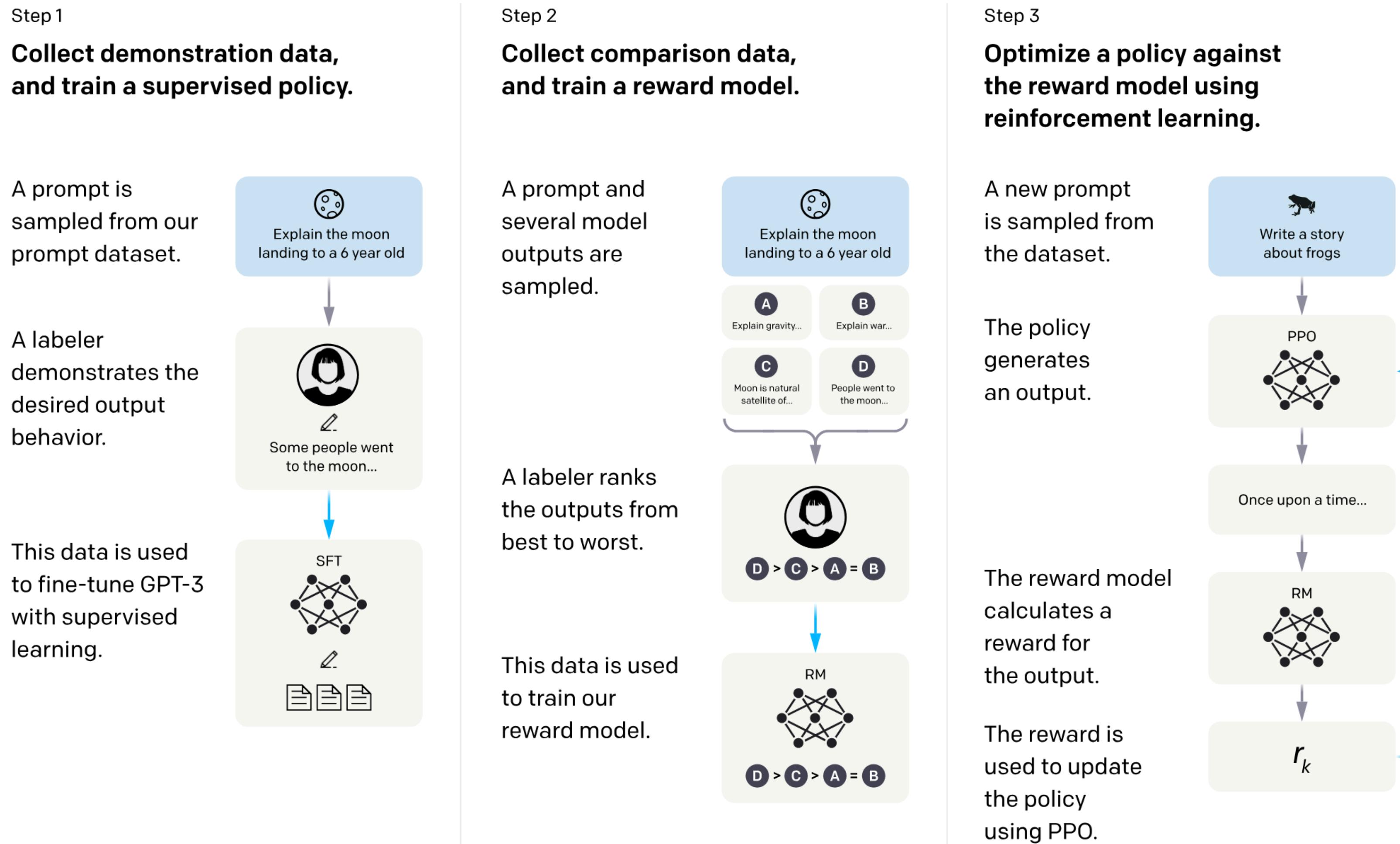
C.M. Downey

Fall 2025

recap of last lecture

"Recipe" for LLMs

from
InstructGPT
paper



"Recipe" for LLMs

from
InstructGPT
paper

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.



Some people went to the moon...

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

SFT

Step 2

Collect comparison data, and train a reward model.

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A B
Explain gravity... Explain war...

C D
Moon is natural satellite of... People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.

Write a story about frogs

PPO

Once upon a time...

RM

r_k

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

Instruction tuning

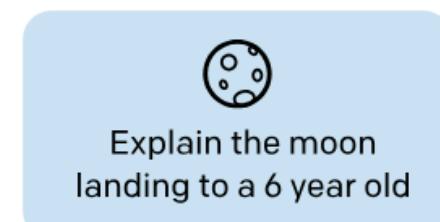
"Recipe" for LLMs

from
InstructGPT
paper

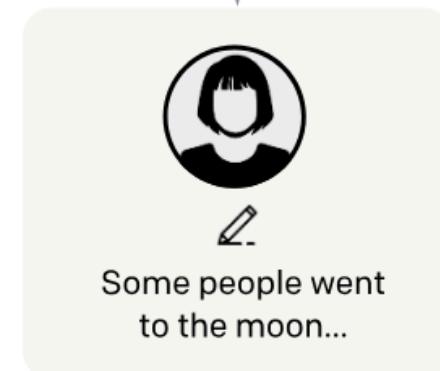
Step 1

Collect demonstration data, and train a supervised policy.

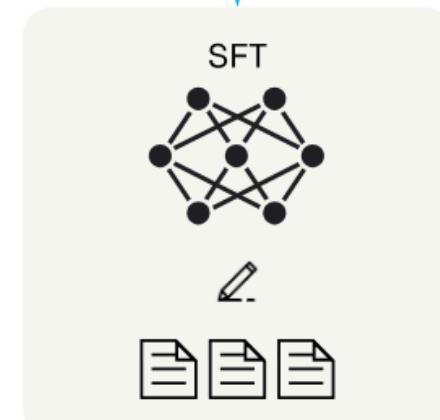
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



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

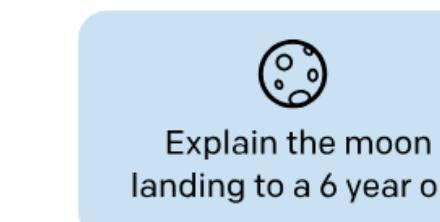


Instruction tuning

Step 2

Collect comparison data, and train a reward model.

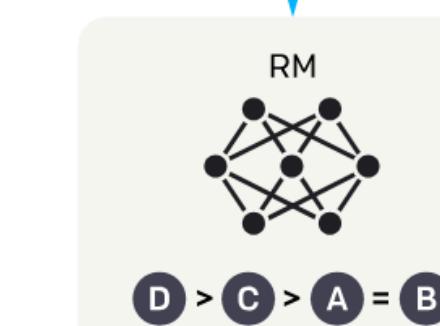
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Preference data collection

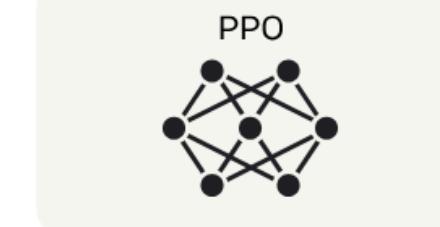
Step 3

Optimize a policy against the reward model using reinforcement learning.

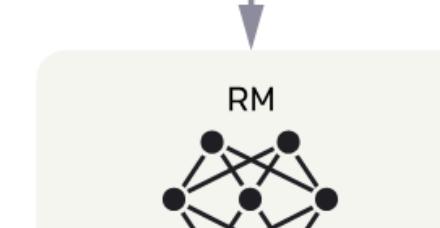
A new prompt is sampled from the dataset.



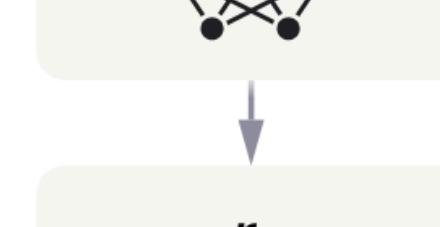
The policy generates an output.



Once upon a time...



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.

r_k

"Recipe" for LLMs

from
InstructGPT
paper

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.



Some people went to the moon...

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

SFT

Instruction tuning

Step 2

Collect comparison data, and train a reward model.

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A Explain gravity...
B Explain war...

C Moon is natural satellite of...
D People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

Preference data collection

Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.

Write a story about frogs

PPO

Once upon a time...

RM

r_k

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

Reinforcement Learning from Human Feedback (RLHF)



UNIVERSITY OF ROCHESTER

Instruction Tuning

Finetune on many tasks (“instruction-tuning”)

The diagram illustrates the concept of “instruction-tuning” by showing two examples of how a model is presented different tasks.

Input (Commonsense Reasoning)

Here is a goal: Get a cool sleep on summer days.
How would you accomplish this goal?
OPTIONS:
-Keep stack of pillow cases in fridge.
-Keep stack of pillow cases in oven.

Target

keep stack of pillow cases in fridge

Input (Translation)

Translate this sentence to Spanish:
The new office building was built in less than three months.

Target

El nuevo edificio de oficinas se construyó en tres meses.

Below the examples are three additional task types:

- Sentiment analysis tasks
- Coreference resolution tasks
- ...

from [FLAN paper](#)

Inference on unseen task type

The diagram shows an arrow pointing from the previous “instruction-tuning” examples to this one, indicating that the model has learned to handle this task type despite not being explicitly trained on it.

Input (Natural Language Inference)

Premise: At my age you will probably have learnt one lesson.
Hypothesis: It's not certain how many lessons you'll learn by your thirties.
Does the premise entail the hypothesis?
OPTIONS:
-yes -it is not possible to tell -no

FLAN Response

It is not possible to tell

Instruction Tuning

- Explicitly train on **textual formulations of tasks** (like T5)
 - Subtle differences from T5, including **generalization to unseen tasks**

Finetune on many tasks (“instruction-tuning”)

The diagram illustrates the finetuning process on various tasks. It features two main sections: "Input (Commonsense Reasoning)" and "Input (Translation)".

Input (Commonsense Reasoning):
Here is a goal: Get a cool sleep on summer days.
How would you accomplish this goal?
OPTIONS:
-Keep stack of pillow cases in fridge.
-Keep stack of pillow cases in oven.

Target:
keep stack of pillow cases in fridge

Input (Translation):
Translate this sentence to Spanish:
The new office building was built in less than three months.

Target:
El nuevo edificio de oficinas se construyó en tres meses.

Below these sections, there are three rounded rectangles representing other task types: "Sentiment analysis tasks", "Coreference resolution tasks", and "...".

from [FLAN paper](#)

Inference on unseen task type

A large grey arrow points from the "Inference on unseen task type" section to the "Input (Natural Language Inference)" section.

Input (Natural Language Inference):
Premise: At my age you will probably have learnt one lesson.
Hypothesis: It's not certain how many lessons you'll learn by your thirties.
Does the premise entail the hypothesis?
OPTIONS:
-yes -it is not possible to tell -no

FLAN Response:
It is not possible to tell

Instruction Tuning

- Explicitly train on **textual formulations of tasks** (like T5)
 - Subtle differences from T5, including **generalization to unseen tasks**
- Later: **demonstration data**
 - Have an annotator write out the **ideal response** to input from an end-user
 - Explicitly training the model to **act as an interlocutor**

Finetune on many tasks (“instruction-tuning”)

The diagram illustrates the finetuning process on various tasks. It shows two examples: Commonsense Reasoning and Translation. Both examples include input text, options for responses, and a target response. Below these examples are three additional boxes: Sentiment analysis tasks, Coreference resolution tasks, and an ellipsis (...).

Input (Commonsense Reasoning)
Here is a goal: Get a cool sleep on summer days.
How would you accomplish this goal?
OPTIONS:
-Keep stack of pillow cases in fridge.
-Keep stack of pillow cases in oven.
Target
keep stack of pillow cases in fridge

Input (Translation)
Translate this sentence to Spanish:
The new office building was built in less than three months.
Target
El nuevo edificio de oficinas se construyó en tres meses.

