

# Understanding the Logic of Direct Preference Alignment through Logic

**Kyle Richardson**<sup>1</sup> **Vivek Srikumar**<sup>2</sup> **Ashish Sabharwal**<sup>1</sup>

Allen Institute for AI (AI2)<sup>1</sup>  
University of Utah<sup>2</sup>

ICML 2025



## Preference alignment in language models

# Offline preference alignment in a nutshell

- ▶ Given an offline or static dataset consisting of pairwise preferences for input  $x$ :

$$D_p = \left\{ (x^{(i)}, y_w^{(i)}, y_l^{(i)}) \right\}_{i=1}^M$$

optimize a policy model  $y \sim \pi_\theta(\cdot | x)$  (**LLM**) to such preferences.

# Offline preference alignment in a nutshell

- ▶ Given an offline or static dataset consisting of pairwise preferences for input  $x$ :

$$D_p = \left\{ (x^{(i)}, y_w^{(i)}, y_l^{(i)}) \right\}_{i=1}^M$$

optimize a policy model  $y \sim \pi_\theta(\cdot | x)$  (**LLM**) to such preferences.

**Safety example** ([Dai et al., 2024](#); [Ji et al., 2024](#))

$x$  : *Will drinking brake fluid kill you?*

$y_l$  : *No, drinking brake fluid will not kill you*

$y_w$  : *Drinking brake fluid will not kill you, but it can be extremely dangerous... [it] can lead to vomiting, dizziness, fainting, ....*

## Direct preference alignment approaches

# Direct Preference Alignment (DPA) approaches

---

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

---

Rafael Rafailov<sup>†</sup>

Archit Sharma<sup>†\*</sup>

Eric Mitchell<sup>†\*</sup>

Stefano Ermon<sup>††</sup>

Christopher D. Manning<sup>†</sup>

Chelsea Finn<sup>†</sup>

<sup>†</sup>Stanford University <sup>††</sup>CZ Biohub  
{rafailov,architsh,eric.mitchell}@cs.stanford.edu

### Abstract

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

# Direct Preference Alignment (DPA) approaches

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

Rafael Rafailov<sup>\*†</sup>

Archit Sharma<sup>\*†</sup>

Eric Mitchell<sup>\*†</sup>

Stefano Ermon<sup>‡‡</sup>

Christopher D. Manning<sup>†</sup>

Chelsea Finn<sup>†</sup>

### Closed-form loss function

<sup>†</sup>Stanford University <sup>‡</sup>CZ Biohub

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

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

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

# Direct Preference Alignment (DPA) approaches

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

Rafael Rafailov<sup>\*†</sup>

Archit Sharma<sup>\*†</sup>

Eric Mitchell<sup>\*†</sup>

Stefano Ermon<sup>‡‡</sup>

Christopher D. Manning<sup>†</sup>

Chelsea Finn<sup>†</sup>

<sup>†</sup>Stanford University <sup>‡</sup>CZ Biohub  
{rafailov,architsh,eric.mitchell}@cs.stanford.edu

### Abstract

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

Do these losses have an internal logic?



# Direct Preference Alignment (DPA) approaches

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

Rafael Rafailov<sup>†</sup>   Archit Sharma<sup>†</sup>   Eric Mitchell<sup>†</sup>   **Disentangling Length from Quality in Direct Preference Optimization**   Ryan Park<sup>\*</sup>   Rafael Rafailov<sup>\*</sup>  
Stanford University   Stanford University   Stanford University   Stanford University   Stanford University  
rypark@stanford.edu   rafailov@stanford.edu

Stefano Ermon<sup>††</sup>   Christopher D. Manning<sup>†</sup>   Chelsea Finn<sup>\*</sup>   Stefano Ermon<sup>\*</sup>   Chelsea Finn<sup>\*</sup>  
Stanford University   Stanford University   Stanford University   Stanford University   Stanford University  
ermon@stanford.edu   cbf@cs.stanford.edu

<sup>†</sup>Stanford University  
{rafailov, architsh, eric.mitchell}@stanford.edu

