

# Unpacking DPO and PPO: Disentangling Best Practices for Learning from Preference Feedback

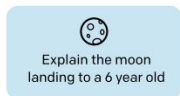
Geyang Guo and Duong Minh Le

# RLHF

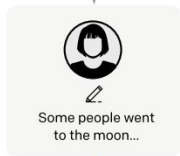
Step 1

**Collect demonstration data,  
and train a supervised policy.**

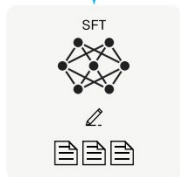
A prompt is  
sampled from our  
prompt dataset.



A labeler  
demonstrates the  
desired output  
behavior.



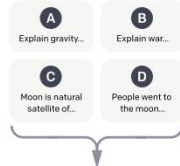
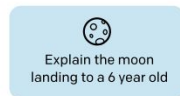
This data is used  
to fine-tune GPT-3  
with supervised  
learning.



Step 2

**Collect comparison data,  
and train a reward model.**

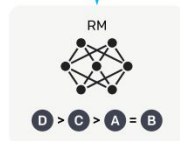
A prompt and  
several model  
outputs are  
sampled.



A labeler ranks  
the outputs from  
best to worst.



This data is used  
to train our  
reward model.



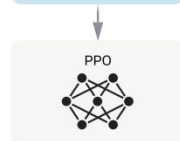
Step 3

**Optimize a policy against  
the reward model using  
reinforcement learning.**

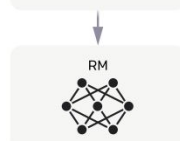
A new prompt  
is sampled from  
the dataset.



The policy  
generates  
an output.



The reward model  
calculates a  
reward for  
the output.



The reward is  
used to update  
the policy  
using PPO.

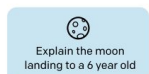


# RLHF

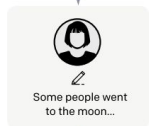
Step 1

**Collect demonstration data, and train a supervised policy.**

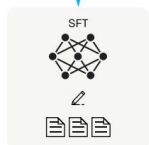
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



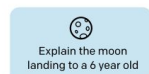
This data is used to fine-tune GPT-3 with supervised learning.



Step 2

**Collect comparison data, and train a reward model.**

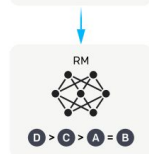
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



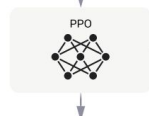
Step 3

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.



The policy generates an output.



Once upon a time...

The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.

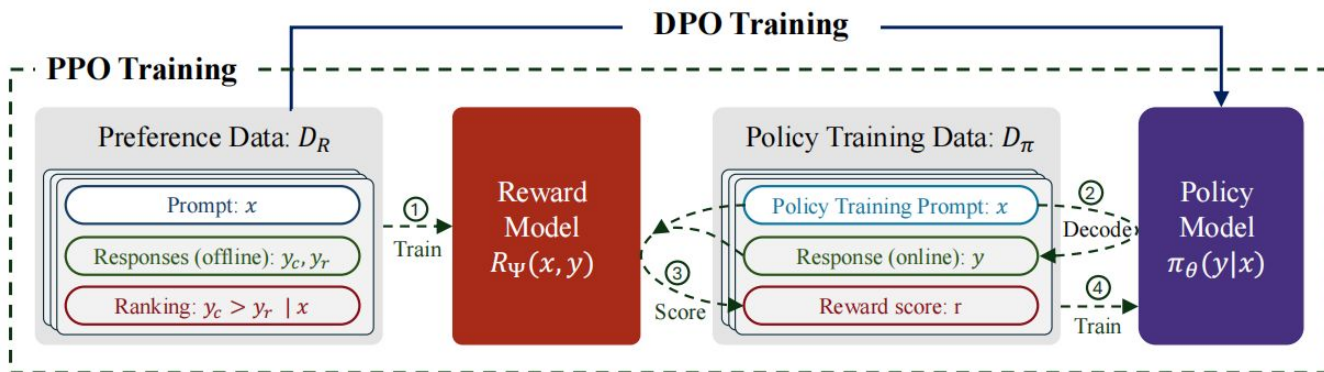


## Limitations?

1. Complex: involves training multiple LMs
2. Significant computational costs: need to sample from LM policy in the loop of training
3. Hard to implement: many hyper-parameters

# PPO -> DPO

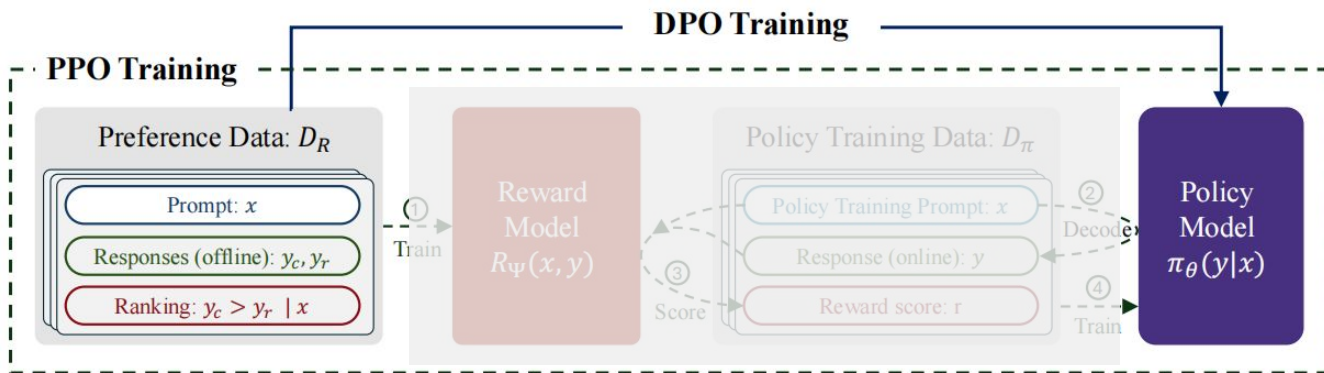
Intuition?



