



Reinforcement Learning in Post-Training

Bridging the Gap in News Production: Aligning AI Models

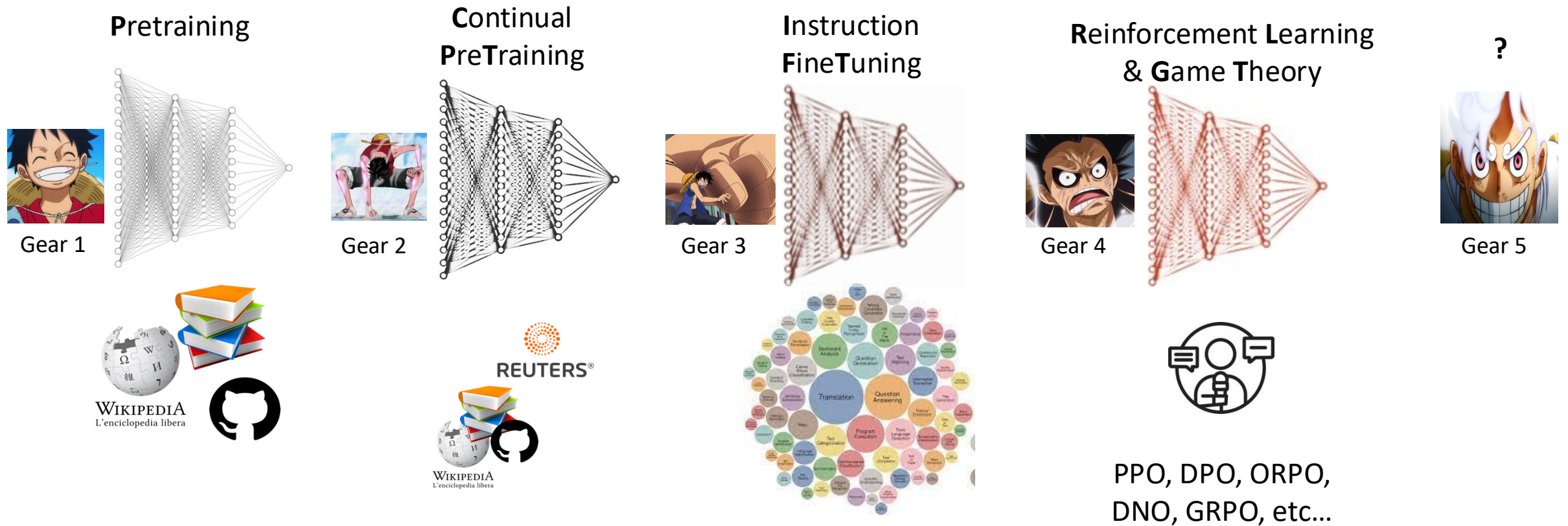
GIOFRÉ Daniele

Contents

- LLM Training – Overview
- RL in Post-Training Alignment
 - InstructGPT – PPO
 - DPO
 - RPO
 - PPO vs DPO/RPO
 - OnlineDPO
 - Reward Models vs LLM-as-Judge

LLM Training - Overview

The Gear of Training



RL in Post-Training Alignment

The Gear 4

Making a LLM better at following instructions and at reasoning:

- Evolution from InstructGPT to modern approaches
- Focus on PPO and DPO methods
- Bypass DPO limitations
- Make the bridge between PPO and DPO via Online DPO

InstructGPT and PPO Framework

Schema

- Three-stage process:
 - Supervised fine-tuning (SFT)
 - Reward Modeling (RM)
 - PPO optimization
- Key innovation:
Using human feedback for alignment

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

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

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

r_k

OUYANG, Long, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 2022, 35: 27730

InstructGPT and PPO Framework

How it Works

- Policy optimization using RL
- Iterative process:
 - Sample responses from policy
 - Evaluate with reward model
 - Update policy with PPO loss
- Ensures controlled policy updates

$$\text{objective}(\phi) = \mathbb{E}_{(x,y) \sim D_{\pi_{\phi}^{RL}}} \left[r_{\theta}(x, y) - \beta \log \frac{\pi_{\phi}^{RL}(y|x)}{\pi_{SFT}(y|x)} \right] + \gamma \mathbb{E}_{x \sim D_{\text{pretrain}}} \log(\pi_{\phi}^{RL}(x))$$