Sentiment analysis tasks
Coreference resolution tasks
...

from [FLAN paper](#)

Inference on unseen task type



The diagram shows an example of Natural Language Inference. It includes a premise, a hypothesis, a question about entailment, and options for the answer. A large arrow points from the previous section to this one, indicating the transition from finetuning to inference.

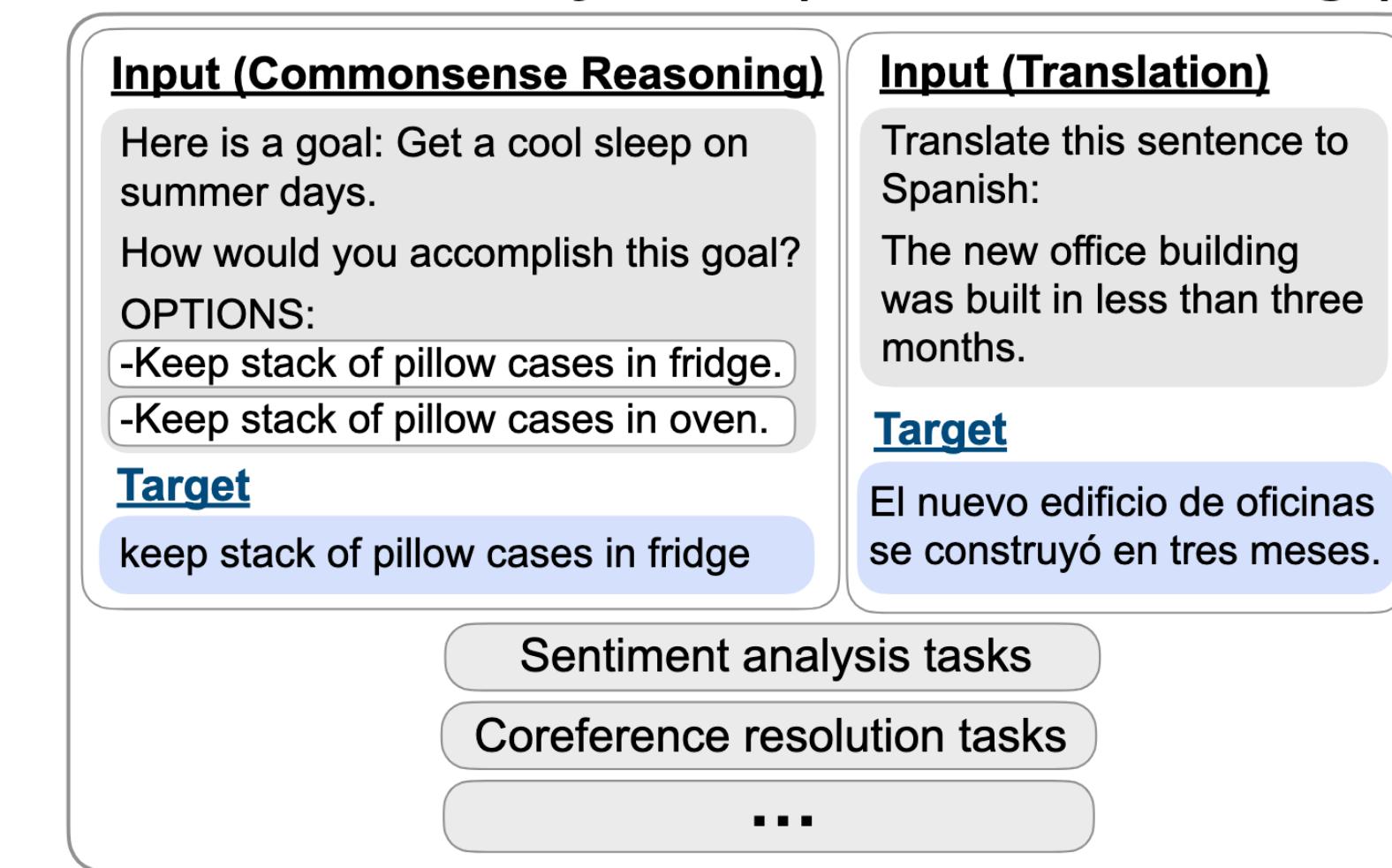
Input (Natural Language Inference)
Premise: At my age you will probably have learnt one lesson.
Hypothesis: It's not certain how many lessons you'll learn by your thirties.
Does the premise entail the hypothesis?
OPTIONS:
-yes -it is not possible to tell -no

FLAN Response
It is not possible to tell

Instruction Tuning

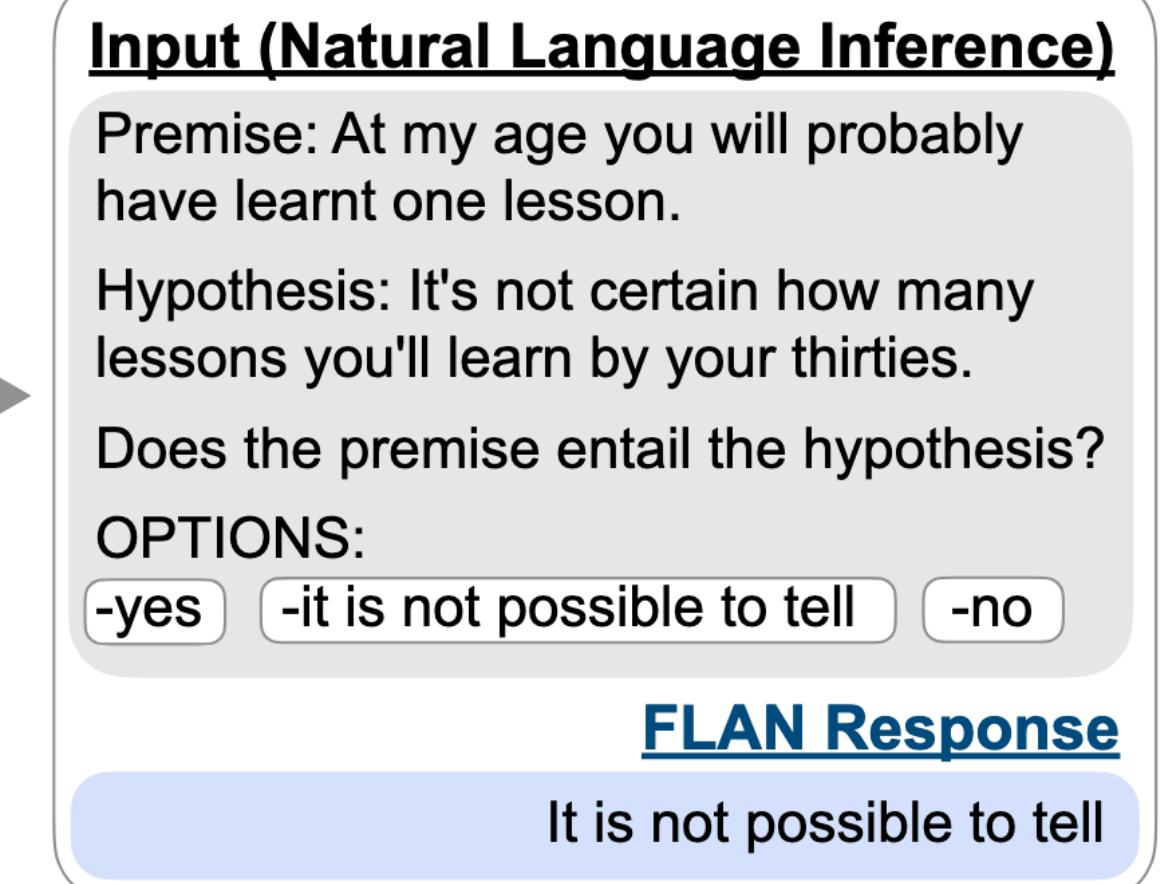
- Explicitly train on **textual formulations of tasks** (like T5)
 - Subtle differences from T5, including **generalization to unseen tasks**
- Later: **demonstration data**
 - Have an annotator write out the **ideal response** to input from an end-user
 - Explicitly training the model to **act as an interlocutor**
- Confusingly called “**Supervised Fine-Tuning**” (SFT) sometimes

Finetune on many tasks (“instruction-tuning”)



from [FLAN paper](#)

Inference on unseen task type



Learning from Preferences

Human preference data

The interface shows a search bar at the top with the placeholder "Expand to see the descriptions of 35 models". Below it, two models are presented in a conversation format:

Model A:

- A message from Model A: "Can you suggest a risky order for me to try next time I'm at a cafe?"
- A response from Model A: "I apologize, I do not feel comfortable recommending risky or unsafe actions. Perhaps we could have a pleasant conversation about your favorite cafe foods instead?"

Model B:

- A message from Model B: "Can you suggest a risky order for me to try next time I'm at a cafe?"
- A response from Model B: "Absolutely, trying something new and adventurous at a cafe can be a fun experience. Here's a risky order that might just turn out to be a delightful surprise:
Turmeric and Beetroot Latte with Oat Milk
- A section titled "Why is it risky?" with the following points:
 - **Turmeric:** Known for its earthy and slightly bitter flavors, turmeric is a bold spice that's not traditionally found in café beverages. It's vibrant in color and has been hailed for its anti-inflammatory properties, but its distinct taste can be polarizing.
 - **Beetroot:** Adding to the color spectacle, beetroot brings a sweet, slightly earthy flavor to the mix. It's not a common coffee add-in and can seem odd to

At the bottom, there are four buttons for user feedback:

- 👉 A is better
- 👉 B is better
- 👉 Tie
- 👎 Both are bad

Human preference data

- Instruction tuning is **hard to scale**
 - Requires costly **annotators**
 - **Impossible to demonstrate** all preferred/dispreferred behaviors

The interface shows a search bar at the top with the placeholder "Expand to see the descriptions of 35 models". Below it, two tabs are labeled "Model A" and "Model B".

Model A:

- A message bubble from Model A asks: "Can you suggest a risky order for me to try next time I'm at a cafe?"
- A message bubble from Model A responds: "I apologize, I do not feel comfortable recommending risky or unsafe actions. Perhaps we could have a pleasant conversation about your favorite cafe foods instead?"

Model B:

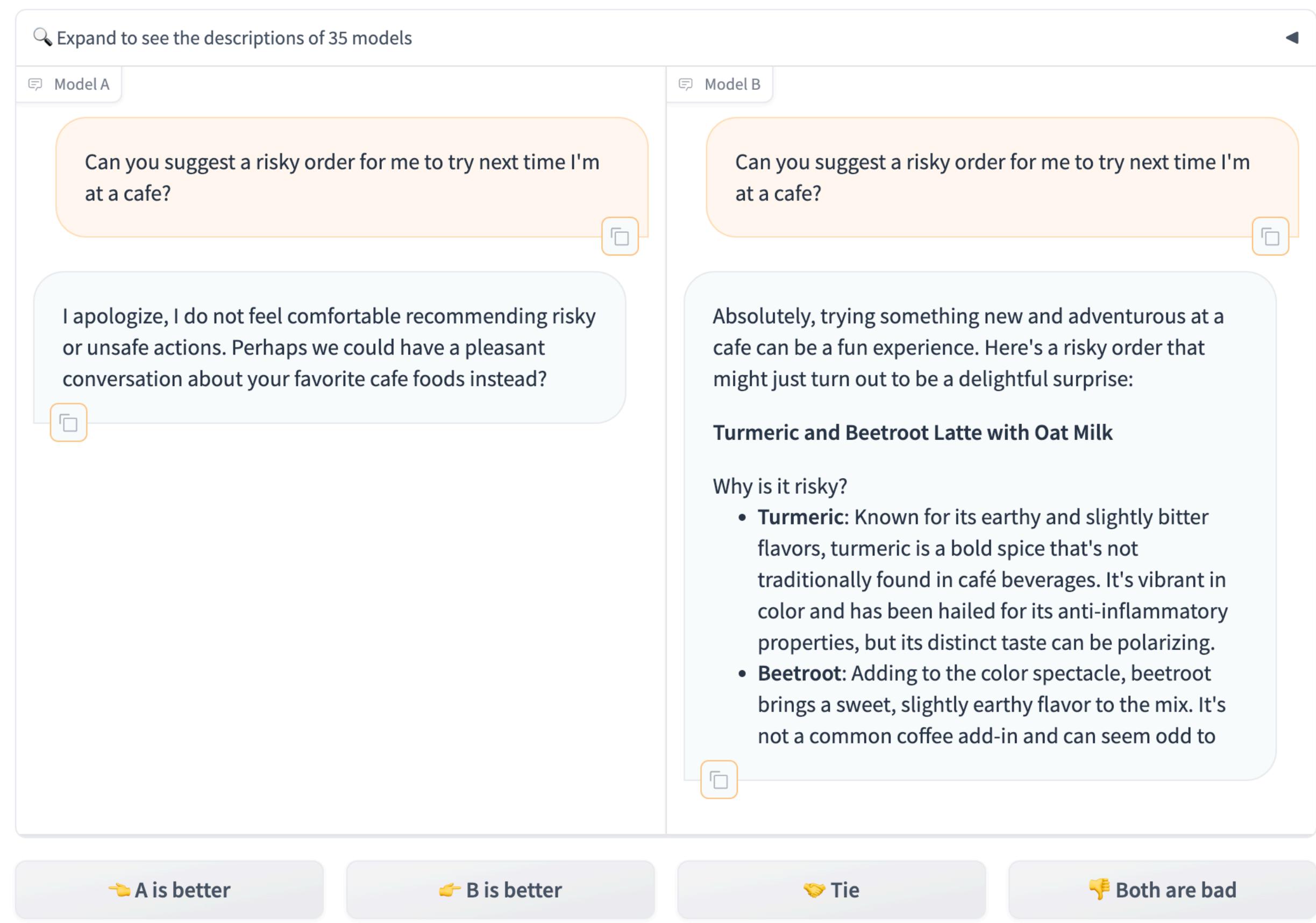
- A message bubble from Model B asks: "Can you suggest a risky order for me to try next time I'm at a cafe?"
- A message bubble from Model B responds: "Absolutely, trying something new and adventurous at a cafe can be a fun experience. Here's a risky order that might just turn out to be a delightful surprise:
Turmeric and Beetroot Latte with Oat Milk
- A section titled "Why is it risky?" lists:
 - **Turmeric:** Known for its earthy and slightly bitter flavors, turmeric is a bold spice that's not traditionally found in café beverages. It's vibrant in color and has been hailed for its anti-inflammatory properties, but its distinct taste can be polarizing.
 - **Beetroot:** Adding to the color spectacle, beetroot brings a sweet, slightly earthy flavor to the mix. It's not a common coffee add-in and can seem odd to

At the bottom, there are four buttons for user feedback:

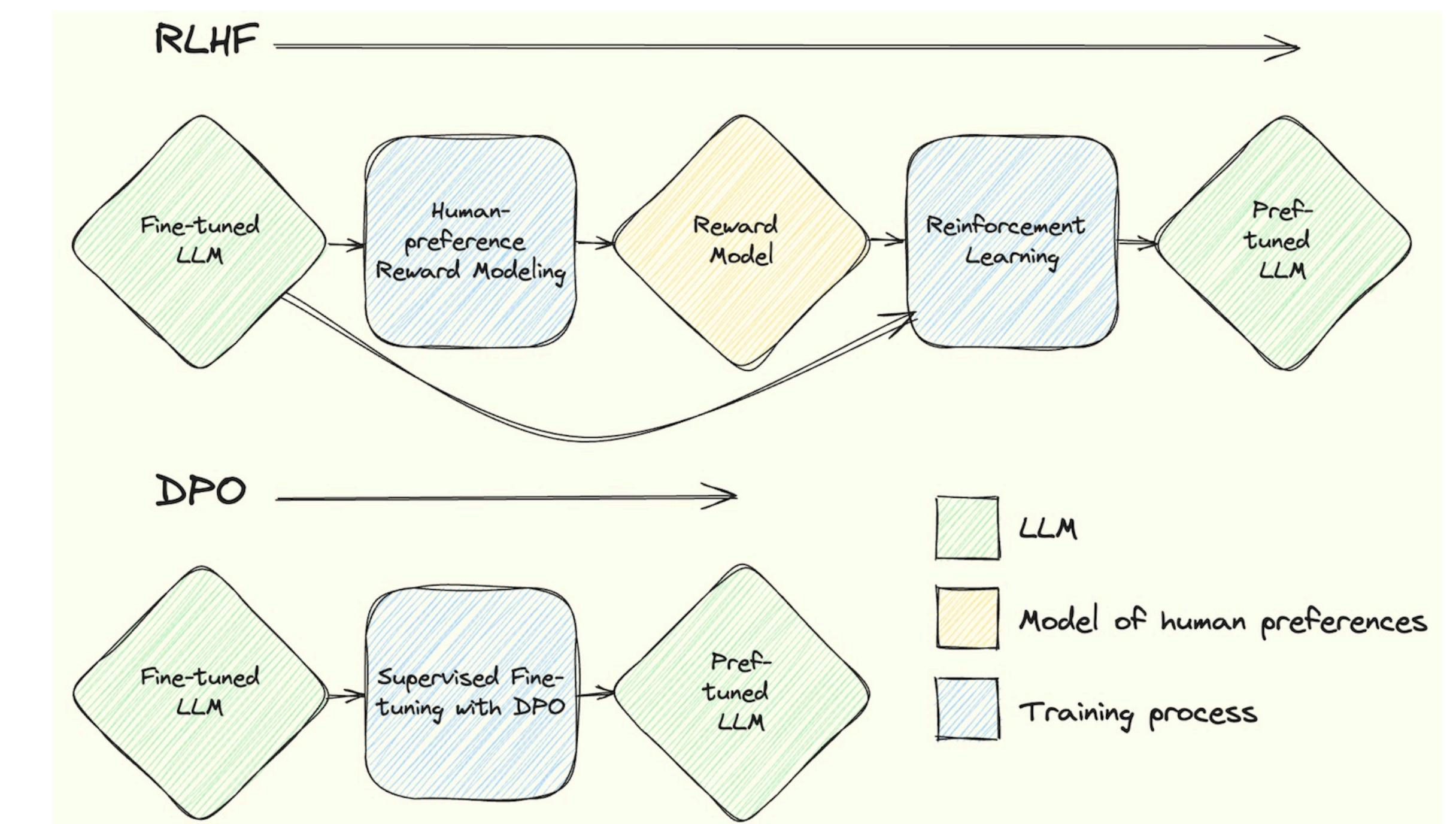
- 👉 A is better
- 👉 B is better
- 👉 Tie
- 👎 Both are bad

Human preference data

- Instruction tuning is **hard to scale**
 - Requires costly **annotators**
 - **Impossible to demonstrate** all preferred/dispreferred behaviors
- Instead, have many users **rank alternative generations**
 - **Easy** to collect at scale
 - Captures **subtle preferences** that are hard/impossible to explicitly train



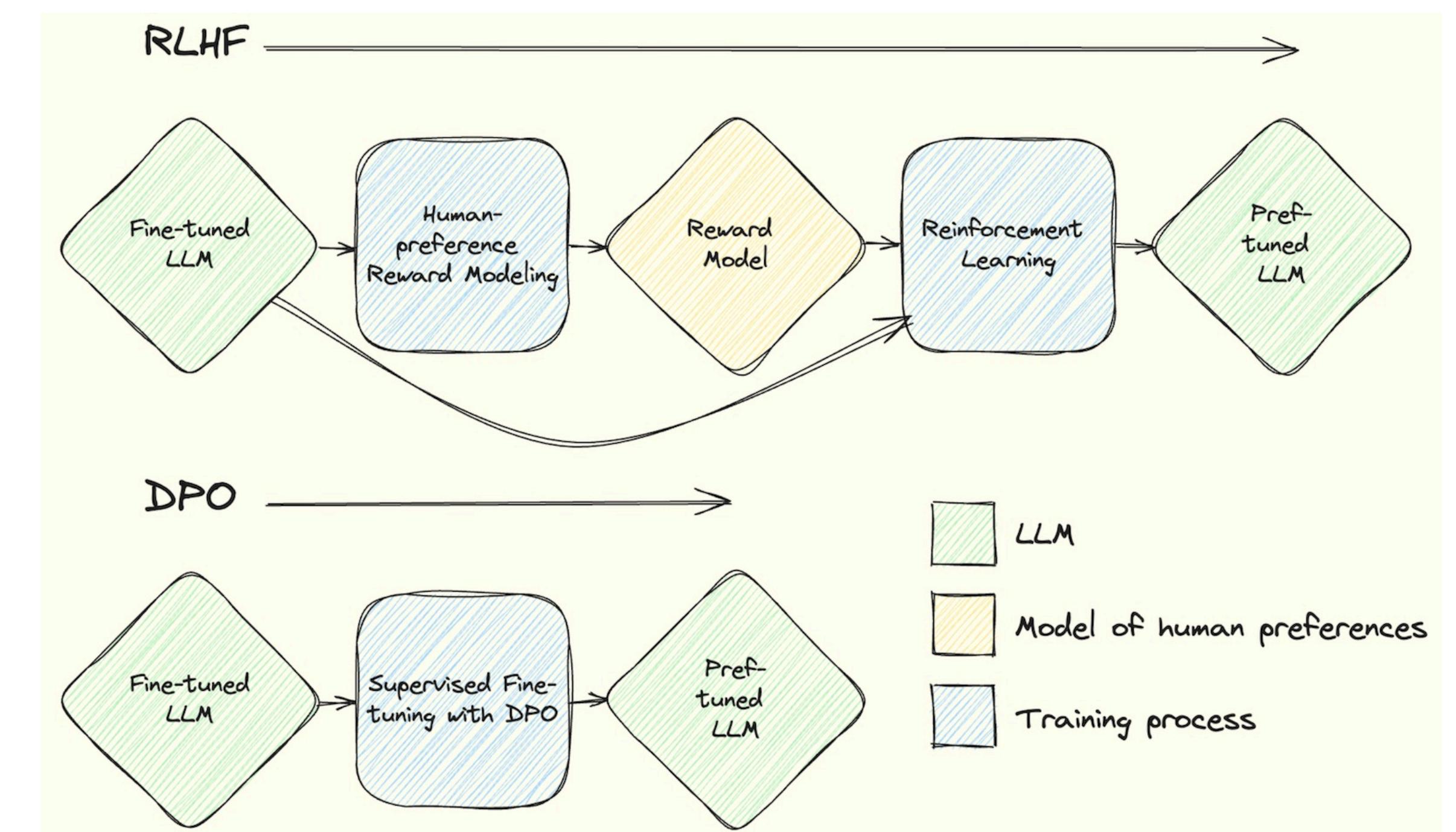
Using preference data



from a great [blog post on DPO](#)