Directly optimize a language model to align to human preferences, without explicit reward modeling or reinforcement learning

# PPO -> DPO

Intuition?



Directly optimize a language model to align to human preferences, without explicit reward modeling or reinforcement learning

# DPO derivation

RLHF objective:  $\max_{\pi} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(y|x)} [r(x, y)] - \beta \mathbb{D}_{\text{KL}}(\pi(\cdot | x) || \pi_{\text{ref}}(\cdot | x))$

get high reward

stay close to reference model

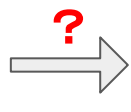
Goal: implicitly optimizes the same objective as existing RLHF algorithms

# DPO derivation

RLHF objective:  $\max_{\pi} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(y|x)} [r(x, y)] - \beta \mathbb{D}_{\text{KL}}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))$

get high reward Stay close to reference model

Goal: implicitly optimizes the same objective as existing RLHF algorithms

  $\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$

# Step 1: Got optimal policy

RLHF objective:  $\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}} [r(x, y)] - \beta \text{KL} [\pi_{\theta}(y|x), \pi_{\theta_{\text{old}}}(y|x)]$

get high reward                      Stay close to reference model

KL function  $\longleftrightarrow$

$$\min_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi_{\theta}(\cdot|x)} \left[ \log \frac{\pi_{\theta}(y|x)}{\pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)} \right]$$



# Step 1: Got optimal policy

RLHF objective:  $\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}} [r(x, y)] - \beta \text{KL} [\pi_{\theta}(y|x), \pi_{\theta_{\text{old}}}(y|x)]$

get high reward                      Stay close to reference model

KL function  $\longleftrightarrow$   $\min_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi_{\theta}(\cdot|x)} \left[ \log \frac{\pi_{\theta}(y|x)}{\pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)} \right]$

Denote  $Z(x) = \sum_y \pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$  and  $\pi^*(y|x) = \frac{1}{Z(x)} \pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$

(1)  $\pi^*(y|x) > 0 (\forall y > 0)$ ; (2)  $\sum_y \pi^*(y|x) = 1$

# Step 1: Got optimal policy

RLHF objective:  $\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}} [r(x, y)] - \beta \text{KL} [\pi_{\theta}(y|x), \pi_{\theta_{\text{old}}}(y|x)]$

$\underbrace{\hspace{10em}}_{\text{get high reward}} \quad \underbrace{\hspace{10em}}_{\text{Stay close to reference model}}$

KL function  $\longleftrightarrow$   $\min_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi_{\theta}(\cdot|x)} \left[ \log \frac{\pi_{\theta}(y|x)}{\pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)} \right]$

Denote  $Z(x) = \sum_y \pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$  and  $\pi^*(y|x) = \frac{1}{Z(x)} \pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$

(1)  $\pi^*(y|x) > 0 (\forall y > 0)$ ; (2)  $\sum_y \pi^*(y|x) = 1$

$\longleftrightarrow \min_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}} [\text{KL} [\pi_{\theta}(y|x), \pi^*(y|x)] - \log Z(x)]$

# Step 1: Got optimal policy

RLHF objective:  $\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}} [r(x, y)] - \beta \text{KL} [\pi_{\theta}(y|x), \pi_{\theta_{\text{old}}}(y|x)]$

get high reward                      Stay close to reference model

KL function  $\longleftrightarrow$   $\min_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi_{\theta}(\cdot|x)} \left[ \log \frac{\pi_{\theta}(y|x)}{\pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)} \right]$

Denote  $Z(x) = \sum_y \pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$  and  $\pi^*(y|x) = \frac{1}{Z(x)} \pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$

(1)  $\pi^*(y|x) > 0 (\forall y > 0)$ ; (2)  $\sum_y \pi^*(y|x) = 1$

$\longleftrightarrow$   $\min_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}} [\text{KL} [\pi_{\theta}(y|x), \pi^*(y|x)] - \log Z(x)]$

Optimal policy  $\longrightarrow$   $\pi_r(y|x) = \pi^*(y|x) = \frac{1}{Z(x)} \pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$

## Step 2: Write any reward function as function of optimal policy

Optimal policy:  $\pi_r(y|x) = \pi^*(y|x) = \frac{1}{Z(x)} \pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$

$$\longleftrightarrow^{\log} r(x, y) = \beta \log\left(\frac{\pi_r(y|x)}{\pi_{\theta_{\text{old}}}(y|x)}\right) + \beta \log(Z(x))$$

## Step 2: Write any reward function as function of optimal policy

Optimal policy:  $\pi_r(y|x) = \pi^*(y|x) = \frac{1}{Z(x)} \pi_{\theta_{\text{old}}}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$

$$\begin{array}{c} \text{log} \\ \longleftrightarrow \end{array} r(x, y) = \beta \log \left( \frac{\pi_r(y|x)}{\pi_{\theta_{\text{old}}}(y|x)} \right) + \beta \log (Z(x))$$

Ratio is positive if the optimal policy likes response **more** than reference model, negative if policy likes response **less** than reference model.

And the **more** the policy favors this response, the **higher** this ratio will be.

