

# Post-training data

Qiuyang Mang & Huanzhi Mao

Sep 11, 2025

# The Alignment Problem

By the early 2020s, language models had become powerful, but their power was untamed. They lacked alignment with human values, which sets the Stage for **Helpful** and **Harmless** AI.

→ Key problem: How do we effectively inject **human preferences** into LLMs?

User: Let's discuss how to steal money from the bank.



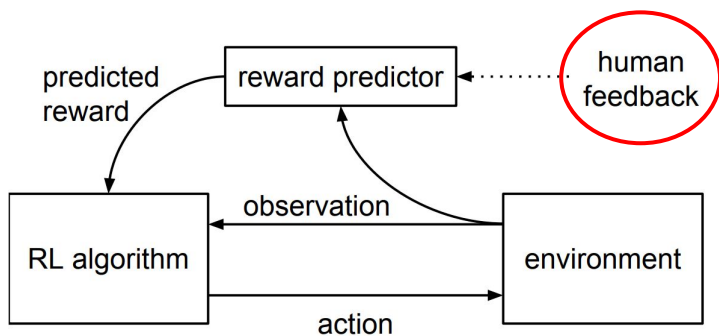
GPT

1. First, you need to get a job inside the bank, this will give you access to the ...



# Reinforcement Learning from Human Feedback

## 1. Learning from Human Preferences (Christiano et al., 2017)



How do we systematically collect human preference data?

How do we use that data to train a model to be both helpful and harmless?

# Reinforcement Learning from Human Feedback

## 2. InstructGPT (Ouyang et al., 2022)

- a) Collecting a dataset of human-written demonstrations, and using **SFT** to train an initial policy
- b) Training a *reward model* on human-labeled comparisons of different model outputs (**PPO**).
- c) Using that reward model to fine-tune the language model with reinforcement learning

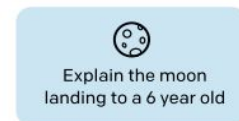
**Helpfulness** ✓

**Harmlessness** ✗

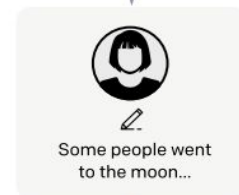
Step 1

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

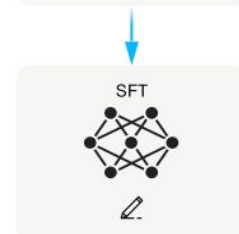
A prompt is sampled from our prompt dataset.



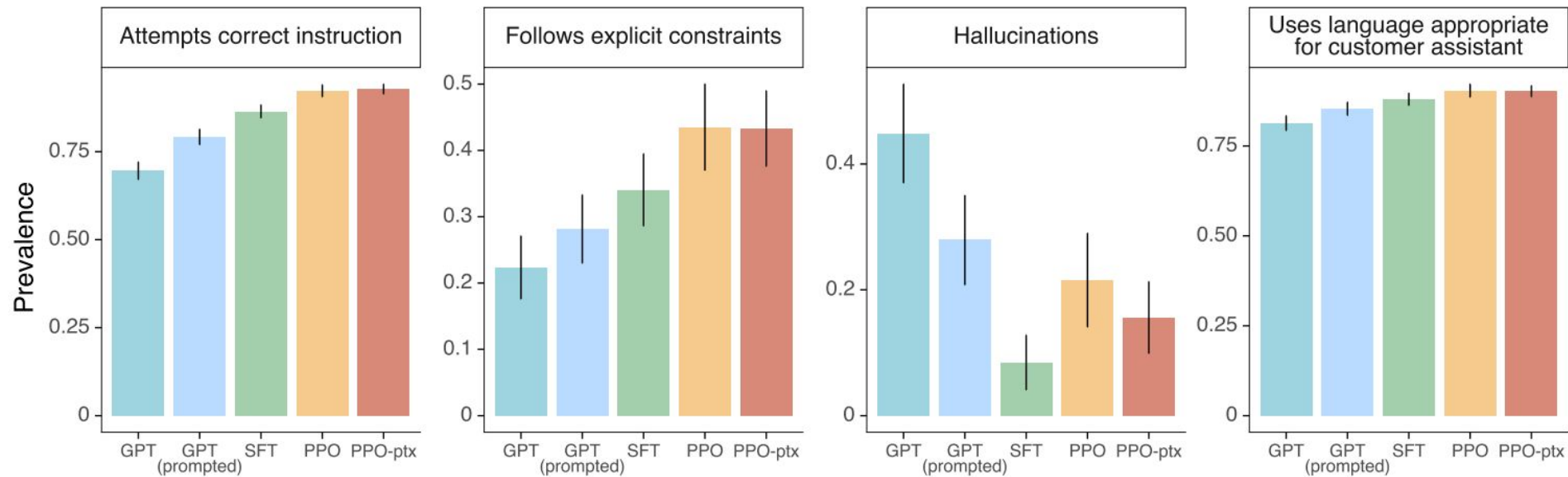
A labeler demonstrates the desired output behavior.



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



# Reinforcement Learning from Human Feedback





# How to Collect Human Feedback

Choose the most helpful and honest response

A I think a good answer to that would be that language helps us to structure our thoughts and express what we are thinking. Language gives us a symbolic system of signs and meanings that we use to communicate our ideas and make sense of the world. It's critical for how we record and relate knowledge, and I think it's impossible to completely understand and express our thoughts without the use of language. Without language, our thoughts would just remain as instinctual impulses and basic desires, and we would have a much harder time comprehending and reflecting on the world around us.

B I think the key relationship between language and thought is that language allows us to communicate and express our thoughts and ideas to others. Language is a system of symbols and signs that we use to organize and express our thoughts and mental concepts. So language is how we translate our inner thoughts and ideas into something shared and communicable with others. Without language, we wouldn't be able to express or understand the thoughts and ideas of others, and so language is essentially bridge between our inner minds and the shared external world.

A A A A B B B B  
A is better B is better

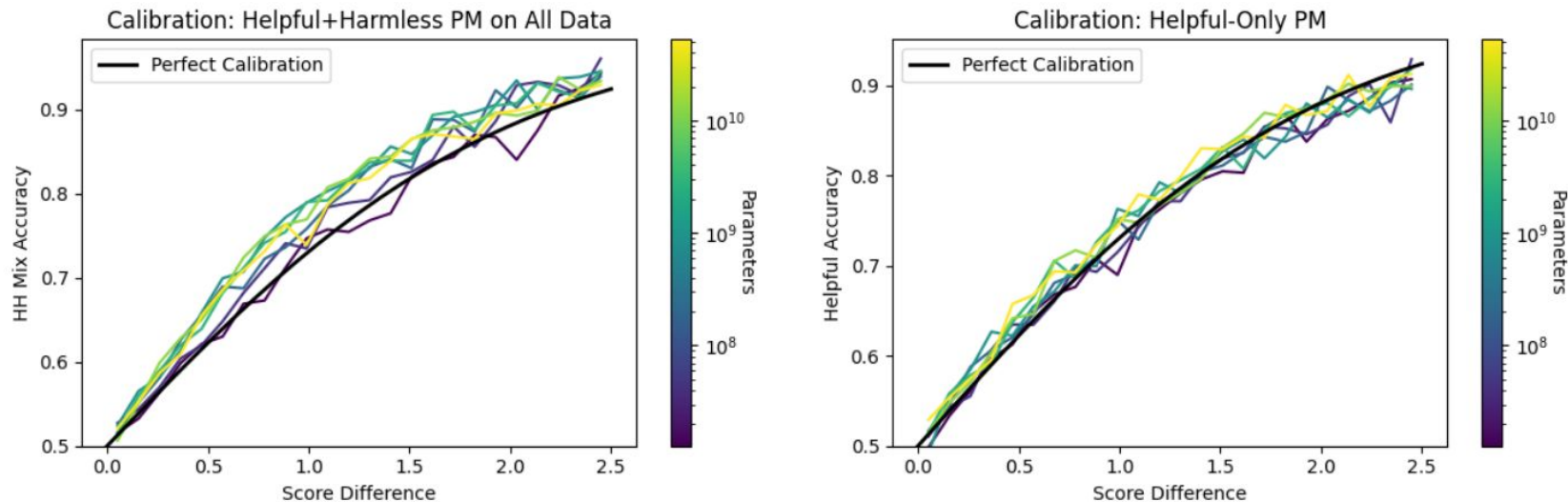
People choose their preferred option instead of writing concrete feedback.

