

# 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>

March 2025



# General purpose large language models (LLMs)



user



ChatGPT

Response

User

What can you tell me about the seminar on formal languages and neural networks (FLaNN) that's held online? Please be brief but tell me some example topics they might discuss.

Assistant

The Seminar on Formal Languages and Neural Networks (FLaNN) is an online event that focuses on the intersection of formal language theory and neural network research. Topics they might discuss include:

- Applications of formal languages in neural network architectures
- Neural network models for language processing
- Algorithmic learning theory related to formal languages
- Automata theory and neural computation
- Symbolic and subsymbolic integration in AI systems

# General purpose large language models (LLMs)



user



ChatGPT

Response

User

What can you tell me about the seminar on formal languages and neural networks (FLaNN) that's held online? Please be brief but tell me some example topics they might discuss.

Assistant

The Seminar on Formal Languages and Neural Networks (FLaNN) is an online event that focuses on the intersection of formal language theory and neural network research. Topics they might discuss include:

- Applications of formal languages in neural network architectures
- Neural network models for language processing
- Algorithmic learning theory related to formal languages
- Automata theory and neural computation
- Symbolic and subsymbolic integration in AI systems

- ▶ We routinely use LLMs for writing code, integrate within complex systems, use for tasks like scientific discovery.

# General purpose large language models (LLMs)



User



User

What can you tell me about the seminar on formal languages and neural networks (FLaNN) that's held online? Please be brief but tell me some example topics they might discuss.



ChatGPT



Response

Assistant

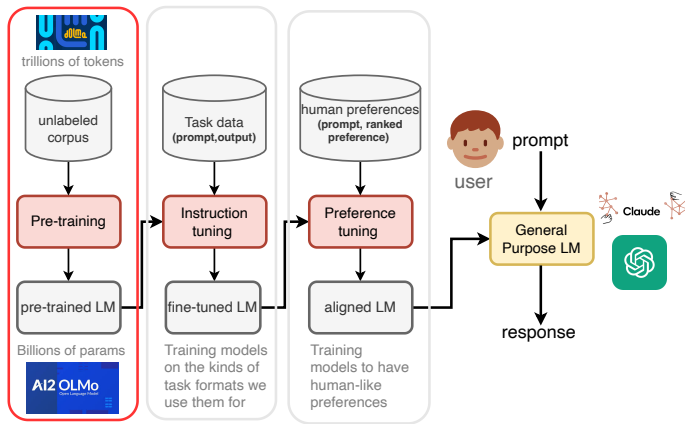
The Seminar on Formal Languages and Neural Networks (FLaNN) is an online event that focuses on the intersection of formal language theory and neural network research. Topics they might discuss include:

## Models have far exceeded expectations

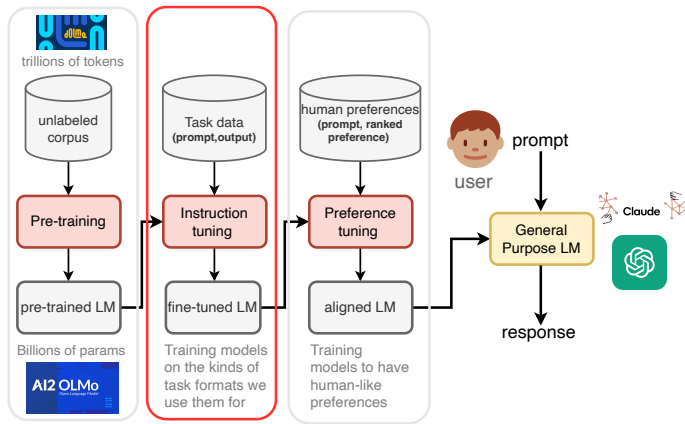
- Neural network models for language processing
- Algorithmic learning theory related to formal languages
- Automata theory and neural computation
- Symbolic and subsymbolic integration in AI systems

- We routinely use LLMs for writing code, integrate within complex systems, use for tasks like scientific discovery.

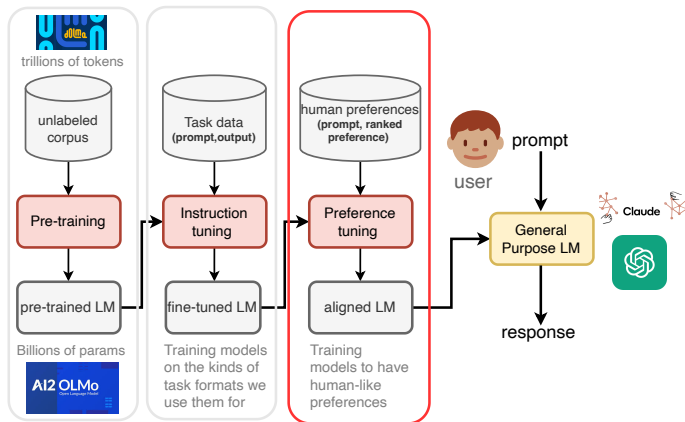
# How do we get to general purpose LLMs? reminder



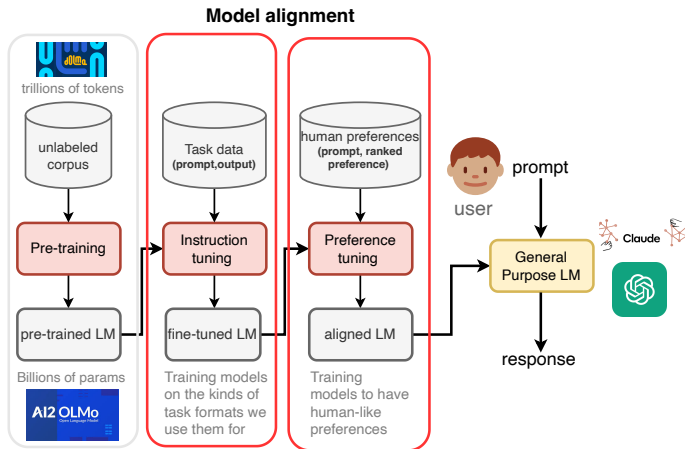
# How do we get to general purpose LLMs? reminder