# Step 3: Loss function

A loss function on  
reward functions:

Derived from the Bradley-Terry model of human preferences:

$$\mathcal{L}_R(r, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r(x, y_w) - r(x, y_l))]$$

+

Transformation between  
reward functions and  
policies:

=

A loss function on  
policy:

# Step 3: Loss function

A loss function on  
reward functions:

Derived from the Bradley-Terry model of human preferences:

$$\mathcal{L}_R(r, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r(x, y_w) - r(x, y_l))]$$

+

Transformation between  
reward functions and  
policies:

$$r(x, y) = \beta \log \left( \frac{\pi_r(y|x)}{\pi_{\theta_{\text{old}}}(y|x)} \right) + \beta \log (Z(x))$$

=

A loss function on  
policy:

# Step 3: Loss function

A loss function on  
reward functions:

Derived from the Bradley-Terry model of human preferences:

$$\mathcal{L}_R(r, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r(x, y_w) - r(x, y_l))]$$

+

Transformation between  
reward functions and  
policies:

$$r(x, y) = \beta \log \left( \frac{\pi_r(y|x)}{\pi_{\theta_{\text{old}}}(y|x)} \right) + \beta \log (Z(x))$$

=

A loss function on  
policy:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$$



# Step 3: Loss function

A loss function on  
reward functions:

Derived from the Bradley-Terry model of human preferences:

$$\mathcal{L}_R(r, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r(x, y_w) - r(x, y_l))]$$

+

Transformation between  
reward functions and  
policies:

$$r(x, y) = \beta \log \left( \frac{\pi_r(y|x)}{\pi_{\theta_{\text{old}}}(y|x)} \right) + \beta \log (Z(x))$$

=

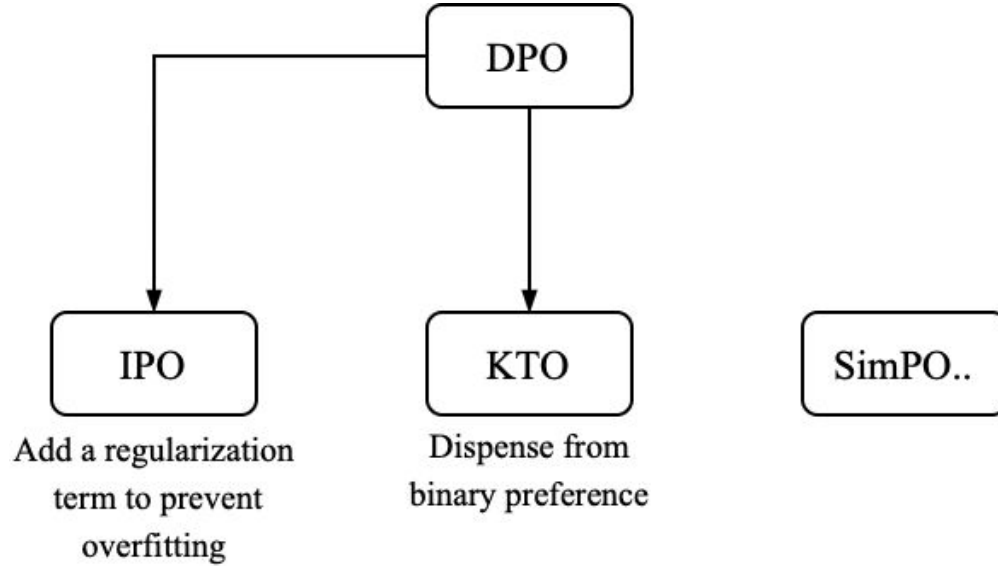
A loss function on  
policy:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$$

reward of preferred  
response

reward of dispreferred  
response

# Explorations beyond DPO



# PPO vs DPO

Contents	Compute and Speed	Exploration and Quality
<b>PPO</b>	<p>More complicated</p> <ol style="list-style-type: none"><li>1. Additional training of reward model and value model</li><li>2. Decode online responses during policy training</li></ol>	Trains on online data generated by the current policy
<b>DPO</b>	More efficient, stable	Trains on pre-generated offline data, thus limit exploration

# A Recipe for Learning from Preferences

Preference data

Preference Learning  
Algorithm

Reward model

Policy training  
prompt

# Experiments: Benchmarks

---

<b>Factuality</b>	<b>Reasoning</b>	<b>Coding</b>	<b>Truthfulness</b>	<b>Safety</b>	<b>Inst. Foll.</b>
-------------------	------------------	---------------	---------------------	---------------	--------------------

---

*MMLU*

*GSM8k*  
*Big Bench Hard*

*HumanEval+*  
*MBPP+*

*TruthfulQA*

*ToxiGen*  
*XSTest*

*AlpacaEval 1&2*  
*IFEval*

# Experiments: Preference Data

Source		# Samples	Factuality	Reasoning	Coding	Truthfulness	Safety	Inst. Following	Average
-	Llama 2 base	-	52.0	37.0	30.7	32.7	32.7	-	-
-	TÜLU 2 (SFT)	-	55.4	47.8	45.1	56.6	91.8	44.2	56.8
Web	SHP-2	500,000	55.4	47.7	40.3	62.2	90.4	45.6	56.9
	StackExchange	500,000	55.7	46.8	39.6	67.4	92.6	44.6	57.8
Human	PRM800k	6,949	55.3	49.7	46.6	54.7	91.9	43.4	56.9
	Chatbot Arena (2023)	20,465	55.4	50.2	45.9	58.5	67.3	50.8	54.7
	Chatbot Arena (2024)	34,269	55.7	50.4	37.7	56.7	58.1	50.7	51.5
	AlpacaF. Human Pref	9,686	55.3	47.6	43.3	56.1	90.7	44.5	56.2
	HH-RLHF	158,530	54.7	46.0	43.6	65.6	93.1	45.4	58.1
	HelpSteer	9,270	55.2	48.2	46.5	60.3	92.5	45.2	58.0
Synthetic	AlpacaF. GPT-4 Pref	19,465	55.3	49.1	43.4	57.7	89.5	46.3	56.9
	Capybara 7k	7,563	55.2	46.4	46.4	57.5	91.5	46.1	57.2
	Orca Pairs	12,859	55.5	46.8	46.0	57.9	90.5	46.2	57.2
	Nectar	180,099	55.3	47.8	43.2	68.2	93.1	47.8	59.2
	UltraF. (overall)	60,908	55.6	48.8	46.5	67.6	92.1	51.1	60.3
	UltraF. (fine-grained)	60,908	55.3	50.9	45.9	69.3	91.9	52.8	61.0

Table 1: **Preference data:** Performance of TÜLU 2 13B models trained on various preference datasets using DPO. **Blue** indicates improvements over the SFT baseline, **orange** degradations. Overall, synthetic data works best. DPO training improves truthfulness and instruction-following most, with limited to no improvements in factuality and reasoning.