Step 3

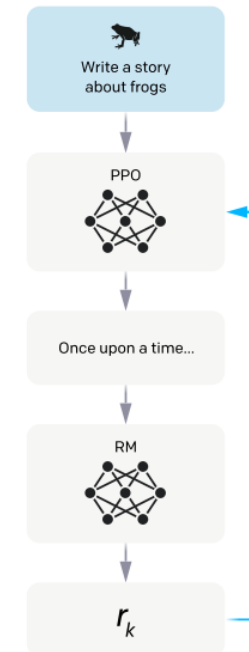
Optimize a policy against the reward model using reinforcement learning.

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.



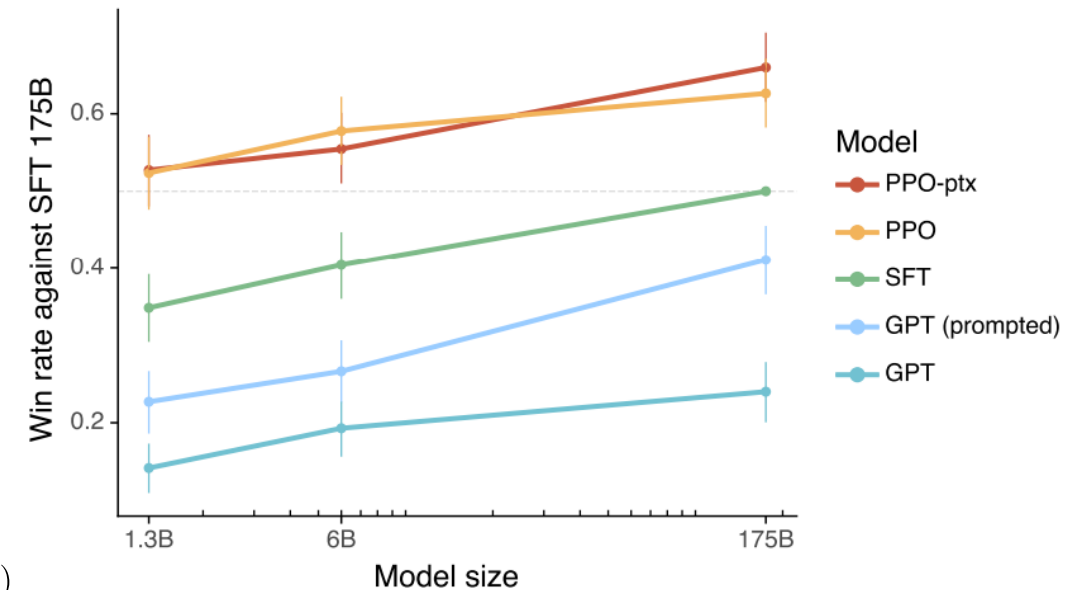
OUYANG, Long, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 2022, 35: 27730

InstructGPT and PPO Framework

Does it Work?

- Policy optimization using RL
- Iterative process:
 - Sample responses from policy
 - Evaluate with reward model
 - Update policy with PPO loss
- Ensures controlled policy updates

$$\text{objective}(\phi) = \mathbb{E}_{(x,y) \sim D_{\pi_{\phi}^{RL}}} \left[r_{\theta}(x, y) - \beta \log \frac{\pi_{\phi}^{RL}(y|x)}{\pi_{SFT}(y|x)} \right] + \gamma \mathbb{E}_{x \sim D_{\text{pretrain}}} \log(\pi_{\phi}^{RL}(x))$$



OUYANG, Long, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 2022, 35: 27730

InstructGPT and PPO Framework

Challenges and Limitations

- Resource intensive: 3 models in memory
 - SFT model
 - Current policy
 - Reward model
- Requires on-policy data collection
- Complex training dynamics – 3 steps
- High computational cost

$$\text{objective}(\phi) = \mathbb{E}_{(x,y) \sim D_{\pi_{\phi}^{RL}}} \left[\underbrace{r_{\theta}(x, y)}_{\text{purple}} - \beta \log \underbrace{\frac{\pi_{\phi}^{RL}(y|x)}{\pi^{SFT}(y|x)}}_{\text{blue}} \right] + \gamma \mathbb{E}_{x \sim D_{\text{pretrain}}} \log(\underbrace{\pi_{\phi}^{RL}(x)}_{\text{orange}})$$