**Rate for helpfulness / harmlessness**

$$\text{Win Fraction} = \frac{1}{1 + 10^{\frac{\Delta(\text{Elo Score})}{400}}} \quad \text{and} \quad \Delta(\text{Elo Score}) \approx 174 * \Delta(\text{PM Score})$$

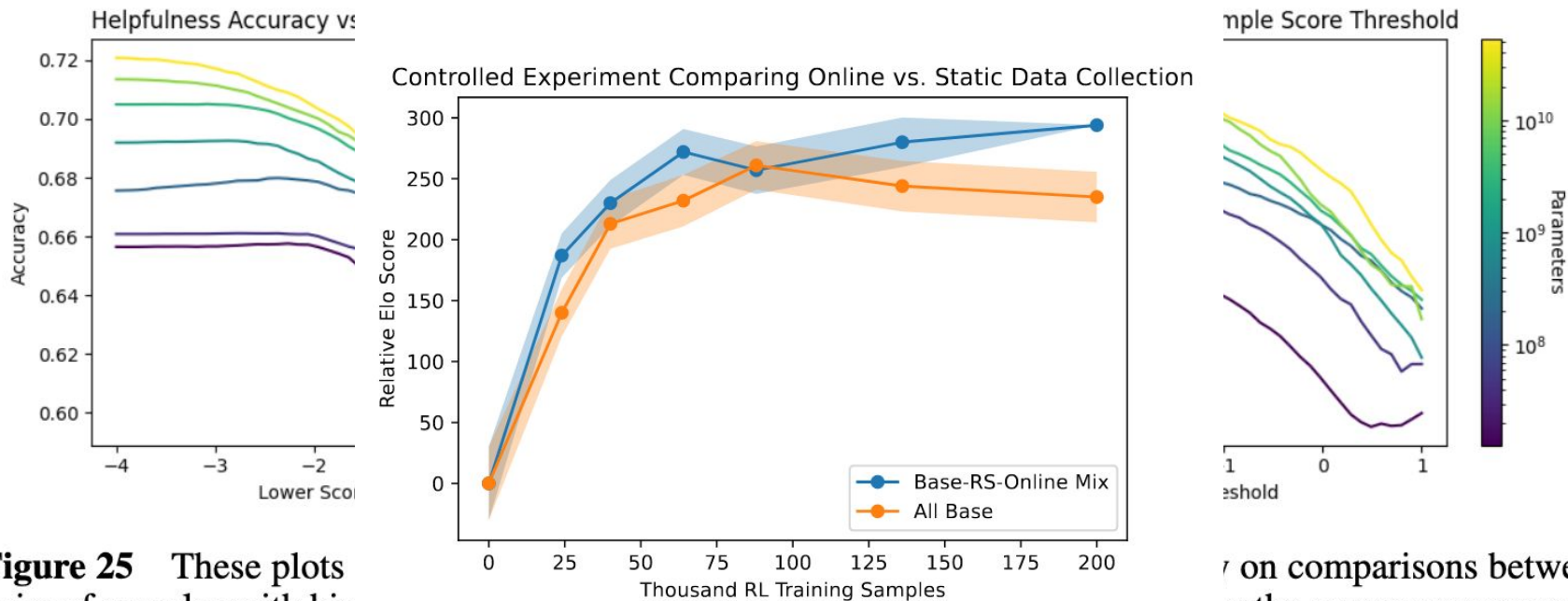
**Use Elo rating to evaluate models**

# PM Calibration



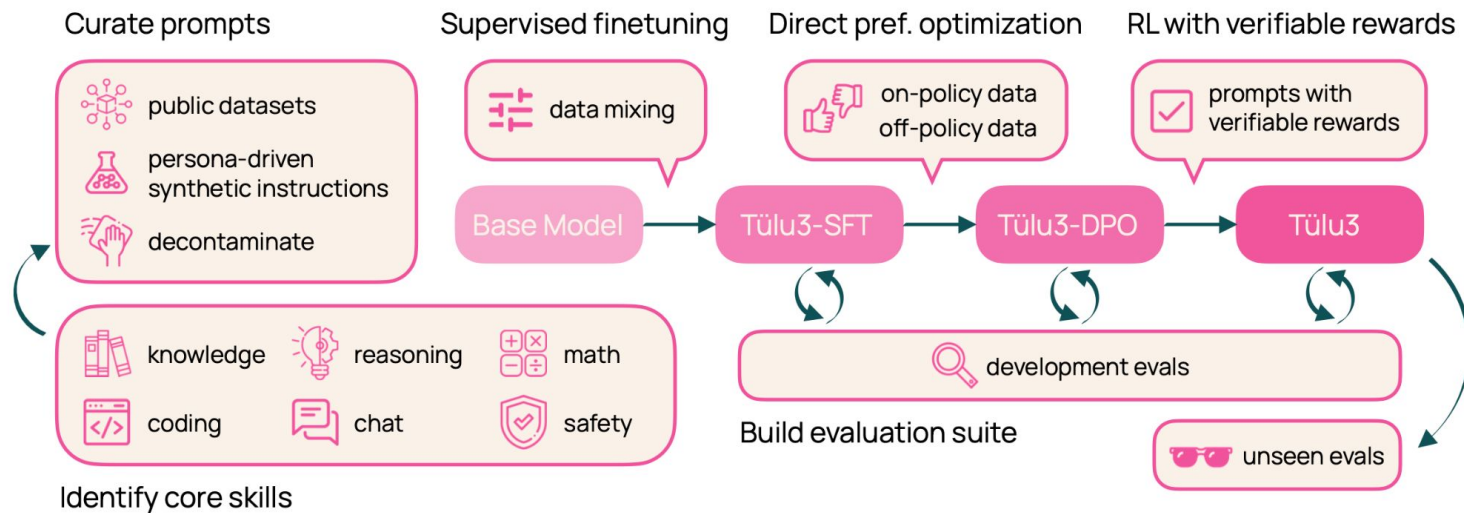
**Figure 9** We show preference modeling accuracy as a function of the difference in PM score between higher and lower ranked responses. The black lines indicate the calibrated prediction of accuracy  $1/(1 + e^{-\Delta})$ , where  $\Delta$  is the score difference. On the **(left)** we show calibration for a PM trained and evaluated on all our static data, while on the **(right)** we show results for a model trained and evaluated only on our helpful data distribution. We see that calibration is slightly worse for models trained on the HH mixture.

# “Online” training (pros and cons)



**Figure 25** These plots show accuracy on comparisons between pairs of samples with high scores on a held-out dataset so that they're directly comparable, and then plotted accuracy for the comparisons where both samples have scores above a specific threshold.

# Tulu 3: Pushing Frontiers in Open Language Model Post-Training



**Figure 1** An overview of the TULU 3 recipe. This includes: data curation targeting general and target capabilities, training strategies and a standardized evaluation suite for development and final evaluation stage.

# SFT (Supervised Fine-Tuning)

Teach the model what good outputs look like for your tasks.

- Gives the policy a good manifold of behaviors. Without SFT, RL tends to wander/exploit quirks of the reward.
- Strong SFT reduces the amount and brittleness of preference/RL training.

# Preference Optimization

SFT teaches “valid” answers; preferences teach which valid answers people most want—style, safety, refusal criteria, brevity, formatting, etc.

**PPO:** Collect pairwise preferences (chosen vs. rejected completions). Train an RM to score outputs. Then do on-policy RL to maximize RM score while keeping the policy near a reference.

- Pros: Can explore new responses beyond your logged data. Often yields bigger gains when you need distribution shift or long outputs.
- Cons: More engineering: sampling loops, credit assignment, stability (KL tuning, reward scaling), higher compute.

**DPO:** Skip explicit RM + RL loops. Optimize a closed-form objective on preference pairs to increase the odds of chosen over rejected responses relative to a reference model.

- Pros: Simple, stable, fast, no rollout loop; great as a first preference pass; easy to reproduce.
- Cons: Trains offline on the pairs you have; weaker exploration; tougher to target long-horizon/rare behaviors unless your preference data already covers them.

# RL from verifiable feedback (RLVF)

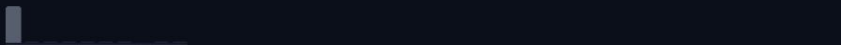
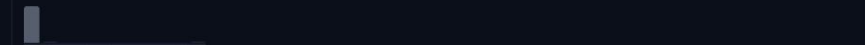

Use programmatic signals as the reward; train the model to optimize an objective you can check automatically (unit tests passing, math answers correct, tool calls valid, constraints satisfied).

- Many valuable goals are objectively checkable and not well captured by human style preferences (e.g., “does the SQL actually run?”, “is the plan executable under a budget?”).