# How do we get to general purpose LLMs? reminder

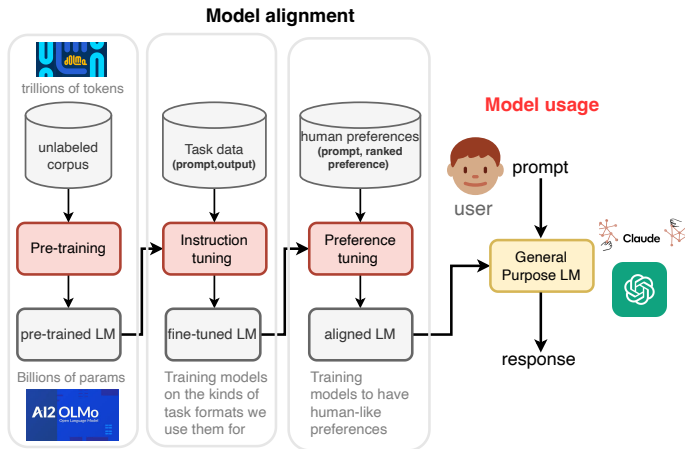


# How do we get to general purpose LLMs? reminder

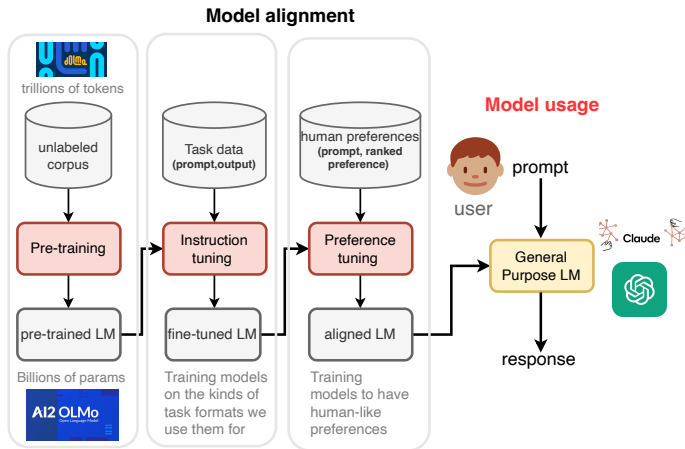




# How do we get to general purpose LLMs? reminder

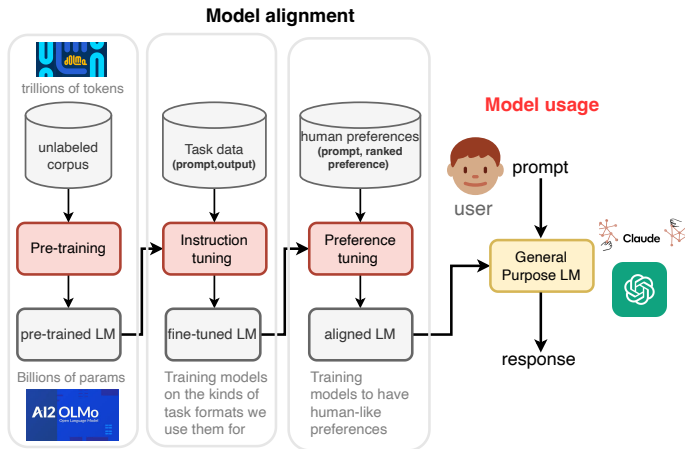


# How do we get to general purpose LLMs? reminder



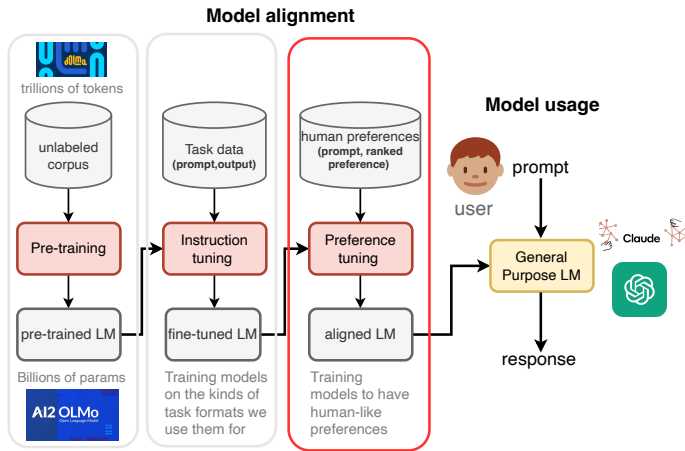
- **Dilemma:** we know vanishingly little about commercial models, models and datasets in general are huge, opaque.

# How do we get to general purpose LLMs? reminder



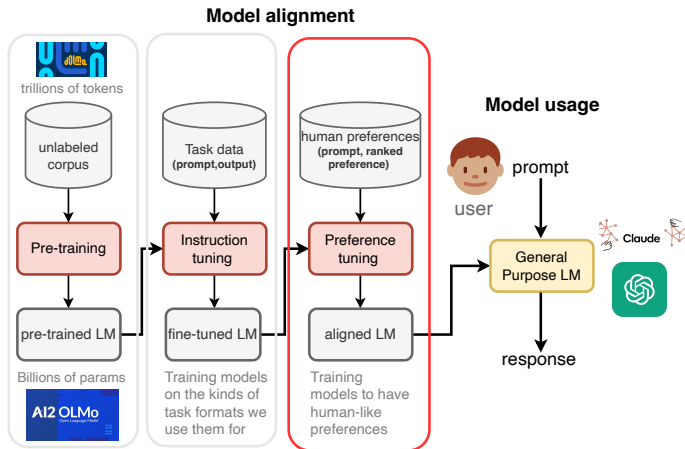
An obvious problem for safety and applications, but also for deciding what research to do, how to innovate.

# Modeling the formal semantics of LLM algorithms



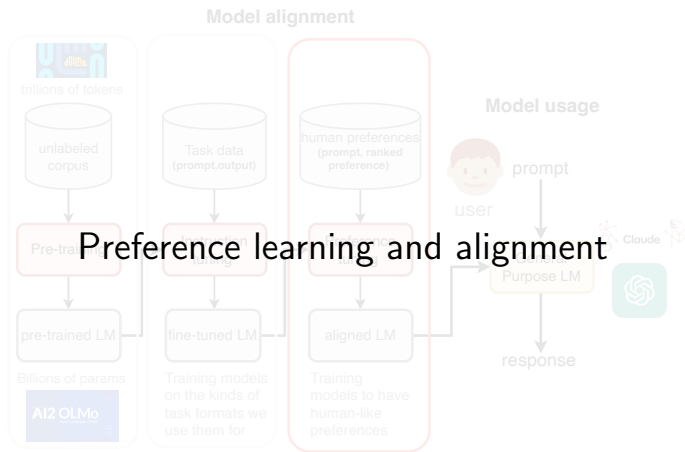
**Today:** can we formally characterize the semantics of preference tuning and alignment? Both for understanding and innovation; **armchair NLP**.

# Modeling the formal semantics of LLM algorithms



**Questions:** What do we do when we tune models to preferences? Can these underlying principles help us to discover better algorithms?

# Modeling the formal semantics of LLM algorithms



**Questions:** What do we do when we tune models to preferences? Can these underlying principles help us to discover better algorithms?

# 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 (DPA) algorithms


- ▶ Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

# Direct preference alignment (DPA) algorithms

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):


$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

  
convex loss function

# Direct preference alignment (DPA) algorithms

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

  
model quantity

# Direct preference alignment (DPA) algorithms

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**Examples:** DPO (Rafailov et al., 2024)

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


# Direct preference alignment (DPA) algorithms

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**Examples:** DPO (Rafailov et al., 2024)

$$\mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ -\log \sigma \left( \beta \cdot \underbrace{\left[ \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right]}_{\text{log ratio difference } \rho_\theta} \right) \right].$$

  
logistic log loss  $f$


# Direct preference alignment (DPA) algorithms

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**Examples:** DPO (Rafailov et al., 2024)

$$\mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ -\log \sigma \left( \beta \cdot \underbrace{\left[ \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right]}_{\text{log ratio difference } \rho_\theta} \right) \right].$$

  
logistic log loss  $f$

As a discrete reasoning problem: reasoning about relationships between our policy model  $\pi_\theta$  and a reference model  $\pi_{\text{ref}}$ .


# Direct preference alignment (DPA) algorithms

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**Examples:** DPO (Rafailov et al., 2024)

$$\mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ -\log \sigma \left( \beta \cdot \underbrace{\left[ \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right]}_{\text{log ratio difference } \rho_\theta} \right) \right].$$

  
logistic log loss  $f$

**Question:** What kind of discrete reasoning problem does  $\rho_\theta$  encode? E.g., if expressed as a symbolic expression.

# Direct preference alignment (DPA) algorithms

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**Examples:** (Azar et al., 2023; Zhao et al., 2022)

|      | $f(\rho_\theta, \beta) =$            | $\rho_\theta$   | properties        |
|------|--------------------------------------|---|-------------------|
| DPO  | $-\log \sigma(\beta \rho_\theta)$    | $\log \frac{\pi_\theta(y_w x)}{\pi_{\text{ref}}(y_w x)} - \log \frac{\pi_\theta(y_l x)}{\pi_{\text{ref}}(y_l x)}$ | logistic log loss |
| IPO  | $(\rho_\theta - \frac{1}{2\beta})^2$ |   | squared loss      |
| SLiC | $\max(0, \beta - \rho_\theta)$       | $\log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)}$  | hinge loss        |



# Direct preference alignment (DPA) algorithms

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**Examples:** (Azar et al., 2023; Zhao et al., 2022)

|      | $f(\rho_\theta, \beta) =$            | $\rho_\theta$   | properties        |
|------|--------------------------------------|---|-------------------|
| DPO  | $-\log \sigma(\beta \rho_\theta)$    | $\log \frac{\pi_\theta(y_w x)}{\pi_{\text{ref}}(y_w x)} - \log \frac{\pi_\theta(y_l x)}{\pi_{\text{ref}}(y_l x)}$ | logistic log loss |
| IPO  | $(\rho_\theta - \frac{1}{2\beta})^2$ |   | squared loss      |
| SlIC | $\max(0, \beta - \rho_\theta)$       | $\log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)}$  | hinge loss        |

Same question: What kind of discrete reasoning problems do SlIC and IPO involve? How are they related?

# Coming up with a new preference loss

- ▶ Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**The procedure:** Select a convex loss function  $f$ , define some model quantity  $\rho_\theta$ , experiment.

# Coming up with a new preference loss

- ▶ Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**The procedure:** Select a convex loss function  $f$ , define some model quantity  $\rho_\theta$ , experiment.

1. (theory) *What is the right  $f$  to use?* Theoretical limitations and properties of Bradley-Terry model or other variants.

# Coming up with a new preference loss

- ▶ Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**The procedure:** Select a convex loss function  $f$ , define some model quantity  $\rho_\theta$ , experiment.

1. (theory) *What is the right  $f$  to use?* Theoretical limitations and properties of Bradley-Terry model or other variants.
2. (empirical) *Can we devise novel algorithmic variants of DPO and  $\rho_\theta$ ?* Find the next best preference algorithm?

# Coming up with a new preference loss

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

**The procedure:** Select a convex loss function  $f$ , define some model quantity  $\rho_\theta$ , experiment.

- (theory) *What is the right  $f$  to use?* Theoretical limitations and properties of Bradley-Terry model or other variants.
- (empirical) *Can we devise novel algorithmic variants of DPO and  $\rho_\theta$ ?* Find the next best preference algorithm?