OUYANG, Long, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 2022, 35: 27730

Direct Preference Optimization (DPO)

- Single-stage optimization: “from reward functions to optimal policies”
- Directly learns from preference data:
- No explicit reward model needed, but implicit rewards: $\hat{r}_\theta(x, y) = \beta \log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)}$
- Loss based on a Berry-Terry pair-wise modelling:

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

RAFAILOV, Rafael, et al. Direct Preference Optimization. Advances in Neural Information Processing Systems, 2024, 36.

Direct Preference Optimization (DPO)

Limitations

- Training instability if $\hat{r}_{\theta,c}(x, y) \simeq \hat{r}_{\theta,r}(x, y)$ weak learning
- Overoptimization in RL
- Unclear separation between chosen and rejected
- Offline training – off-policy issue
- Double samples for same micro batch – pair-wise comparison
- Dataset quality dependency

RAFAILOV, Rafael, et al. Direct Preference Optimization. Advances in Neural Information Processing Systems, 2024, 36.

Direct Preference Optimization (DPO)

Limitations

- Training instability if $\hat{r}_{\theta,c}(x, y) \simeq \hat{r}_{\theta,r}(x, y)$ weak learning
- Overoptimization in RL
- Unclear separation between chosen and rejected
- Offline training – all off-policy training issues
- ~~Double samples for same micro batch – pair-wise comparison~~
- ~~Dataset quality dependency~~

RAFAILOV, Rafael, et al. Direct Preference Optimization. Advances in Neural Information Processing Systems, 2024, 36.

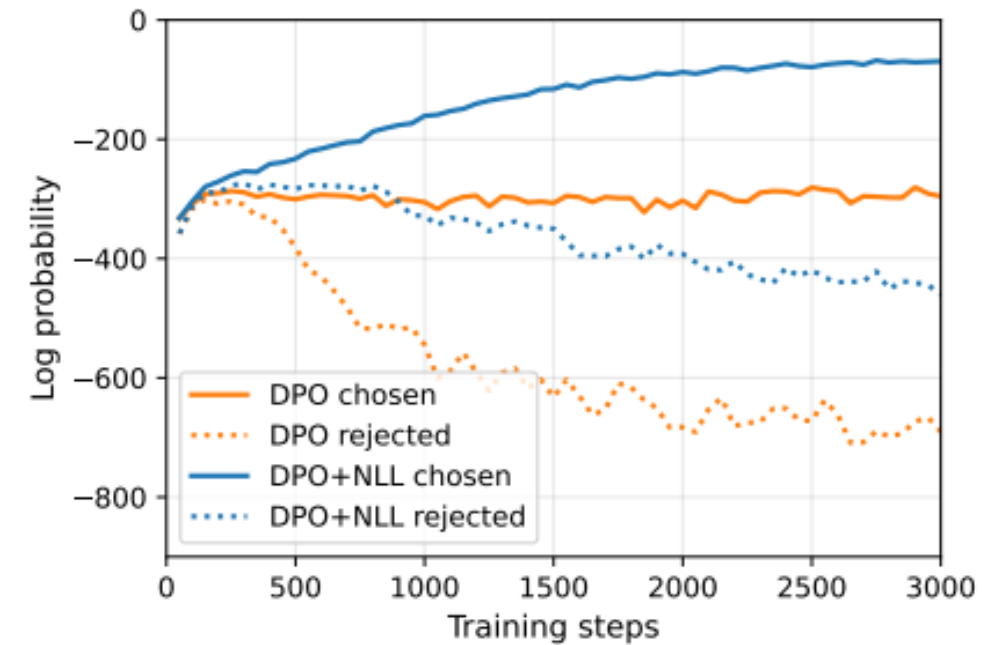
Regularized Preference Optimization (RPO)

SFT Loss is Implicitly an Adversarial Regularizer

RPO objective = Preference optimization loss + Imitation (SFT) loss

$$\mathcal{L}_{\text{DPO+NLL}} = \mathcal{L}_{\text{DPO}}(c_i^c, y_i^c, c_i^r, y_i^r | x_i) + \alpha \mathcal{L}_{\text{NLL}}(c_i^c, y_i^c | x_i)$$

Win rate (%)	RPO (beta)	Ref. (beta)	DPO (beta)
RPO (beta)	50.0	79.0	56.0
Ref. (beta)	21.0	50.0	22.7
DPO (beta)	44.0	77.3	50.0



(a) ARC-Challenge

LIU, Zhihan, et al. Provably mitigating overoptimization in RLHF. *arXiv preprint arXiv:2405.16436*, 2024.