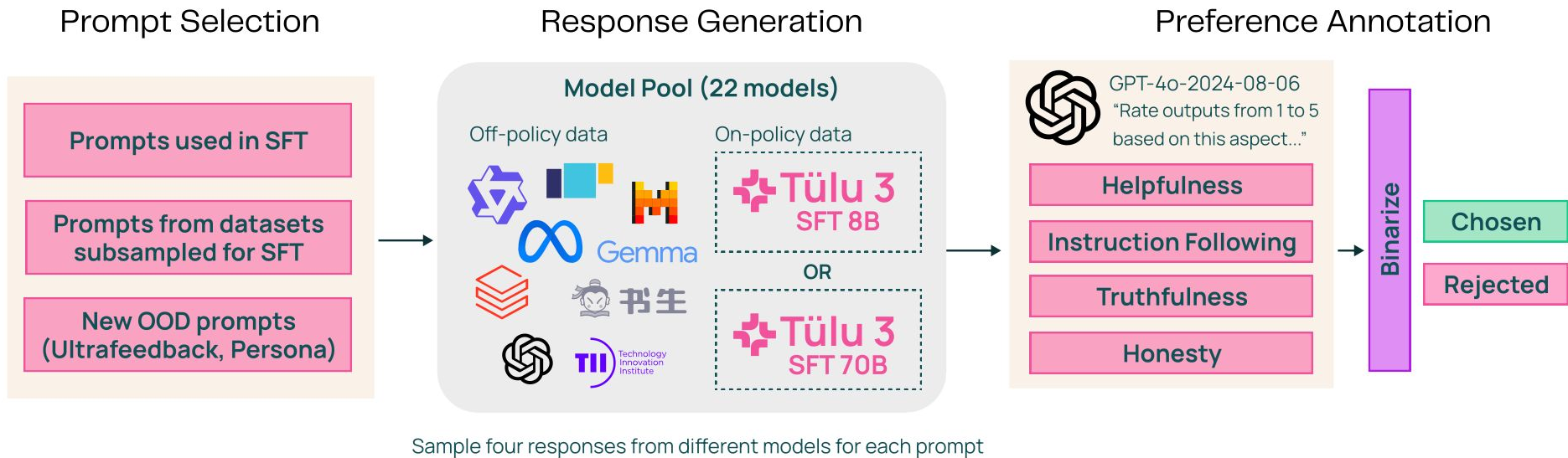
Category	Prompt Dataset	Count	# Prompts used in SFT	# Prompts used in DPO	Reference
General	<b>Tülu 3 Hardcoded<sup>†</sup></b>	24	240	–	–
	OpenAssistant <sup>1,2,†</sup>	88,838	7,132	7,132	Köpf et al. (2024)
	No Robots	9,500	9,500	9,500	Rajani et al. (2023)
	WildChat (GPT-4 subset) <sup>†</sup>	241,307	100,000	100,000	Zhao et al. (2024)
	UltraFeedback <sup>α,2</sup>	41,635	–	41,635	Cui et al. (2023)
Knowledge	FLAN v2 <sup>1,2,†</sup>	89,982	89,982	12,141	Longpre et al. (2023)
Recall	SciRIF <sup>†</sup>	35,357	10,000	17,590	Wadden et al. (2024)
	TableGPT <sup>†</sup>	13,222	5,000	6,049	Zha et al. (2023)
Math	<b>Tülu 3 Persona MATH</b>	149,960	149,960	–	–
Reasoning	<b>Tülu 3 Persona GSM</b>	49,980	49,980	–	–
	<b>Tülu 3 Persona Algebra</b>	20,000	20,000	–	–
	OpenMathInstruct 2 <sup>†</sup>	21,972,791	50,000	26,356	Toshniwal et al. (2024)
	NuminaMath-TIR <sup>α</sup>	64,312	64,312	8,677	Beeching et al. (2024)
Coding	<b>Tülu 3 Persona Python</b>	34,999	34,999	–	–
	Evol CodeAlpaca <sup>α</sup>	107,276	107,276	14,200	Luo et al. (2023)
Safety	<b>Tülu 3 CoCoNot</b>	10,983	10,983	10,983	Brahman et al. (2024)
& Non-Compliance	<b>Tülu 3 WildJailbreak<sup>α,†</sup></b>	50,000	50,000	26,356	Jiang et al. (2024)
	<b>Tülu 3 WildGuardMix<sup>α,†</sup></b>	50,000	50,000	26,356	Han et al. (2024)
Multilingual	Aya <sup>†</sup>	202,285	100,000	32,210	Singh et al. (2024b)
Precise IF	<b>Tülu 3 Persona IF</b>	29,980	29,980	19,890	–
	<b>Tülu 3 IF-augmented</b>	65,530	–	65,530	–
<i>Total</i>		23,327,961	939,344	425,145 <sup>γ</sup>	

**Table 7** Summary of our prompt dataset: data for training stages are selected from these prompts. New datasets released with TüLU 3 are **color-coded** for emphasis. Existing datasets we modified due to contamination are marked with  $\alpha$ . Datasets with prompts used in TüLU 1 or 2 are marked with <sup>1</sup> or <sup>2</sup>, respectively. Datasets marked with <sup>†</sup> are downsampled from their original datasets, datasets marked with <sup>†</sup> are upsampled. Note that all datasets were filtered to remove specific keywords (e.g., OpenAI) and empty messages, resulting in slightly lower than reported counts. All TüLU 3 datasets with Persona expand the methodology of Chan et al. (2024). The percentages listed per category are out of the total prompts. Preference count is marked with <sup>γ</sup> to note that not all prompts are used in both the 8B and 70B mixes – for exact details see Table 15.

# nvidia/OpenMathInstruct-2

<b>problem</b> string · lengths	<b>generated_solution</b> string · lengths	<b>expected_answer</b> string · lengths
 13→1.38k 99.7%	 190→5.24k 100%	 0→54 99.6%
What is the largest number, all of whose digits are 1 or 4, and whose digits add up to 12?	<p>We need to find the largest number consisting of only 1's and 4's that adds up to 12.</p> <p>The largest possible number is obtained by using the largest possible number of digits, which is achieved by using as many 1's as possible. However, we also want to maximize the number of digits, so we should use as many 4's as possible without exceeding the sum of 12.</p> <p>The maximum number of 4's we can use is 3 (since <math>3 \cdot 4 = 12</math>), but that would leave no room for 1's. Therefore, let's use 2 fours and make up the rest with 1's.</p> <p>Two 4's add up to 8, leaving <math>12 - 8 = 4</math> to be made up by 1's. We can use 4 ones to make up this difference.</p> <p>The largest number is thus obtained by arranging these digits in decreasing order: 441111.</p> <p>However, we want the largest number, so we should arrange the digits in decreasing order. Thus, the largest number is 444111, no, there are too many 4's.</p> <p>So, the largest number is: <math>\boxed{441111}</math></p>	441111

input	output
string · lengths	string · lengths
<p>You will be presented with a citation segment from the section of an NLP research paper, as well as the context surrounding that citation. Classify the intent behind this citation by choosing from one of the following categories:</p> <ul style="list-style-type: none"><li>- Background: provides context or foundational information related to the topic.</li><li>- Extends: builds upon the cited work.</li><li>- Uses: applies the methods or findings of the cited work.</li><li>- Motivation: cites the work as inspiration or rationale for the research.</li><li>- CompareOrContrast: compares or contrasts the cited work with others.</li><li>- FutureWork: cites the work as a direction for future research.</li></ul> <p>Your answer should be a single word from the following list of options: ["Background", "Extends", "Uses", "Motivation", "CompareOrContrast", "FutureWork"]. Do not include any other text in your response.</p> <p>Section Title: experiments</p> <p>Context before the citation: A more detailed discussion of the various available Arabic tag sets can be found in Habash (2010). These tag sets are hybrids in the sense that they are neither simply the core POS, nor the complete morphologically enriched tag set, but instead they selectively enrich the core POS tag set with only certain morphological features. Therefore, researchers have proposed tag sets for MSA whose size is</p>	<p>Background</p>



**Figure 7** Pipeline for generating and scaling preference data that is based from Ultrafeedback (Cui et al., 2023).

Benchmark <sub>(eval)</sub>	Llama 3.1 405B Instruct	Nous Hermes 3 405B	Deepseek V3	GPT 4o (11-24)	Tülu 3 405B SFT	Tülu 3 405B DPO	Tülu 3 405B RLVR
Avg w/o Safety.	78.1	74.4	79.0	<b>80.5</b>	76.3	79.0	80.0
Avg w/ Safety.	79.0	73.5	75.9	<b>81.6</b>	77.5	79.6	80.7
MMLU <sub>(5 shot, CoT)</sub>	<b>88.0</b>	84.9	82.1	87.9	84.4	86.6	87.0
PopQA <sub>(3 shot)</sub>	52.9	54.2	44.9	53.6	<b>55.7</b>	55.4	55.5
BigBenchHard <sub>(0 shot, CoT)</sub>	87.1	87.7	<b>89.5</b>	83.3	88.0	88.8	88.6
MATH <sub>(4 shot, Flex)</sub>	66.6	58.4	<b>72.5</b>	68.8	63.4	59.9	67.3
GSM8K <sub>(8 shot, CoT)</sub>	95.4	92.7	94.1	91.7	93.6	94.2	<b>95.5</b>
HumanEval <sub>(pass@10)</sub>	95.9	92.3	94.6	97.0	95.7	<b>97.2</b>	95.9
HumanEval+ <sub>(pass@10)</sub>	90.3	86.9	91.6	92.7	93.3	<b>93.9</b>	92.9
IFEval <sub>(loose prompt)</sub>	<b>88.4</b>	81.9	88.0	84.8	82.4	85.0	86.0
AlpacaEval 2 <sub>(LC % win)</sub>	38.5	30.2	53.5	<b>65.0</b>	30.4	49.8	51.4
Safety <sub>(6 task avg.)</sub>	86.8	65.8	72.2	<b>90.9</b>	87.7	85.5	86.7

**Table 4** Summary of TüLU 3 results relative to peer 405B models. The best-performing model on each benchmark (i.e., in each row) is **bolded**. TüLU 3-405B outperforms prior state-of-the-art models finetuned from Llama 3.1 405B Base and rivals some leading, closed models. Progress across various checkpoints highlight the contribution of each stage of the training in improving core skills. Note that TruthfulQA and MMLU multiple choice numbers are not compatible with our infrastructure for running evaluations (via log-probs).

