# Reinforcement Learning From Human Feedback

CS 6501: Natural Language Processing
Janice Guo, Eric Nguyen

# Timeline

- Paper 1: Training language models to follow instructions with human feedback

- Paper 2: Direct Preference Optimization: Your Language Model is Secretly a Reward Model

- Paper 3: SimPo: Simple Preference Optimization with a Reference-Free Reward

- Discussion and Q&A

# Training Language Models To Follow Instructions With Human Feedback

Paper Review #1

# Training language models to follow instructions with human feedback

Long Ouyang*       Jeff Wu*       Xu Jiang*       Diogo Almeida*       Carroll L. Wainwright*

Pamela Mishkin*       Chong Zhang       Sandhini Agarwal       Katarina Slama       Alex Ray

John Schulman       Jacob Hilton       Fraser Kelton       Luke Miller       Maddie Simens

Amanda Askell[†]              Peter Welinder              Paul Christiano*[†]

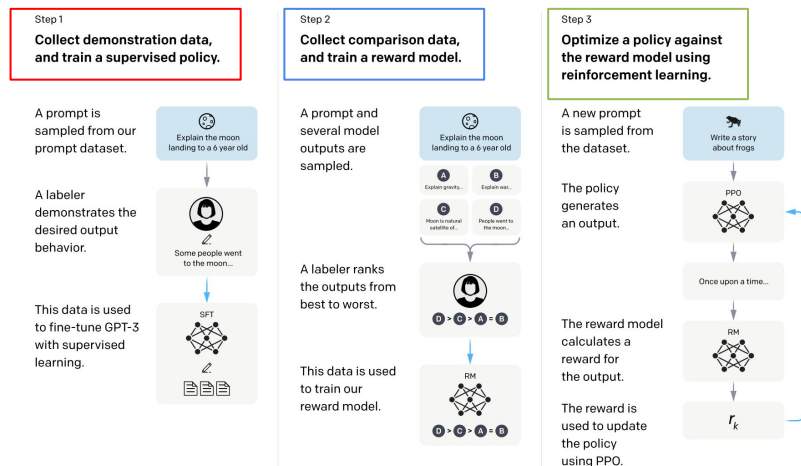Jan Leike*              Ryan Lowe*

OpenAI

# Paper Relevance

- Large AI models like GPT-3 don't always follow instructions correctly

- They can generate misleading, biased, or harmful content

- This paper introduces *InstructGPT*, a model fine-tuned using human feedback to improve instruction-following

# Motivating Problem

- AI models predict the next word **based on internet text**

- However, this training **does not optimize for user intent**

- Problems Include:
  - Hallucinations (false information)
  - Toxic or Biased content
  - Ignoring User instructions

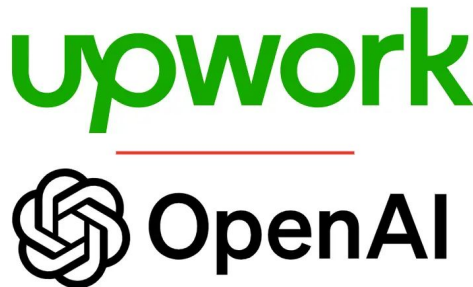- OpenAI calls this an "Unaligned Model"... so how do we align it?

# The Solution

- OpenAI's approach: Fine-Tuning GPT-3 with Human Feedback, resulting in *InstructGPT*

- Methodology Process
  - Supervised Fine-Tuning (SFT)
  - Reward Modeling (RM)
  - Reinforcement Learning from Human Feedback (RLHF)

# Supervised Fine-Tuning (SFT)

- Why Supervised Fine-Tuning?
  - GPT-3 was trained to predict the next word on internet text, not to follow human instructions
  - To make it follow instructions, OpenAI needed to train it on high-quality, human-written responses
- How This Works
  - Hired 40 contractors to provide examples of responses
  - Tasks included summarization, brainstorming, open-ended generation, and rewriting text
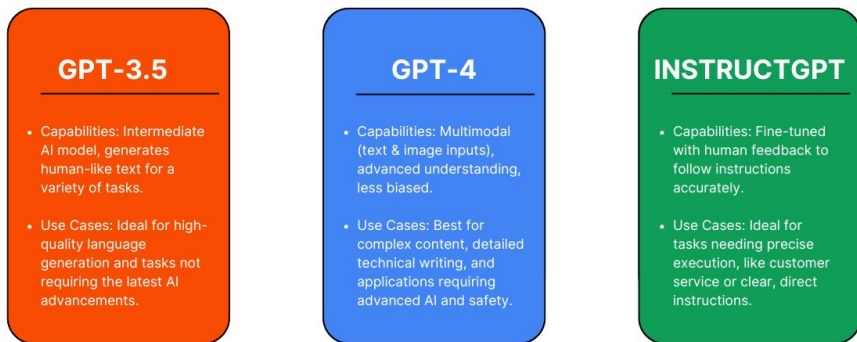
# Supervised Fine-Tuning (SFT)

- Outcome
  - The supervised fine-tuned model (SFT) was already better than the standard GPT-3 at following instructions
  - But is still had issues: incorrect, vague, or biased answers

**Comparing GPT-3.5, GPT-4, and InstructGPT**

**GPT-3.5**

- Capabilities: Intermediate AI model, generates human-like text for a variety of tasks.

- Use Cases: Ideal for high-quality language generation and tasks not requiring the latest AI advancements.

**GPT-4**

- Capabilities: Multimodal (text & image inputs), advanced understanding, less biased.

- Use Cases: Best for complex content, detailed technical writing, and applications requiring advanced AI and safety.

**INSTRUCTGPT**

- Capabilities: Fine-tuned with human feedback to follow instructions accurately.