**our work:** How do we define new  $\rho_\theta$ s and what is the size and structure of this space?

## Coming up with a new preference loss

- Recent direct preference alignment (DPA) approaches assume a closed-form loss function, takes the form (Tang et al., 2024):

$$\ell_{\text{DPA}}(\theta, D) := \mathbb{E}_{(x, y_w, y_l) \sim D_p} \left[ f(\rho_\theta, \beta) \right]$$

The procedure: Select a convex loss function  $f$ , define some model quantity  $\rho_\theta$ , experiment.

- (theory) *What is the right  $f$  to use?* Theoretical limitations and properties of Bradley-Terry model or other variants.
- (empirical) *Can we devise novel algorithmic variants of DPA and  $\rho_\theta$ ?*  
Find the next best preference algorithm?

**our work:** How do we define new  $\rho_\theta$ s and what is the size and structure of this space?

# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

| Approach                 | Loss equation $f(\rho_\theta, \beta)$   |
|--------------------------|---|
| SliC (Zhao et al., 2022) | $\max \left( 0, \beta - \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$ |



# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

| Approach                 | Loss equation $f(\rho_\theta, \beta)$  |
|--------------------------|--|
| SliC (Zhao et al., 2022) | $\max \left( 0, \beta - \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$    |
| CPD (Xu et al., 2024)    | $-\log \sigma \left( \beta \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$ |

# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

| Approach                 | Loss equation $f(\rho_\theta, \beta)$  |
|--------------------------|--|
| SliC (Zhao et al., 2022) | $\max \left( 0, \beta - \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$    |
| CPD (Xu et al., 2024)    | $-\log \sigma \left( \beta \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$ |

# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

| Approach                 | Loss equation $f(\rho_\theta, \beta)$  |
|--------------------------|--|
| SliC (Zhao et al., 2022) | $\max \left( 0, \beta - \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$    |
| CPD (Xu et al., 2024)    | $-\log \sigma \left( \beta \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$ |

Remove optimization details, log ratio  $s_\theta(y_w, y_l) := \log \frac{\pi_\theta(y_w|x)}{\pi_\theta(y_l|x)}$ ,  
regularizer terms, length normalization, **core loss equations**  $\rho_\theta$ .

# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

| Approach                 | Loss equation $f(\rho_\theta, \beta)$  |
|--------------------------|--|
| SliC (Zhao et al., 2022) | $\max \left( 0, \beta - \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$          |
| CPD (Xu et al., 2024)    | $-\log \sigma \left( \beta \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$       |
| RRHF (Yuan et al., 2024) | $\max \left( 0, -\log \frac{\pi_\theta(y_w x)^{ y_w }}{\pi_\theta(y_l x)^{ y_l }} \right)$ |

**note:** This is a messy area, idiosyncratic.

Remove optimization details, log ratio  $s_\theta(y_w, y_l) := \log \frac{\pi_\theta(y_w|x)}{\pi_\theta(y_l|x)}$ ,  
regularizer terms, length normalization, **core loss equations**  $\rho_\theta$ .

# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

| Approach                 | Loss equation $f(\rho_\theta, \beta)$  |
|--------------------------|--|
| SliC (Zhao et al., 2022) | $\max \left( 0, \beta - \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$          |
| CPD (Xu et al., 2024)    | $-\log \sigma \left( \beta \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$       |
| RRHF (Yuan et al., 2024) | $\max \left( 0, -\log \frac{\pi_\theta(y_w x)^{ y_w }}{\pi_\theta(y_l x)^{ y_l }} \right)$ |

Remove optimization details, log ratio  $s_\theta(y_w, y_l) := \log \frac{\pi_\theta(y_w|x)}{\pi_\theta(y_l|x)}$ ,  
regularizer terms, length normalization, **core loss equations**  $\rho_\theta$ .

# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

| Approach                 | Loss equation $f(\rho_\theta, \beta)$  |
|--------------------------|--|
| SliC (Zhao et al., 2022) | $\max \left( 0, \beta - \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$          |
| CPO (Xu et al., 2024)    | $-\log \sigma \left( \beta \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$       |
| RRHF (Yuan et al., 2024) | $\max \left( 0, -\log \frac{\pi_\theta(y_w x)^{ y_w }}{\pi_\theta(y_l x)^{ y_l }} \right)$ |

**observation:** same log ratios keep coming up

Remove optimization details, log ratio  $s_\theta(y_w, y_l) := \log \frac{\pi_\theta(y_w|x)}{\pi_\theta(y_l|x)}$ ,  
regularizer terms, length normalization, **core loss equations**  $\rho_\theta$ .

# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

| Approach                 | Loss equation $f(\rho_\theta, \beta)$  |
|--------------------------|--|
| SliC (Zhao et al., 2022) | $\max \left( 0, \beta - \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$          |
| CPD (Xu et al., 2024)    | $-\log \sigma \left( \beta \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$       |
| PRHE (Yuan et al., 2024) | $\max \left( 0, -\log \frac{\pi_\theta(y_w x)^{ y_w }}{\pi_\theta(y_l x)^{ y_l }} \right)$ |

**question:** What do these log ratios mean semantically?

Remove optimization details, log ratio  $s_\theta(y_w, y_l) := \log \frac{\pi_\theta(y_w|x)}{\pi_\theta(y_l|x)}$ ,  
regularizer terms, length normalization, **core loss equations**  $\rho_\theta$ .

# A look at the structure of DPO loss functions

- ▶ Looking at different preference approaches through time.

| Approach                 | Loss equation $f(\rho_\theta, \beta)$  |
|--------------------------|--|
| SliC (Zhao et al., 2022) | $\max \left( 0, \beta - \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$          |
| CPD (Xu et al., 2024)    | $-\log \sigma \left( \beta \log \frac{\pi_\theta(y_w x)}{\pi_\theta(y_l x)} \right)$       |
| RRHF (Yuan et al., 2024) | $\max \left( 0, -\log \frac{\pi_\theta(y_w x)^{ y_w }}{\pi_\theta(y_l x)^{ y_l }} \right)$ |

Adding a reference model, involves adding an additional term:

$$\text{DPO} \quad -\log \sigma \left( \beta \cdot \left[ \underbrace{\log \frac{\pi_\theta(y_w | x)}{\pi_\theta(y_l | x)}}_{s_\theta(y_w, y_l)} - \underbrace{\log \frac{\pi_{\text{ref}}(y_w | x)}{\pi_{\text{ref}}(y_l | x)}}_{\text{reference policy}} \right] \right)$$



# The recipe for creating new variants of DPO

DPO

$$-\log \sigma \left( \beta \cdot \left[ \underbrace{\log \frac{\pi_{\theta}(y_w | x)}{\pi_{\theta}(y_l | x)}}_{s_{\theta}(y_w, y_l)} - \underbrace{\log \frac{\pi_{\text{ref}}(y_w | x)}{\pi_{\text{ref}}(y_l | x)}}_{s_{\text{ref}}(y_w, y_l)} \right] \right)$$

# The recipe for creating new variants of DPO

$$\text{DPO} \quad -\log \sigma \left( \beta \cdot \left[ \underbrace{\log \frac{\pi_{\theta}(y_w | x)}{\pi_{\theta}(y_l | x)}}_{s_{\theta}(y_w, y_l)} - \underbrace{\log \frac{\pi_{\text{ref}}(y_w | x)}{\pi_{\text{ref}}(y_l | x)}}_{s_{\text{ref}}(y_w, y_l)} \right] \right)$$

- **Type 1:** Add additional terms to loss.

| Variant                   | Core Loss equation $\rho_{\theta}$   |
|---------------------------|--|
| R-DPO (Park et al., 2024) | $s_{\theta}(y_w, y_l) - s_{\text{ref}}(y_w, y_l) + \gamma_{\text{length}}$                                     |
| ODPO (Amini et al., 2024) | $s_{\theta}(y_w, y_l) - s_{\text{ref}}(y_w, y_l) - \gamma_{\text{offset}}$                                     |
| DPOP (Pal et al., 2024)   | $s_{\theta}(y_w, y_l) - s_{\text{ref}}(y_w, y_l) - \log \frac{\pi_{\text{ref}}(x, y_w)}{\pi_{\theta}(x, y_w)}$ |

# The recipe for creating new variants of DPO

$$\text{DPO} \quad -\log \sigma \left( \beta \cdot \left[ \underbrace{\log \frac{\pi_{\theta}(y_w | x)}{\pi_{\theta}(y_l | x)}}_{s_{\theta}(y_w, y_l)} - \underbrace{\log \frac{\pi_{\text{ref}}(y_w | x)}{\pi_{\text{ref}}(y_l | x)}}_{s_{\text{ref}}(y_w, y_l)} \right] \right)$$