# Experiments: Preference Data

Source		# Samples	Factuality	Reasoning	Coding	Truthfulness	Safety	Inst. Following	Average
-	Llama 2 base	-	52.0	37.0	30.7	32.7	32.7	-	-
-	TÜLU 2 (SFT)	-	55.4	47.8	45.1	56.6	91.8	44.2	56.8
Web	SHP-2	500,000	55.4	47.7	40.3	62.2	90.4	45.6	56.9
	StackExchange	500,000	55.7	46.8	39.6	67.4	92.6	44.6	57.8
Human	PRM800k	6,949	55.3	49.7	46.6	54.7	91.9	43.4	56.9
	Chatbot Arena (2023)	20,465	55.4	50.2	45.9	58.5	67.3	50.8	54.7
	Chatbot Arena (2024)	34,269	55.7	50.4	37.7	56.7	58.1	50.7	51.5
	AlpacaF. Human Pref	9,686	55.3	47.6	43.3	56.1	90.7	44.5	56.2
	HH-RLHF	158,530	54.7	46.0	43.6	65.6	93.1	45.4	58.1
	HelpSteer	9,270	55.2	48.2	46.5	60.3	92.5	45.2	58.0
Synthetic	AlpacaF. GPT-4 Pref	19,465	55.3	49.1	43.4	57.7	89.5	46.3	56.9
	Capybara 7k	7,563	55.2	46.4	46.4	57.5	91.5	46.1	57.2
	Orca Pairs	12,859	55.5	46.8	46.0	57.9	90.5	46.2	57.2
	Nectar	180,099	55.3	47.8	43.2	68.2	93.1	47.8	59.2
	UltraF. (overall)	60,908	55.6	48.8	46.5	67.6	92.1	51.1	60.3
	UltraF. (fine-grained)	60,908	55.3	50.9	45.9	69.3	91.9	52.8	61.0

Table 1: **Preference data:** Performance of TÜLU 2 13B models trained on various preference datasets using DPO. **Blue** indicates improvements over the SFT baseline, **orange** degradations. Overall, synthetic data works best. DPO training improves truthfulness and instruction-following most, with limited to no improvements in factuality and reasoning.