PANG, Richard Yuanzhe, et al. Iterative reasoning preference optimization. *arXiv preprint arXiv:2404.19733*, 2024.

Regularized Preference Optimization (RPO)

By passing some limits of DPO

- ~~Training instability if $\hat{r}_{\theta,c}(x, y) \simeq \hat{r}_{\theta,r}(x, y)$~~ ~~weak learning~~
- ~~Overoptimization in RL~~
- Unclear separation between chosen and rejected [tackled in *]
- Offline training – all off-policy training issues

$$\mathcal{L}_{\text{DPO+NLL}} = \mathcal{L}_{\text{DPO}}(c_i^c, y_i^c, c_i^r, y_i^r | x_i) + \alpha \mathcal{L}_{\text{NLL}}(c_i^c, y_i^c | x_i)$$

LIU, Zhihan, et al. Provably mitigating overoptimization in RLHF. *arXiv preprint arXiv:2405.16436*, 2024.

* ADLER, Bo, et al. Nemotron-4 340B Technical Report. *arXiv preprint arXiv:2406.11704*, 2024.

PPO vs DPO Comparison

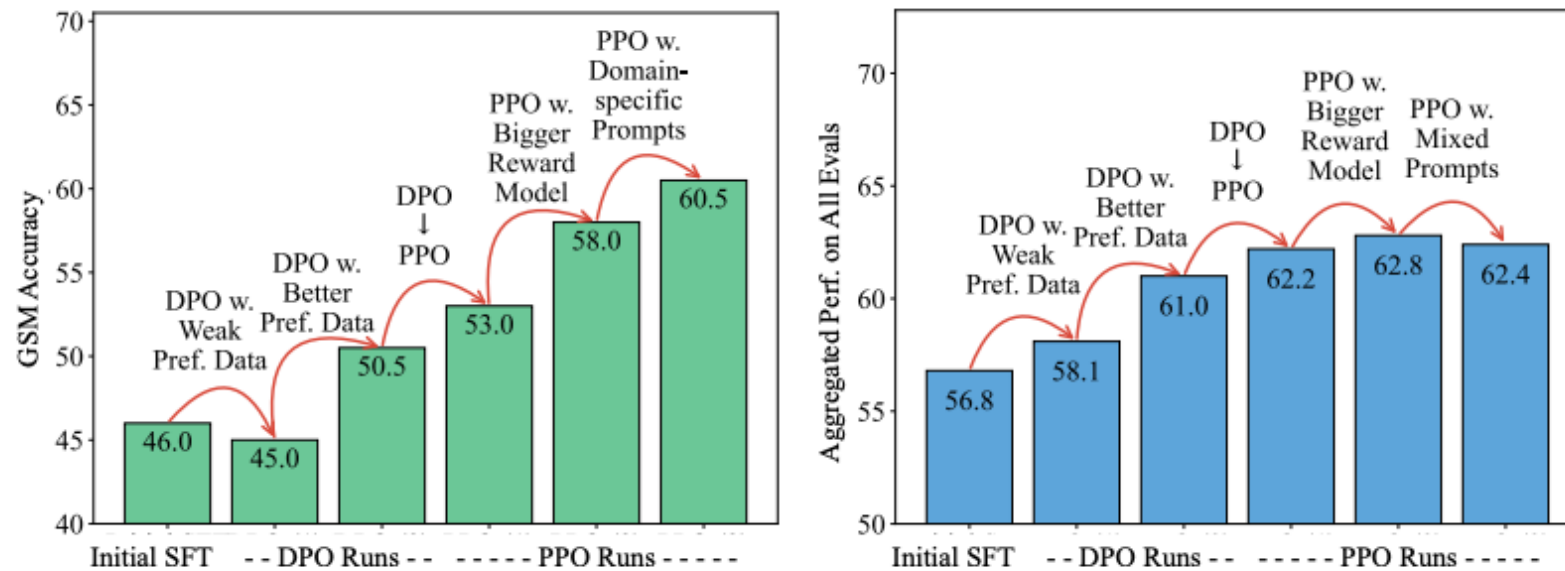


Figure 1: Performance improvements resulted by changing different components in the preference training of TüLU. Left: Accuracy on GSM [9], for testing math capabilities. Right: Overall performance, aggregated over the 11 benchmarks described in §2.2.