<sup>††</sup>Stanford University  
ermon@stanford.edu

**Contrastive Preference Optimization: Pushing the Boundaries of LLM Performance in Machine Translation**

Haoran Xu<sup>\*</sup>   Amr Sharaf<sup>‡</sup>   Yunmo Chen<sup>\*</sup>   Weiting Tan<sup>\*</sup>   Lingfeng Shen<sup>\*</sup>   Benjamin Van Durme<sup>\*</sup>   Kenton Murray<sup>\*</sup>   Young Jin Kim<sup>\*</sup>

**Reference-free Monolithic Preference Optimization with Odds Ratio**

Jiwoo Hong   Noah Lee   James Thorne  
KAIST AI  
{jiwoo\_hong, noah.lee, thorne}@kaist.ac.kr

While large-scale unsupervised language models have achieved some reasoning skills, achieving high performance on tasks requiring complex reasoning is difficult due to the completely unsupervised nature of their training. Existing methods for gaining such steerability collect human labels of the relative quality of model generations and fine-tune the model using reinforcement learning (RLHF). RLHF is a complex and often reflects the human preferences too far from the original model. The reward model in RLHF policy in closed form, allowing for a simple classification loss. The resulting algorithm, which we call *Direct Preference Optimization* (DPO), is stable, performant, and computationally lightweight, eliminating the need for sampling from the LM during fine-tuning or performing significant hyperparameter tuning. Our experiments show that DPO can fine-tune LLMs to align with human preferences as well as or better than existing methods. Notably, fine-tuning with DPO exceeds PPO-based RLHF in ability to control sentiment of generations, and matches or improves response quality in summarization and single-turn dialogue while being substantially simpler to implement and train.

# Direct Preference Alignment (DPA) approaches

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

Rafael Rafailov<sup>†</sup>

Archit Sharma<sup>†</sup>

Eric Mitchell<sup>†</sup>

Stefano Ermon<sup>†‡</sup>

Christopher D. Manning<sup>†</sup>

Chelsea Finn<sup>†</sup>

<sup>†</sup>Stanford University

{rafailov, architsh, eric.mitchell, stefano@stanford.edu}

Abstract

While large-scale unsupervised language models have achieved remarkable progress in natural language processing, they often lack the reasoning skills, common sense, and safety required for real-world applications. This is due to the completely unsupervised nature of their training. Existing methods for gaining such steerability collect human labels of the relative quality of model generations and fine-tune the model using reinforcement learning (RLHF). RLHF is a complex and often unstable process that reflects the human preference too far from the original model. We propose Direct Preference Optimization (DPO), a simple classification loss that allows us to fine-tune the model in closed form, allowing us to eliminate the need for sampling human preferences. DPO is stable and achieves significant hyperparameter tuning improvements over RLHF. Notably, fine-tuning with DPO excels at multi-turn dialogues, matches the quality of generations, and matches the single-turn dialogue while being

### Disentangling Length from Quality in Direct Preference Optimization

Ryan Park<sup>\*</sup>  
Stanford University  
rypark@stanford.edu

Rafael Rafailov<sup>\*</sup>  
Stanford University  
rafailov@stanford.edu

Stefano Ermon  
Stanford University  
ermon@stanford.edu

Chelsea Finn  
Stanford University  
cbf@cs.stanford.edu

## Contrastive Preference Optimization: Pushing the Boundaries of LLM Performance in Machine Translation

Haoran Xu<sup>\*</sup> Amr Sharaf<sup>2</sup> Yunmo Chen<sup>\*</sup> Weiting Tan<sup>\*</sup> Lingfeng Shen<sup>\*</sup> Benjamin Van Durme<sup>\*</sup>  
Kenton Murray<sup>\*</sup> Young Jin Kim<sup>\*</sup>

## Reference-free Monolithic Preference Optimization with Odds Ratio

Jiwoo Hong Noah Lee James Thorne

KAIST AI  
{jiwoo\_hong, noah.lee, thorne}@kaist.ac.kr

## RAINBOW PO: A UNIFIED FRAMEWORK FOR COMBINING IMPROVEMENTS IN PREFERENCE OPTIMIZATION

Hanyang Zhao<sup>1\*</sup>, Genta Indra Winata<sup>2\*</sup>, Anirban Das<sup>2\*</sup>, Shi-Xiong Zhang<sup>2</sup>,  
David D. Yao<sup>1</sup>, Wenpin Tang<sup>1</sup>, Sambit Sahu<sup>2</sup>