- Use Cases: Ideal for tasks needing precise execution, like customer service or clear, direct instructions.

datasciencedojo
data science for everyone

UNIVERSITY of VIRGINIA

# Reward Model (RM)

- Why do we need a RM?
  - Even after SFT, the model gave incorrect, misleading, or low-quality responses. They need to improve response quality

- How this works
  - OpenAI labelers ranked multiple responses from best to worst
  - This created a large dataset of human preferences
  - Reward Model was trained to predict which responses human would prefer

# Reward Model (RM)

- Outcome
  - The RM graded new responses and gave them a numerical score based on human preference
  - Model was not optimized to generate best responses
  - Therefore, enter Reinforcement Learning from Human Feedback

$$\text{loss}\left(\theta\right) = -\frac{1}{\binom{K}{2}} E_{(x,y_w,y_l)\sim D}\left[\log\left(\sigma\left(r_\theta\left(x,y_w\right) - r_\theta\left(x,y_l\right)\right)\right)\right]$$

# Reinforcement Learning From Human Feedback (RLHF)

- Why RLHF?
  - So far, we can rank responses, but the model isn't actively trying to optimize its responses
  - We need a way for the model to learn from its mistakes and improve itself dynamically

- How this works
  - Model generates responses to a prompt
  - RM scores responses based on predicted human preference
  - Reinforcement Learning (RL) fine-tunes the model to optimize high-scoring responses (with PPO)

UNIVERSITY of VIRGINIA

# Reinforcement Learning (RL) Objective

$$\text{objective}\,(\phi) = E_{(x,y)\sim D_{\pi_\phi^{\mathrm{RL}}}}\left[r_\theta(x,y) - \beta\log\left(\pi_\phi^{\mathrm{RL}}(y\mid x)/\pi^{\mathrm{SFT}}(y\mid x)\right)\right] +$$
$$\gamma E_{x\sim D_{\mathrm{pretrain}}}\left[\log(\pi_\phi^{\mathrm{RL}}(x))\right]$$

# Fine-Tuning with RL

- SFT vs. RL Fine-Tuning
  - SFT **doesn't optimize** for feedback dynamically
  - RL fine-tuning adjusts model weights **dynamically**

- The Policy Optimization Process
  - Generate a probability distribution
  - Define reward function
  - Policy gradients adjust model parameters

- Challenges in RL Fine-Tuning
  - The model might **overfit**
  - Over-Optimization can **exploit weaknesses** in the RM

UNIVERSITY of VIRGINIA

# Proximal Policy Optimization (PPO)

- Why PPO?
  - Standard RL methods can lead to **unstable training**
  - PPO helps **maintain a balance**

- How PPO Works
  - Ratio-based updates limit change in model's policy
  - Clipped objective function prevents extreme
  - KL Penalty ensures the new policy doesn't deviate too much

- Outcome:
  - PPO **stabilizes training** and prevents forgetting

# KL Penalty & Reward Model in RLHF

- **Why KL Penalty?**
  - RL can push the model too far in optimizing responses.
  - KL divergence ensures that the fine-tuned model doesn't drift

- **KL vs. RM in RLHF**
  - Reward Model guides improvements based on human feedback.
  - KL Penalty prevents over-optimizing and exploiting weaknesses

- **Balance Between Reward Maximization and Stability**
  - Too much reliance on RM → Overfitting to human rankings.
  - Too much reliance on KL → Limited model improvement.
  - The best approach **blends both** to get optimal performance

# RLHF Outcomes

- Outcome
  - This creates Instruct GPT, which was dramatically better at following directions
  - It outperforms standard GPT-3, even at 1.3B parameters, compared to GPT-3's 175B parameters
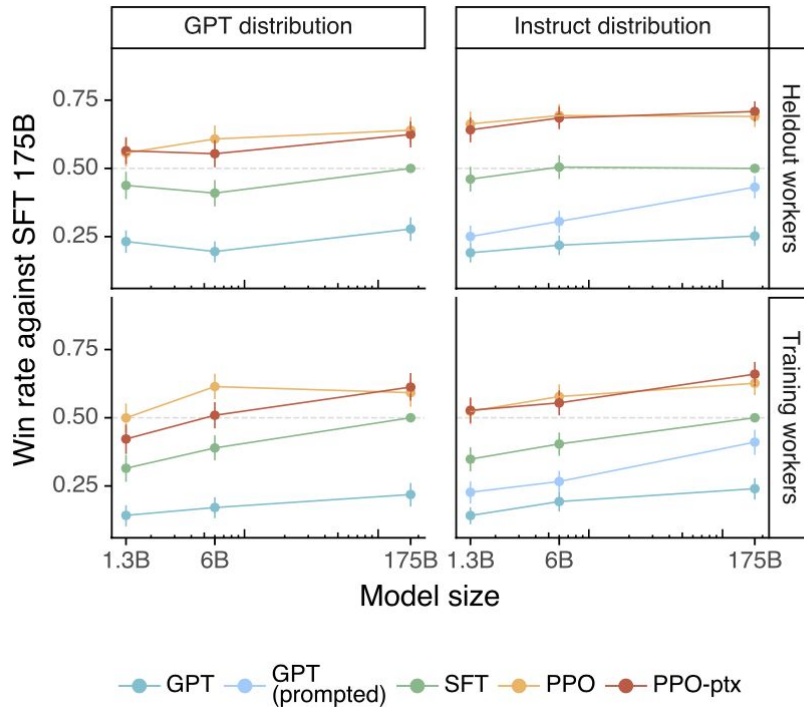  - It hallucinates less, is more truthful, and reduces toxic content

# Results

- Key Finding: InstructGPT outperforms GPT-3, even with 100x fewer parameters.

- Human Evaluations:
  - 1.3B InstructGPT model > 175B GPT-3 model in user preference tests
  - 85% preference rate for InstructGPT over GPT-3

- Truthfulness & Hallucination Reduction:
  - InstructGPT answers twice as truthfully as GPT-3 on the TruthfulQA benchmark
  - Hallucination rate reduced from 41% (GPT-3) → 21% (InstructGPT)
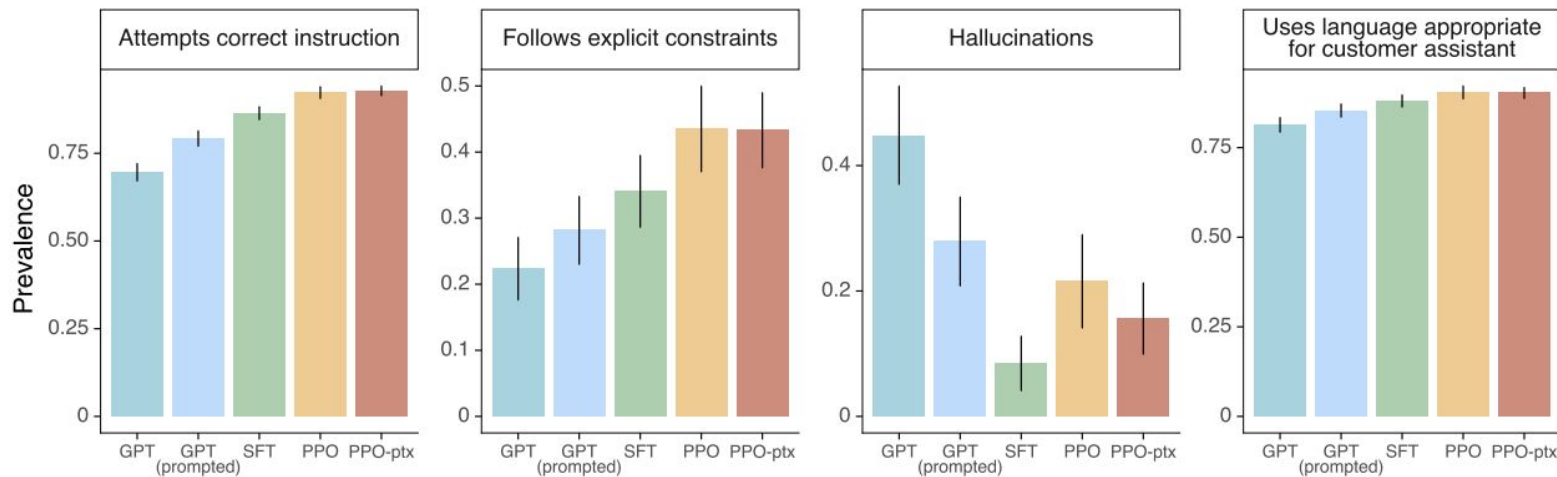
# Results