- **Type 1:** Add additional terms to loss.

| Variant                   | Core Loss equation $\rho_{\theta}$   |
|---------------------------|--|
| R-DPO (Park et al., 2024) | $s_{\theta}(y_w, y_l) - s_{\text{ref}}(y_w, y_l) + \gamma_{\text{length}}$                                     |
| ODPO (Amini et al., 2024) | $s_{\theta}(y_w, y_l) - s_{\text{ref}}(y_w, y_l) - \gamma_{\text{offset}}$                                     |
| DPOP (Pal et al., 2024)   | $s_{\theta}(y_w, y_l) - s_{\text{ref}}(y_w, y_l) - \log \frac{\pi_{\text{ref}}(x, y_w)}{\pi_{\theta}(x, y_w)}$ |

- **Type 2:** Change or re-parameterize reference ratio.

| Variant                   | Core Loss equation $\rho_{\theta}$  |
|---------------------------|---|
| ORPO (Hong et al., 2024)  | $s_{\theta}(y_w, y_l) - \log \frac{1 - \pi_{\theta}(y_w   x)}{1 - \pi_{\theta}(y_l   x)}$ |
| SimPO (Meng et al., 2024) | $s_{\theta}(y_w, y_l) - \gamma_{\text{penalty}}$  |

# The recipe for creating new variants of DPO

$$\text{DPO} \quad -\log \sigma \left( \beta \cdot \left[ \underbrace{\log \frac{\pi_{\theta}(y_w | x)}{\pi_{\theta}(y_l | x)}}_{s_{\theta}(y_w, y_l)} - \underbrace{\log \frac{\pi_{\text{ref}}(y_w | x)}{\pi_{\text{ref}}(y_l | x)}}_{s_{\text{ref}}(y_w, y_l)} \right] \right)$$

- **Type 1:** Add additional terms to loss.

**question:** What happens when we add new terms?

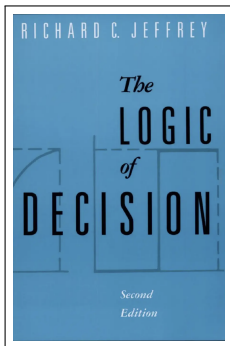
| Variant                   | Core Loss equation $\rho_{\theta}$   |
|---------------------------|--|
| ODPO (Amini et al., 2024) | $s_{\theta}(y_w, y_l) - s_{\text{ref}}(y_w, y_l) - \gamma_{\text{offset}}$                                     |
| DPOP (Pal et al., 2024)   | $s_{\theta}(y_w, y_l) - s_{\text{ref}}(y_w, y_l) - \log \frac{\pi_{\text{ref}}(x, y_w)}{\pi_{\theta}(x, y_w)}$ |

- **Type 2:** Change or re-parameterize reference ratio.

| Variant                   | Core Loss equation $\rho_{\theta}$  |
|---------------------------|---|
| ORPO (Hong et al., 2024)  | $s_{\theta}(y_w, y_l) - \log \frac{1 - \pi_{\theta}(y_w   x)}{1 - \pi_{\theta}(y_l   x)}$ |
| SimPO (Meng et al., 2024) | $s_{\theta}(y_w, y_l) - \gamma_{\text{penalty}}$  |

# Haven't these semantic questions been looked at before?

**Analytic philosophy:** Much work on the semantics of pairwise preference, rich languages for expressing ideas.



(Jeffrey, 1965)

| THE STATUS OF VARIOUS PREFERENCE PRINCIPLES  |            |                  |        |                  |                  |            |
|--|------------|------------------|--------|------------------|------------------|------------|
| Preference Principle   | Von Wright | Chisholm<br>Sosa | Martin | $P^\#$           | $P^\star$        | $P^\omega$ |
| 1. $pPq \rightarrow \sim(qPp)$   | ✓          | ✓                | ✓      | +                | +                | +          |
| 2. $(pPq \ \& \ qPr) \rightarrow pPr$  | ✓          | ✓                | ✓      | +                | +                | +          |
| 3. $pPq \rightarrow \sim qP\sim p$   |            | x                | ✓      | (+) <sup>1</sup> | +                | +          |
| 4. $\sim qP\sim p \rightarrow pPq$   |            | x                | ✓      | (+) <sup>1</sup> | +                | +          |
| 5. $pPq \rightarrow (p \ \& \ \sim q) P(\sim p \ \& \ q)$  | ✓          | x                |        | +                | +                | +          |
| 6. $(p \ \& \ \sim q) P(\sim p \ \& \ q) \rightarrow pPq$  | ✓          | x                |        | +                | +                | +          |
| 7. $[\sim(pP\sim p) \ \& \ \sim(\sim pPp)] \ \& \ \sim(qP\sim q) \ \& \ \sim(\sim qPq) \rightarrow [\sim(pPq) \ \& \ \sim(qPp)]$ | ✓          | ✓                |        | +                | +                | +          |
| 8. $[\sim(qP\sim q) \ \& \ \sim(\sim qPq) \ \& \ pPq] \rightarrow pP\sim p$  |            | ✓                |        | +                | +                | —          |
| 9. $[\sim(qP\sim q) \ \& \ \sim(\sim qPq) \ \& \ qP\sim p] \rightarrow pP\sim p$   |            | ✓                |        | +                | +                | —          |
| 10. $pPq \rightarrow [(p \ \& \ r) P(q \ \& \ r) \ \& \ (p \ \& \ \sim r) P(q \ \& \ \sim r)]$                                   | ✓          |                  |        | —                | —                | +          |
| 11. $[(p \ \& \ r) P(q \ \& \ r) \ \& \ (p \ \& \ \sim r) P(q \ \& \ \sim r)] \rightarrow pPq$                                   | ✓          |                  |        | (+) <sup>2</sup> | (+) <sup>3</sup> | +          |
| 12. $[\sim(pPq) \ \& \ \sim(qPr)] \rightarrow \sim(pPr)$   |            | ✓                |        | +                | +                | —          |
| 13. $(pPr \vee qPr) \rightarrow (p \vee q) Pr$   |            |                  | ✓      | —                | —                | —          |
| 14. $(p \vee q) Pr \rightarrow [pPr \ \& \ qPr]$   | ✓          |                  |        | —                | —                | —          |
| 15. $[pPr \ \& \ qPr] \rightarrow (p \vee q) Pr$   | ✓          |                  |        | —                | —                | —          |
| 16. $(p \vee q) Pr \rightarrow (pPr \vee qPr)$   |            |                  |        | —                | —                | —          |
| 17. $pP(q \vee r) \rightarrow (pPq \ \vee \ pPr)$  |            |                  | ✓      | —                | —                | —          |
| 18. $(pPq \ \& \ pPr) \rightarrow pP(q \vee r)$  |            |                  |        | —                | —                | —          |
| 19. $(pPr \ \& \ qPr) \rightarrow (p \ \& \ q) Pr$   |            |                  |        | —                | —                | —          |

Semantic foundations for the logic of preference Rescher (1967)

# The language of machine learning

## Loss functions

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

# The language of machine learning

## Loss functions

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

- Frustration: the language of machine learning is not very rich, hard to express complex ideas, come up with improved algorithms, barrier.

# The language of machine learning

## Loss functions

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



Specification or theory of preference?

- Frustration: the language of machine learning is not very rich, hard to express complex ideas, come up with improved algorithms, barrier.



# The language of machine learning

## Loss functions

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



Specification or theory of preference?

- **Frustration:** the language of machine learning is not very rich, hard to express complex ideas, come up with improved algorithms, barrier.

**Broader goal:** High-level modeling languages for specifying and better understanding LLMs and their algorithms.

# The language of machine learning

## Loss functions

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

## Formalization of preference losses

- Frustration: the language of machine learning is not very rich, hard to express complex ideas, come up with improved algorithms, barrier.