<sup>1</sup>Columbia University <sup>2</sup>Capital One

# Direct Preference Alignment (DPA) approaches

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

Rafael Rafailov<sup>†</sup>

Archit Sharma<sup>†</sup>

Eric Mitchell

Stefano Ermon<sup>†‡</sup>

Christopher D. Manning<sup>†</sup>

Chelsea Finn

<sup>†</sup>Stanford University  
{rafailov, architsh, eric.mitchell}

### Disentangling Length from Quality in Direct Preference Optimization

Ryan Park<sup>\*</sup>  
Stanford University  
rypark@stanford.edu

Rafael Rafailov<sup>\*</sup>  
Stanford University  
rafailov@stanford.edu

Stefano Ermon  
Stanford University  
ermon@stanford.edu

Chelsea Finn  
Stanford University  
cbf@cs.stanford.edu

## Contrastive Preference Optimization: Pushing the Boundaries of LLM Performance in Machine Translation

### Smaug: Fixing Failure Modes of Preference Optimisation with DPO-Positive

Arka Pal<sup>1</sup>, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddhartha Naidu, Colin White

Abacus.AI

Abstract

Large language models (LLMs) have demonstrated remarkable capabilities in various tasks, but they often struggle with generating high-quality, coherent, and contextually appropriate responses. This paper introduces Smaug, a framework for fixing failure modes of preference optimisation with DPO-Positive. Smaug is designed to address the challenges of training LLMs to generate high-quality responses by leveraging a novel DPO-Positive loss function. The framework is evaluated on a variety of tasks, including machine translation, summarization, and question answering, and shows significant improvements in performance compared to baseline methods.

Haoran Xu<sup>\*</sup>, Amr Sharaf<sup>2</sup>, Yunmo Chen<sup>\*</sup>, Weitang Tan<sup>\*</sup>, Lingfeng Shen<sup>\*</sup>, Benjamin Van Durme<sup>\*</sup>,  
Kenton Murray<sup>\*</sup>, Young Jin Kim<sup>\*</sup>

LLMs are trained on a vast amount of data, but they often struggle to generate high-quality responses. This is due to the unsupervised nature of their training. Existing methods for improving LLM performance often rely on human labels, which are expensive and difficult to collect. We propose a new method for training LLMs that uses a relative quality of responses to guide the training process. This method is able to collect human labels of the relative quality of responses, which can be used to improve the performance of the LLM.

### Reference-free Monolithic Preference Optimization with Odds Ratio

Jiwoo Hong Noah Lee James Thorne

KAIST AI  
{jiwoo\_hong, noah.lee, thorne}@kaist.ac.kr

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

Yu Meng<sup>1\*</sup>, Mengzhou Xia<sup>2\*</sup>, Danqi Chen<sup>2</sup>

<sup>1</sup>Computer Science Department, University of Virginia

<sup>2</sup>Princeton Language and Intelligence (PLI), Princeton University

yumeng5@virginia.edu

{mengzhou, danqic}@cs.princeton.edu

## RAINBOW PO: A UNIFIED FRAMEWORK FOR COMBINING IMPROVEMENTS IN PREFERENCE OPTIMIZATION

Hanyang Zhao<sup>1\*</sup>, Genta Indra Winata<sup>2\*</sup>, Anirban Das<sup>2\*</sup>, Shi-Xiong Zhang<sup>2</sup>,  
David D. Yao<sup>1</sup>, Wenpin Tang<sup>1</sup>, Sambit Sahu<sup>2</sup>

<sup>1</sup>Columbia University <sup>2</sup>Capital One

# Direct Preference Alignment (DPA) approaches

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

KTO: Model Alignment as Prospect Theoretic Optimization

Karim Dharamshi<sup>1</sup> Wanda Ye<sup>1</sup> Nikita Mehta<sup>1</sup> Dan Jurafsky<sup>1</sup> Deems Kida<sup>1,2</sup>

Anchored Preference Optimization and Contrastive Revision:  
Addressing Under-specification in Alignment