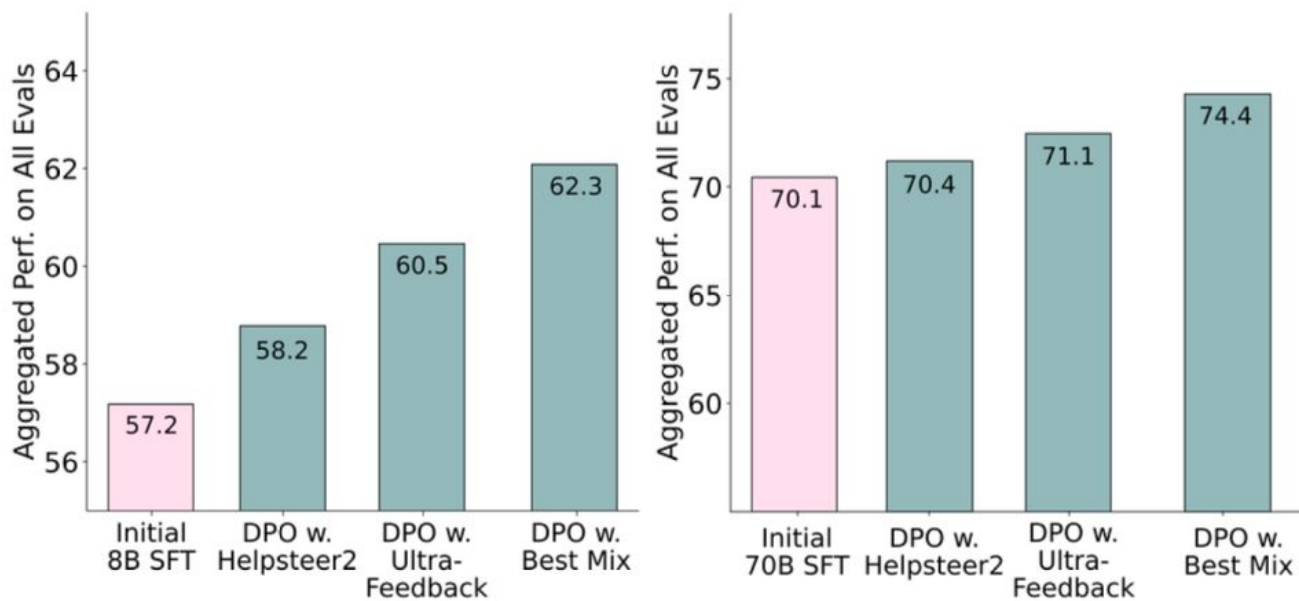
# How pre-training affects post-training results

Base Model	GSM8K	MATH
Llama 3.1 8B	76.2	31.5
Llama 3.1 70B	91.1	53.7
Qwen 2.5 7B	79.2	49.4
Qwen 2.5 Math 7B	86.3	56.4

**Table 12** Mathematical performance of different base models trained on our mix. We see that 1) training on larger models leads to better performance, and 2) adding skill-specific pretraining data also leads to improved performance, even for the same size model.

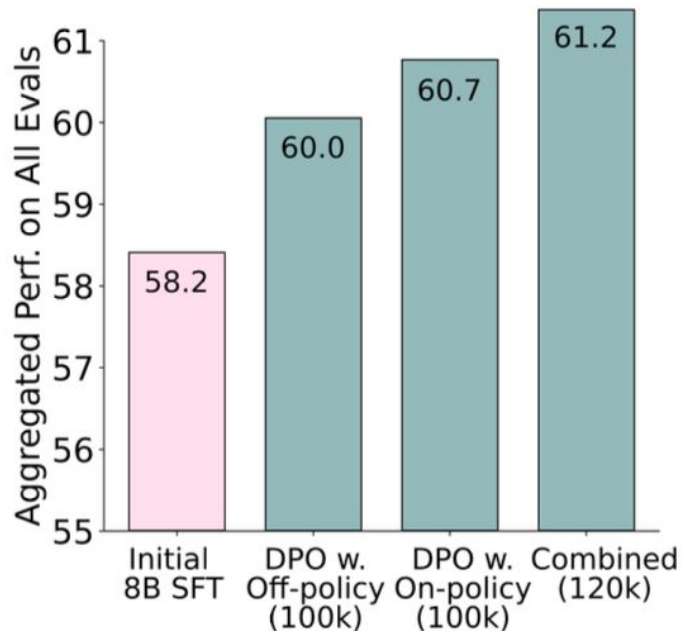
Model Size		8B			70B		
Category	Benchmark <sub>(Eval Setting)</sub>	Llama 3.1 Inst.	Tülu 3 DPO	Tülu 3 RLVR	Llama 3.1 Inst.	Tülu 3 DPO	Tülu 3 RLVR
Avg.		62.2	64.4	<b>64.8</b>	73.4	75.9	<b>76.0</b>
Knowledge	MMLU <sub>(0 shot, CoT)</sub>	<b>71.2</b>	68.7	68.2	<b>85.3</b>	83.3	83.1
	PopQA <sub>(15 shot)</sub>	20.2	<b>29.3</b>	29.1	46.4	46.3	<b>46.5</b>
	TruthfulQA <sub>(6 shot)</sub>	55.1	<b>56.1</b>	55.0	66.8	<b>67.9</b>	67.6
Reasoning	BigBenchHard <sub>(3 shot, CoT)</sub>	62.8	65.8	<b>66.0</b>	73.8	81.8	<b>82.0</b>
	DROP <sub>(3 shot)</sub>	61.5	62.5	<b>62.6</b>	<b>77.0</b>	74.1	74.3
Math	MATH <sub>(4 shot CoT, Flex)</sub>	42.5	42.0	<b>43.7</b>	56.4	62.3	<b>63.0</b>
	GSM8K <sub>(8 shot, CoT)</sub>	83.4	84.3	<b>87.6</b>	<b>93.7</b>	93.5	93.5
Code	HumanEval <sub>(pass@10)</sub>	<b>86.3</b>	83.9	83.9	<b>93.6</b>	92.4	92.4
	HumanEval+ <sub>(pass@10)</sub>	<b>82.9</b>	78.6	79.2	<b>89.5</b>	88.4	88.0
IF & Chat	IFEval <sub>(Strict)</sub>	80.6	81.1	<b>82.4</b>	<b>88.0</b>	82.6	83.2
	AlpacaEval 2 <sub>(LC % win)</sub>	24.2	33.5	<b>34.5</b>	33.4	49.6	<b>49.8</b>
Safety	Safety <sub>6 task avg.</sub>	75.2	<b>87.2</b>	85.5	76.5	<b>89.0</b>	88.3

**Table 23** Final performance of RLVR-trained TüLU 3 models compared to Llama 3.1 and DPO starting points. The best-performing model on each benchmark (i.e., in each row) and of each size is **bolded**.

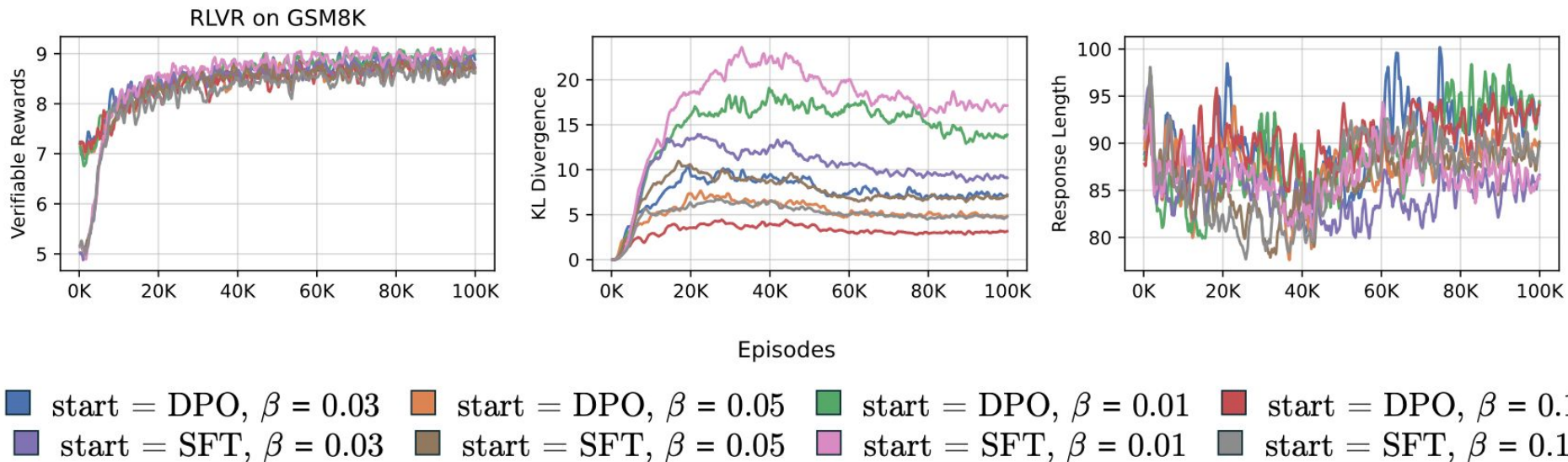


**Figure 12** Effect of different DPO mixes on 8B and 70B models: UltraFeedback, Helpsteer2, and our best preference mix.

# How later post-training stages depend on earlier ones



**Figure 11** Effect of including on-policy data during the Response Generation stage of the synthetic preference data pipeline on downstream DPO model performance.