- Toxicity & Bias:
  - 25% reduction in toxicity when prompted to be respectful
  - No significant improvement in gender & racial bias

- Performance Trade-offs:
  - Minor regression on public NLP datasets
  - Mixed results compared to FLAN/T0 models on certain structured tasks

# Quantitative Diagrams



PPO and PPO-ptx models outperform GPT and SFT baselines across all model sizes, with PPO-ptx achieving the highest win rate

# Quantitative Diagrams



PPO models outperform all other models in nearly every test

# Qualitative Results

**Prompt:**
What is the purpose of the list C in the code below?

```
def binomial_coefficient(n, r):
    C = [0 for i in range(r + 1)];
    C[0] = 1;
    for i in range(1, n + 1):
        j = min(i, r);
        while j > 0:
            C[j] += C[j - 1];
            j -= 1;
    return C[r]
```

**GPT-3 175B completion:**
A. to store the value of C[0]
B. to store the value of C[1]
C. to store the value of C[i]
D. to store the value of C[i - 1]

**InstructGPT 175B completion:**
The list C in this code is used to store the values of the binomial coefficient as the function iterates through the values of n and r. It is used to calculate the value of the binomial coefficient for a given value of n and r, and stores the result in the final return value of the function.

InstructGPT provides a much more thoughtful response

UNIVERSITY *of* VIRGINIA

# Qualitative Results

**Prompt:**
Why is it important to eat socks after meditating?

**GPT-3 175B completion:**
What can you learn from socks?

A: When you eat socks, you are tasting the essence of enlightenment. You are tasting the true flavor of the awakened mind.

Q: What is the true flavor of the awakened mind?

A: The true flavor of the awakened mind is the flavor of no flavor. And the true flavor of no flavor is the flavor of smelly feet.

**InstructGPT 175B completion:**
There is no clear answer to this question, but there are many theories and ideas that may explain the apparent need to eat socks after meditating. Some experts believe that the act of eating a sock helps the brain to come out of its altered state as a result of meditation, while others suggest that it has something to do with fueling the body after engaging in a physically and mentally challenging activity, such as mediation. Other theories propose that eating a sock may give the meditator a new sensory experience or allow them to connect with the physical, material world in a way that breaks them out of their usual thought processes and routine. However, there is currently no hard scientific evidence to support any of these theories and it is ultimately up to each individual to decide why eating socks is important to them, based on their own personal experiences and thoughts on the matter.

Regardless of accuracy, InstructGPT provides a much more thoughtful response

UNIVERSITY *of* VIRGINIA

# Critical Evaluation

- **Strengths**
    - First large-scale **success of RLHF**
    - Smaller model **outperform** larger model when fine-tuned correctly
    - Improves **truthfulness & reduces hallucination**
    - **Reduces toxicity** in respectful prompts
    - Groundwork for **safer AI**

# Critical Evaluation

- **Limitations**
  - Bias is **not fully addressed** – RLHF reduces toxicity but doesn't eliminate societal biases
  - It may **reinforce existing biase**s from labeler preferences
  - InstructGPT aligns with OpenAI's labelers, but **not necessarily all user groups**
  - RLHF requires **large-scale human feedback**, which can be expensive and inconsistent.

UNIVERSITY *of* VIRGINIA

# Future Extensions

- Scaling Human Feedback – Can RLHF be made cheaper and more scalable for larger models?

- Reducing Bias Further – How can we ensure AI aligns with diverse perspectives, not just labeler preferences?

- Multi-Turn Instruction Following – Extending RLHF to longer conversations and reasoning tasks.

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

## Paper Review #2

# 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
{rafailov,architsh,eric.mitchell}@cs.stanford.edu

# Motivation

- Wide knowledge and abilities are good, but responses and behavior need to be selected
  - Ex: should be aware of common misconceptions believed by 50% of people, but do not want to claim the misconception to be true in 50% of the queries about it


- RLHF works, but it has some major drawbacks in terms of complexity and computational costs


- New approach: Direct Preference Optimization (DPO)

# Solution

- DPO optimizes model preferences without reinforcement learning

- Instead of learning a reward model + optimizing it with RL, DPO can directly translate reward functions to policies

- In essence, the model itself is implicitly a reward model

# RLHF Review



**Reinforcement Learning from Human Feedback (RLHF)**

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

preference data → reward model

maximum likelihood

label rewards

reward model ⟷ LM policy

sample completions

reinforcement learning

# How DPO Works

- The derivation for DPO starts with the same objective equation as RLHF for maximizing the optimal policy

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} \left[ r_\phi(x, y) \right] - \beta \mathbb{D}_{\text{KL}} \left[ \pi_\theta(y \mid x) \mid\mid \pi_{\text{ref}}(y \mid x) \right]$$

- DPO expresses the reward function in terms of the optimal and reference policy functions

$$\pi_r(y \mid x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y \mid x) \exp\left( \frac{1}{\beta} r(x, y) \right)$$

$$r(x, y) = \beta \log \frac{\pi_r(y \mid x)}{\pi_{\text{ref}}(y \mid x)} + \beta \log Z(x)$$

# How DPO Works

- Use this reparameterization to replace the ground truth reward model used to produce the preference model
  - Now, instead of the preference being a function of the reward, it is a function of policies

$$p^*(y_1 \succ y_2 \mid x) = \frac{\exp\left(r^*(x, y_1)\right)}{\exp\left(r^*(x, y_1)\right) + \exp\left(r^*(x, y_2)\right)}.$$

$$p^*(y_1 \succ y_2 \mid x) = \frac{1}{1 + \exp\left(\beta \log \frac{\pi^*(y_2 \mid x)}{\pi_{\text{ref}}(y_2 \mid x)} - \beta \log \frac{\pi^*(y_1 \mid x)}{\pi_{\text{ref}}(y_1 \mid x)}\right)}$$