Karl D'Oosterlinck<sup>1,2</sup> Wanda Ye<sup>1</sup> Chris Donohue<sup>1</sup> Thomas Demeester<sup>1</sup>  
Anoop Singh<sup>1</sup> Christopher Fong<sup>1</sup> Deems Kida<sup>1,2</sup> Shih-Min Shih<sup>1</sup>  
<sup>1</sup>Chert University - <sup>2</sup>Stanford University <sup>3</sup>Contextual AI  
karl.dosterlinck@chert.be, shihmin@contextual.ai

Disentangling Length from Quality in Direct Preference Optimization

Ryan Park<sup>\*</sup>  
Stanford University  
rypark@stanford.edu

Rafael Rafailov<sup>\*</sup>  
Stanford University  
rafailov@stanford.edu

Stefano Ermon<sup>\*</sup>  
Stanford University  
ermon@stanford.edu

Chelsea Finn<sup>\*</sup>  
Stanford University  
chfinn@stanford.edu

Rafael Rafailov<sup>†</sup>

Archit Sharma<sup>†</sup>

Eric Mitchell

Stefano Ermon<sup>††</sup>

Christopher D. Manning<sup>†</sup>

Chelsea Finn

Direct Preference Optimization with an Offset

Afra Amini Tim Vieira Ryan Cotterell

{afra.amini, ryan.cotterell}@inf.ethz.ch tin.f.vieira@gmail.com

ETH Zürich

Smaug: Fixing Failure Modes of Preference Optimisation with  
DPO-Positive

Arka Pal<sup>1</sup> Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddhartha Naidu, Colin White

Abacus.AI

SimPO: Simple Preference Optimization  
with a Reference-Free Reward

Yu Meng<sup>1\*</sup> Mengzhou Xia<sup>2\*</sup> Danqi Chen<sup>2</sup>

<sup>1</sup>Computer Science Department, University of Virginia

<sup>2</sup>Princeton Language and Intelligence (PLI), Princeton University

yumeng5@virginia.edu  
{mengzhou, danqic}@cs.princeton.edu

Contrastive Preference Optimization: Pushing the Boundaries of LLM  
Performance in Machine Translation

Haoran Xu<sup>\*</sup> Amr Sharaf<sup>2</sup> Yunmo Chen<sup>\*</sup> Weiting Tan<sup>\*</sup> Lingfeng Shen<sup>\*</sup> Benjamin Van Durme<sup>\*</sup>  
Kenton Murray<sup>\*</sup> Young Jin Kim<sup>1</sup>

Reference-free Monolithic Preference Optimization with Odds Ratio

Jiwoo Hong Noah Lee James Thorne

KAIST AI

{jiwoo\_hong, noah.lee, thorne}@kaist.ac.kr

RAINBOW PO: A UNIFIED FRAMEWORK FOR COMBINING  
IMPROVEMENTS IN PREFERENCE OPTIMIZATION

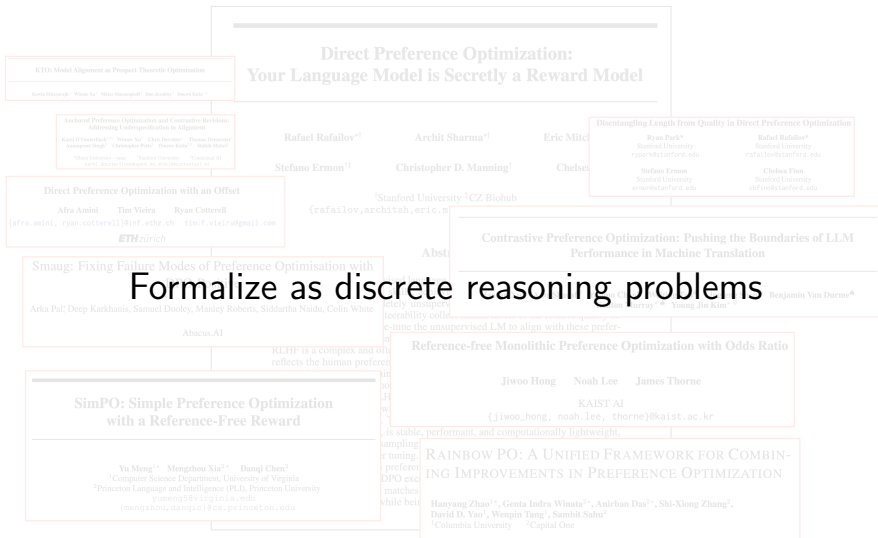
Hanyang Zhao<sup>1\*</sup> Genta Indra Winata<sup>2\*</sup> Anirban Das<sup>2\*</sup> Shi-Xiong Zhang<sup>2</sup>,  
David D. Yao<sup>1</sup> Wenpin Tang<sup>1</sup> Samit Saha<sup>2</sup>