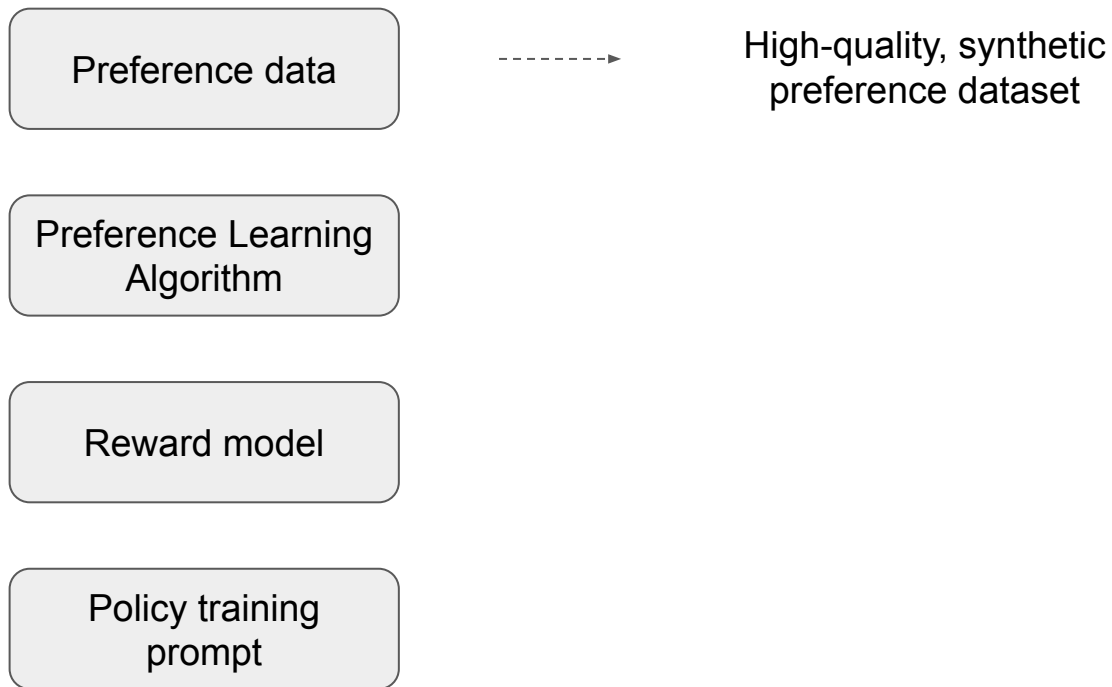
# Experiments: Preference Data

Source		# Samples	Factuality	Reasoning	Coding	Truthfulness	Safety	Inst. Following	Average
-	Llama 2 base	-	52.0	37.0	30.7	32.7	32.7	-	-
-	TÜLU 2 (SFT)	-	55.4	47.8	45.1	56.6	91.8	44.2	56.8
Web	SHP-2	500,000	55.4	47.7	40.3	62.2	90.4	45.6	56.9
	StackExchange	500,000	55.7	46.8	39.6	67.4	92.6	44.6	57.8
Human	PRM800k	6,949	55.3	49.7	46.6	54.7	91.9	43.4	56.9
	Chatbot Arena (2023)	20,465	55.4	50.2	45.9	58.5	67.3	50.8	54.7
	Chatbot Arena (2024)	34,269	55.7	50.4	37.7	56.7	58.1	50.7	51.5
	AlpacaF. Human Pref	9,686	55.3	47.6	43.3	56.1	90.7	44.5	56.2
	HH-RLHF	158,530	54.7	46.0	43.6	65.6	93.1	45.4	58.1
	HelpSteer	9,270	55.2	48.2	46.5	60.3	92.5	45.2	58.0
Synthetic	AlpacaF. GPT-4 Pref	19,465	55.3	49.1	43.4	57.7	89.5	46.3	56.9
	Capybara 7k	7,563	55.2	46.4	46.4	57.5	91.5	46.1	57.2
	Orca Pairs	12,859	55.5	46.8	46.0	57.9	90.5	46.2	57.2
	Nectar	180,099	55.3	47.8	43.2	68.2	93.1	47.8	59.2
	UltraF. (overall)	60,908	55.6	48.8	46.5	67.6	92.1	51.1	60.3
	UltraF. (fine-grained)	60,908	55.3	50.9	45.9	69.3	91.9	52.8	61.0

- Synthetic data with *per-aspect annotations* performs best (*i.e.*, *UltraF.*)
- Per-aspect annotations (*i.e.*, *UltraF*, *HelpSteer*): Datasets collected by first getting per-aspect annotations (e.g., helpfulness, harmlessness) then averaging



# A Recipe for Learning from Preferences



# Experiments: Preference Learning Algorithm (DPO vs. PPO)

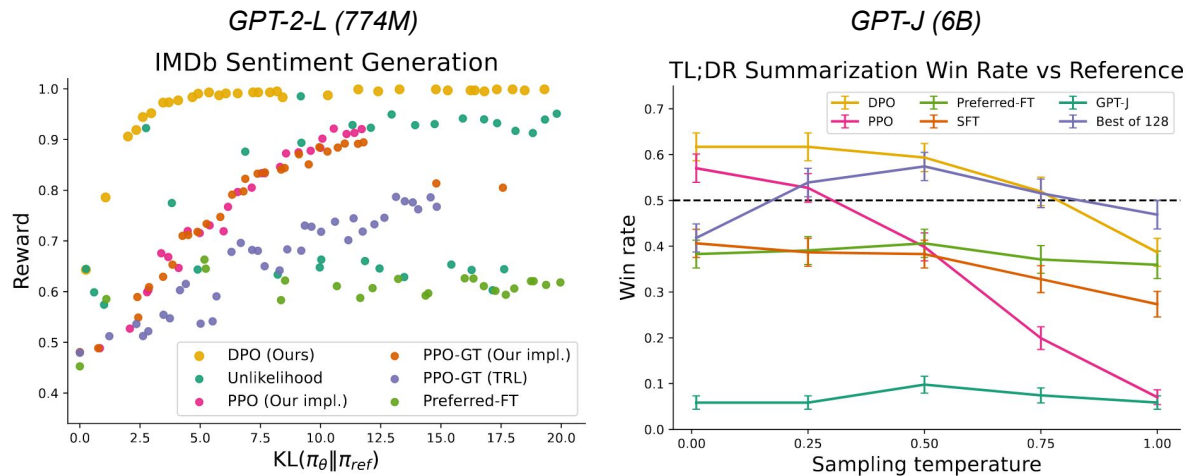


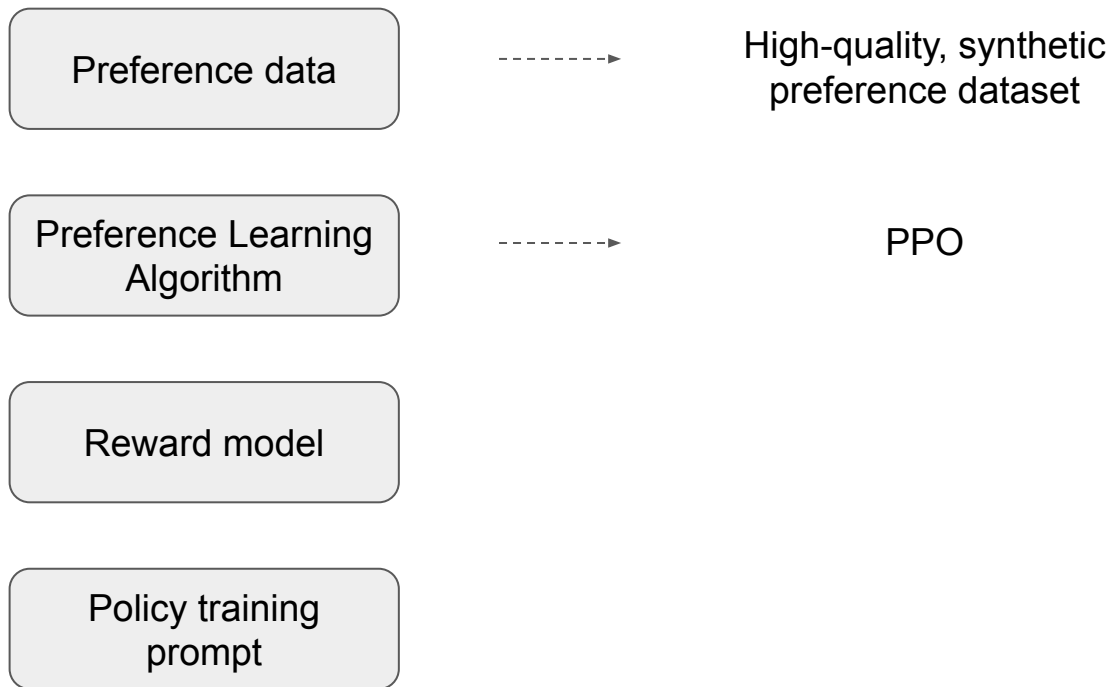
Figure 2: **Left.** The frontier of expected reward vs KL to the reference policy. **DPO provides the highest expected reward for all KL values**, demonstrating the quality of the optimization. **Right.** TL;DR summarization win rates vs. human-written summaries, using GPT-4 as evaluator. **DPO exceeds PPO’s best-case performance on summarization**, while being more robust to changes in the sampling temperature.