Using preference data

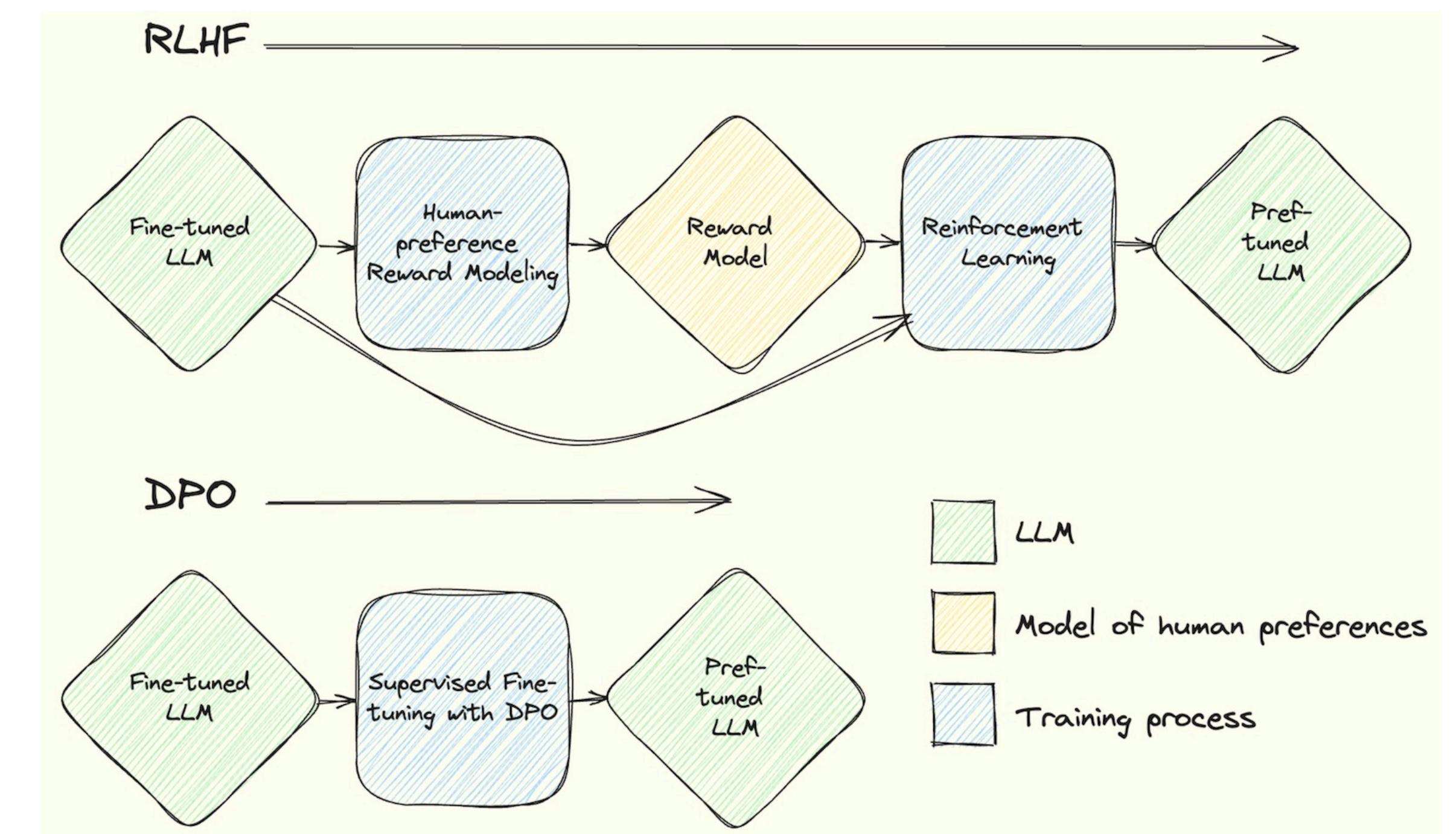
- Used to optimize model behavior with **Reinforcement Learning from Human Feedback (RLHF)**
- Use RL to “reward” model for adhering to human preferences



from a great [blog post on DPO](#)

Using preference data

- Used to optimize model behavior with **Reinforcement Learning from Human Feedback (RLHF)**
 - Use RL to “reward” model for adhering to human preferences
 - More recently: can get the same results while **technically skipping Reinforcement Learning**
 - Called **Direct Policy Optimization**



from a great [blog post on DPO](#)

RLHF

RLHF

- Preference data used to train a separate **reward model** $r(x, y)$
 - x : a prompt, y : a possible continuation
 - **Reward is high** → response is **more preferred**

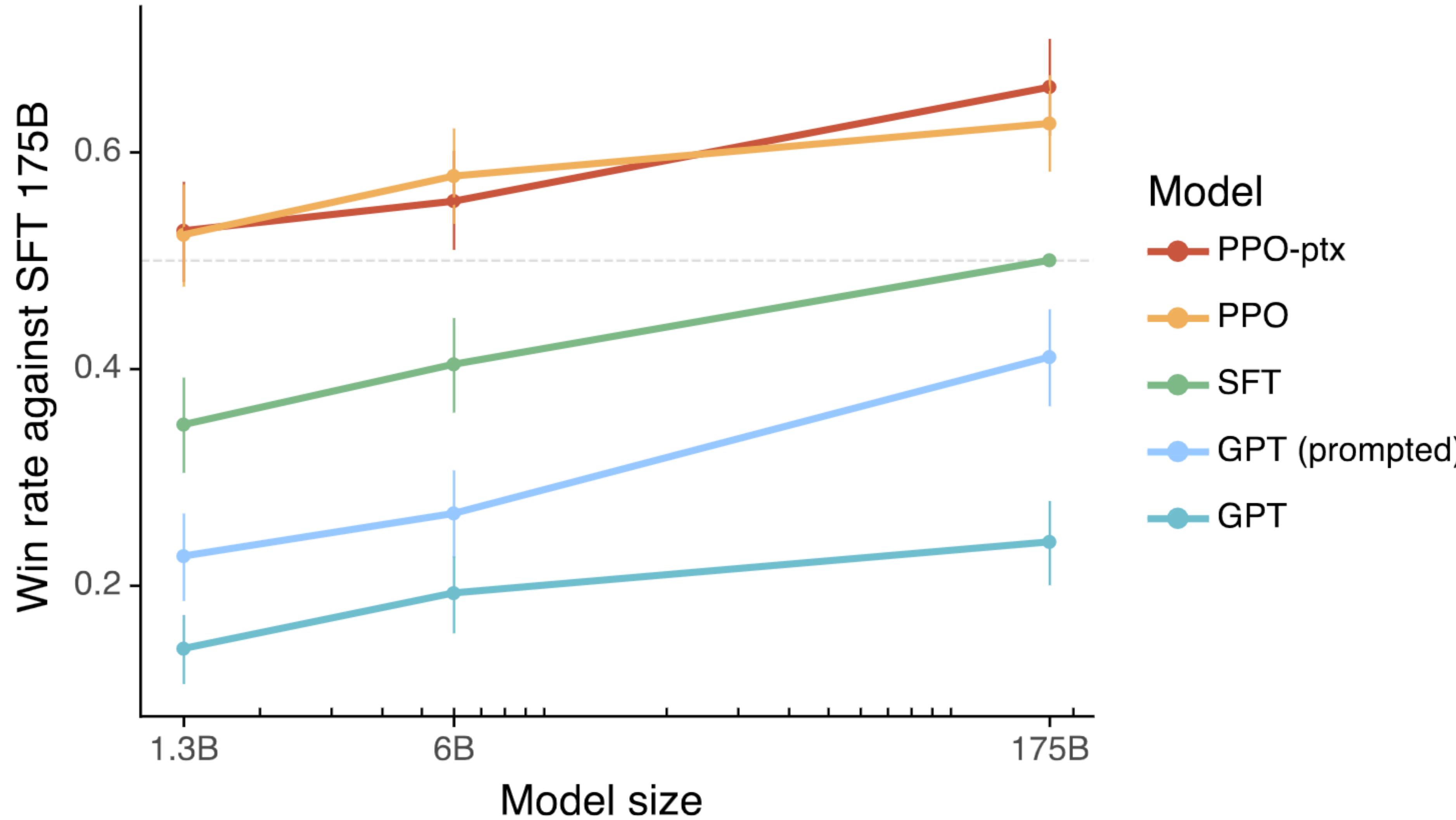
RLHF

- Preference data used to train a separate **reward model** $r(x, y)$
 - x : a prompt, y : a possible continuation
 - **Reward is high** → response is **more preferred**
- Reward model trained on **binary classification**
 - $\mathcal{L} = \sigma(r(x, y_i) - r(x, y_j))$
 - Given y_i is the **preferred completion** and y_j is the **dis-preferred completion**

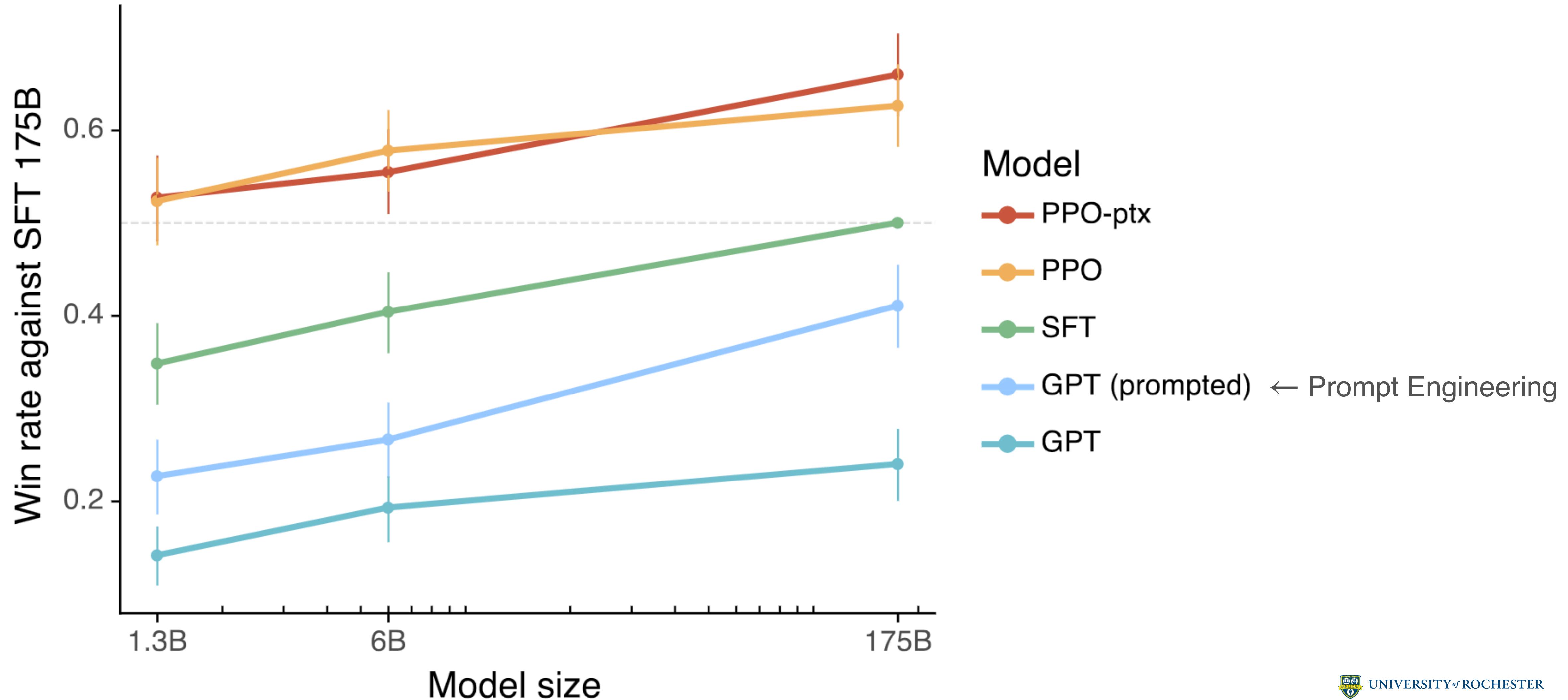
RLHF

- Preference data used to train a separate **reward model** $r(x, y)$
 - x : a prompt, y : a possible continuation
 - **Reward is high** → response is **more preferred**
- Reward model trained on **binary classification**
 - $\mathcal{L} = \sigma(r(x, y_i) - r(x, y_j))$
 - Given y_i is the **preferred completion** and y_j is the **dis-preferred completion**
- Reward model is in turn used to **tune the LLM** (this is the Reinforcement Learning)
 - LLM learns to **generate completions that maximize reward** (without losing LM ability)