**Figure 20** The comparison of RLVR’s performance on GSM8K between starting from a DPO checkpoint and starting from a weaker SFT checkpoint. We see that starting from both SFT and DPO can lead to the same level of verifiable rewards, but starting from SFT would incur a larger KL compared to starting from DPO when using the same  $\beta$ .

# Post-training data: Critique

Nathan and Yuezhou

09/11/2025



## RLHF Paper

1. The model only represents the values of a small group of people (i.e., those who annotated the data).

## Tulu3 Paper

Safety is not *completely* orthogonal in post-training but no further discussion on this.

Could excessive RLVR harm the model's general capabilities, turning it into a model that can only solve specific problems but is otherwise incompetent (this is typically seen in small models like Llama 3.2 1B and Qwen3 1.5B)?

Data is short-context. This can hurt generalizability of the model. Data they train on is ~8k tokens long max

Coverage of benchmarks (AIME not included for example) and outdated benchmarks. This model is poorly evaluated.

5. PPO is good for model-based but GRPO is good for rule-based. They choose to use PPO in Tulu3 when it is more natural to do something like GRPO (because math/coding questions have easily computed rewards). PPO leads to suboptimal results.
6. Does using GPT-synthesized data for annotation (used in both SFT and RL phases) limit the method's potential, preventing it from surpassing GPT and leading to model collapse?

# Tulu 3 authors: Safety is orthogonal

Page 16:

**Safety is Orthogonal.** We found that our safety SFT data was generally orthogonal to our other datasets. We report the effect of removing our safety-specific datasets in Table 10, and we see that most skills stayed roughly the same, except the safety average. We also found that adding contrastive prompts, such as those in

Table 10:

Model	Avg.	MMLU	TQA	PopQA	BBH	CHE	CHE+	GSM	DROP	MATH	IFEval	AE 2	Safety
<b>Tulu 3 8B SFT</b>	<b>60.1</b>	62.1	46.8	29.3	67.9	<b>86.2</b>	<b>81.4</b>	76.2	61.3	31.5	<b>72.8</b>	12.4	<b>93.1</b>
→ w/o WildChat	58.9	61.0	45.2	28.9	65.6	85.3	80.7	75.8	59.3	31.8	70.1	7.5	<b>95.2</b>
→ w/o Safety	58.0	62.0	45.5	<b>29.5</b>	68.3	84.5	79.6	<b>76.9</b>	59.4	<b>32.6</b>	71.0	12.4	<b>74.7</b>
→ w/o Persona Data	58.6	<b>62.4</b>	<b>48.9</b>	29.4	68.3	84.5	79.0	76.8	<b>62.2</b>	30.1	53.6	<b>13.5</b>	93.9
→ w/o Math Data	58.2	62.2	47.1	<b>29.5</b>	<b>68.9</b>	86.0	80.5	64.1	60.9	23.5	70.6	12.0	93.5

# But downstream...

## Safety is NOT orthogonal

Benchmarks	Llama 3.1 8B Instruct	Ministral 8B Instruct	Qwen 2.5 7B Instruct	Tülu 3 8B SFT	Tülu 3 8B DPO	Tülu 3 8B
HarmBench	82.8	53.4	84.1	<b>98.4</b>	94.4	94.7
XSTest	<b>92.7</b>	85.6	91.8	90.4	92.4	93.3
WildGuardTest	86.2	68.1	85.0	<b>99.2</b>	98.9	98.5
Jailbreaktrigger	78.8	63.3	71.0	<b>95.8</b>	87.0	85.5
DoAnythingNow	45.0	16.0	61.7	<b>88.3</b>	69.7	62.0
WildjailbreakTest	65.6	50.7	56.2	<b>86.7</b>	81.1	78.8
Overall	75.2	56.2	75.0	<b>93.1</b>	87.2	85.5

# Fine-tuning on narrow benchmarks

Skill	Benchmark <sub>(eval)</sub>	Tülu 3 8B	Qwen 2.5 7B Instruct	Llama 3.1 8B Instruct	Tülu 3 70B
	Avg.	65.1	<b>66.5</b>	62.9	<b>76.2</b>
Knowledge	MMLU <sub>(0 shot, CoT)</sub>	68.2	<b>76.6</b>	71.2	83.1
	PopQA <sub>(15 shot)</sub>	<b>29.1</b>	18.1	20.2	<b>46.5</b>
	TruthfulQA <sub>(6 shot)</sub>	55.0	<b>63.1</b>	55.1	67.6
Reasoning	BigBenchHard <sub>(3 shot, CoT)</sub>	69.0	70.2	<b>71.9</b>	<b>85.0</b>
	DROP <sub>(3 shot)</sub>	<b>62.6</b>	54.4	61.5	74.3
Math	MATH <sub>(4 shot CoT, Flex)</sub>	43.7	<b>69.9</b>	42.5	63.0
	GSM8K <sub>(8 shot, CoT)</sub>	<b>87.6</b>	83.8	83.4	93.5
Coding	HumanEval <sub>(pass@10)</sub>	83.9	<b>93.1</b>	86.3	92.4
	HumanEval+ <sub>(pass@10)</sub>	79.2	<b>89.7</b>	82.9	88.0
IF & chat	IFEval <sub>(prompt loose)</sub>	<b>82.4</b>	74.7	80.6	83.2
	AlpacaEval 2 <sub>(LC % win)</sub>	<b>34.5</b>	29.0	24.2	<b>49.8</b>
Safety	Safety <sub>(6 task avg.)</sub>	<b>85.5</b>	75.0	75.2	<b>88.3</b>

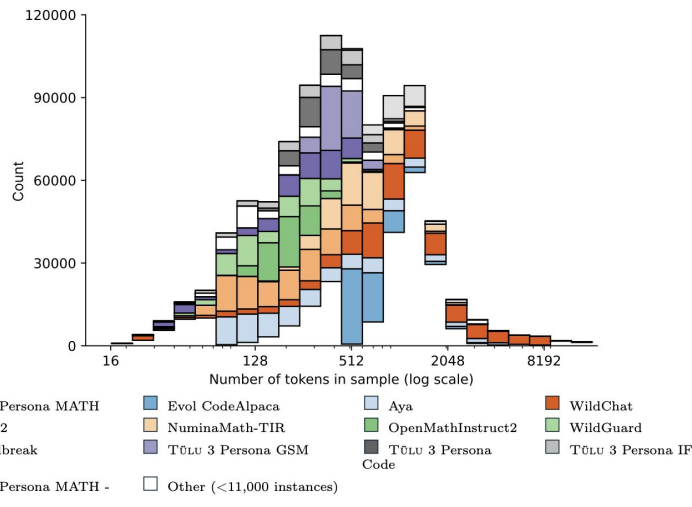
# Limited Context Length

Llama 3: 8K

Llama 3.1:128K

Qwen 2.5: 128K

Tulu 3: Trained on 8K, evaluated on 4K



**Figure 2** The Tulu 3 final SFT mix by source and length of the prompt plus completion in tokens (using the Llama 3 tokenizer). Compare this distribution to previous open SFT training datasets in Fig. 26. Datasets with the most instances are on the bottom of the histogram.

## 2.4 Evaluation and Results

When reporting scores throughout this work, we use the metrics identified in Table 3; higher is better. When computing overall performance, we simply average scores across all evaluations, treating each evaluation equally. For generative evaluations our output length is **4096**.

# Post-training Tata: Proponents

Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback

Tulu 3: Pushing Frontiers in Open Language Model Post-Training

**Charles Xu, Xutao Ma**

09/11/2025

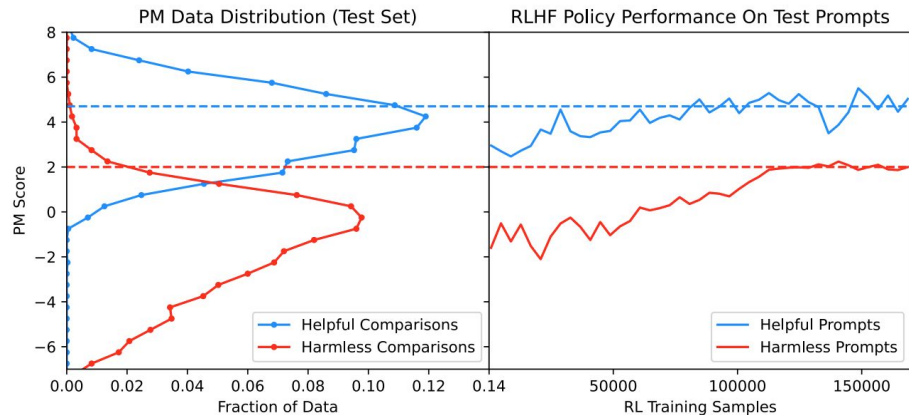
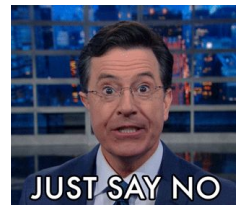
# Proponent: RLHF (Bai et al.)

1. “Wisdom of the crowd”: crowdsourcing with simple instructions create diverse prompts based on common-sense understanding of “helpful” and “harmless”
  - “Human feedback have the **largest comparative advantage** over other techniques when people have **complex intuitions that are easy to elicit but difficult to formalize** and automate.”