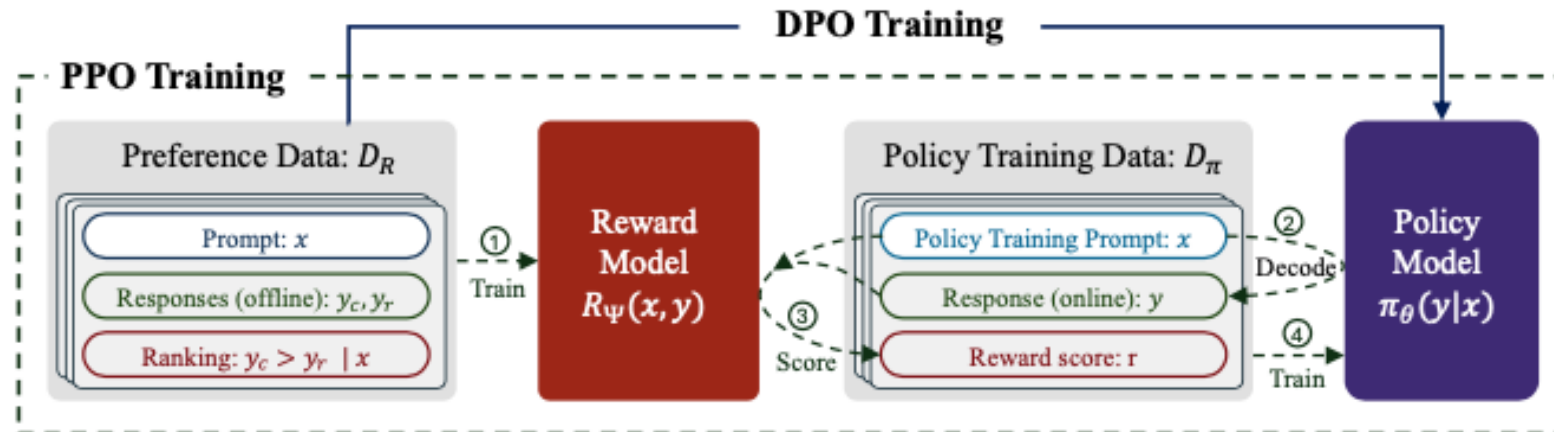
PPO vs DPO Comparison

PPO

- Online generation – Online Learning
- Requires reward model
- More stable training and slightly better metrics
- Higher computational cost