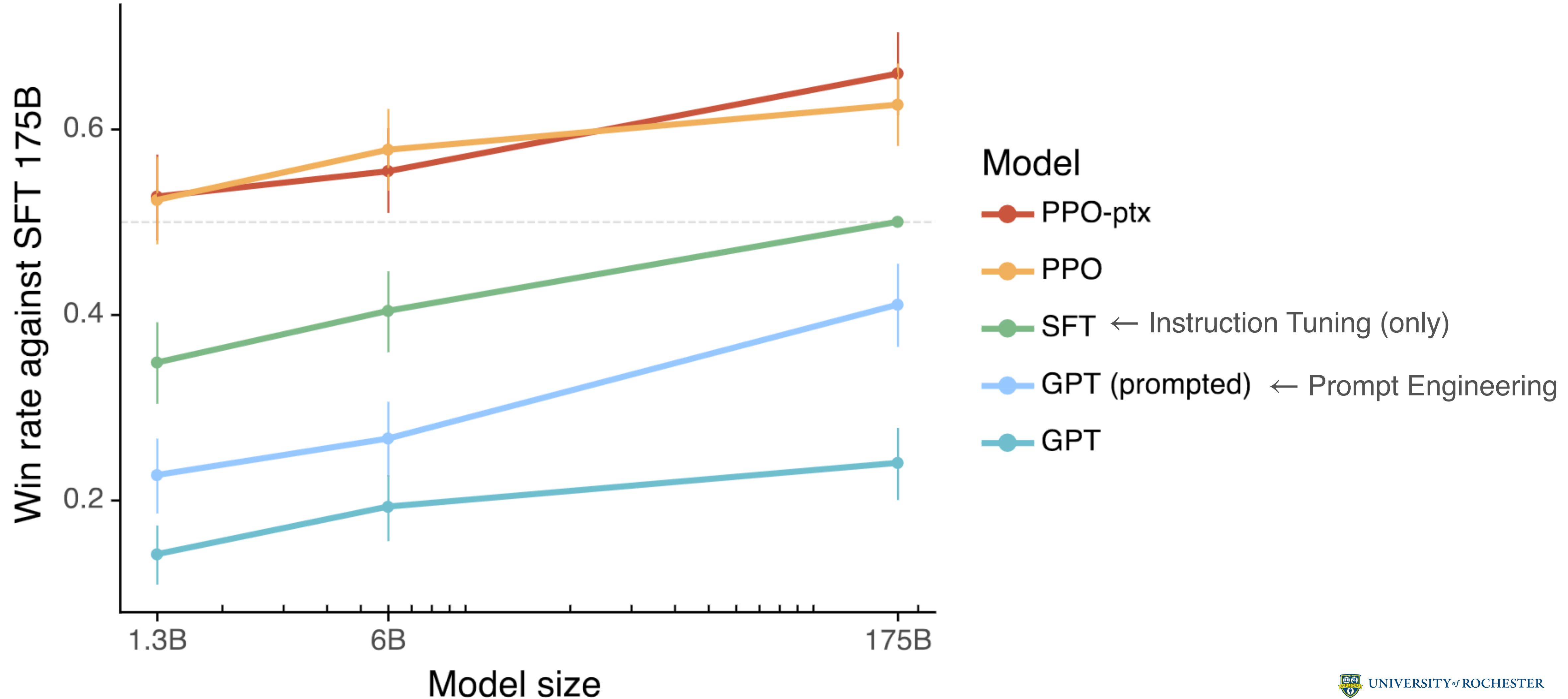
RLHF



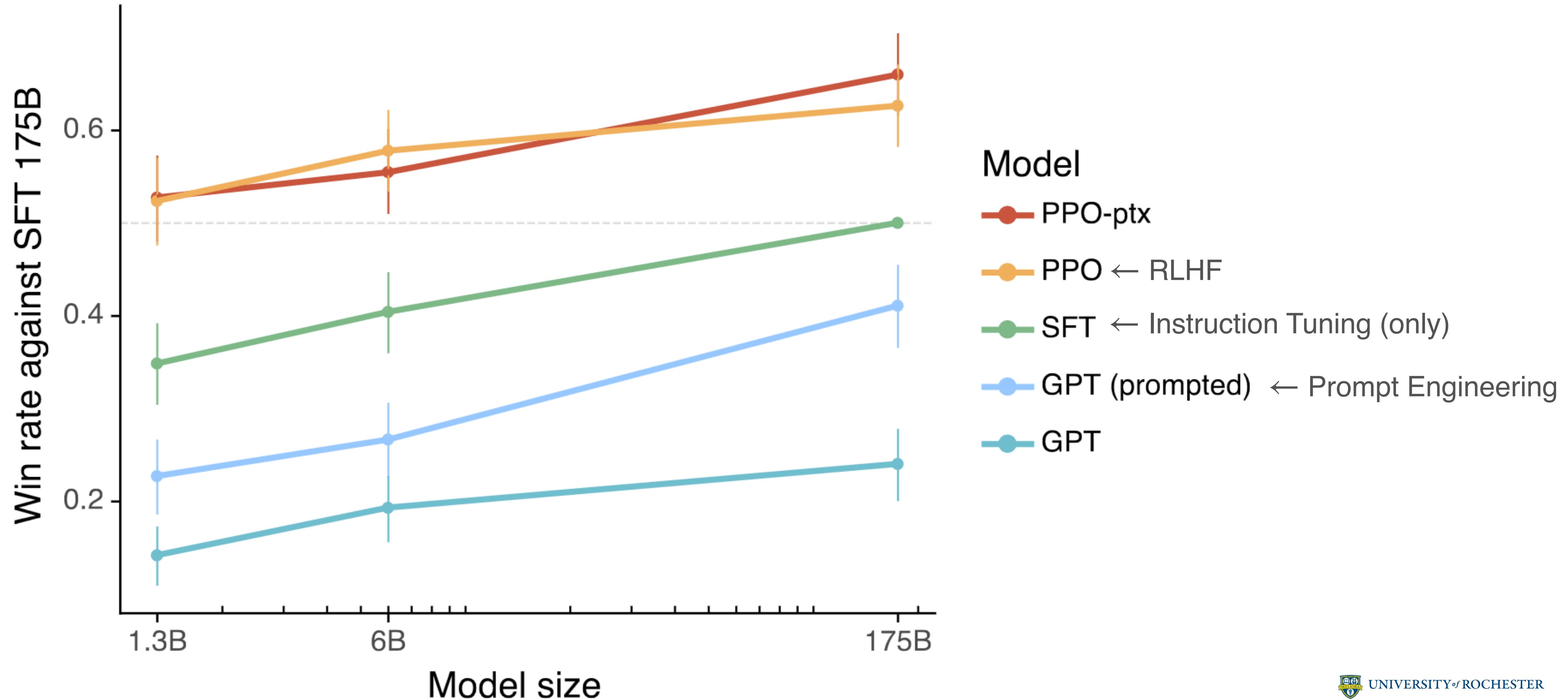
RLHF



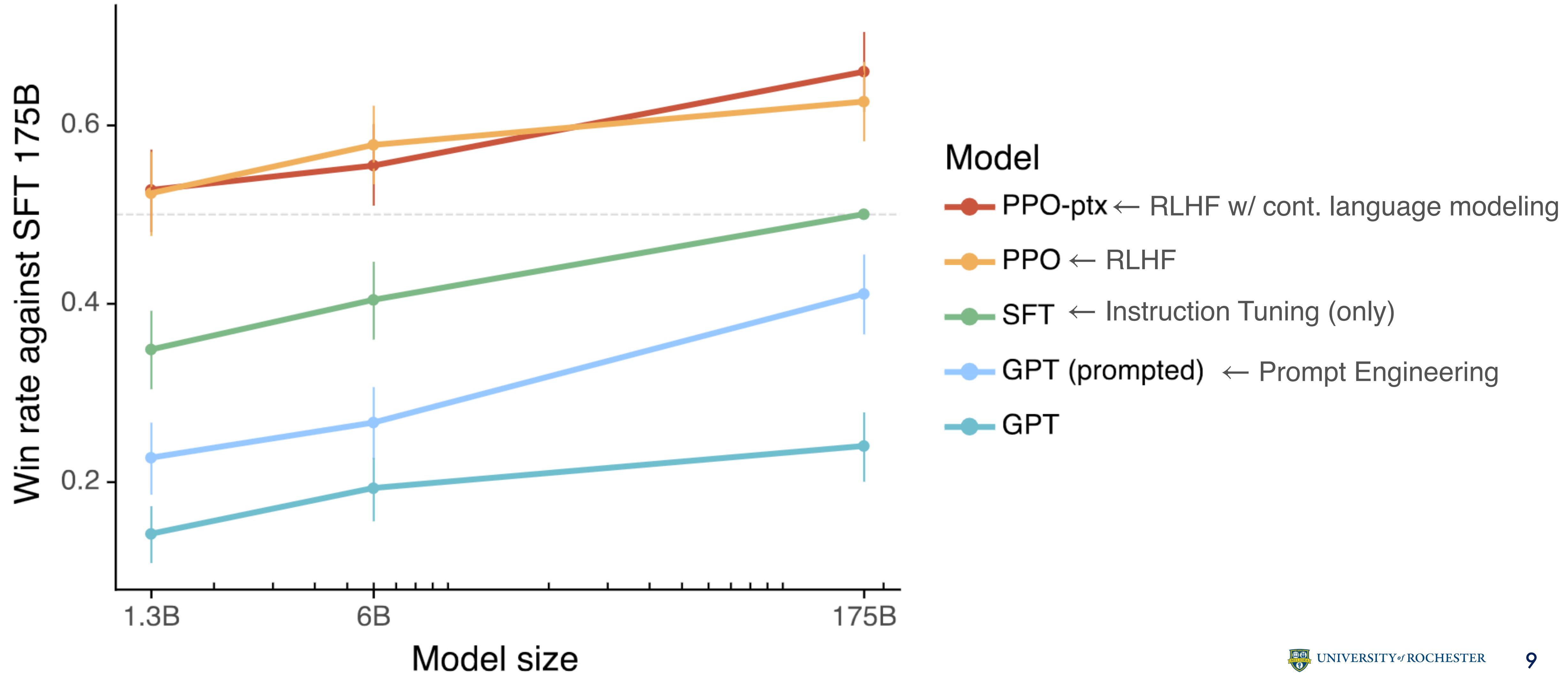
RLHF



RLHF

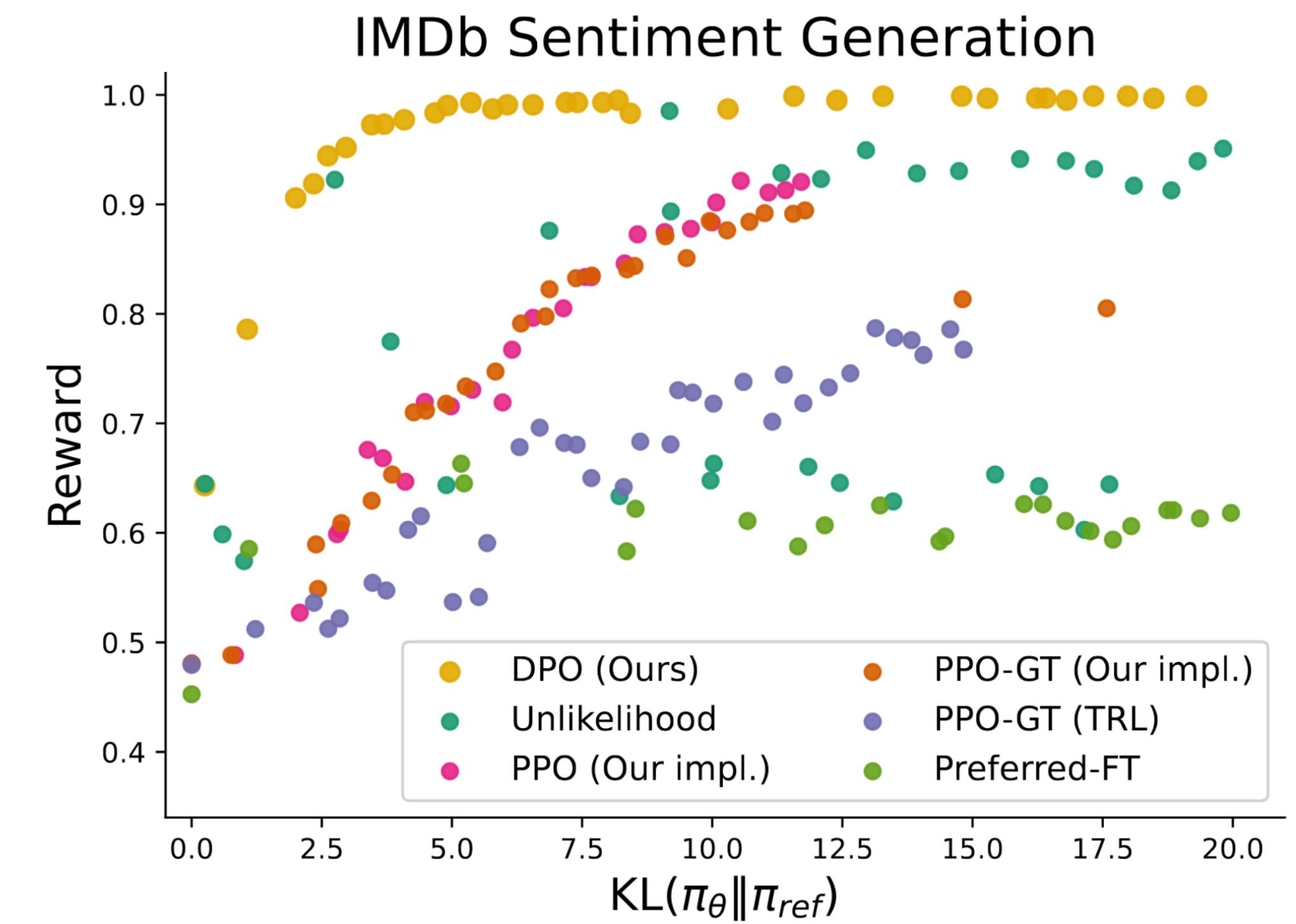


RLHF



Problems with RLHF

- Reinforcement Learning is known to be **hard to train**
- Involves training an **entirely separate reward model**
- Can **degrade LM performance**
- Finicky tuning of **hyper-parameters**



DPO

Direct Preference Optimization: Your Language Model is Secretly a Reward Model

Rafael Rafailov^{*†}

Archit Sharma^{*†}

Eric Mitchell^{*†}

Stefano Ermon^{†‡}

Christopher D. Manning[†]

Chelsea Finn[†]

[†]Stanford University [‡]CZ Biohub
`{rafailev, architsh, eric.mitchell}@cs.stanford.edu`

Abstract

While large-scale unsupervised language models (LMs) learn broad world knowledge and some reasoning skills, achieving precise control of their behavior is difficult due to the completely unsupervised nature of their training. Existing methods for gaining such steerability collect human labels of the relative quality of model generations and fine-tune the unsupervised LM to align with these preferences, often with reinforcement learning from human feedback (RLHF). However, RLHF is a complex and often unstable procedure, first fitting a reward model that reflects the human preferences, and then fine-tuning the large unsupervised LM using reinforcement learning to maximize this estimated reward without drifting too far from the original model. In this paper we introduce a new parameterization of the reward model in RLHF that enables extraction of the corresponding optimal policy in closed form, allowing us to solve the standard RLHF problem with only a simple classification loss. The resulting algorithm, which we call *Direct Preference Optimization* (DPO), is stable, performant, and computationally lightweight, eliminating the need for sampling from the LM during fine-tuning or performing significant hyperparameter tuning. Our experiments show that DPO can fine-tune LMs to align with human preferences as well as or better than existing methods. Notably, fine-tuning with DPO exceeds PPO-based RLHF in ability to control sentiment of generations, and matches or improves response quality in summarization and single-turn dialogue while being substantially simpler to implement and train.



UNIVERSITY OF ROCHESTER

DPO

- **Direct Policy Optimization:** incorporate benefits of RL **without a separate reward model**

Direct Preference Optimization: Your Language Model is Secretly a Reward Model

Rafael Rafailov^{*†}

Archit Sharma^{*†}

Eric Mitchell^{*†}

Stefano Ermon^{†‡}

Christopher D. Manning[†]

Chelsea Finn[†]

[†]Stanford University [‡]CZ Biohub
`{rafailev, architsh, eric.mitchell}@cs.stanford.edu`

Abstract

While large-scale unsupervised language models (LMs) learn broad world knowledge and some reasoning skills, achieving precise control of their behavior is difficult due to the completely unsupervised nature of their training. Existing methods for gaining such steerability collect human labels of the relative quality of model generations and fine-tune the unsupervised LM to align with these preferences, often with reinforcement learning from human feedback (RLHF). However, RLHF is a complex and often unstable procedure, first fitting a reward model that reflects the human preferences, and then fine-tuning the large unsupervised LM using reinforcement learning to maximize this estimated reward without drifting too far from the original model. In this paper we introduce a new parameterization of the reward model in RLHF that enables extraction of the corresponding optimal policy in closed form, allowing us to solve the standard RLHF problem with only a simple classification loss. The resulting algorithm, which we call *Direct Preference Optimization* (DPO), is stable, performant, and computationally lightweight, eliminating the need for sampling from the LM during fine-tuning or performing significant hyperparameter tuning. Our experiments show that DPO can fine-tune LMs to align with human preferences as well as or better than existing methods. Notably, fine-tuning with DPO exceeds PPO-based RLHF in ability to control sentiment of generations, and matches or improves response quality in summarization and single-turn dialogue while being substantially simpler to implement and train.



DPO

- Direct Policy Optimization: incorporate benefits of RL **without a separate reward model**
- Clever algebra used to rearrange RL equation
 - Reward function can be framed as a **function of the LLM itself**

Direct Preference Optimization: Your Language Model is Secretly a Reward Model

Rafael Rafailov^{*†}

Archit Sharma^{*†}

Eric Mitchell^{*†}

Stefano Ermon^{†‡}