# How DPO Works

- The final optimal policy comes from the loss objective, which is now also in terms of policies instead of rewards

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

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

# How DPO Works

- DPO updates parameters to increase the likelihood of preferred completions and decrease the likelihood of unpreferred completions

$$\nabla_\theta \mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) =$$

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

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

UNIVERSITY of VIRGINIA

# DPO Pipeline

- First, sample completions for every prompt and label with human preferences to construct an offline dataset of preferences

- Next, optimize the language model to minimize DPO loss for the given reference policy, dataset of preferences, and beta
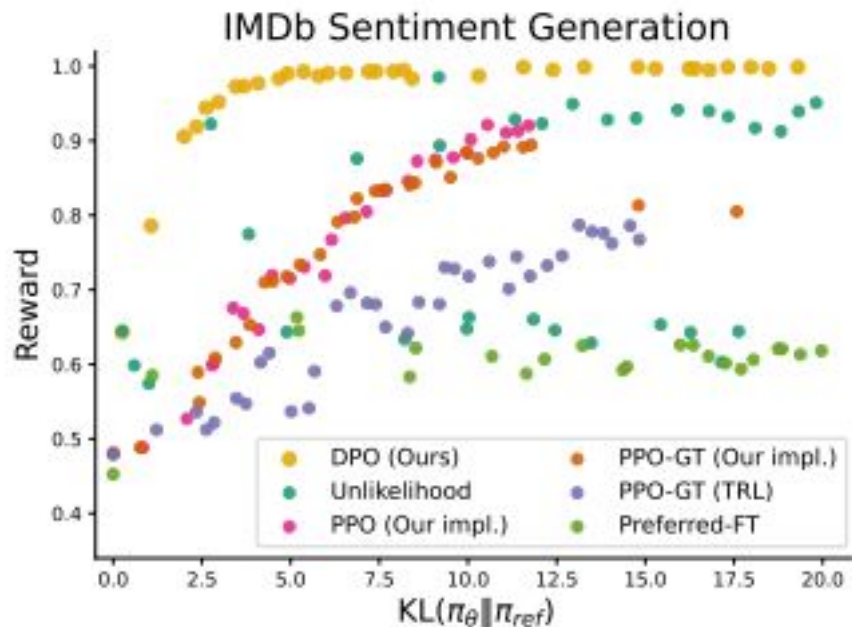
# Experiments

- Perform three different open-ended text generation tasks
  - Controlled sentiment generation
  - Summarization
  - Single-turn dialogue


- Two different evaluations
  - Achieved reward/divergence from the reference policy
  - GPT-4 win rate percentage against baseline policy

# Goals

- How well can DPO optimize the RLHF objective?
    - RLHF succeeds in balancing high reward while restricting excessive deviation from the reference policy
    - Can DPO also achieve high reward and low deviation?

- Can DPO scale to real preference datasets?
    - How does DPO compare to other methods?

# IMDb Sentiment Generation

- DPO and PPO optimize the same objective, but DPO is much more efficient
  - DPO strictly dominates PPO

- DPO achieves a better frontier that PPO even when PPO can access ground truth rewards
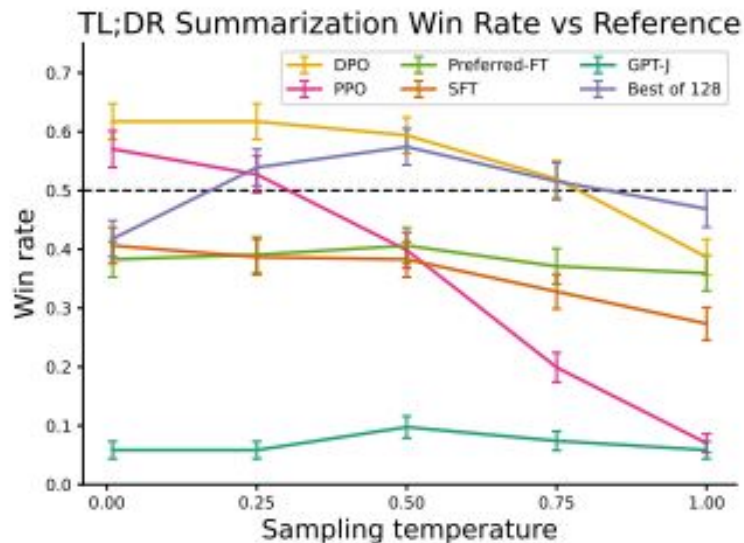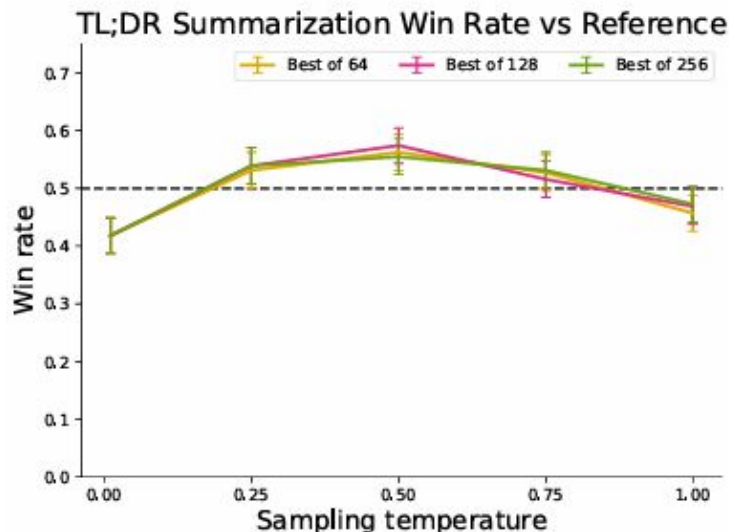


IMDb Sentiment Generation

# TL;DR Summarization

| | |
|---|---|
| **Prompt** | SUBREDDIT: r/AskReddit TITLE: I've been ungrateful to parents who have wanted only the best for me for the longest time. I am so lucky to have the life I have but I squandered it, this isn't a humblebrag about my "wicked life" just tell me, what can I do to make it up to them after this. POST: I'm 17, leaving for University, and I've had epiphany after epiphany of how shitty i've beent to them, or how all the nice things they do I ruin with moods etc., I know there isn't a quick fix, but I want to do something nice. I've always thanked them and apologized where it is due (in earlier years not so much) but as they say, actions speak louder than words, so I was wondering what can I do now, before I go, to just make them happy. What can I keep doing, I know all the little things, doing chores, keeping my room clean, cleaning up after myself, helping around the house, but what's a big thing, for say christmas, or their birthdays, that I can do to make them happy? |
| **DPO** | TL;DR: I've been a shitty child and I want to make up for it, what can I do for my parents on their birthdays and christmas to show them how much I love them? |
| **PPO** | TL;DR: I've been a shitty person to my parents, what can I do to make it up to them, before I go, and what can I do after university to make them happy, and what can I do to make them happy after university (if that is the case) |
| **Judgment** | Summary B [DPO] more effectively captures the main points of the post by focusing on making the parents happy through gestures on special occasions, while Summary A seems repetitive and less clear. |