# Experiments: Preference Learning Algorithm (DPO vs. PPO)

Data / Model	Alg.	Factuality	Reasoning	Coding	Truthfulness	Safety	Inst. Foll.	Average
Llama 2 base	-	52.0	37.0	30.7	32.7	32.7	-	-
TÜLU 2 (SFT)	-	55.4	47.8	45.1	56.6	91.8	44.2	56.8
StackExchange	DPO	<b>55.3</b>	<b>47.8</b>	42.4	<b>56.2</b>	92.0	46.7	56.7
	PPO	55.1	<b>47.8</b>	<b>46.4</b>	54.2	<b>92.6</b>	<b>47.4</b>	<b>57.3</b>
ChatArena (2023)	DPO	<b>55.4</b>	<b>50.2</b>	45.9	<b>58.5</b>	67.3	<b>50.8</b>	54.7
	PPO	55.2	49.2	<b>46.4</b>	55.8	<b>79.4</b>	49.7	<b>55.9</b>
HH-RLHF	DPO	<b>55.2</b>	47.6	44.2	<b>60.0</b>	<b>93.4</b>	46.6	57.8
	PPO	54.9	<b>48.6</b>	<b>45.9</b>	58.0	92.8	<b>47.0</b>	<b>57.9</b>
Nectar	DPO	<b>55.6</b>	45.8	39.0	<b>68.1</b>	<b>93.3</b>	<b>48.4</b>	58.4
	PPO	55.2	<b>51.2</b>	<b>45.6</b>	60.1	92.6	47.4	<b>58.7</b>
UltraFeedback (FG)	DPO	55.3	50.9	45.9	69.3	<b>91.9</b>	52.8	61.0
	PPO	<b>56.0</b>	<b>52.0</b>	<b>47.7</b>	<b>71.5</b>	91.8	<b>54.4</b>	<b>62.2</b>
Avg. $\Delta$ b/w PPO & DPO		-0.1	+1.3	+2.9	-2.5	+2.3	+0.1	+0.7

Table 2: **DPO vs PPO**: Average performance of 13B models trained using DPO and PPO across different datasets, along with the performance difference between DPO and PPO ( $\Delta$ ). Blue indicates improvements over the SFT baseline, orange degradations. All datasets are downsampled to 60,908 examples (except ChatArena, which is made up of 20,465 responses). **PPO outperforms DPO by an average of 0.7 points, where most improvements are in reasoning, coding, and chat capabilities.**

# A Recipe for Learning from Preferences



# Experiments: Reward Models

- Scaling up the training data for RM:
  - **Mix RM:** Construct a data mixture of the top performing preference datasets (i.e., UltraFeedback, HelpSteer, Nectar, StackExchange, HH-RLHF, PRM800k)
  - **UltraF. RM:** Reward model trained only on UltraFeedback
  
- Scaling up the reward model size: 13B and 70B

# Experiments: Reward Models

Reward Model	Direct Eval.	
	RewardBench Score	Best-of-N over SFT Avg. Perf. ( $\Delta$ )
13B UltraF. RM	61.0	56.9 (+5.8)
13B Mix RM	<b>79.8</b>	58.3 (+7.3)
70B UltraF. RM	73.6	<b>61.1 (+10.3)</b>
70B Mix RM	73.9	60.6 (+9.5)

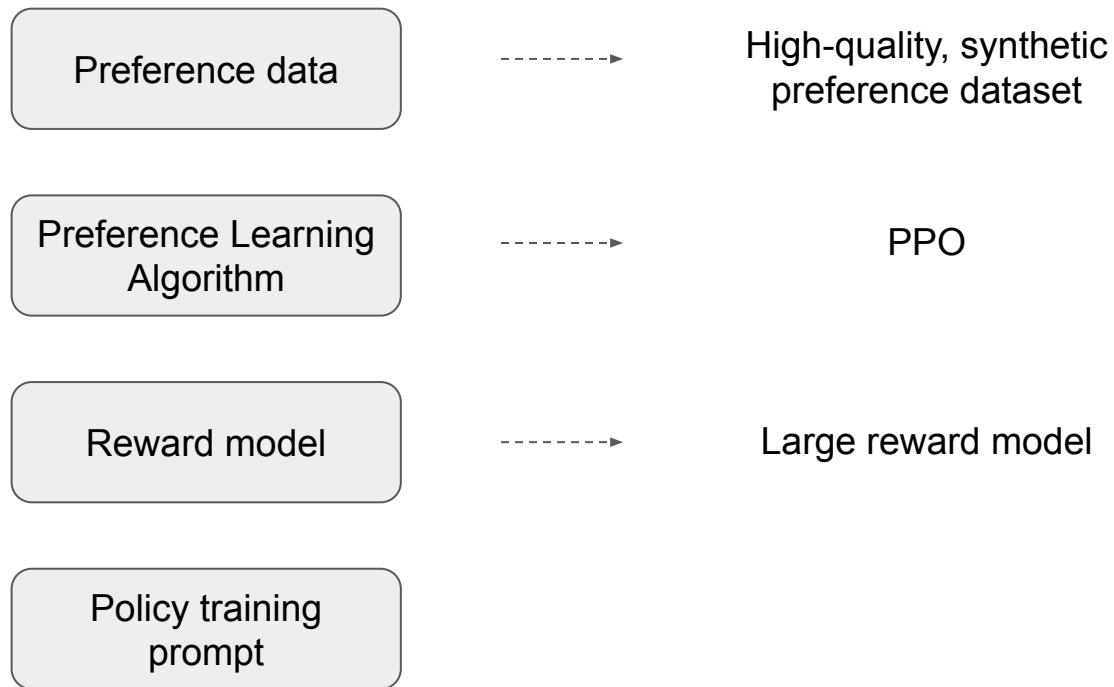
- Best-of-N:
  - Sample 16 responses for each evaluation task
  - Pick the top-scoring response according to the RM as the final output
- RewardBench: evaluating if the relative scores given by reward models match a test set of chosen-rejected pairs from diverse sources