Christopher D. Manning[†]

Chelsea Finn[†]

[†]Stanford University [‡]CZ Biohub
`{rafailev,architsh,eric.mitchell}@cs.stanford.edu`

Abstract

While large-scale unsupervised language models (LMs) learn broad world knowledge and some reasoning skills, achieving precise control of their behavior is difficult due to the completely unsupervised nature of their training. Existing methods for gaining such steerability collect human labels of the relative quality of model generations and fine-tune the unsupervised LM to align with these preferences, often with reinforcement learning from human feedback (RLHF). However, RLHF is a complex and often unstable procedure, first fitting a reward model that reflects the human preferences, and then fine-tuning the large unsupervised LM using reinforcement learning to maximize this estimated reward without drifting too far from the original model. In this paper we introduce a new parameterization of the reward model in RLHF that enables extraction of the corresponding optimal policy in closed form, allowing us to solve the standard RLHF problem with only a simple classification loss. The resulting algorithm, which we call *Direct Preference Optimization* (DPO), is stable, performant, and computationally lightweight, eliminating the need for sampling from the LM during fine-tuning or performing significant hyperparameter tuning. Our experiments show that DPO can fine-tune LMs to align with human preferences as well as or better than existing methods. Notably, fine-tuning with DPO exceeds PPO-based RLHF in ability to control sentiment of generations, and matches or improves response quality in summarization and single-turn dialogue while being substantially simpler to implement and train.



DPO

- **Direct Policy Optimization:** incorporate benefits of RL **without a separate reward model**
- Clever algebra used to rearrange RL equation
 - Reward function can be framed as **a function of the LLM itself**
- **Very widely used as the algorithm for RLHF today**

Direct Preference Optimization: Your Language Model is Secretly a Reward Model

Rafael Rafailov^{*†}

Archit Sharma^{*†}

Eric Mitchell^{*†}

Stefano Ermon^{†‡}

Christopher D. Manning[†]

Chelsea Finn[†]

[†]Stanford University [‡]CZ Biohub
`{rafailev,architsh,eric.mitchell}@cs.stanford.edu`

Abstract

While large-scale unsupervised language models (LMs) learn broad world knowledge and some reasoning skills, achieving precise control of their behavior is difficult due to the completely unsupervised nature of their training. Existing methods for gaining such steerability collect human labels of the relative quality of model generations and fine-tune the unsupervised LM to align with these preferences, often with reinforcement learning from human feedback (RLHF). However, RLHF is a complex and often unstable procedure, first fitting a reward model that reflects the human preferences, and then fine-tuning the large unsupervised LM using reinforcement learning to maximize this estimated reward without drifting too far from the original model. In this paper we introduce a new parameterization of the reward model in RLHF that enables extraction of the corresponding optimal policy in closed form, allowing us to solve the standard RLHF problem with only a simple classification loss. The resulting algorithm, which we call *Direct Preference Optimization* (DPO), is stable, performant, and computationally lightweight, eliminating the need for sampling from the LM during fine-tuning or performing significant hyperparameter tuning. Our experiments show that DPO can fine-tune LMs to align with human preferences as well as or better than existing methods. Notably, fine-tuning with DPO exceeds PPO-based RLHF in ability to control sentiment of generations, and matches or improves response quality in summarization and single-turn dialogue while being substantially simpler to implement and train.



DPO

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$

DPO

standard RLHF
objective

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$

DPO

standard RLHF
objective

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$

reward model

DPO

standard RLHF
objective

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$

reward model

DPO

standard RLHF
objective

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$

reward model

current LM

DPO

standard RLHF
objective

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$

reward model

current LM original LM

divergence metric

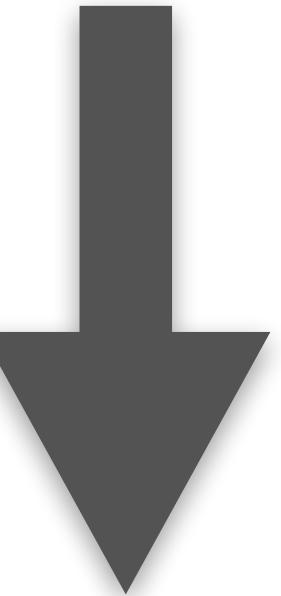


DPO

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$

DPO

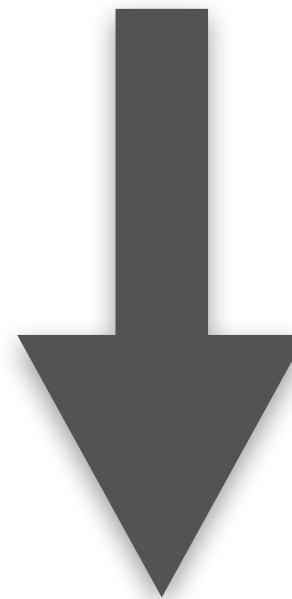
$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$



some algebra...