**Broader goal:** High-level modeling languages for specifying and better understanding LLMs and their algorithms.

# Preference learning as a discrete reasoning problem

## Loss Function

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

# Preference learning as a discrete reasoning problem

## Loss Function

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

Two models, four predictions

# Preference learning as a discrete reasoning problem

## Loss Function

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

Two models, four predictions

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

# Preference learning as a discrete reasoning problem

## 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?

# Preference learning as a discrete reasoning problem

## 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?

# Preference learning as a discrete reasoning problem

## 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?



# Preference learning as a discrete reasoning problem

## 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.

# Preference learning as a discrete reasoning problem

## 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.

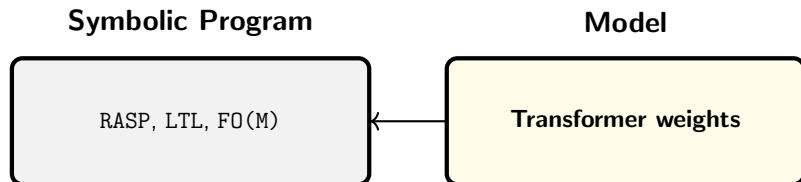
# Distilling LLMs to symbolic programs in general

**Model**

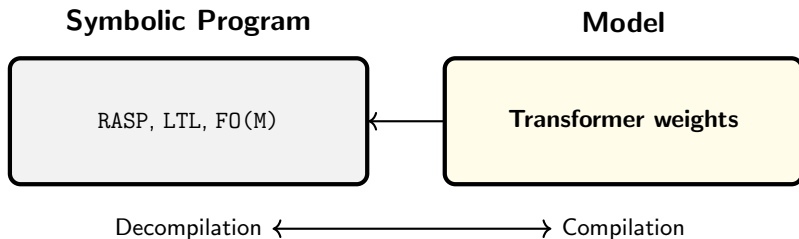


**Transformer weights**

# Distilling LLMs to symbolic programs in general



# Distilling LLMs to symbolic programs in general



- We know what the *target languages* are (Weiss et al., 2021; Merrill and Sabharwal, 2023; Yang and Chiang, 2024), how to compile, decompile (Friedman et al., 2023).

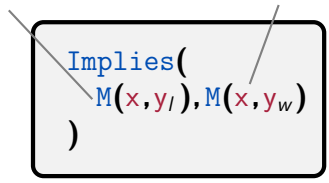
What do these programs tell us?

```
Implies(  
  M(x, yl), M(x, yw)  
)
```

# What do these programs tell us?

Model predicts loser

Model predicts winner



A diagram showing a logical formula `Implies(M(x, y_l), M(x, y_w))` enclosed in a rounded rectangle. Two lines point from external text to parts of the formula: one from 'Model predicts loser' to `M(x, y_l)`, and another from 'Model predicts winner' to `M(x, y_w)`. The word `Implies` is blue, while the `M` and `x` in both arguments are blue, and `y_l` and `y_w` are red.

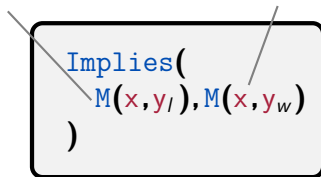
```
Implies(  
  M(x, y_l), M(x, y_w)  
)
```

**Conceptually:** Model predications are logical propositions, Boolean variables inside of formulas, weighted by prediction probability.

# What do these programs tell us?

Model predicts loser

Model predicts winner



$$w(\mathbf{M}(\mathbf{x}, \mathbf{y})) = \pi_{\mathbf{M}}(\mathbf{y} \mid \mathbf{x})$$

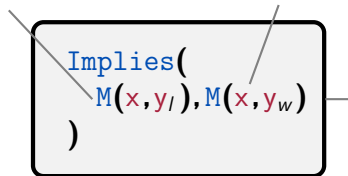
**Conceptually:** Model predications are logical propositions, Boolean variables inside of formulas, weighted by prediction probability.



# What do these programs tell us?

Model predicts loser

Model predicts winner



A diagram consisting of a light gray rounded rectangle with a black border. Inside the rectangle, the text `Implies(` is in blue, followed by `M(x, yl), M(x, yw)` in blue with red subscripts, and a closing parenthesis `)` in blue. A line from the text 'Model predicts loser' points to the `M(x, yl)` part. Another line from the text 'Model predicts winner' points to the `M(x, yw)` part. A horizontal line extends from the right side of the rectangle towards the explanatory text.

```
Implies(  
  M(x, yl), M(x, yw)  
)
```

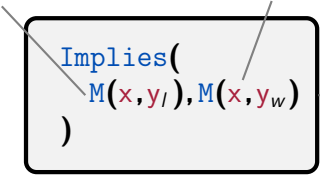
*Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.*

**Conceptually:** Predictions are connected through Boolean operators, express constraints on predictions;  $\rho_\theta$  as formulas.

# What do these programs tell us?

Model predicts loser

Model predicts winner



```
Implies(  
  M(x, yl), M(x, yw)  
)
```

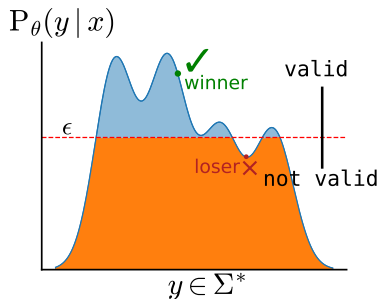
*Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.*

**Assumption:** Every loss function has an internal logic that can be expressed in this way, we want to uncover that logic.

# What do these programs tell us?

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

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

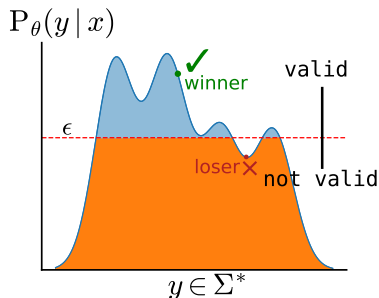


**Assumption:** Every loss function has an internal logic that can be expressed in this way, we want to uncover that logic.

# What do these programs tell us?

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

```
And(  
  M(x, y_w),  
  Not(M(x, y_l)))
```

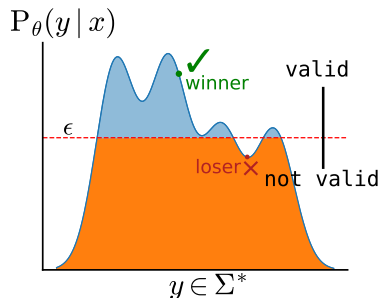


**Assumption:** Every loss function has an internal logic that can be expressed in this way, we want to uncover that logic.

# What do these programs tell us?

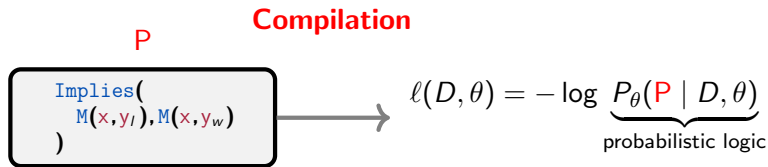
```
Implies(  
  M(x, yl), M(x, yw)  
)
```

```
And(  
  M(x, yw),  
  Not(M(x, yl)))
```



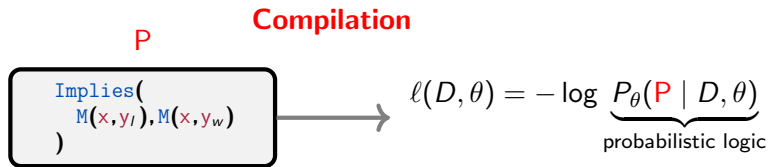
**Observation:** The second program is more strict than the first, involves semantic entailment.

# What do these programs tell us?



*Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.*

# What do these programs tell us?

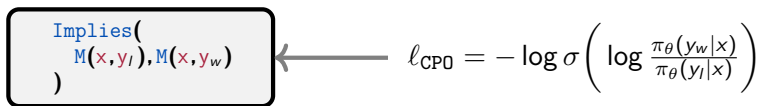


*Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.*

**What we did:** defined a novel probabilistic logic for preference modeling,  
**note:** logic useful not only for learning and loss.

# What do these programs tell us?

P **Decompilation**



*Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.*



# What do these programs tell us?

**P**      **Decompilation**

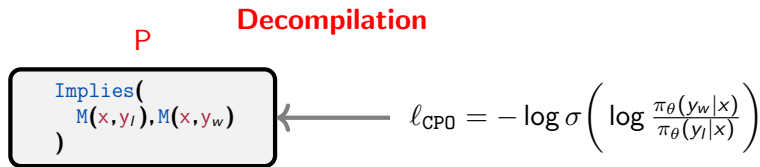
```
Implies(  
  M(x, yl), M(x, yw)  
)
```

$$\leftarrow \ell_{\text{CPO}} = -\log \sigma \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\theta}(y_l|x)} \right)$$

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

$$\underbrace{\ell_{\text{CPO}}(D, \theta) = -\log P_{\theta}(\mathbf{P} \mid D, \theta)}_{\text{correctness property}}$$

# What do these programs tell us?



Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

$$\underbrace{\ell_{\text{CPD}}(D, \theta) = -\log P_{\theta}(\mathbf{P} \mid D, \theta)}_{\text{correctness property}}$$

**The second thing we did:** Defined a mechanical procedure for decompilation, proved its correctness, invariance to choice of  $f$ .