DPO/RPO

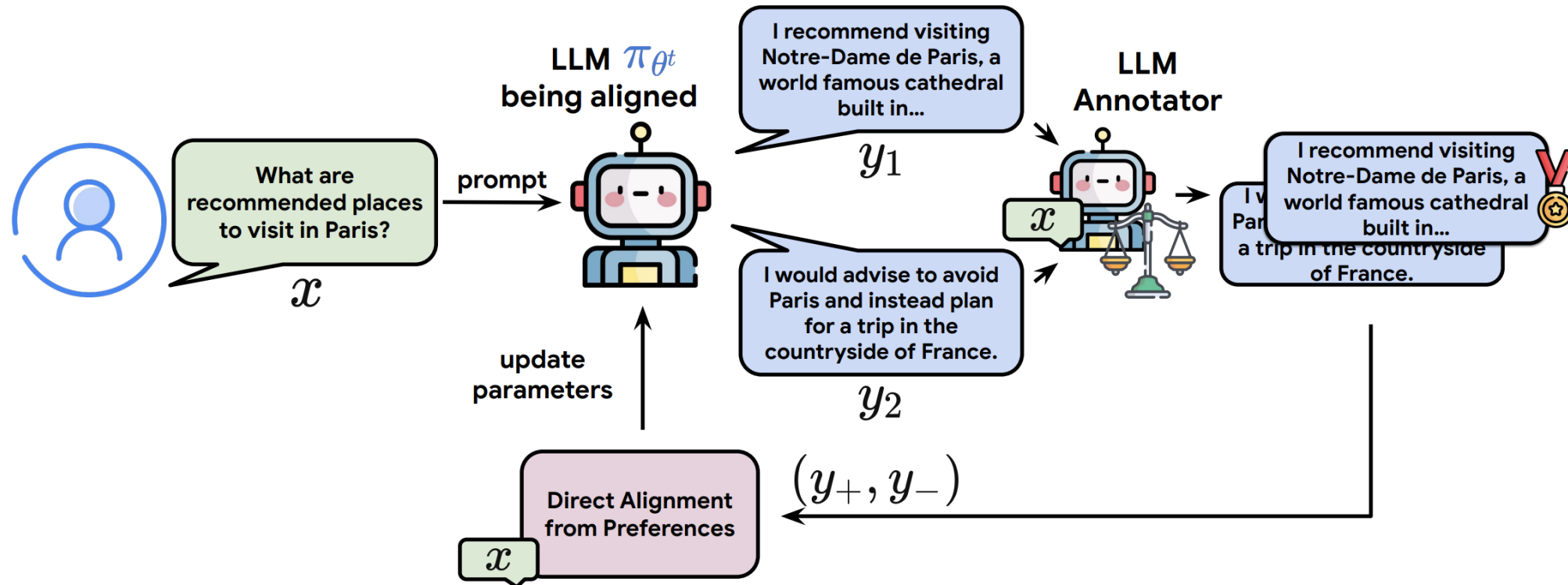
- Offline data allowed – Offline Learning
- No reward model needed
- Simpler implementation and negative samples as zero
- More efficient and fast training



IVISON, Hamish, et al. Unpacking DPO and PPO. arXiv preprint arXiv:2406.09279, 2024.