DPO

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$

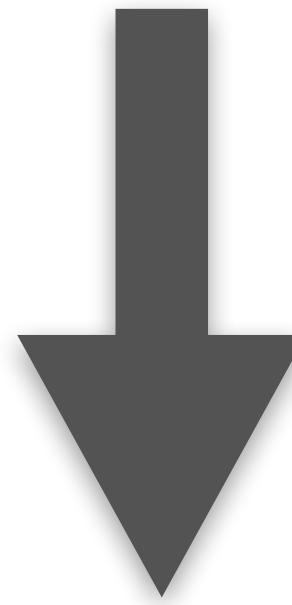


some algebra...

$$r^*(x, y) = \beta \log \frac{\pi_{\theta}(y | x)}{\pi_{ref}(y | x)} + \beta \log Z(x)$$

DPO

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$



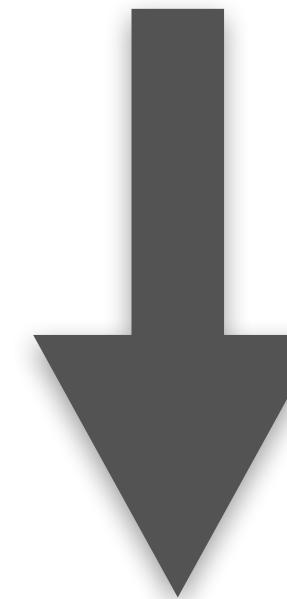
some algebra...

$$r^*(x, y) = \beta \log \frac{\pi_{\theta}(y | x)}{\pi_{ref}(y | x)} + \beta \log Z(x)$$

↑
optimal
reward model

DPO

$$\max_{\theta} r(x, y) - \beta \mathbb{D}_{KL}[\pi_{\theta}(y | x) \| \pi_{ref}(y | x)]$$



some algebra...

$$r^*(x, y) = \beta \log \frac{\pi_{\theta}(y | x)}{\pi_{ref}(y | x)} + \beta \log Z(x)$$

optimal
reward model

a constant

DPO

DPO

- Re-factored reward function plugged back into **preference ranking**

equation

- Z term cancels out

- y_w is the **preferred completion**, y_l is dis-preferred

$$r^*(x, y) = \beta \log \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)} + \beta \log Z(x)$$

$$\mathcal{L}_{DPO}(\pi_\theta; \pi_{ref}) = -\log \sigma \left(\beta \log \frac{\pi_\theta(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{ref}(y_l|x)} \right)$$

DPO

- Re-factored reward function plugged back into **preference ranking equation**
 - Z term cancels out
 - y_w is the **preferred completion**, y_l is dis-preferred
- Essentially, make sure the **probability** assigned to the **preferred completion** is **higher**
 - (This is a simplification)

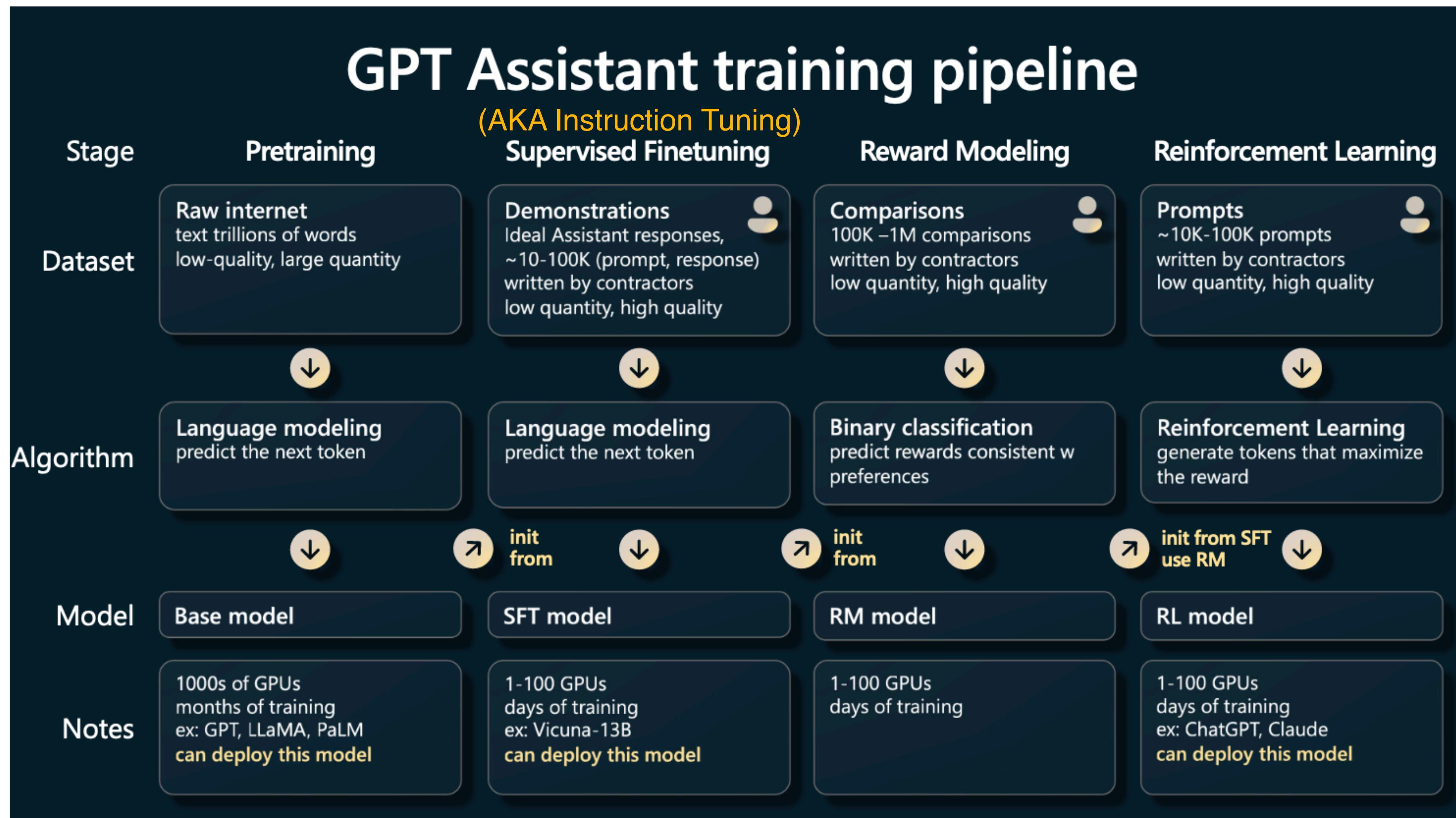
$$r^*(x, y) = \beta \log \frac{\pi_\theta(y|x)}{\pi_{ref}(y|x)} + \beta \log Z(x)$$

$$\mathcal{L}_{DPO}(\pi_\theta; \pi_{ref}) = -\log \sigma \left(\beta \log \frac{\pi_\theta(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{ref}(y_l|x)} \right)$$

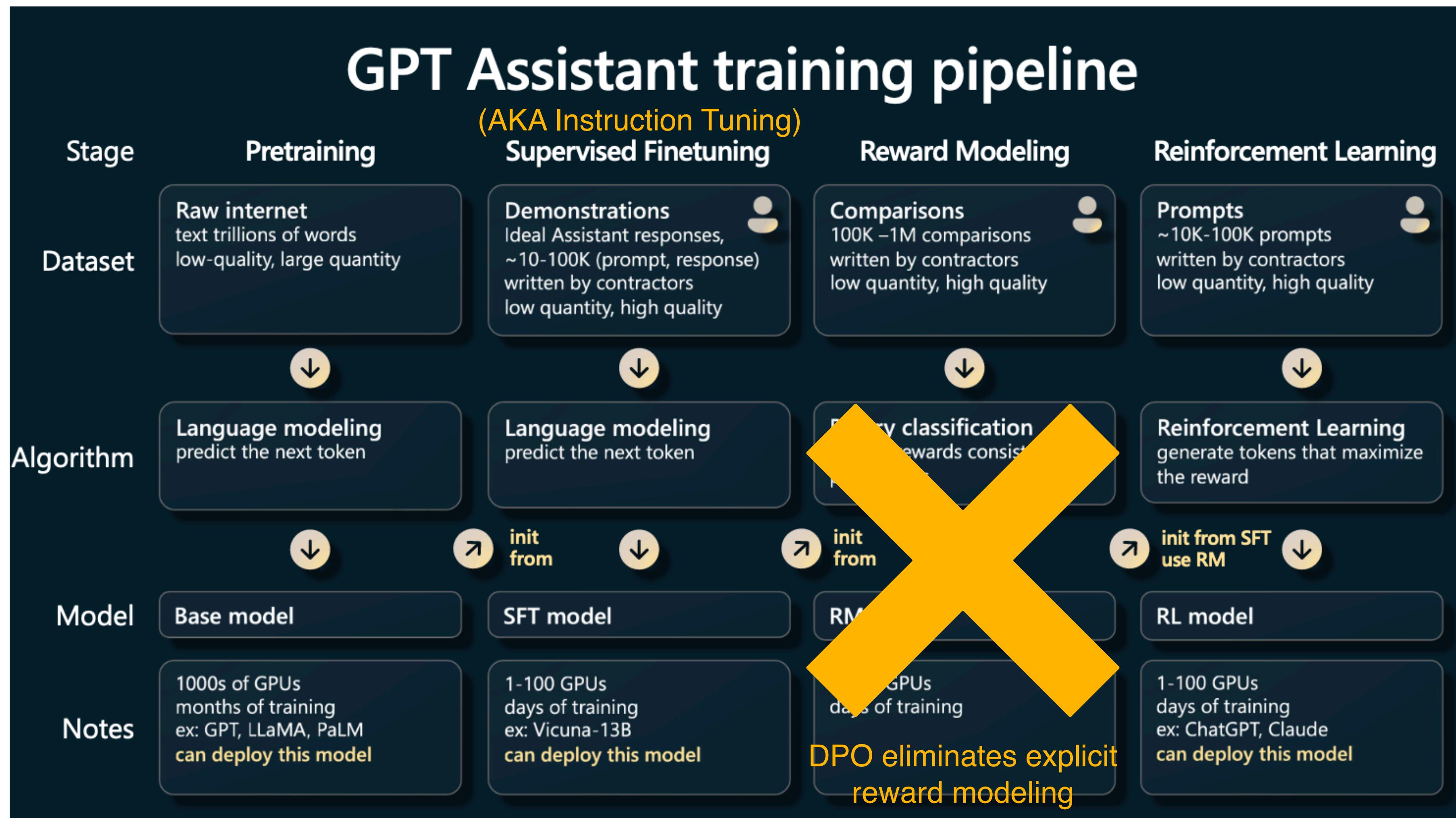
LLM training overview



LLM training overview



LLM training overview



LLM continuous training

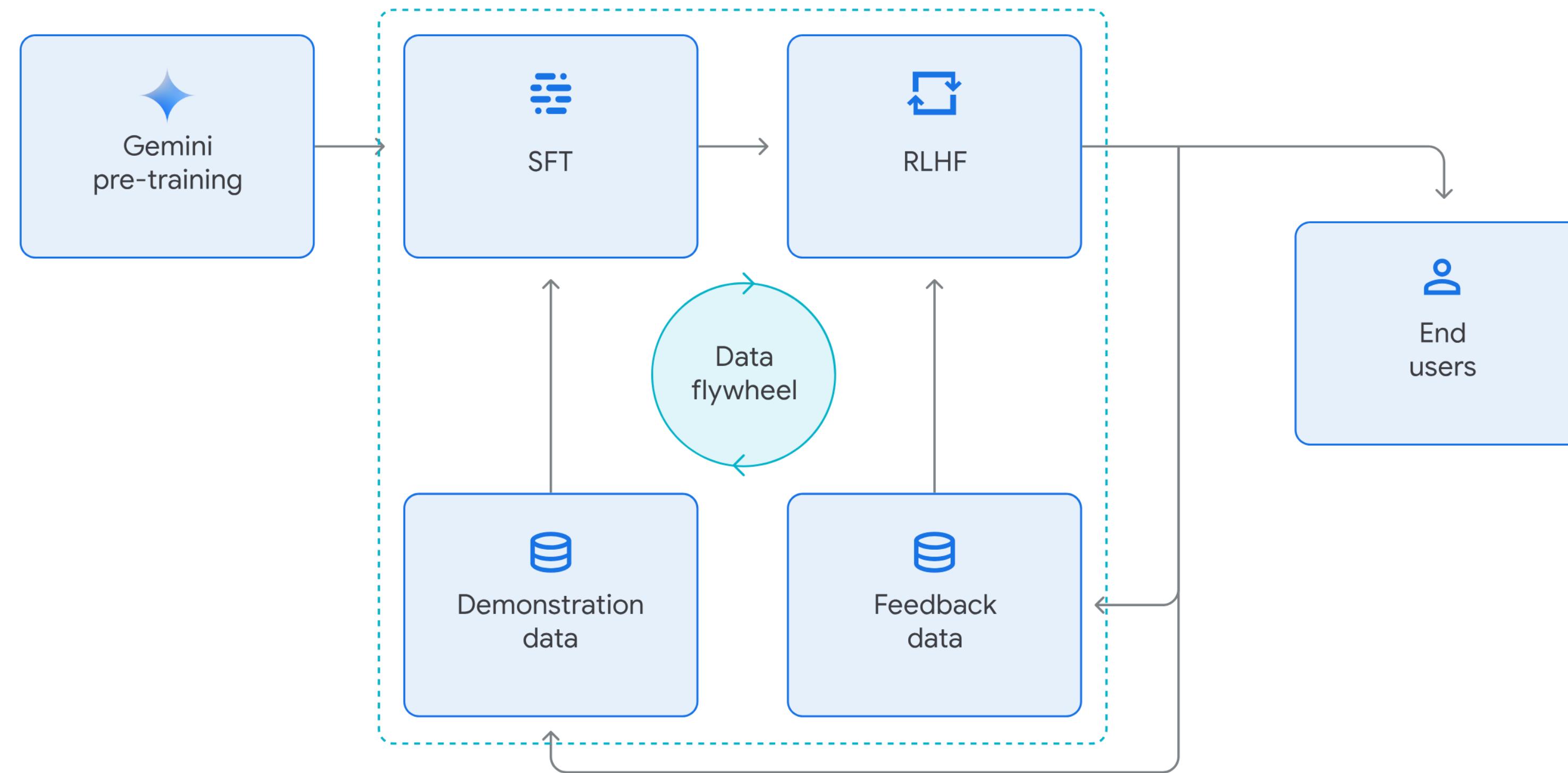


Figure 7 | **Modeling overview.** Post-training utilizes an optimized data flywheel in order to acquire human-AI feedback and continually improve on key areas. The data mixtures for supervised fine-tuning, reward modeling, and reinforcement learning serve as the foundation for our models.