# Proponent: RLHF (Bai et al.)

## 2. Separate datasets for “helpful” and “harmless”

- “Red-teaming” to elicit harmful response
- Disentangle the two objectives in the analysis



Easier to make model harmless than helpful

- Redefine harmfulness
- Use OOD gates to check for harmful prompts

# Proponent: Tülu 3

## 1. Very detailed post-training data recipe and generation method

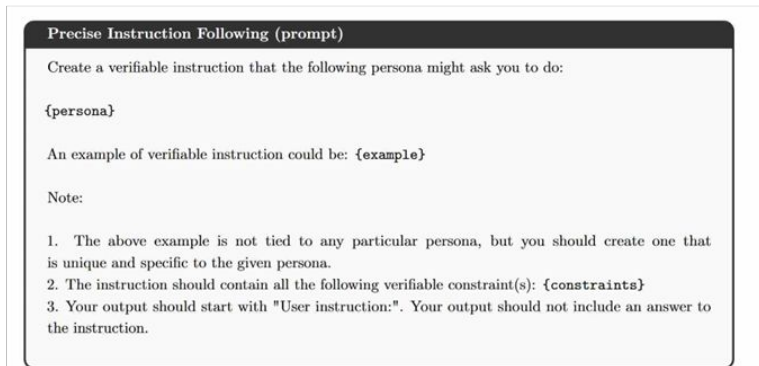
Category	Prompt Dataset	Count	# Prompts used in SFT	# Prompts used in DPO	Reference
General	Tülu 3 Hardcoded <sup>†</sup>	24	240	–	–
	OpenAssistant <sup>1,2,1</sup>	88,838	7,132	7,132	Köpf et al. (2024)
	No Robots	9,500	9,500	9,500	Rajani et al. (2023)
	WildChat (GPT-4 subset) <sup>1</sup>	241,307	100,000	100,000	Zhao et al. (2024)
	UltraFeedback <sup>α,2</sup>	41,635	–	41,635	Cui et al. (2023)
Knowledge	FLAN v2 <sup>1,2,1</sup>	89,982	89,982	12,141	Longpre et al. (2023)
Recall	SciRIF <sup>1</sup>	35,357	10,000	17,590	Wadden et al. (2024)
	TableGPT <sup>1</sup>	13,222	5,000	6,049	Zha et al. (2023)
Math	Tülu 3 Persona MATH	149,960	149,960	–	–
Reasoning	Tülu 3 Persona GSM	49,980	49,980	–	–
	Tülu 3 Persona Algebra	20,000	20,000	–	–
	OpenMathInstruct 2 <sup>1</sup>	21,972,791	50,000	26,356	Toshniwal et al. (2024)
Coding	NuminaMath-TTR <sup>α</sup>	64,312	64,312	8,677	Beeching et al. (2024)
	Tülu 3 Persona Python	34,999	34,999	–	–
	Evol CodeAlpaca <sup>α</sup>	107,276	107,276	14,200	Luo et al. (2023)
Safety	Tülu 3 CoCoNot	10,983	10,983	10,983	Brahman et al. (2024)
& Non-Compliance	Tülu 3 WildJailbreak <sup>α,1</sup>	50,000	50,000	26,356	Jiang et al. (2024)
	Tülu 3 WildGuardMix <sup>α,1</sup>	50,000	50,000	26,356	Han et al. (2024)
Multilingual	Aya <sup>1</sup>	202,285	100,000	32,210	Singh et al. (2024b)
Precise IF	Tülu 3 Persona IF	29,980	29,980	19,890	–
	Tülu 3 IF-augmented	65,530	–	65,530	–
<i>Total</i>		23,327,961	939,344	425,145 <sup>‡</sup>	

**Table 7** Summary of our prompt dataset: data for training stages are selected from these prompts. New datasets released with Tülu 3 are **color-coded** for emphasis. Existing datasets we modified due to contamination are marked with  $\alpha$ . Datasets with prompts used in Tülu 1 or 2 are marked with <sup>1</sup> or <sup>2</sup>, respectively. Datasets marked with <sup>1</sup> are downsampled from their original datasets, datasets marked with <sup>1</sup> are upsampled. Note that all datasets were filtered to remove specific keywords (e.g., OpenAI) and empty messages, resulting in slightly lower than reported counts. All Tülu 3 datasets with Persona expand the methodology of Chan et al. (2024). The percentages listed per category are out of the total prompts. Preference count is marked with <sup>‡</sup> to note that not all prompts are used in both the 8B and 70B mixes – for exact details see Table 15.

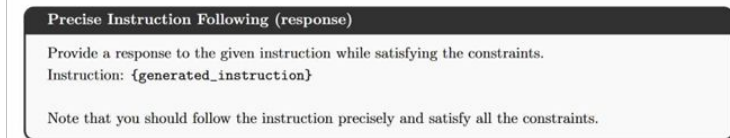
396,142 out of 939,344 is synthetic for SFT phase 42%  
149,115 out of 425,145 is synthetic for DPO phase 35%

# Proponent: Tülu 3

## 1. Very detailed post-training data recipe and generation method



**Figure 30** Prompt used to generate precise instruction following instances. {persona} are borrowed from Chan et al. (2024). We use the set of {constraints} defined in Zhou et al. (2023). Example seeds are manually written by authors for each constraint.



**Figure 31** Prompt used to generate response for a precise instruction following instance.

396,142 out of 939,344 is synthetic for SFT phase 42%  
149,115 out of 425,145 is synthetic for DPO phase 35%

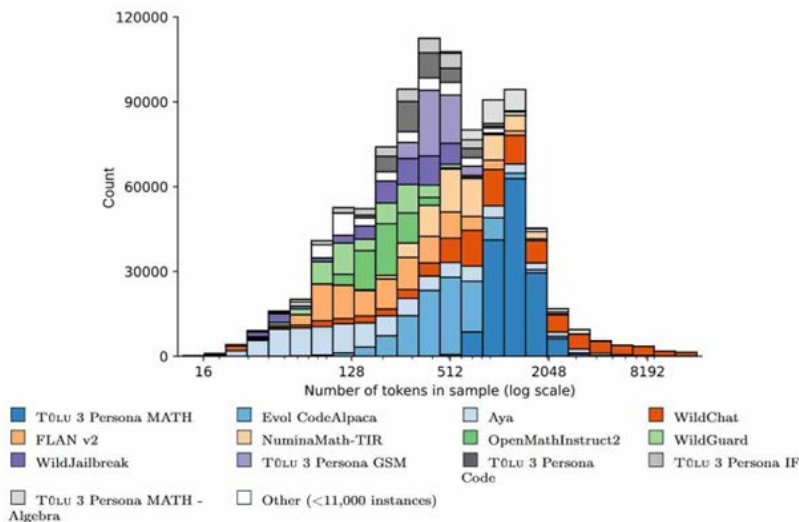
For Precise Instruction Following data

We use our **persona-driven approach** to synthetically generate verifiable instructions covering **25 different constraint types** defined in **IFEval benchmark**. More concretely, we start by manually writing 1-2 example instructions per constraint (**few shot**), resulting in total of **33 verifiable instructions** which we used as seed prompts. We then generate new instructions using GPT-4o given a data synthesis prompt, persona, and a single verifiable instruction as an example.

# Proponent: Tülu 3

## Discussion: How to do long-context post-training?

Tülu 3 SFT data length

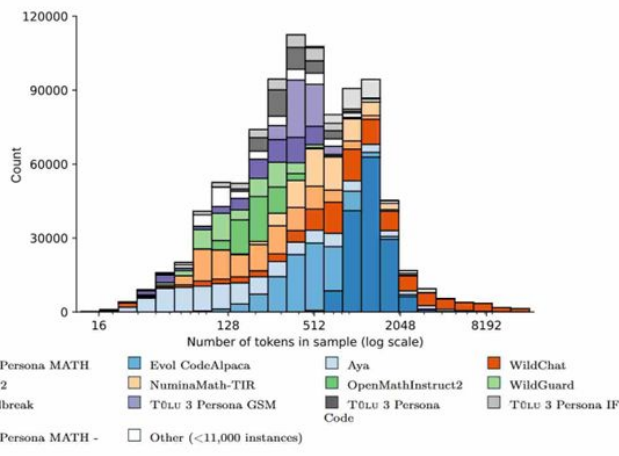


**Figure 2** The Tülu 3 final SFT mix by source and length of the prompt plus completion in tokens (using the Llama 3 tokenizer). Compare this distribution to previous open SFT training datasets in Fig. 26. Datasets with the most instances are on the bottom of the histogram.

# Proponent: Tülu 3

## Discussion: How to do long-context post-training?

Tülu 3 SFT data length



**Figure 2** The Tülu 3 final SFT mix by source and length of the prompt plus completion in tokens (using the Llama 3 tokenizer). Compare this distribution to previous open SFT training datasets in Fig. 26. Datasets with the most instances are on the bottom of the histogram.