# Experiments: Reward Models

Reward Model	Direct Eval.		PPO Training Perf. (w. UltraF. prompts)		
	RewardBench Score	Best-of-N over SFT Avg. Perf. ( $\Delta$ )	GSM Acc.	AlpacaEval2 winrate	Avg. on All Evals.
13B UltraF. RM	61.0	56.9 (+5.8)	53.0	26.1	62.2
13B Mix RM	<b>79.8</b>	58.3 (+7.3)	51.0	25.7	61.6
70B UltraF. RM	73.6	<b>61.1 (+10.3)</b>	<b>58.0</b>	26.7	<b>62.8</b>
70B Mix RM	73.9	60.6 (+9.5)	51.5	<b>31.6</b>	61.8

- It is difficult to translate improvements in reward models to the underlying policy
- Increasing scale and data improves reward models, but these only minimally impact the average downstream performance

# A Recipe for Learning from Preferences





# Experiments: Policy Training Prompts

- UltraF. Prompts: 20 random prompts from UltraFeedback
- Mined Math: math-related prompts from varied dataset
- GSM train: prompts from GASM train set

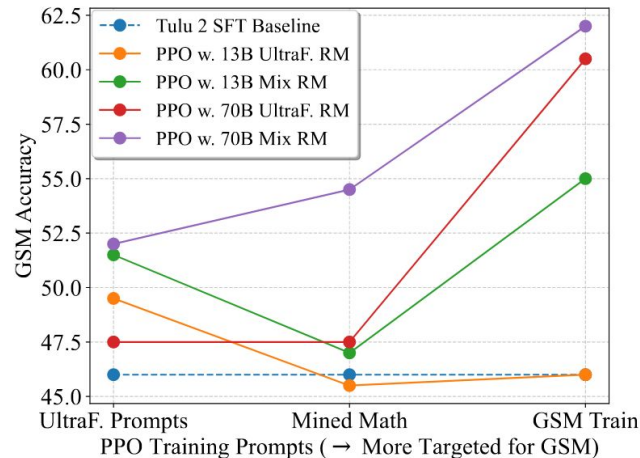


Figure 3: **Policy training prompt math evaluation:** Performance of models trained on 20K prompts from varying sources using PPO and evaluated on GSM. Training with larger RMs trained on more data benefits more from in-domain prompts (i.e., prompts directly from the GSM train set), while weaker RMs struggle to generalize beyond their training prompts.

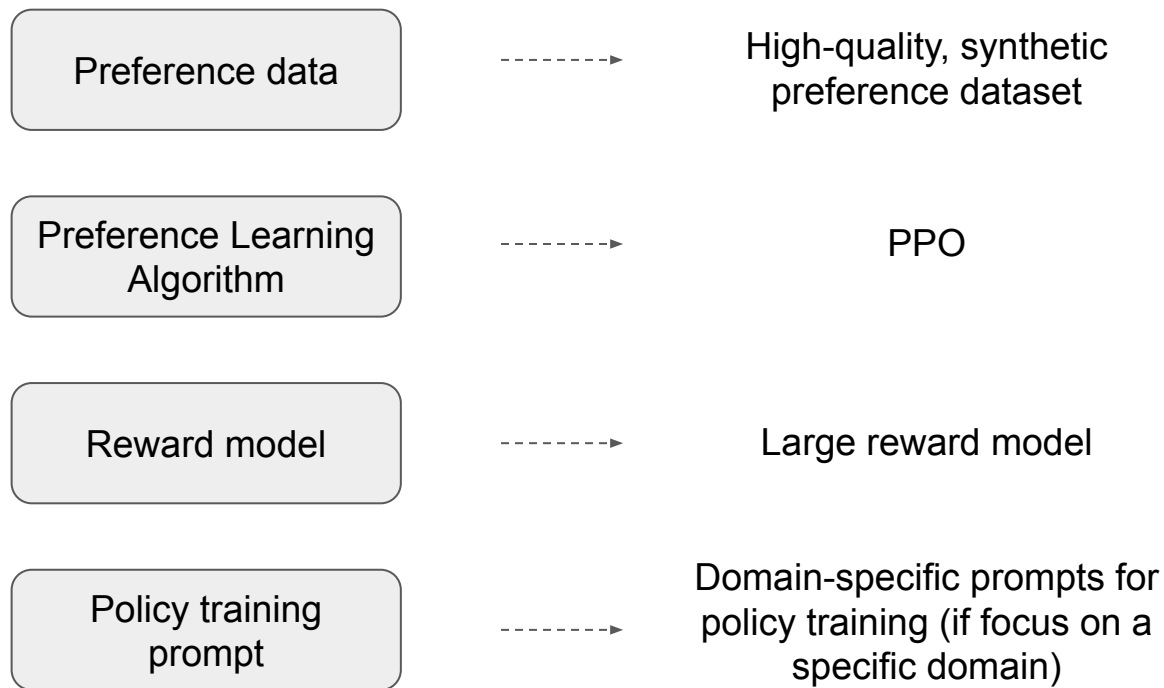
# Experiments: Policy Training Prompts

- **Mixed**: mix math- and code-related prompts with UltraFeedback prompts, then downsample to the same size of UF
- **UF**: UltraFeedback