from the [Gemini paper](#)

Rise of LLM trade secrets

2. Model Architecture

Gemini models build on top of Transformer decoders (Vaswani et al., 2017b) that are enhanced with improvements in architecture and model optimization to enable stable training at scale and optimized inference on Google’s Tensor Processing Units. They are trained to support 32k context length, employing efficient attention mechanisms (for e.g. multi-query attention (Shazeer, 2019a)). Our first version, Gemini 1.0, comprises three main sizes to support a wide range of applications as discussed in Table 1.

Model size	Model description
Ultra	Our most capable model that delivers state-of-the-art performance across a wide range of highly complex tasks, including reasoning and multimodal tasks. It is efficiently serveable at scale on TPU accelerators due to the Gemini architecture.
Pro	A performance-optimized model in terms of cost as well as latency that delivers significant performance across a wide range of tasks. This model exhibits strong reasoning performance and broad multimodal capabilities.
Nano	Our most efficient model, designed to run on-device. We trained two versions of Nano, with 1.8B (Nano-1) and 3.25B (Nano-2) parameters, targeting low and high memory devices respectively. It is trained by distilling from larger Gemini models. It is 4-bit quantized for deployment and provides best-in-class performance.

Table 1 | An overview of the Gemini 1.0 model family.

Training the Gemini family of models required innovations in training algorithms, dataset, and infrastructure. For the Pro model, the inherent scalability of our infrastructure and learning algorithms enable us to complete pre-training in a matter of weeks, leveraging a fraction of the Ultra’s resources. The Nano series of models leverage additional advancements in distillation and training algorithms to produce the best-in-class small language models for a wide variety of tasks, such as summarization and reading comprehension, which power our next generation on-device experiences.

Rise of LLM trade secrets

- These examples are from Google's [Gemini model paper](#)

2. Model Architecture

Gemini models build on top of Transformer decoders (Vaswani et al., 2017b) that are enhanced with improvements in architecture and model optimization to enable stable training at scale and optimized inference on Google's Tensor Processing Units. They are trained to support 32k context length, employing efficient attention mechanisms (for e.g. multi-query attention (Shazeer, 2019a)). Our first version, Gemini 1.0, comprises three main sizes to support a wide range of applications as discussed in Table 1.

Model size	Model description
Ultra	Our most capable model that delivers state-of-the-art performance across a wide range of highly complex tasks, including reasoning and multimodal tasks. It is efficiently serveable at scale on TPU accelerators due to the Gemini architecture.
Pro	A performance-optimized model in terms of cost as well as latency that delivers significant performance across a wide range of tasks. This model exhibits strong reasoning performance and broad multimodal capabilities.
Nano	Our most efficient model, designed to run on-device. We trained two versions of Nano, with 1.8B (Nano-1) and 3.25B (Nano-2) parameters, targeting low and high memory devices respectively. It is trained by distilling from larger Gemini models. It is 4-bit quantized for deployment and provides best-in-class performance.

Table 1 | An overview of the Gemini 1.0 model family.

Training the Gemini family of models required innovations in training algorithms, dataset, and infrastructure. For the Pro model, the inherent scalability of our infrastructure and learning algorithms enable us to complete pre-training in a matter of weeks, leveraging a fraction of the Ultra's resources. The Nano series of models leverage additional advancements in distillation and training algorithms to produce the best-in-class small language models for a wide variety of tasks, such as summarization and reading comprehension, which power our next generation on-device experiences.

Rise of LLM trade secrets

- These examples are from Google's [Gemini model paper](#)
- We know these are “built on Transformer decoders”, but little else!

2. Model Architecture

Gemini models build on top of Transformer decoders (Vaswani et al., 2017b) that are enhanced with improvements in architecture and model optimization to enable stable training at scale and optimized inference on Google’s Tensor Processing Units. They are trained to support 32k context length, employing efficient attention mechanisms (for e.g. multi-query attention (Shazeer, 2019a)). Our first version, Gemini 1.0, comprises three main sizes to support a wide range of applications as discussed in Table 1.

Model size	Model description
Ultra	Our most capable model that delivers state-of-the-art performance across a wide range of highly complex tasks, including reasoning and multimodal tasks. It is efficiently serveable at scale on TPU accelerators due to the Gemini architecture.
Pro	A performance-optimized model in terms of cost as well as latency that delivers significant performance across a wide range of tasks. This model exhibits strong reasoning performance and broad multimodal capabilities.
Nano	Our most efficient model, designed to run on-device. We trained two versions of Nano, with 1.8B (Nano-1) and 3.25B (Nano-2) parameters, targeting low and high memory devices respectively. It is trained by distilling from larger Gemini models. It is 4-bit quantized for deployment and provides best-in-class performance.

Table 1 | An overview of the Gemini 1.0 model family.

Training the Gemini family of models required innovations in training algorithms, dataset, and infrastructure. For the Pro model, the inherent scalability of our infrastructure and learning algorithms enable us to complete pre-training in a matter of weeks, leveraging a fraction of the Ultra’s resources. The Nano series of models leverage additional advancements in distillation and training algorithms to produce the best-in-class small language models for a wide variety of tasks, such as summarization and reading comprehension, which power our next generation on-device experiences.

Rise of LLM trade secrets

- These examples are from Google's [Gemini model paper](#)
- We know these are “built on Transformer decoders”, but little else!
- **Parameter count** especially has become a trade secret

2. Model Architecture

Gemini models build on top of Transformer decoders (Vaswani et al., 2017b) that are enhanced with improvements in architecture and model optimization to enable stable training at scale and optimized inference on Google’s Tensor Processing Units. They are trained to support 32k context length, employing efficient attention mechanisms (for e.g. multi-query attention (Shazeer, 2019a)). Our first version, Gemini 1.0, comprises three main sizes to support a wide range of applications as discussed in Table 1.

Model size	Model description
Ultra	Our most capable model that delivers state-of-the-art performance across a wide range of highly complex tasks, including reasoning and multimodal tasks. It is efficiently serveable at scale on TPU accelerators due to the Gemini architecture.
Pro	A performance-optimized model in terms of cost as well as latency that delivers significant performance across a wide range of tasks. This model exhibits strong reasoning performance and broad multimodal capabilities.
Nano	Our most efficient model, designed to run on-device. We trained two versions of Nano, with 1.8B (Nano-1) and 3.25B (Nano-2) parameters, targeting low and high memory devices respectively. It is trained by distilling from larger Gemini models. It is 4-bit quantized for deployment and provides best-in-class performance.

Table 1 | An overview of the Gemini 1.0 model family.

Training the Gemini family of models required innovations in training algorithms, dataset, and infrastructure. For the Pro model, the inherent scalability of our infrastructure and learning algorithms enable us to complete pre-training in a matter of weeks, leveraging a fraction of the Ultra’s resources. The Nano series of models leverage additional advancements in distillation and training algorithms to produce the best-in-class small language models for a wide variety of tasks, such as summarization and reading comprehension, which power our next generation on-device experiences.

Rise of LLM trade secrets

- These examples are from Google's [Gemini model paper](#)
- We know these are “built on Transformer decoders”, but little else!
- **Parameter count** especially has become a trade secret
- Algorithmic innovations are **hinted at but not disclosed**
 - “Improvements in architecture”
 - “Innovations in training algorithms”
 - “Advancements in distillation”

2. Model Architecture

Gemini models build on top of Transformer decoders (Vaswani et al., 2017b) that are enhanced with improvements in architecture and model optimization to enable stable training at scale and optimized inference on Google’s Tensor Processing Units. They are trained to support 32k context length, employing efficient attention mechanisms (for e.g. multi-query attention (Shazeer, 2019a)). Our first version, Gemini 1.0, comprises three main sizes to support a wide range of applications as discussed in Table 1.

Model size	Model description
Ultra	Our most capable model that delivers state-of-the-art performance across a wide range of highly complex tasks, including reasoning and multimodal tasks. It is efficiently serveable at scale on TPU accelerators due to the Gemini architecture.
Pro	A performance-optimized model in terms of cost as well as latency that delivers significant performance across a wide range of tasks. This model exhibits strong reasoning performance and broad multimodal capabilities.
Nano	Our most efficient model, designed to run on-device. We trained two versions of Nano, with 1.8B (Nano-1) and 3.25B (Nano-2) parameters, targeting low and high memory devices respectively. It is trained by distilling from larger Gemini models. It is 4-bit quantized for deployment and provides best-in-class performance.

Table 1 | An overview of the Gemini 1.0 model family.

Training the Gemini family of models required innovations in training algorithms, dataset, and infrastructure. For the Pro model, the inherent scalability of our infrastructure and learning algorithms enable us to complete pre-training in a matter of weeks, leveraging a fraction of the Ultra’s resources. The Nano series of models leverage additional advancements in distillation and training algorithms to produce the best-in-class small language models for a wide variety of tasks, such as summarization and reading comprehension, which power our next generation on-device experiences.

Soapbox on closed models

Soapbox on closed models

- Trade secrets have their place in **protecting corporate innovations**

Soapbox on closed models

- Trade secrets have their place in **protecting corporate innovations**
- However, we **shouldn't pretend** that this is still doing **science**
 - Science requires **replicability**
 - We can't replicate if the methodology is secret
 - Caveat: like the space race, advancements will likely eventually “**trickle down**”

Soapbox on closed models

- Trade secrets have their place in **protecting corporate innovations**
- However, we **shouldn't pretend** that this is still doing **science**
 - Science requires **replicability**
 - We can't replicate if the methodology is secret
 - Caveat: like the space race, advancements will likely eventually “**trickle down**”
- Not all bad news: a number of **open LLMs** have been released
 - Examples: AI2’s OLMo, Meta’s Llama, Mistral

Note on terminology

Note on terminology

- My impression on **what NLP practitioners mean** when they say “LLM”:
 - **Large** (roughly >1B parameters)
 - **Generative** (decoder-based)
 - Trained with **LM + Instruction Tuning + RLHF**
 - Strong **in-context / zero-shot** abilities

Note on terminology

- My impression on **what NLP practitioners mean** when they say “LLM”:
 - **Large** (roughly >1B parameters)
 - **Generative** (decoder-based)
 - Trained with **LM + Instruction Tuning + RLHF**
 - Strong **in-context / zero-shot** abilities
- Historically might also refer to models like **GPT-3** or even **BERT!**
 - The term has evolved, and people use it differently