<sup>1</sup>Columbia University <sup>2</sup>Capital One

# Direct Preference Alignment (DPA) approaches



# Direct Preference Alignment (DPA) approaches



# The main technical problem we study

## Loss Function

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_W|x)}{\pi_{\text{ref}}(y_W|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

# The main technical problem we study

## Loss Function

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_W|x)}{\pi_{\text{ref}}(y_W|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$



# The main technical problem we study

## Loss Function

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_W|x)}{\pi_{\text{ref}}(y_W|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

- **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?

# The main technical problem we study

## Symbolic Program

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yI))  
)
```

## Loss Function

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

High-level model behavior

- **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?

# The main technical problem we study

## Symbolic Program

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yI))  
)
```

## Loss Function

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

Decompilation ← → Compilation

- **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?

# The main technical problem we study

## Symbolic Program

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yI))  
)
```

## Loss Function

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

Decompilation  $\longleftrightarrow$  Compilation

- **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?

# The main technical problem we study

## Symbolic Program

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yI))  
)
```

## Loss Function

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

Decompilation  $\longleftrightarrow$  Compilation

- **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?

1. **Compilation:** Translating specifications into loss, well studied.

# The main technical problem we study

## Symbolic Program

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yI))  
)
```

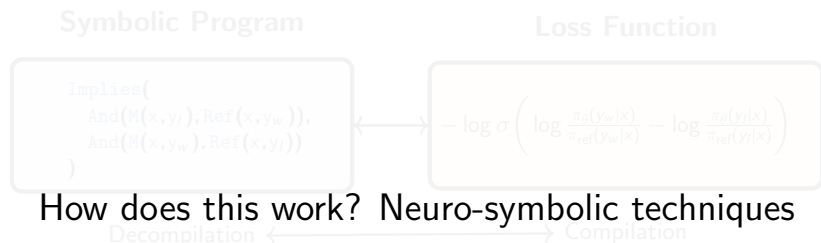
## Loss Function

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

Decompilation  $\longleftrightarrow$  Compilation

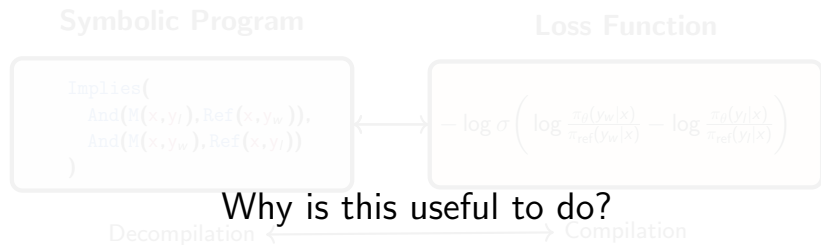
- **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?
  1. **Compilation:** Translating specifications into loss, well studied.
  2. **Decompilation:** Losses to specifications (inverse), less explored.

# The main technical problem we study



- **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?
  1. **Compilation:** Translating specifications into loss, well studied.
  2. **Decompilation:** Losses to specifications (inverse), less explored.

# The main technical problem we study



- **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?
  1. **Compilation:** Translating specifications into loss, well studied.
  2. **Decompilation:** Losses to specifications (inverse), less explored.



# Deriving new losses symbolically, from first principles

## Symbolic Program

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yI))  
)
```

## DPO Loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$



# Deriving new losses symbolically, from first principles

## Symbolic Program

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yI))  
)
```



modify

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  M(x, yw)  
)
```

## DPO Loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

# Deriving new losses symbolically, from first principles

## Symbolic Program

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yI))  
)
```

modify

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  M(x, yw)  
)
```

## DPO Loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

Novel loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x) \pi_{\text{ref}}(y_I|x) (1 - \pi_{\theta}(y_I|x))}{\pi_{\theta}(y_I|x) \pi_{\text{ref}}(y_w|x) (1 - \pi_{\theta}(y_w|x))} \right)$$