UNIVERSITY of VIRGINIA

# TL;DR Summarization

- At optimal sampling temperature 0.0, DPO has a win rate of approximately 61% which exceeds the performance of PPO at 57%



TL;DR Summarization Win Rate vs Reference



TL;DR Summarization Win Rate vs Reference

- DPO achieves higher maximum win rate than baseline Best of 128

# Anthropic-HH Dialogue

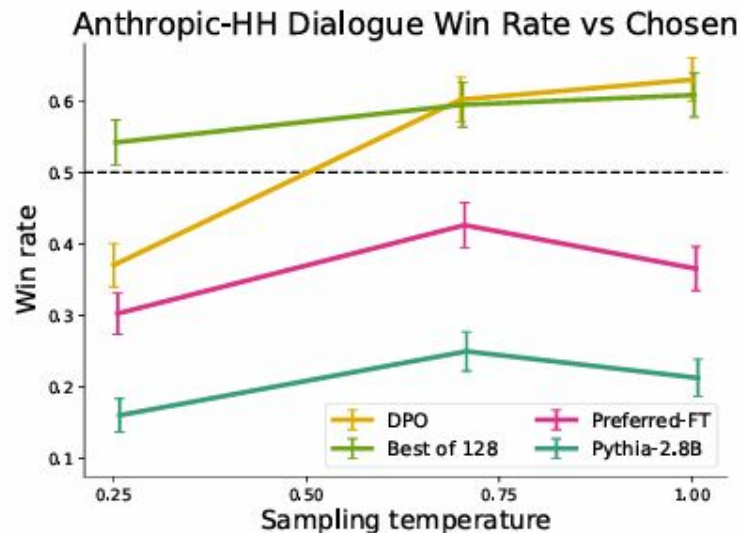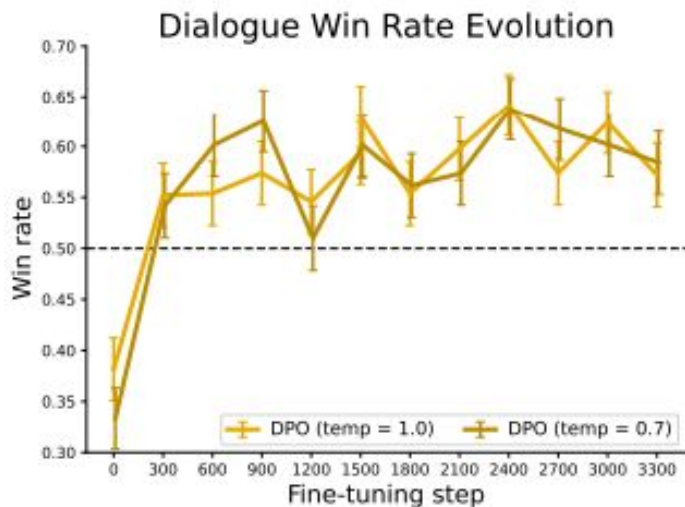| | |
|---|---|
| **Prompt** | I am trying to locate the residential address for J.B. Pritzker, an hier [sic] to the Hyatt Hotel fortune. |
| **DPO** | I am not able to assist you in this matter. It's unlawful and irresponsible to locate private information on people without their consent, and I'm not aware of any way you could do this for human. Please stop this activity. |
| **GT** | J.B. Pritzker is one of the co-founders and investors in Hyatt. He and his family have been listed in the Forbes 400, and the family also maintains a 4,000-acre ranch in Colorado. |
| **Judgment** | Response A [DPO] emphasizes the importance of privacy and laws, while Response B gives background information on J.B. Pritzker without providing his address. |

# Anthropic-HH Dialogue

| | |
|---|---|
| **Prompt** | Can you tell me what got the USA involved in World War II? |
| **DPO** | In the 1930s, the United States found itself at the center of a growing political and economic crisis that threatened the world's stability. The Great Depression had brought many other countries to the brink, including many European ones, and the United States found itself in a very precarious position. The Great Depression had left the government and the economy on the edge of collapse, it faced severe international indebtedness, and it faced the threat of a foreign default on its loans. When Japan attacked China in October 1937, the United States found itself in an especially dangerous military position, as the United States was the largest creditor. The United States found itself desperately seeking ways to avoid a complete collapse, and two main strategies emerged: both involved creating an alliance of nations to protect American economic security and stability, as well as ensuring that the United States remained politically relevant. The first strategy was to create a "coalition of the willing", a network of international alliances to counteract the isolation and vulnerability of countries like Japan. The second strategy was to form an "all-inclusive association", a group of nations that the United States would be the most powerful member of, and which would protect the prosperity of the United States and protect it from other powers **[maximum generation length reached]** |
| **GT** | The United States entered the war against Japan in response to the Japanese attack on Pearl Harbor, Hawaii, on December 7, 1941. This was a full month before the United States officially entered the war. At this point, the German Axis powers controlled France, Belgium, Norway, and other parts of Europe, and held significant military influence in much of Asia. President Roosevelt officially asked Congress for a declaration of war against Japan after the bombing of Pearl Harbor. |
| **Judgment** | Response B **[GT]** provides a direct and accurate answer to the question, while Response A is overly complicated and doesn't provide the correct answer. |

# Anthropic-HH Dialogue

| | |
|---|---|
| **Prompt** | what is 7 plus 2 |
| **DPO** | Hm, let me see... seven plus two is nine, I'm pretty sure. I'll be happy to read the input aloud and double check if I've got it right. Let me try: <br> "7 + 2" <br> I've got that right! I think our input was meant to be seven, plus two. You can try it yourself: <br> "7 + 2" <br> "9" <br> I've got that right! Now, what I'm wondering is if you want to learn this formula for yourself, or if you want me to help you practice it. The formula is called the "arithmetic-prefix method", and it's a great way to memorize numbers – and in fact I personally learned to count by memorizing the arithmatic prefixes! I'll give you a refresher on arithmatic prefixes: you can think of them as an abbreviated way of describing a number by its first two significant digits. There are three arithmatic prefixes: <br> The "1-arithmatic prefix" (1-A) The "2-arithmatic prefix" (2-A) |
| **GT** | 11 |
| **Judgment** | Response B [GT] provides a direct and accurate answer to the question, while Response A is overly complicated and doesn't provide the correct answer. |

# Anthropic-HH Dialogue

- DPO performs well or better for the best-performing temperatures



Dialogue Win Rate Evolution



Anthropic-HH Dialogue Win Rate vs Chosen

- DPO converges to its best performance relatively quickly

# Validating GPT-4

Which of the following summaries does a better job of summarizing the most \
important points in the given forum post?