# Illustration of approach and results

**Input Loss**  $\ell_{\text{ORPO}}$

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

# Illustration of approach and results

**Input Loss**  $\ell_{\text{ORPO}}$

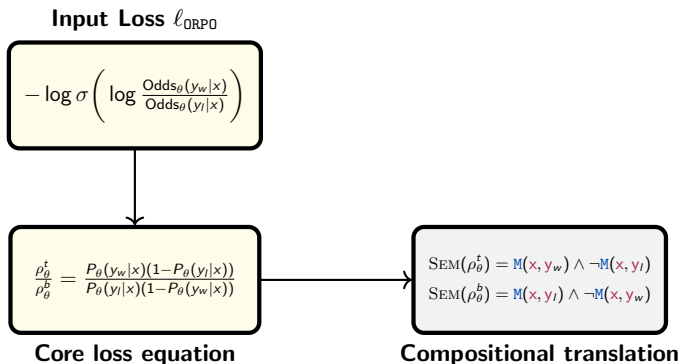
$$-\log \sigma \left( \log \frac{\text{Odds}_{\theta}(y_w|x)}{\text{Odds}_{\theta}(y_l|x)} \right)$$



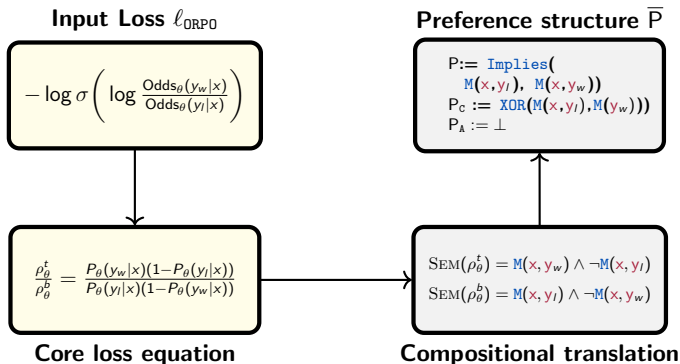
$$\frac{\rho_{\theta}^t}{\rho_{\theta}^b} = \frac{P_{\theta}(y_w|x)(1-P_{\theta}(y_l|x))}{P_{\theta}(y_l|x)(1-P_{\theta}(y_w|x))}$$

**Core loss equation**

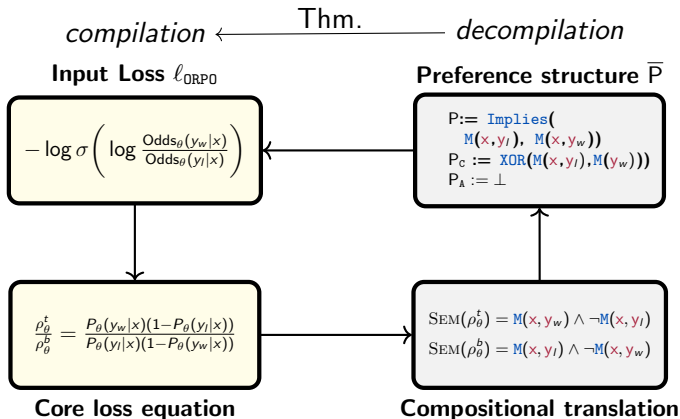
# Illustration of approach and results



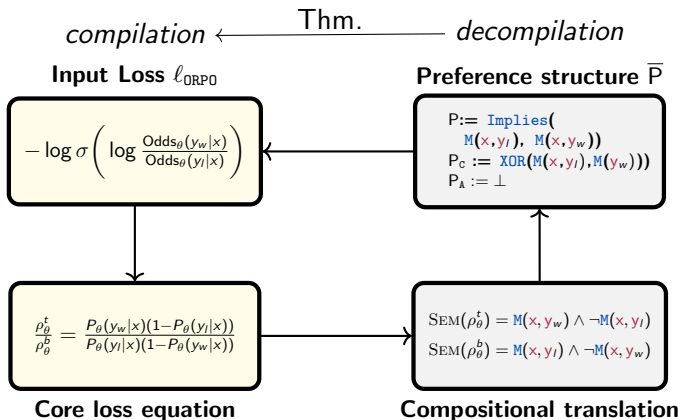
# Illustration of approach and results



# Illustration of approach and results



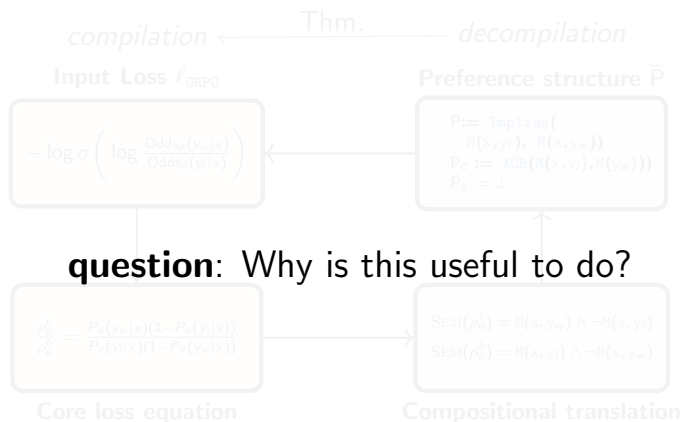
# Illustration of approach and results



- **Preference structure**, a core construct in our logic, encoding for preference losses, has a natural Boolean interpretation.

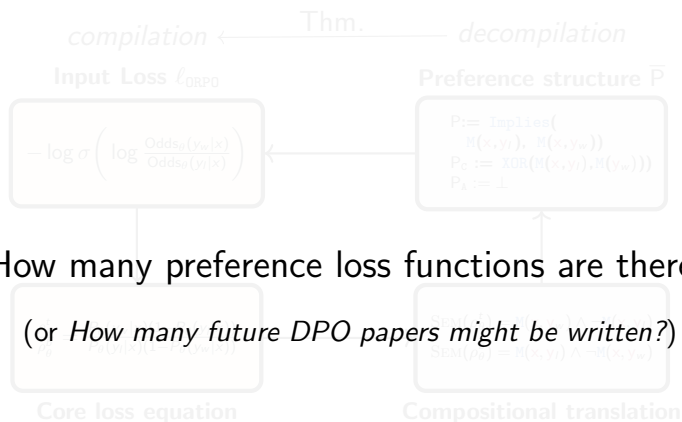


# Illustration of approach and results



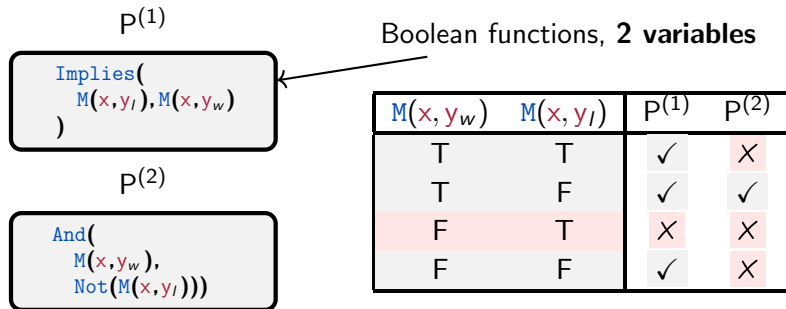
- **Preference structure**, a core construct in our logic, encoding for preference losses, has a natural Boolean interpretation.

# Illustration of approach and results

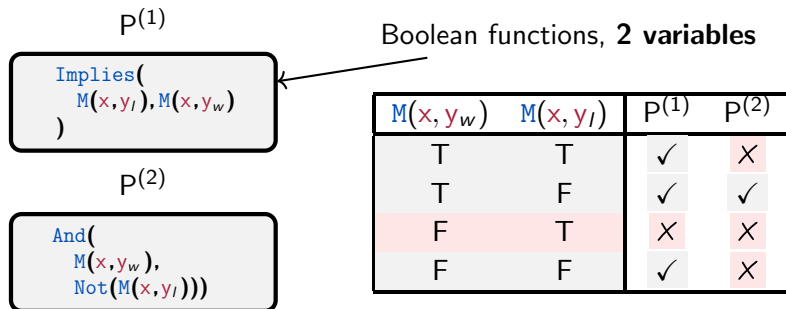


- **Preference structure**, a core construct in our logic, encoding for preference losses, has a natural Boolean interpretation.

## Why is this useful? understanding the space

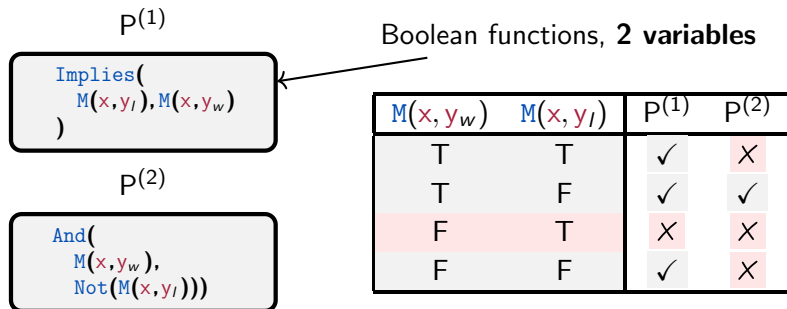


## Why is this useful? understanding the space



- ▶ Every program (in our logic) is pair of Boolean functions (in  $n$  variables),  
corr. to ✓ and X, leads to  $4^{2^n}$  possible loss functions.