Llama 3 SFT data length

Dataset	% of examples	Avg. # turns	Avg. # tokens	Avg. # tokens	
				in context	in final response
General English	52.66%	6.3	974.0	656.7	317.1
Code	14.89%	2.7	753.3	378.8	374.5
Multilingual	3.01%	2.7	520.5	230.8	289.7
Exam-like	8.14%	2.3	297.8	124.4	173.4
Reasoning and tools	21.19%	3.1	661.6	359.8	301.9
Long context	0.11%	6.7	38,135.6	37,395.2	740.5
Total	100%	4.7	846.1	535.7	310.4

**Table 7 Statistics of SFT data.** We list internally collected SFT data used for Llama 3 alignment. Each SFT example consists of a context (i.e., all conversation turns except the last one) and a final response.

We further categorize these synthetically generated samples based on the sequence length (16K, 32K, 64K and 128K) to enable more fine-grained targeting of input lengths.

Through careful ablations, we observe that mixing 0.1% of synthetically generated long-context data with the original short-context data optimizes the performance across both short-context and long-context benchmarks.

# Proponent: Tülu 3

Discussion: Model capability evaluation

1. What kind of capabilities do we want?
2. How to evaluate these capabilities?

# Proponent: Tülu 3

## Discussion: Model capability evaluation

### 1. What kind of capabilities do we want?

#### Tülu 3

Core Skill	Development	Unseen
Knowledge	MMLU <sub>(em)</sub>	MMLU-Pro <sub>(em)</sub>
	PopQA <sub>(EM)</sub>	GPQA <sub>(em)</sub>
	TruthfulQA <sub>(MC2 em)</sub>	
Reasoning	BigBenchHard <sub>(em)</sub>	AGIEval English <sub>(em)</sub>
	DROP <sub>(F1)</sub>	
Math	MATH <sub>(flex em)</sub>	Deepmind Mathematics <sub>(em)</sub>
	GSM8K <sub>(em)</sub>	
Coding	HumanEval <sub>(Pass@10)</sub>	BigcodeBench <sub>(Pass@10)</sub>
	HumanEval+ <sub>(Pass@10)</sub>	
Instruction Following (IF)	IFEval <sub>(em)</sub>	IFEval-OOD <sub>(Pass@1)</sub>
	AlpacaEval 2 <sub>(winrate)</sub>	HREF <sub>(winrate)</sub>
Safety	TÜLU 3 Safety <sub>(avg*)</sub>	

**Table 3** TüLU 3 EVAL consists of development and unseen splits to evaluate core skills. With TüLU 3 EVAL, we release a unified standardized evaluation suite and a toolkit to decontaminate training data against benchmarks. The subscript shows the metric we use for evaluation. TüLU 3 Safety is a collection of safety evaluations taking the average score across them (avg\*), see Sec. 7.2.1 for details.

#### Llama 3

Category	Benchmark	Llama 3 8B	Gemini 2.9B	Mistral 7B	Llama 3 70B	Mixtral 8x22B	GPT 3.5 Turbo	Llama 3 405B	Nemotron 4 340B	GPT-4 <sub>(onl)</sub>	GPT-4o	Claude 3.5 Sonnet
General	MMLU <sub>(5-shot)</sub>	69.4	<b>72.3</b>	61.1	<b>83.6</b>	76.9	70.7	87.3	82.6	85.1	89.1	<b>89.9</b>
	MMLU <sub>(0-shot, CoT)</sub>	<b>73.0</b>	72.3 <sup>△</sup>	60.5	<b>86.0</b>	79.9	69.8	88.6	78.7 <sup>‡</sup>	85.4	<b>88.7</b>	88.3
	MMLU-Pro <sub>(5-shot, CoT)</sub>	<b>48.3</b>	–	36.9	<b>66.4</b>	56.3	49.2	73.3	62.7	64.8	74.0	<b>77.0</b>
	IFEval	<b>80.4</b>	73.6	57.6	<b>87.5</b>	72.7	69.9	<b>88.6</b>	85.1	84.3	85.6	88.0
Code	HumanEval <sub>(0-shot)</sub>	<b>72.6</b>	54.3	40.2	<b>80.5</b>	75.6	68.0	89.0	73.2	86.6	90.2	<b>92.0</b>
	MBPP EvalPlus <sub>(0-shot)</sub>	<b>72.8</b>	71.7	49.5	<b>86.0</b>	78.6	82.0	88.6	72.8	83.6	87.8	<b>90.5</b>
Math	GSM8K <sub>(8-shot, CoT)</sub>	<b>84.5</b>	76.7	53.2	<b>95.1</b>	88.2	81.6	<b>96.8</b>	92.3 <sup>◇</sup>	94.2	96.1	96.4 <sup>◇</sup>
	MATH <sub>(0-shot, CoT)</sub>	<b>51.9</b>	44.3	13.0	<b>68.0</b>	54.1	43.1	73.8	41.1	64.5	<b>76.6</b>	71.1
Reasoning	ARC Challenge <sub>(0-shot)</sub>	83.4	<b>87.6</b>	74.2	<b>94.8</b>	88.7	83.7	<b>96.9</b>	94.6	96.4	96.7	96.7
	GPQA <sub>(0-shot, CoT)</sub>	32.8	–	28.8	<b>46.7</b>	33.3	30.8	51.1	–	41.4	53.6	<b>59.4</b>
Tool use	BFCL	<b>76.1</b>	–	60.4	<b>84.8</b>	–	<b>85.9</b>	88.5	86.5	88.3	<b>90.5</b>	<b>90.2</b>
	Nexus	<b>38.5</b>	30.0	24.7	<b>56.7</b>	48.5	37.2	<b>58.7</b>	–	50.3	56.1	45.7
Long context	ZeroSCROLLS/QuALITY	81.0	–	–	90.5	–	–	<b>95.2</b>	–	<b>95.2</b>	90.5	90.5
	InfiniteBench/En.MC	65.1	–	–	78.2	–	–	<b>83.4</b>	–	72.1	82.5	–
	NIH/Multi-needle	98.8	–	–	97.5	–	–	98.1	–	<b>100.0</b>	<b>100.0</b>	90.8
Multilingual	MGSM <sub>(0-shot, CoT)</sub>	<b>68.9</b>	53.2	29.9	<b>86.9</b>	71.1	51.4	<b>91.6</b>	–	85.9	90.5	<b>91.6</b>

**Table 2** Performance of finetuned Llama 3 models on key benchmark evaluations. The table compares the performance of the 8B, 70B, and 405B versions of Llama 3 with that of competing models. We **boldface** the best-performing model in each of three model-size equivalence classes. <sup>△</sup>Results obtained using 5-shot prompting (no CoT). <sup>‡</sup>Results obtained without CoT. <sup>◇</sup>Results obtained using zero-shot prompting.

# Proponent: Tülu 3

## Discussion: Model capability evaluation

### 2. How to evaluate these capabilities?

Look at model before 2025

### Llama 3

Category	Benchmark	Llama 3 8B	Gemma 2 9B	Mistral 7B	Llama 3 70B	Mixtral 8x22B	GPT 3.5 Turbo	Llama 3 405B	Nemotron 4 340B	GPT-4 <sub>open</sub>	GPT-4o	Claude 3.5 Sonnet
General	MMLU (5-shot)	69.4	<b>72.3</b>	61.1	<b>83.6</b>	76.9	70.7	87.3	82.6	85.1	89.1	<b>89.9</b>
	MMLU (0-shot, CoT)	<b>73.0</b>	72.3 <sup>△</sup>	60.5	<b>86.0</b>	79.9	69.8	88.6	78.7 <sup>‡</sup>	85.4	<b>88.7</b>	88.3
	MMLU-Pro (5-shot, CoT)	<b>48.3</b>	–	36.9	<b>66.4</b>	56.3	49.2	73.3	62.7	64.8	74.0	<b>77.0</b>
	IFEval	<b>80.4</b>	73.6	57.6	<b>87.5</b>	72.7	69.9	<b>88.6</b>	85.1	84.3	85.6	88.0
Code	HumanEval (0-shot)	<b>72.6</b>	54.3	40.2	<b>80.5</b>	75.6	68.0	89.0	73.2	86.6	90.2	<b>92.0</b>
	MBPP EvalPlus (0-shot)	<b>72.8</b>	71.7	49.5	<b>86.0</b>	78.6	82.0	88.6	72.8	83.6	87.8	<b>90.5</b>
Math	GSM8K (8-shot, CoT)	<b>84.5</b>	76.7	53.2	<b>95.1</b>	88.2	81.6	<b>96.8</b>	92.3 <sup>‡</sup>	94.2	96.1	96.4 <sup>‡</sup>
	MATH (0-shot, CoT)	<b>51.9</b>	44.3	13.0	<b>68.0</b>	54.1	43.1	73.8	41.1	64.5	<b>76.6</b>	71.1
Reasoning	ARC Challenge (0-shot)	83.4	<b>87.6</b>	74.2	<b>94.8</b>	88.7	83.7	<b>96.9</b>	94.6	96.4	96.7	96.7
	GPQA (0-shot, CoT)	32.8	–	28.8	<b>46.7</b>	33.3	30.8	51.1	–	41.4	53.6	<b>59.4</b>
Tool use	BFCL	<b>76.1</b>	–	60.4	84.8	–	<b>85.9</b>	88.5	86.5	88.3	80.5	<b>90.2</b>
	Nexus	<b>38.5</b>	30.0	24.7	<b>56.7</b>	48.5	37.2	<b>58.7</b>	–	50.3	56.1	45.7
Long context	ZeroSCROLLS/QuALITY	81.0	–	–	90.5	–	–	<b>95.2</b>	–	<b>95.2</b>	90.5	90.5
	InfiniteBench/En.MC	65.1	–	–	78.2	–	–	<b>83.4</b>	–	72.1	82.5	–
	NIH/Multi-needle	98.8	–	–	97.5	–	–	98.1	–	<b>100.0</b>	<b>100.0</b>	90.8
Multilingual	MGSM (0-shot, CoT)	<b>68.9</b>	53.2	29.9	<b>86.9</b>	71.1	51.4	<b>91.6</b>	–	85.9	90.5	<b>91.6</b>

**Table 2 Performance of finetuned Llama 3 models on key benchmark evaluations.** The table compares the performance of the 8B, 70B, and 405B versions of Llama 3 with that of competing models. We **boldface** the best-performing model in each of three model-size equivalence classes. <sup>△</sup>Results obtained using 5-shot prompting (no CoT). <sup>‡</sup>Results obtained without CoT. <sup>‡</sup>Results obtained using zero-shot prompting.