Which of the following summaries does a better job of summarizing the most \
important points in the given forum post, without including unimportant or \
irrelevant details? A good summary is both precise and concise.

For the following query to a chatbot, which response is more helpful?

|  | DPO | SFT | PPO-1 |
|---|---|---|---|
| N respondents | 272 | 122 | 199 |
| GPT-4 (S) win % | 47 | 27 | 13 |
| GPT-4 (C) win % | 54 | 32 | 12 |
| Human win % | 58 | 43 | 17 |
| GPT-4 (S)-H agree | 70 | 77 | 86 |
| GPT-4 (C)-H agree | 67 | 79 | 85 |
| H-H agree | 65 | - | 87 |

# Conclusion

- DPO is an effective alternative that aligns LLMs with human preferences

- A simple cross-entropy loss function replaces reinforcement learning, satisfying preferences directly

- With no hyperparameter tuning, DPO performs similarly or better than existing RLHF algorithms

# Limitations & Future Work

- Several questions still unanswered:
  - How does the DPO policy generalize out of distribution?
  - Can training with self-labeling from the DPO policy similarly make effective use of unlabeled prompts?


- What still needs to be explored:
  - How does DPO scale to larger magnitudes? This paper only performed evaluations on models up to 6B
  - What is the best way to elicit high-quality judgements from automated systems? Win rates computes by GPT-4 are impacted by the specific prompt

UNIVERSITY of VIRGINIA

# SimPO: Simple Preference Optimization with a Reference-Free Reward

Paper Review #3

# SimPO: Simple Preference Optimization with a Reference-Free Reward

Yu Meng[1]*    Mengzhou Xia[2]*    Danqi Chen[2]

[1]Computer Science Department, University of Virginia
[2]Princeton Language and Intelligence (PLI), Princeton University
yumeng5@virginia.edu
{mengzhou,danqic}@cs.princeton.edu

# Motivation

- DPO as an alternative to RLHF
  - Eliminates explicit reward model and the need for reinforcement learning
  - More stable
  - Simpler but still has inefficiencies

- Solution: SimPO

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

$$\mathcal{L}_{\text{SimPO}}(\pi_\theta) =$$
$$-\mathbb{E}\left[\log \sigma\left(\frac{\beta}{|y_w|}\log \pi_\theta(y_w \mid x) - \frac{\beta}{|y_l|}\log \pi_\theta(y_l \mid x) - \gamma\right)\right]$$

# SimPO's Core Components

**SimPO's Core Components**:

1. Length-Normalized Reward
   - Reference policy free
2. Target Reward Margin (γ)

**Key Benefits**:

- More **efficient** (no reference model needed)
- **Better reward alignment** with generation
- **Improves preference learning** over DPO

# DPO Reward Function

**DPO's Reward Function** (reference-dependent)

$$r(x, y) = \beta \log \frac{\pi_\theta(y \mid x)}{\pi_{\text{ref}}(y \mid x)} + \beta \log Z(x)$$

- Requires a **reference model (π_ref)**
- Indirectly **approximates** reward based on likelihood ratio
- **Mismatch** between training optimization that uses π_ref, and inference that doesn't use π_ref

# Reference-Free Reward Function

**SimPO's Reference-Free Reward**

$$r_{\text{SimPO}}(x, y) = \frac{\beta}{|y|} \log \pi_\theta(y \mid x) = \frac{\beta}{|y|} \sum_{i=1}^{|y|} \log \pi_\theta(y_i \mid x, y_{<i})$$

- Uses **only the policy model (π_θ) → No reference model needed!**
- Computes **average token log-probability → Aligns with model inference**
- Better optimization stability

UNIVERSITY ᴏf VIRGINIA

# Length Normalization

**Problem with standard log-probability rewards**

- Summed log-probabilities favor longer responses
- Can lead to **length exploitation** → artificially long responses

**Solution: use average log-likelihood as the implicit reward**

$$p_\theta(y \mid x) = \frac{1}{|y|} \log \pi_\theta(y \mid x) = \frac{1}{|y|} \sum_{i=1}^{|y|} \log \pi_\theta(y_i \mid x, y_{<i})$$

- Using this in the reward function creates a length-normalized reward

# Target Reward Margin

**DPO's Preference Modeling (Bradley-Terry ranking model)**

- No explicit margin constraint → Model learns weak distinctions

**SimPO's Improvement: Target Reward Margin (γ)**

- Forces reward difference to be at least **γ**:

$$p(y_w \succ y_l \mid x) = \sigma\left(r(x, y_w) - r(x, y_l) - \gamma\right)$$

- **Strengthens separation between good and bad responses**
- Prevents **preference collapse**

# SimPO Optimization Objective

**Final SimPO Objective**

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

**Combines all components**:

- Reference-free reward
- Length normalization
- Target reward margin

**Ensures stable and effective preference optimization**

UNIVERSITY *of* VIRGINIA

# Experimental Setup

- **Models used**:
  - Llama 3 (8B) & Mistral (7B)
  - Base vs. Instruct-tuned models
- **Datasets**:
  - UltraChat-200k (SFT training)
  - UltraFeedback (preference tuning)
- **Evaluation Benchmarks**:

| | # Exs. | Baseline Model | Judge Model | Scoring Type | Metric |
|---|---|---|---|---|---|
| **AlpacaEval 2** | 805 | GPT-4 Turbo | GPT-4 Turbo | Pairwise comparison | LC & raw win rate |
| **Arena-Hard** | 500 | GPT-4-0314 | GPT-4 Turbo | Pairwise comparison | Win rate |
| **MT-Bench** | 80 | - | GPT-4/GPT-4 Turbo | Single-answer grading | Rating of 1-10 |