## Why is this useful? understanding the space



**no reference:** 256 losses

- ▶ Every program (in our logic) is pair of Boolean functions (in  $n$  variables), corr. to ✓ and X, leads to  $4^{2^n}$  possible loss functions.

# Loss functions as truth tables

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

**4 variables**

| Ref(x,y <sub>w</sub> ) | M(x,y <sub>I</sub> ) | Ref(x,y <sub>I</sub> ) | M(x,y <sub>w</sub> ) |
|------------------------|----------------------|------------------------|----------------------|
| F                      | F                    | F                      | F                    |
| F                      | F                    | F                      | T                    |
| F                      | F                    | T                      | F                    |
| F                      | F                    | T                      | T                    |
| F                      | T                    | F                      | F                    |
| F                      | T                    | F                      | T                    |
| F                      | T                    | T                      | F                    |
| F                      | T                    | T                      | T                    |
| T                      | F                    | F                      | F                    |
| T                      | F                    | F                      | T                    |
| T                      | F                    | T                      | F                    |
| T                      | F                    | T                      | T                    |
| T                      | T                    | F                      | F                    |
| T                      | T                    | F                      | T                    |
| T                      | T                    | T                      | F                    |
| T                      | T                    | T                      | T                    |

**w/ reference:** 4,294,967,296 losses

## Loss functions as truth tables

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

4 variables

**answer:** loads.

| Ref(x,y <sub>w</sub> ) | M(x,y <sub>I</sub> ) | Ref(x,y <sub>I</sub> ) | M(x,y <sub>w</sub> ) |
|------------------------|----------------------|------------------------|----------------------|
| F                      | F                    | F                      | F                    |
| F                      | F                    | F                      | T                    |
| F                      | F                    | T                      | F                    |
| F                      | F                    | T                      | T                    |
| F                      | T                    | F                      | F                    |
| F                      | T                    | F                      | T                    |
| F                      | T                    | T                      | F                    |
| F                      | T                    | T                      | T                    |
| T                      | F                    | F                      | F                    |
| T                      | F                    | F                      | T                    |
| T                      | F                    | T                      | F                    |
| T                      | F                    | T                      | T                    |
| T                      | T                    | F                      | F                    |
| T                      | T                    | F                      | T                    |
| T                      | T                    | T                      | F                    |
| T                      | T                    | T                      | T                    |

w/ reference: 4,294,967,296 losses

# Loss functions as truth tables

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

**question:** How are losses related to one another?

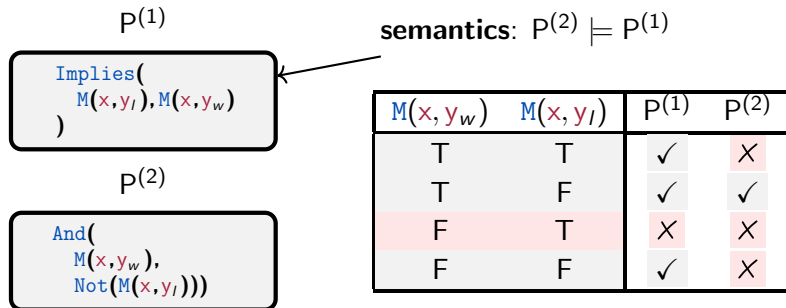
4 variables

| Ref(x,y <sub>w</sub> ) | M(x,y <sub>I</sub> ) | Ref(x,y <sub>I</sub> ) | M(x,y <sub>w</sub> ) |
|------------------------|----------------------|------------------------|----------------------|
| F                      | F                    | F                      | F                    |
| F                      | F                    | F                      | T                    |
| F                      | F                    | T                      | F                    |
| F                      | F                    | T                      | T                    |
| F                      | T                    | F                      | F                    |
| F                      | T                    | F                      | T                    |
| F                      | T                    | T                      | F                    |
| F                      | T                    | T                      | T                    |
| T                      | F                    | F                      | F                    |
| T                      | F                    | F                      | T                    |
| T                      | F                    | T                      | F                    |
| T                      | F                    | T                      | T                    |
| T                      | T                    | F                      | F                    |
| T                      | T                    | F                      | T                    |
| T                      | T                    | T                      | F                    |
| T                      | T                    | T                      | T                    |

w/ reference: 4,294,967,296 losses

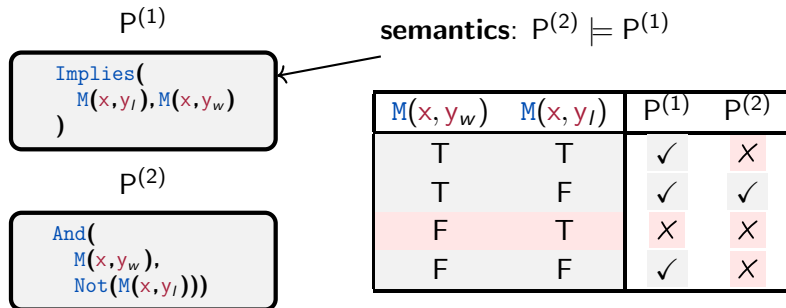


## Why is this useful? understanding the structure



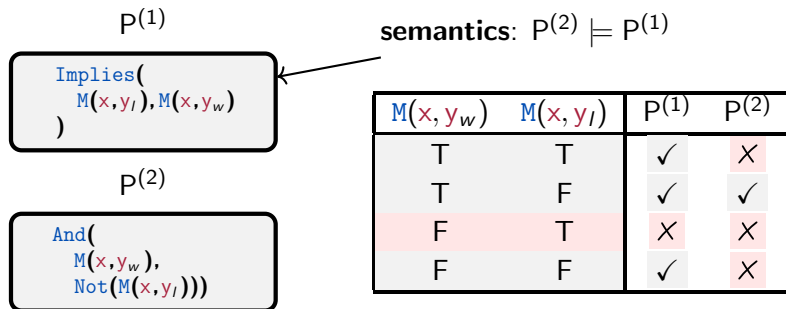
**Proposition** (Xu et al., 2018): Loss behavior is monotonic w.r.t semantic entailment: if  $P^{(2)} \models P^{(1)}$  then  $\ell(D, \theta, P^{(2)}) \geq \ell(D, \theta, P^{(1)})$ .

# Why is this useful? understanding the structure



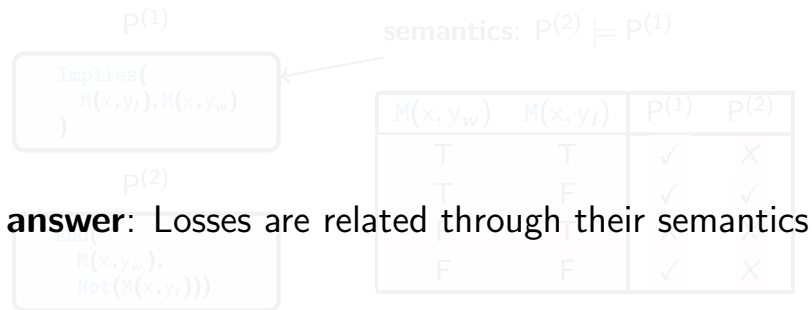
**Proposition** (Xu et al., 2018): Loss is equivalent under semantic equivalence: If  $P^{(2)} \equiv P^{(1)}$  then  $\ell(D, \theta, P^{(2)}) = \ell(D, \theta, P^{(1)})$ .

# Why is this useful? understanding the structure



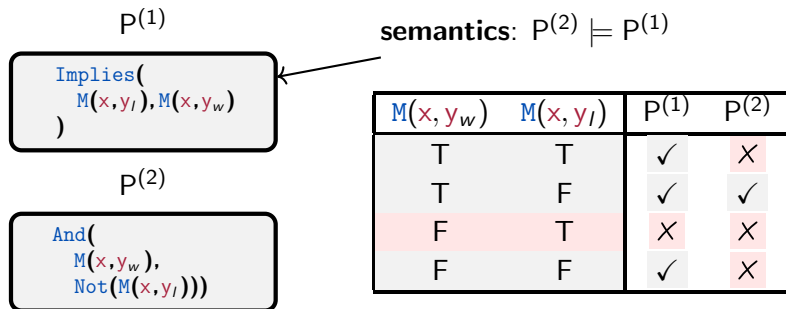
**Theorem:**  $\ell(D, \theta, P^{(2)}) > \ell(D, \theta, P^{(1)})$  (the loss of  $P^{(1)}$  is contained in the loss of  $P^{(2)}$ ).