Online DPO

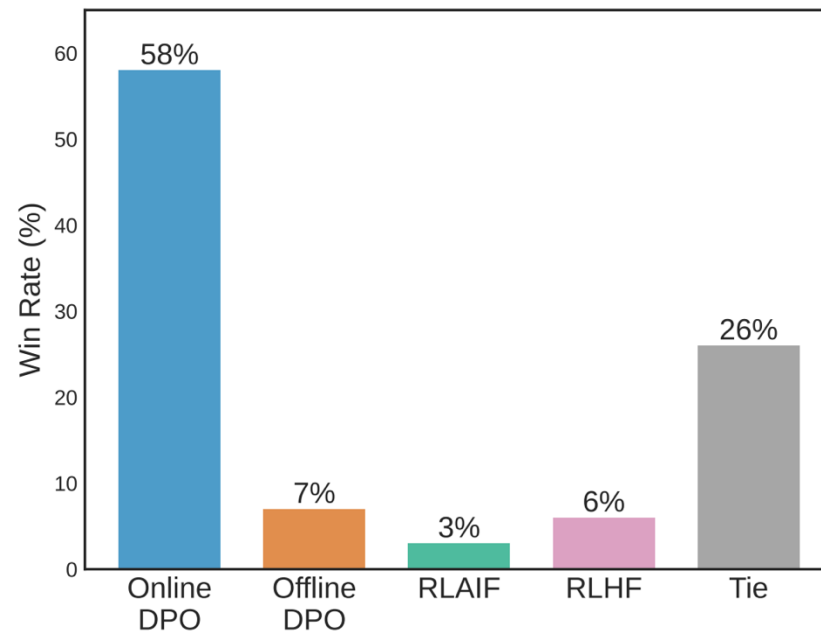
Bridges gap with PPO



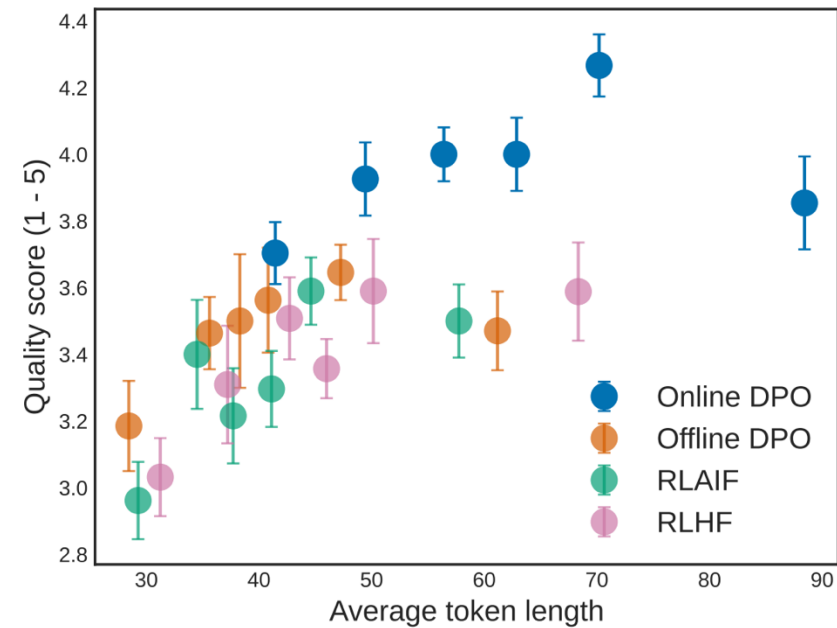
GUO, Shangmin, et al. Direct language model alignment from online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.

Online DPO

Bridges gap with PPO



(a) Fraction of responses preferred by humans



(b) Quality against length of responses

GUO, Shangmin, et al. Direct language model alignment from online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.

Online DPO

Reward Model or LLM-as-Judge?



▲	Model	▲	Model Type	▲	Score	▲	Chat	▲	Chat Hard	▲	Safety	▲	Reasoning	▲
1	infly/INF-ORM-Llama3.1-70B		Seq. Classifier		95.1		96.6		91.0		93.6		99.1	
2	nicolinho/ORM-Gemma-2-27B		Seq. Classifier		94.4		96.6		90.1		92.7		98.3	
3	Skywork/Skywork-Reward-Gemma-2-27B-v0.2		Seq. Classifier		94.3		96.1		89.9		93.0		98.1	
4	nvidia/Llama-3.1-Nemotron-70B-Reward *		Custom Classifier		94.1		97.5		85.7		95.1		98.1	
5	Skywork/Skywork-Reward-Gemma-2-27B ⚠		Seq. Classifier		93.8		95.8		91.4		91.9		96.1	
6	SF-Foundation/TextEval-Llama3.1-70B * ⚠		Generative		93.5		94.1		90.1		93.2		96.4	
7	meta-metrics/MetaMetrics-RM-v1.0		Custom Classifier		93.4		98.3		86.4		90.8		98.2	
8	Skywork/Skywork-Critic-Llama-3.1-70B ⚠		Generative		93.3		96.6		87.9		93.1		95.5	
9	nicolinho/ORM-Llama3.1-8B-v2		Seq. Classifier		93.1		96.4		86.8		92.6		96.8	
10	Skywork/Skywork-Reward-Llama-3.1-8B-v0.2		Seq. Classifier		93.1		94.7		88.4		92.7		96.7	

RewardBench: <https://huggingface.co/spaces/allenai/reward-bench>

We are hiring!

Join our Labs and Foundational Research team!

Hiring in **EMEA** (**Switzerland** and UK) & **North America** for :

- Applied Scientists (Interns, ICs, Team Leads)
- Research Scientists (Interns, ICs, Team Leads)
- Research Engineers (Interns, ICs)



Appendix



Direct Preference Optimization (DPO)

Digging into the loss via its Gradient

- Single-stage optimization
- Directly learns from preference data
- No explicit reward model needed, but implicit rewards
- Mathematical formulation:

$$\hat{r}_\theta(x, y) = \beta \log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)}$$

$$\nabla_\theta \mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\beta \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\underbrace{\sigma(\hat{r}_\theta(x, y_l) - \hat{r}_\theta(x, y_w))}_{\text{higher weight when reward estimate is wrong}} \left[\underbrace{\nabla_\theta \log \pi(y_w | x)}_{\text{increase likelihood of } y_w} - \underbrace{\nabla_\theta \log \pi(y_l | x)}_{\text{decrease likelihood of } y_l} \right] \right]$$

RAFAILOV, Rafael, et al. Direct Preference Optimization. Advances in Neural Information Processing Systems, 2024, 36.