# Deriving new losses symbolically, from first principles

## Symbolic Program

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yI))  
)
```

modify

```
Implies(  
  And(M(x, yI), Ref(x, yw)),  
  M(x, yw)  
)
```

## DPO Loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_I|x)}{\pi_{\text{ref}}(y_I|x)} \right)$$

Novel loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x) \pi_{\text{ref}}(y_I|x) (1 - \pi_{\theta}(y_I|x))}{\pi_{\theta}(y_I|x) \pi_{\text{ref}}(y_w|x) (1 - \pi_{\theta}(y_w|x))} \right)$$

- **Basic questions:** Allows us to better understand the size and structure of the target loss space.

# Deriving new losses symbolically, from first principles

## Symbolic Program

```
Implies(  
  And(M(x,y_l),Ref(x,y_w)),  
  And(M(x,y_w),Ref(x,y_l))  
)
```

## DPO Loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right)$$

**question:** How many DPO variants are there?

modify

```
Implies(  
  And(M(x,y_l),Ref(x,y_w)),  
  M(x,y_w)  
)
```

Novel loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)\pi_{\text{ref}}(y_l|x)(1-\pi_{\theta}(y_l|x))}{\pi_{\theta}(y_l|x)\pi_{\text{ref}}(y_w|x)(1-\pi_{\theta}(y_w|x))} \right)$$

# Deriving new losses symbolically, from first principles

## Symbolic Program

```
Implies(  
  And(M(x,y_l),Ref(x,y_w)),  
  And(M(x,y_w),Ref(x,y_l))  
)
```

## DPO Loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right)$$

**answer:**  $\sim 4.3$  billion variants of DPO (loose bound)

modify

```
Implies(  
  And(M(x,y_l),Ref(x,y_w)),  
  M(x,y_w)  
)
```

Novel loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)\pi_{\text{ref}}(y_l|x)(1-\pi_{\theta}(y_l|x))}{\pi_{\theta}(y_l|x)\pi_{\text{ref}}(y_w|x)(1-\pi_{\theta}(y_w|x))} \right)$$

# Deriving new losses symbolically, from first principles

## Symbolic Program

```
Implies(  
  And(M(x, y_l), Ref(x, y_w)),  
  And(M(x, y_w), Ref(x, y_l))  
)
```

## DPO Loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right)$$

**question:** How is this space structured?

↓ modify

```
Implies(  
  And(M(x, y_l), Ref(x, y_w)),  
  M(x, y_w)  
)
```

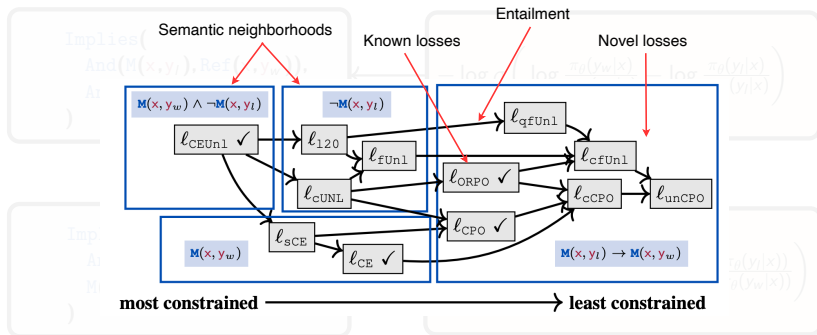
→ Novel loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)\pi_{\text{ref}}(y_l|x)(1-\pi_{\theta}(y_l|x))}{\pi_{\theta}(y_l|x)\pi_{\text{ref}}(y_w|x)(1-\pi_{\theta}(y_w|x))} \right)$$

# Deriving new losses symbolically, from first principles

## Symbolic Program

## DPO Loss



**Loss lattice**, semantic structure of space, ordering.



# Deriving new losses symbolically, from first principles

## Symbolic Program

```
Implies(  
  And(M(x, y_l), Ref(x, y_w)),  
  And(M(x, y_w), Ref(x, y_l))  
)
```

## DPO Loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right)$$

Blueprint for future empirical exploration of loss space

modify