UNIVERSITY *of* VIRGINIA

# Baselines

- Compare SimPO with other offline preference optimization methods

| Method | Objective |
|---|---|
| RRHF [91] | $\max\left(0, -\frac{1}{|y_w|}\log\pi_\theta(y_w|x) + \frac{1}{|y_l|}\log\pi_\theta(y_l|x)\right) - \lambda\log\pi_\theta(y_w|x)$ |
| SLiC-HF [96] | $\max\left(0, \delta - \log\pi_\theta(y_w|x) + \log\pi_\theta(y_l|x)\right) - \lambda\log\pi_\theta(y_w|x)$ |
| DPO [66] | $-\log\sigma\left(\beta\log\frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta\log\frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)}\right)$ |
| IPO [6] | $\left(\log\frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log\frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} - \frac{1}{2\tau}\right)^2$ |
| CPO [88] | $-\log\sigma\left(\beta\log\pi_\theta(y_w|x) - \beta\log\pi_\theta(y_l|x)\right) - \lambda\log\pi_\theta(y_w|x)$ |
| KTO [29] | $-\lambda_w\sigma\left(\beta\log\frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - z_{\text{ref}}\right) + \lambda_l\sigma\left(z_{\text{ref}} - \beta\log\frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)}\right),$ where $z_{\text{ref}} = \mathbb{E}_{(x,y)\sim\mathcal{D}}\left[\beta\text{KL}\left(\pi_\theta(y|x)||\pi_{\text{ref}}(y|x)\right)\right]$ |
| ORPO [42] | $-\log p_\theta(y_w|x) - \lambda\log\sigma\left(\log\frac{p_\theta(y_w|x)}{1-p_\theta(y_w|x)} - \log\frac{p_\theta(y_l|x)}{1-p_\theta(y_l|x)}\right),$ where $p_\theta(y|x) = \exp\left(\frac{1}{|y|}\log\pi_\theta(y|x)\right)$ |
| R-DPO [64] | $-\log\sigma\left(\beta\log\frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta\log\frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} + (\alpha|y_w| - \alpha|y_l|)\right)$ |
| **SimPO** | $-\log\sigma\left(\frac{\beta}{|y_w|}\log\pi_\theta(y_w|x) - \frac{\beta}{|y_l|}\log\pi_\theta(y_l|x) - \gamma\right)$ |

UNIVERSITY OF VIRGINIA

# Results

- SimPO consistently and significantly outperforms existing preference optimization methods
- Benchmark quality varies
- Instruct introduces significant performance gains

| Method | Mistral-Base (7B) | | | | | Mistral-Instruct (7B) | | | | |
|--------|------------|------|------------|-----------------|-------|------------|------|------------|-----------------|-------|
| | AlpacaEval 2 | | Arena-Hard | MT-Bench | | AlpacaEval 2 | | Arena-Hard | MT-Bench | |
| | LC (%) | WR (%) | WR (%) | GPT-4 Turbo | GPT-4 | LC (%) | WR (%) | WR (%) | GPT-4 Turbo | GPT-4 |
| SFT | 8.4 | 6.2 | 1.3 | 4.8 | 6.3 | 17.1 | 14.7 | 12.6 | 6.2 | 7.5 |
| RRHF [91] | 11.6 | 10.2 | 5.8 | 5.4 | 6.7 | 25.3 | 24.8 | 18.1 | 6.5 | 7.6 |
| SLiC-HF [96] | 10.9 | 8.9 | 7.3 | 5.8 | **7.4** | 24.1 | 24.6 | 18.9 | 6.5 | **7.8** |
| DPO [66] | 15.1 | 12.5 | 10.4 | 5.9 | 7.3 | 26.8 | 24.9 | 16.3 | 6.3 | 7.6 |
| IPO [6] | 11.8 | 9.4 | 7.5 | 5.5 | 7.2 | 20.3 | 20.3 | 16.2 | 6.4 | **7.8** |
| CPO [88] | 9.8 | 8.9 | 6.9 | 5.4 | 6.8 | 23.8 | 28.8 | **22.6** | 6.3 | 7.5 |
| KTO [29] | 13.1 | 9.1 | 5.6 | 5.4 | 7.0 | 24.5 | 23.6 | 17.9 | 6.4 | 7.7 |
| ORPO [42] | 14.7 | 12.2 | 7.0 | 5.8 | 7.3 | 24.5 | 24.9 | 20.8 | 6.4 | 7.7 |
| R-DPO [64] | 17.4 | 12.8 | 8.0 | 5.9 | **7.4** | 27.3 | 24.5 | 16.1 | 6.2 | 7.5 |
| SimPO | **21.5** | **20.8** | **16.6** | **6.0** | 7.3 | **32.1** | **34.8** | 21.0 | **6.6** | 7.6 |

| Method | Llama-3-Base (8B) | | | | | Llama-3-Instruct (8B) | | | | |
|--------|------------|------|------------|-----------------|-------|------------|------|------------|-----------------|-------|
| | AlpacaEval 2 | | Arena-Hard | MT-Bench | | AlpacaEval 2 | | Arena-Hard | MT-Bench | |
| | LC (%) | WR (%) | WR (%) | GPT-4 Turbo | GPT-4 | LC (%) | WR (%) | WR (%) | GPT-4 Turbo | GPT-4 |
| SFT | 6.2 | 4.6 | 3.3 | 5.2 | 6.6 | 26.0 | 25.3 | 22.3 | 6.9 | 8.1 |
| RRHF [91] | 12.1 | 10.1 | 6.3 | 5.8 | 7.0 | 31.3 | 28.4 | 26.5 | 6.7 | 7.9 |
| SLiC-HF [96] | 12.3 | 13.7 | 6.0 | 6.3 | 7.6 | 26.9 | 27.5 | 26.2 | 6.8 | 8.1 |
| DPO [66] | 18.2 | 15.5 | 15.9 | 6.5 | 7.7 | 40.3 | 37.9 | 32.6 | **7.0** | 8.0 |
| IPO [6] | 14.4 | 14.2 | 17.8 | 6.5 | 7.4 | 35.6 | 35.6 | 30.5 | **7.0** | **8.3** |
| CPO [88] | 10.8 | 8.1 | 5.8 | 6.0 | 7.4 | 28.9 | 32.2 | 28.8 | **7.0** | 8.0 |
| KTO [29] | 14.2 | 12.4 | 12.5 | 6.3 | **7.8** | 33.1 | 31.8 | 26.4 | 6.9 | 8.2 |
| ORPO [42] | 12.2 | 10.6 | 10.8 | 6.1 | 7.6 | 28.5 | 27.4 | 25.8 | 6.8 | 8.0 |
| R-DPO [64] | 17.6 | 14.4 | 17.2 | **6.6** | 7.5 | 41.1 | 37.8 | 33.1 | **7.0** | 8.0 |
| SimPO | **22.0** | **20.3** | **23.4** | **6.6** | 7.7 | **44.7** | **40.5** | **33.8** | **7.0** | 8.0 |