Regularized Preference Optimization (RPO)

SFT Loss is Implicitly an Adversarial Regularizer

RPO objective = Preference optimization loss + Imitation (SFT) loss

$$\begin{aligned}\mathcal{L}_{\text{DPO+NLL}} &= \mathcal{L}_{\text{DPO}}(c_i^w, y_i^w, c_i^l, y_i^l | x_i) + \alpha \mathcal{L}_{\text{NLL}}(c_i^w, y_i^w | x_i) \\ &= -\log \sigma \left(\beta \log \frac{M_\theta(c_i^w, y_i^w | x_i)}{M_t(c_i^w, y_i^w | x_i)} - \beta \log \frac{M_\theta(c_i^l, y_i^l | x_i)}{M_t(c_i^l, y_i^l | x_i)} \right) - \alpha \frac{\log M_\theta(c_i^w, y_i^w | x_i)}{|c_i^w| + |y_i^w|}.\end{aligned}$$

LIU, Zhihan, et al. Provably mitigating overoptimization in RLHF. *arXiv preprint arXiv:2405.16436*, 2024.

PANG, Richard Yuanzhe, et al. Iterative reasoning preference optimization. *arXiv preprint arXiv:2404.19733*, 2024.

Online DPO

Bridges gap with PPO

Method	Win	Tie	Loss	Quality
TL;DR				
Online DPO	63.74%	28.57%	7.69%	3.95
Offline DPO	7.69%		63.74%	3.46
Helpfulness				
Online DPO	58.60%	21.20%	20.20%	4.08
Offline DPO	20.20%		58.60%	3.44
Harmlessness				
Online DPO	60.26%	35.90%	3.84%	4.41
Offline DPO	3.84%		60.26%	3.57

Table 2: Win/tie/loss rate of DPO with OAIF (online DPO) against vanilla DPO (offline DPO) on the TL;DR, Helpfulness, Harmlessness tasks, along with the quality score of their generations, judged by *human raters*.

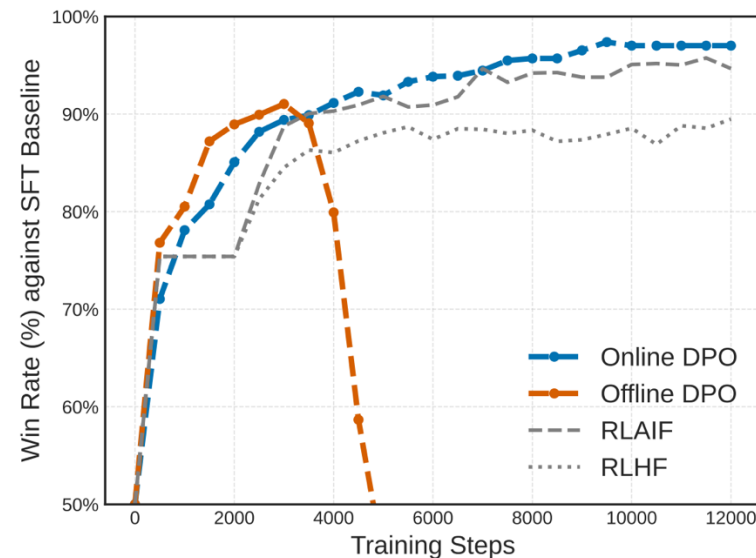


Figure 3: Win rate of DPO with OAIF (online DPO), vanilla DPO (offline DPO), RLAIF, and RLHF against the SFT baseline on the TL;DR task, judged by *Gemini Pro*.

Method	No RM needed	On-policy generation	Online feedback
Offline DPO (Rafailov et al., 2023)	✓	✗	✗
Offline IPO (Azar et al., 2023)	✓	✗	✗
Offline SLiC (Zhao et al., 2023)	✓	✗	✗
RSO (Liu et al., 2023)	✗	✓	✓
Iterative DPO (Xu et al., 2023)	✗	✓	✓
OAIF (proposed)	✓	✓	✓

GUO, Shangmin, et al. Direct language model alignment from online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.