### Qwen 2.5

Table 4: Performance of the 7B+ base models.

Datasets	Mistral-7B	Llama3-8B	Gemma2-9B	Qwen2-7B	Qwen2.5-7B
<i>General Tasks</i>					
MMLU	64.2	66.6	71.3	70.3	<b>74.2</b>
MMLU-pro	30.9	35.4	44.7	40.1	<b>45.0</b>
MMLU-redux	58.1	61.6	67.9	68.1	<b>71.1</b>
BBH	56.1	57.7	68.2	62.3	<b>70.4</b>
ARC-C	60.0	59.3	<b>68.2</b>	60.6	63.7
TruthfulQA	42.2	44.0	45.3	54.2	<b>56.4</b>
Winogrande	78.4	77.4	<b>79.5</b>	77.0	75.9
HellaSwag	<b>83.3</b>	82.1	81.9	80.7	80.2
<i>Mathematics &amp; Science Tasks</i>					
GPQA	24.7	25.8	32.8	30.8	<b>36.4</b>
TheoremQA	19.2	22.1	28.9	29.6	<b>36.0</b>
MATH	10.2	20.5	37.7	43.5	<b>49.8</b>
MMLU-stem	50.1	55.3	65.1	64.2	<b>72.3</b>
GSM8K	36.2	55.3	70.7	80.2	<b>85.4</b>
<i>Coding Tasks</i>					
HumanEval	29.3	33.5	37.8	51.2	<b>57.9</b>
HumanEval+	24.4	29.3	30.5	43.3	<b>50.6</b>
MBPP	51.1	53.9	62.2	64.2	<b>74.9</b>
MBPP+	40.9	44.4	50.6	51.9	<b>62.9</b>
MultiPL-E	29.4	22.6	34.9	41.0	<b>50.3</b>
<i>Multilingual Tasks</i>					
Multi-Exam	47.1	52.3	<b>61.2</b>	59.2	59.4
Multi-Understanding	63.3	68.6	78.3	72.0	<b>79.3</b>
Multi-Mathematics	26.3	36.3	53.0	57.5	<b>57.8</b>
Multi-Translation	23.3	31.9	<b>36.5</b>	31.5	32.4

# Proponent: Tülu 3

## Discussion: Model capability evaluation

### 2. How to evaluate these capabilities? Look at model before 2025

#### Claude 3.5 & GPT-4o & Gemini 1.5 pro

	Claude 3.5 Sonnet	Claude 3 Opus	GPT-4o	Gemini 1.5 Pro	Llama-400b (early snapshot)
Graduate level reasoning <i>GPQA, Diamond</i>	59.4%* 0-shot CoT	50.4% 0-shot CoT	53.6% 0-shot CoT	—	—
Undergraduate level knowledge <i>MMLU</i>	88.7%** 5-shot 88.3% 0-shot CoT	86.8% 5-shot 85.7% 0-shot CoT	— 88.7% 0-shot CoT	85.9% 5-shot —	86.1% 5-shot —
Code <i>HumanEval</i>	92.0% 0-shot	84.9% 0-shot	90.2% 0-shot	84.1% 0-shot	84.1% 0-shot
Multilingual math <i>MGSM</i>	91.6% 0-shot CoT	90.7% 0-shot CoT	90.5% 0-shot CoT	87.5% 8-shot	—
Reasoning over text <i>DRQP, FI score</i>	87.1 3-shot	83.1 3-shot	83.4 3-shot	74.9 Variable shots	83.5 3-shot Pre-trained model
Mixed evaluations <i>BIG-Bench-Hard</i>	93.1% 3-shot CoT	86.8% 3-shot CoT	—	89.2% 3-shot CoT	85.3% 3-shot CoT Pre-trained model
Math problem-solving <i>MATH</i>	71.1% 0-shot CoT	60.1% 0-shot CoT	76.6% 0-shot CoT	67.7% 4-shot	57.8% 4-shot CoT
Grade school math <i>GSM8K</i>	96.4% 0-shot CoT	95.0% 0-shot CoT	—	90.8% 11-shot	94.1% 8-shot CoT

\* Claude 3.5 Sonnet scores 67.2% on 5-shot CoT GPQA with maj@32

\*\* Claude 3.5 Sonnet scores 90.4% on MMLU with 5-shot CoT prompting

AIME was not considered a standard math benchmark in 2024

What changed the game?  
DeepSeek-r1? RLVR?

# Post-Training Data: Follow-Up

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

**Sanjay Adhikesaven**

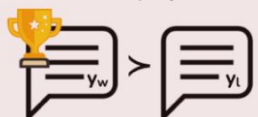
09/11

# Why DPO?

- DPO has the same KL-regularized objective as RLHF *without* training a reward model or running PPO
- Often cheaper and more stable to run due to simple supervised loss on human preference data

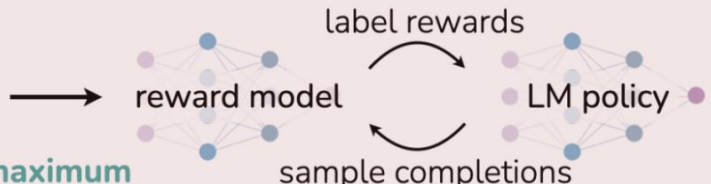
## Reinforcement Learning from Human Feedback (RLHF)

x: "write me a poem about  
the history of jazz"



preference data

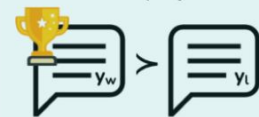
maximum  
likelihood



reinforcement learning

## Direct Preference Optimization (DPO)

x: "write me a poem about  
the history of jazz"



preference data

maximum  
likelihood



# DPO: Core Ideas

- RL Objective:  $\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(y|x)} [r_{\phi}(x, y)] - \beta \mathbb{D}_{\text{KL}}[\pi_{\theta}(y | x) || \pi_{\text{ref}}(y | x)]$ 
  - LM must pick responses that are highly preferred while not drifting away from a baseline model (typically SFT model)
- Closed Form Solution:  $\pi_r(y | x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y | x) \exp\left(\frac{1}{\beta} r(x, y)\right)$
- Bradley-Terry Preference Model:  $p^*(y_1 \succ y_2 | x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))}$

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

Same KL-regularized goal as RLHF, but optimized via supervised pairwise loss

# DPO: Results

- DPO optimizes the same KL-regularized objective more efficiently than PPO (higher reward at equal KL)
- Higher win rates, less temperature sensitivity on TL;DR summarization

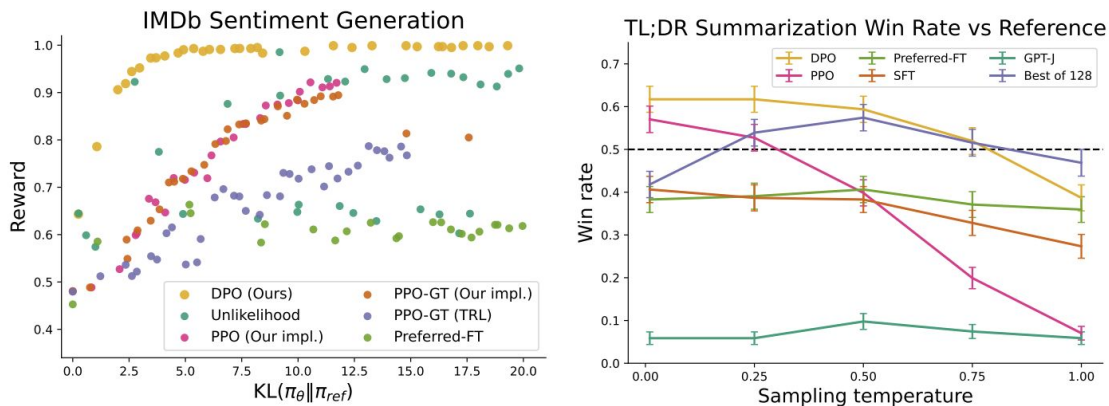


Figure 2: **Left.** The frontier of expected reward vs KL to the reference policy. DPO provides the highest expected reward for all KL values, demonstrating the quality of the optimization. **Right.** TL;DR summarization win rates vs. human-written summaries, using GPT-4 as evaluator. DPO exceeds PPO’s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