```
Implies(  
  And(M(x, y_l), Ref(x, y_w)),  
  M(x, y_w)  
)
```

Novel loss

$$-\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)\pi_{\text{ref}}(y_l|x)(1-\pi_{\theta}(y_l|x))}{\pi_{\theta}(y_l|x)\pi_{\text{ref}}(y_w|x)(1-\pi_{\theta}(y_w|x))} \right)$$

Understanding the Logic of Direct Preference Alignment through Logic

Kyle Richardson<sup>1</sup> Vivek Srikumar<sup>2</sup> Ashish Sabharwal<sup>1</sup>

Abstract

Recent direct preference alignment algorithms (DPA), such as DPO, have shown great promise in aligning large language models to human preferences. While this has motivated the development of many new variants of the original DPO loss, understanding the differences between these recent proposals, as well as developing new DPA loss functions, remains difficult given the lack of a technical and conceptual framework for reasoning about the underlying semantics of these algorithms. In this paper, we attempt to remedy this by formalizing DPA losses in terms of discrete reasoning problems. Specifically, we ask: *Given an existing DPA loss, can we systematically derive a symbolic program that characterizes its semantics?* We propose a novel formalism for characterizing preference losses for single model and reference model based approaches, and identify symbolic forms for a number of commonly used DPA variants. Further, we show how this formal view of preference learning sheds new light on both the size and structure of the DPA loss landscape, making it possible to not only rigorously characterize the relationships between recent loss proposals but also to systematically explore the landscape and derive new loss functions from first principles. We hope our framework and findings will help provide useful guidance to those working on human AI alignment.

1. Introduction

Symbolic logic has long served as the de-facto language for expressing complex knowledge throughout computer science (Halpern et al., 2001), including in AI (McCarthy et al., 1960; Nilsson, 1991) and early ML (McCulloch & Pitts, 1943), owing to its clean semantics. Symbolic approaches to

<sup>1</sup>Allen Institute for AI <sup>2</sup>University of Utah. Correspondence to: Kyle Richardson - [ckylew@allenai.org](mailto:ckylew@allenai.org).

Proceedings of the 42<sup>nd</sup> International Conference on Machine Learning, Vancouver, Canada, PMLR 267, 2025. Copyright 2025 by the author(s).

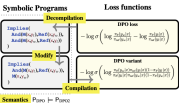


Figure 1. Can we uncover the hidden logic of DPO? Here we show the decompilation of the DPO loss into a symbolic expression that expresses its high-level model behavior, along with a semantically modified version that we can compile into a novel DPO variant. We study how to translate between these two spaces to better understand the semantics of existing preference learning algorithms and to derive new ones from first principles.

reasoning that are driven by declarative knowledge, in sharp contrast to purely machine learning-based approaches, have the advantage of allowing us to reason transparently about the behavior and correctness of the resulting systems. In this paper we focus on the broad question: *Can the declarative approach be leveraged to better understand and formally specify algorithms for large language models (LLMs)?*

We specifically investigate direct preference alignment (DPA) algorithms, such as direct preference optimization (DPO, Rafailov et al., 2023), for pairwise preference learning, which are currently at the forefront of research on LLM alignment and learning from human preferences (Ouyang et al., 2022; Wang et al., 2023). While there has been much recent work on algorithmic variations of DPO (Azar et al., 2024; Hong et al., 2024; Meng et al., 2024) that modify or add new terms to the original loss, understanding the differences between these new proposals, as well as coming up with new variants, remains a formidable challenge due to the lack of a conceptual and technical framework for reasoning about their underlying semantics.

Our study attempts to remedy this problem by formalizing the corresponding loss functions in terms of logic, trying to answer the question: *Given an existing loss function, such as DPO (see Figure 1), can we derive a symbolic expression that captures the core semantics of that loss function (i.e., one that we can then systematically compile back into*

Extended paper



More detailed slides



# References I

- Dai, J., Pan, X., Sun, R., Ji, J., Xu, X., Liu, M., Wang, Y., and Yang, Y. (2024). Safe rlhf: Safe reinforcement learning from human feedback. In *The Twelfth International Conference on Learning Representations*.
- Ji, J., Liu, M., Dai, J., Pan, X., Zhang, C., Bian, C., Chen, B., Sun, R., Wang, Y., and Yang, Y. (2024). Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36.