Reward Model	Prompts	GSM %	Coding	Avg. Across All Evals
Tulu 2 SFT	-	46.0	45.1	56.8
13B UltraF.	UF	53.0	47.7	<b>62.2</b>
13B UltraF.	Mixed	<b>54.5</b>	<b>47.8</b>	61.9
13B Mix	UF	<b>51.0</b>	<b>46.8</b>	<b>61.6</b>
13B Mix	Mixed	50.5	43.8	60.9
70B UltraF.	UF	<b>58.0</b>	47.3	<b>62.8</b>
70B UltraF.	Mixed	56.5	<b>48.4</b>	62.4
70B Mix	UF	51.5	<b>46.1</b>	<b>61.8</b>
70B Mix	Mixed	<b>52.0</b>	44.9	61.1

Table 4: **Policy training prompt overall evaluation:** Performance of PPO policy models trained with the given reward models on 60K prompts from either UltraFeedback or the remixed prompt set that adds additional unlabeled math and coding-related prompts. **Using the remixed prompt set does not improve performance, either on specific evaluations (math, code) or in terms of overall performance.**

# A Recipe for Learning from Preferences



# A Recipe for Learning from Preferences

Model	Factuality	Reasoning	Coding	Truthfulness	Safety	Instr. Foll.	Average
Llama 2 Chat 13B [52]	53.2	24.7	36.9	<b>88.0</b>	91.9	51.2	57.7
Nous Hermes 13B [51]	53.2	43.5	47.7	80.5	43.9	38.7	51.3
Vicuna 1.5 13B [64]	54.5	39.3	38.5	62.8	92.4	45.8	55.6
Llama 2 13B Base	52.0	37.0	30.7	32.7	32.7	-	-
TÜLU 2 13B SFT	55.4	47.8	45.1	56.6	91.8	44.2	56.8
TÜLU 2+DPO 13B	55.3	50.9	45.9	69.3	91.9	52.8	61.0
TÜLU 2+PPO 13B (13B UFRM)	<b>56.0</b>	52.0	47.7	71.5	91.8	54.4	62.2
TÜLU 2+PPO 13B (70B UFRM)	55.4	<b>53.9</b>	47.3	72.3	91.9	<b>55.8</b>	<b>62.8</b>
TÜLU 2+PPO 13B (70B UFRM+MP)	55.3	53.1	<b>48.4</b>	71.0	<b>92.7</b>	54.0	62.4

Table 5: **Putting together a recipe for preference-based learning:** Performance of our best-performing models along with popular open models based on Llama 2 13B. ‘MP’ refers to using the mixed prompt set described in §4. Using PPO with a large reward model performs best overall.

# Summary

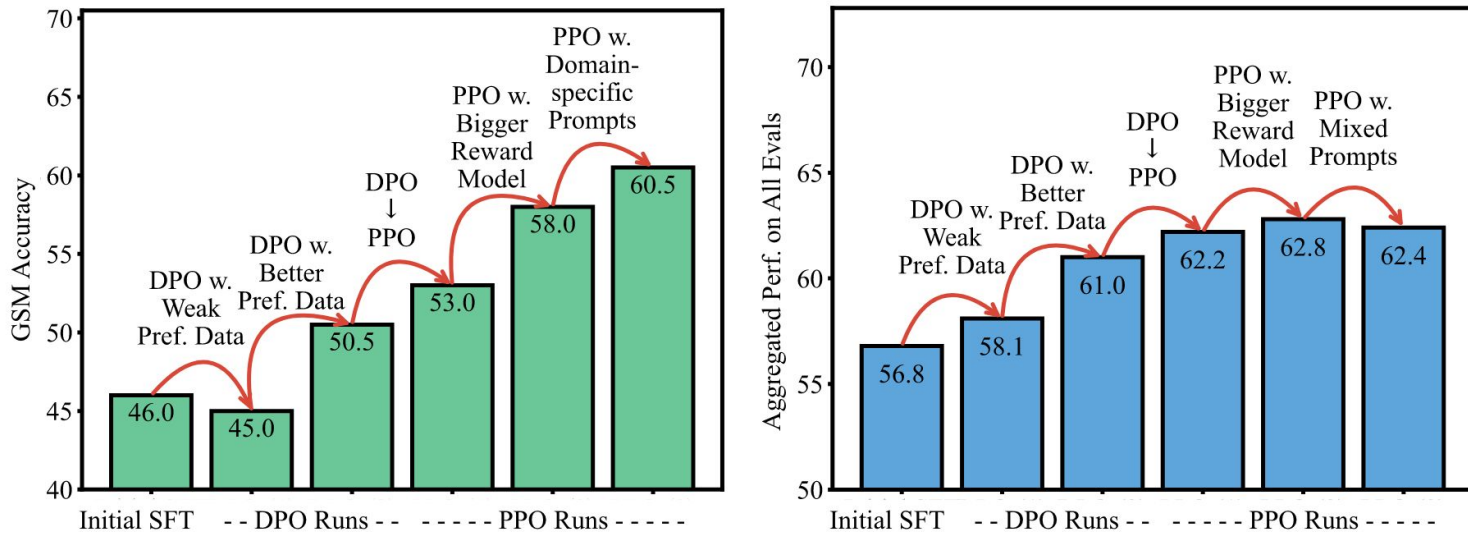


Figure 1: Performance improvements resulted by changing different components in the preference training of TULU. Left: Accuracy on GSM [9], for testing math capabilities. Right: Overall performance, aggregated over the 11 benchmarks described in §2.2.