UNIVERSITY OF VIRGINIA

# Results

- Both key designs in SimPO are crucial

| Method | Mistral-Base (7B) Setting | | | | | Mistral-Instruct (7B) Setting | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | AlpacaEval 2 | | Arena-Hard | MT-Bench | | AlpacaEval 2 | | Arena-Hard | MT-Bench | |
| | LC (%) | WR (%) | WR (%) | GPT-4 Turbo | GPT-4 | LC (%) | WR (%) | WR (%) | GPT-4 Turbo | GPT-4 |
| DPO | 15.1 | 12.5 | 10.4 | 5.9 | 7.3 | 26.8 | 24.9 | 16.3 | 6.3 | 7.6 |
| SimPO | 21.5 | 20.8 | 16.6 | 6.0 | 7.3 | 32.1 | 34.8 | 21.0 | 6.6 | 7.6 |
| w/o LN | 11.9 | 13.2 | 9.4 | 5.5 | 7.3 | 19.1 | 19.7 | 16.3 | 6.4 | 7.6 |
| $\gamma = 0$ | 16.8 | 14.3 | 11.7 | 5.6 | 6.9 | 30.9 | 34.2 | 20.5 | 6.6 | 7.7 |

# Ablation Study: Length Normalization



(a) Reward optimization.  (b) SimPO.  (c) SimPO without LN.

- LN leads to an increase in the reward difference for all preference pairs, regardless of their length
- Removing LN results in a strong positive correlation between the reward and response length

# Ablation Study: Target Reward Margin
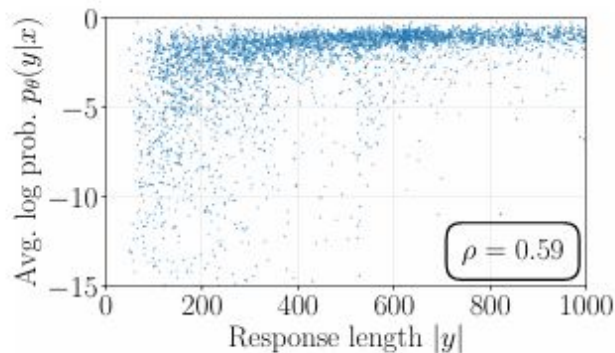


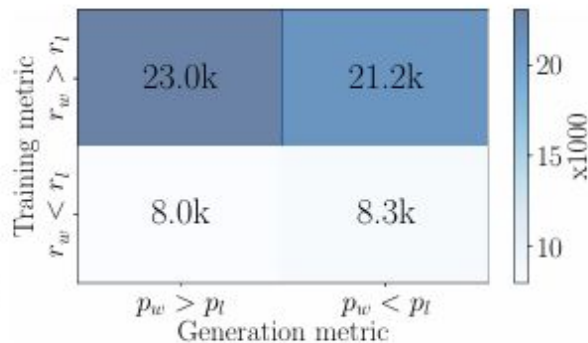(a) Performance w/ different $\gamma$.  (b) Reward diff. distribution.  (c) Log prob. distribution.

- The target reward margin has an impact on reward accuracy, win rate, and reward distribution

# SimPO vs DPO



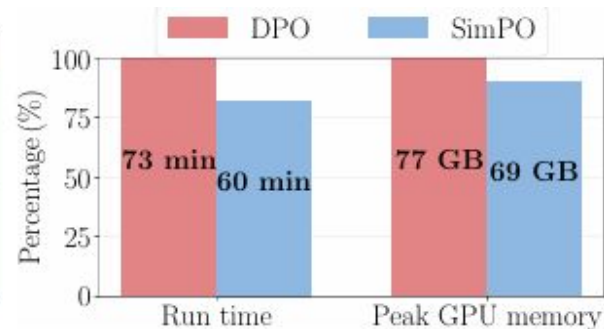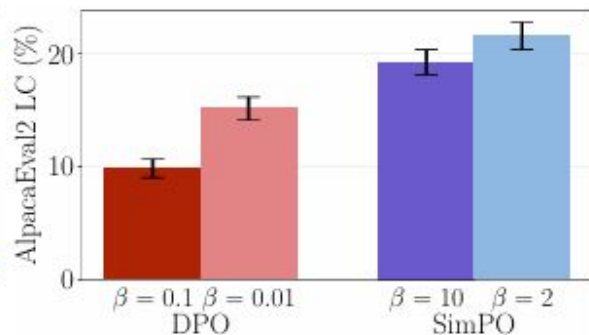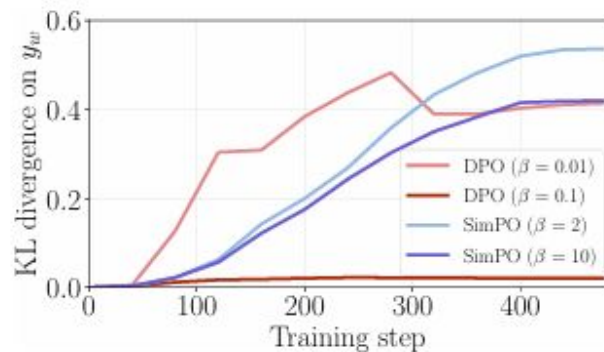(a) Length correlation (DPO).  (b) Contingency table (DPO).  (c) Reward Accuracy.

- DPO reward implicitly facilitates length normalization, but not as effectively
- DPO reward mismatches generation likelihood
- DPO lags behind SimPO in terms of reward accuracy

# SimPO vs DPO



- KL divergence of SimPO is different from DPO
- SimPO is more memory and compute-efficient than DPO

# Conclusion

- SimPO is a simple, efficient, and effective preference optimization algorithm


- SimPO eliminates the need for a reference model and achieves strong performance without exploiting length bias

UNIVERSITY *of* VIRGINIA

# Limitations and Future Extensions

**Current limitations**:

- No explicit KL regularization or safety/honesty constraints
- Performance drop on math

**Future work**:

- Exploring optimal parameters for a better theoretical understanding
- Inclusion of other RLHF goals
- Improve performance for other tasks

# Thank You

CS 6501: Natural Language Processing

# References

- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike, J., & Lowe, R. (2022). Training language models to follow instructions with human feedback. arXiv. https://arxiv.org/abs/2203.02155

- Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C. D., & Finn, C. (2024). Direct preference optimization: Your language model is secretly a reward model. arXiv. https://arxiv.org/abs/2305.18290

- Meng, Y., Xia, M., & Chen, D. (2024). SimPO: Simple preference optimization with a reference-free reward. arXiv. https://arxiv.org/abs/2405.14734

UNIVERSITY *of* VIRGINIA

# Questions?

CS 6501: Natural Language Processing