Why is this useful? understanding the structure



**Theorem:**  $\ell(D, \theta, P^{(2)}) > \ell(D, \theta, P^{(1)})$  (the loss of  $P^{(1)}$  is contained in the loss of  $P^{(2)}$ ).

# Why is this useful? understanding the structure



**Practical strategy:** Start with empirically successful losses, modify semantics (make more or less constrained), then experiment accordingly.

# Why is this useful? understanding the structure

$P(1)$

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

semantics:  $P(2) \models P(1)$

$P(2)$

```
And(  
  M(x, yw),  
  Not(M(x, yI)))
```

| $M(x, y_w)$ | $M(x, y_I)$ | $P(1)$ | $P(2)$ |
|-------------|-------------|--------|--------|
| T           | T           | ✓      | X      |
| T           | F           | ✓      | ✓      |
| F           | T           | X      | X      |
| F           | F           | ✓      | X      |

**questions:** How does our logic work? What do we see?

**Practical strategy:** Start with empirically successful losses, modify semantics (make more or less constrained), then experiment accordingly.

## How does the logic work? **compilation**

| $M(x, y_w)$ | $M(x, y_l)$ | CPO | ORPO | unCPO |
|-------------|-------------|-----|------|-------|
| T           | T           | ✓ X |      | ✓     |
| T           | F           | ✓   | ✓    | ✓     |
| F           | T           | X   | X    | X     |
| F           | F           |     |      | ✓     |

P

```
Implies(  
   $M(x, y_l)$ ,  $M(x, y_w)$   
)
```

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

## How does the logic work? **compilation**

| $M(x, y_w)$ | $M(x, y_l)$ | CPO | ORPO | unCPO |
|-------------|-------------|-----|------|-------|
| T           | T           | ✓ X |      | ✓     |
| T           | F           | ✓   | ✓    | ✓     |
| F           | T           | X   | X    | X     |
| F           | F           |     |      | ✓     |

$$S(i) := \prod_{i \models v} w(v) \cdot \prod_{i \models \neg v} 1 - w(v)$$

P

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

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.



# How does the logic work? **compilation**

|             |             | P   |      |       |
|-------------|-------------|-----|------|-------|
| $M(x, y_w)$ | $M(x, y_l)$ | CPO | ORPO | unCPO |
| T           | T           | ✓ X |      | ✓     |
| T           | F           | ✓   | ✓    | ✓     |
| F           | T           | X   | X    | X     |
| F           | F           |     |      | ✓     |

Implies(  
 $M(x, y_l), M(x, y_w)$   
 )

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

$$S(i) := \prod_{i \models v} w(v) \cdot \prod_{i \models \neg v} 1 - w(v)$$

- Formula probability computed as a weighted count  $\sum \checkmark$  (Chavira and Darwiche, 2008), loss is  $-\log$  (Xu et al., 2018); generalizing:

$$-\log P_{\theta}(P_x) := -\log \sigma \left( \underbrace{\log \frac{\sum \checkmark}{\sum X}}_{\text{symmetric to DPA, } \rho_{\theta}} \right)$$

# How does the logic work? compilation

| $M(x, y_w)$ | $M(x, y_l)$ | CPO | ORPO | unCPO |
|-------------|-------------|-----|------|-------|
| T           | T           | ✓ X |      | ✓     |
| T           | F           | ✓   | ✓    | ✓     |
| F           | T           | X   | X    | X     |
| F           | F           |     |      | ✓     |

P

Implies(  
 $M(x, y_l), M(x, y_w)$   
 )

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

$$S(i) := \prod_{i \models v} w(v) \cdot \prod_{i \models \neg v} 1 - w(v)$$

- Formula probability computed as a weighted count  $\sum \checkmark$  (Chavira and Darwiche, 2008), loss is  $-\log$  (Xu et al., 2018); generalizing:

$$\begin{aligned}
 -\log P_\theta(P_{\text{ORPO}}) &:= -\log \sigma \left( \log \frac{\sum \checkmark}{\sum X} \right) \\
 &= -\log \sigma \left( \underbrace{\log \frac{\pi_\theta(y_w | x)(1 - \pi_\theta(y_l | x))}{\pi_\theta(y_l | x)(1 - \pi_\theta(y_w | x))}}_{\ell_{\text{ORPO}}, P_\theta(P|\text{one hot})} \right)
 \end{aligned}$$

# How does the logic work? compilation

| $M(x, y_w)$ | $M(x, y_l)$ | CP0 | ORP0 | unCP0 |
|-------------|-------------|-----|------|-------|
| T           | T           | ✓ X |      | ✓     |
| T           | F           | ✓   | ✓    | ✓     |
| F           | T           | X   | X    | X     |
| F           | F           |     |      | ✓     |

P

$\text{Implies}($   
 $\quad M(x, y_l), M(x, y_w)$   
 $)$

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

$$S(i) := \prod_{i \models v} w(v) \cdot \prod_{i \models \neg v} 1 - w(v)$$

- Formula probability computed as a weighted count  $\sum \checkmark$  (Chavira and Darwiche, 2008), loss is  $-\log$  (Xu et al., 2018); generalizing:

$$\begin{aligned}
 -\log P_{\theta}(P_{\text{CP0}}) &:= -\log \sigma \left( \log \frac{\sum \checkmark}{\sum X} \right) \\
 &= \underbrace{-\log \sigma \left( \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\theta}(y_l | x)} \right)}_{\ell_{\text{CP0}}, \sim P_{\theta}(P | \text{one true})}
 \end{aligned}$$

## How does the logic work? compilation

| $M(x, y_w)$ | $M(x, y_l)$ | CP0 | ORPO | unCP0 |
|-------------|-------------|-----|------|-------|
| T           | T           | ✓ X |      | ✓     |
| T           | F           | ✓   | ✓    | ✓     |
| F           | T           | X   | X    | X     |
| F           | F           |     |      | ✓     |

$P$

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

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

**observation:** losses differ in conditioning constraints

- ▶ Formula probability computed as a weighted count  $\sum \checkmark$  (Chavira and Darwiche, 2008), loss is  $-\log$  (Xu et al., 2018); generalizing:

$$\begin{aligned} -\log P_{\theta}(P_{\text{CP0}}) &:= -\log \sigma \left( \log \frac{\sum \checkmark}{\sum X} \right) \\ &= \underbrace{-\log \sigma \left( \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\theta}(y_l | x)} \right)}_{\ell_{\text{CP0}}, \sim P_{\theta}(P | \text{one true})} \end{aligned}$$

# How does the logic work? compilation

| $M(x, y_w)$ | $M(x, y_l)$ | CPO | ORPO | unCPO |
|-------------|-------------|-----|------|-------|
| T           | T           | ✓ X |      | ✓     |
| T           | F           | ✓   | ✓    | ✓     |
| F           | T           | X   | X    | X     |
| F           | F           |     |      | ✓     |

P

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

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

$$S(i) := \prod_{i \models v} w(v) \cdot \prod_{i \models \neg v} 1 - w(v)$$

| Loss  | Representation $\bar{P}$  |
|-------|---|
| CE    | $P := M(x, y_w), P_C := \perp$  |
| CEUnl | $P := \text{And}(M(x, y_w), \text{Not}(M(x, y_l)))$<br>$P_C := \perp$   |
| CPO   | $;;$ core semantic formula<br>$P := \text{Implies}(M(x, y_l), M(x, y_w))$<br>$;;$ one-true constraint<br>$P_C := \text{Or}(M(x, y_l), M(x, y_w))$ |
| ORPO  | $P := \text{Implies}(M(x, y_l), M(x, y_w))$<br>$;;$ one-hot constraint<br>$P_C := \text{XOR}(M(x, y_l), M(x, y_w))$                               |

# How does the logic work? compilation

| $M(x, y_w)$ | $M(x, y_l)$ | CPO | ORPO | unCPO |
|-------------|-------------|-----|------|-------|
| T           | T           | ✓ X |      | ✓     |
| T           | F           | ✓   | ✓    | ✓     |
| F           | T           | X   | X    | X     |
| F           | F           |     |      | ✓     |

P

Implies(  
 $M(x, y_l), M(x, y_w)$   
 )

*Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.*

$$S(i) := \prod_{i \models v} w(v) \cdot \prod_{i \models \neg v} 1 - w(v)$$

- Formula probability computed as a weighted count  $\sum \checkmark$  (Chavira and Darwiche, 2008), loss is  $-\log$  (Xu et al., 2018); generalizing:

$$\begin{aligned}
 -\log P_{\theta}(P_{\text{unCPO}}) &:= -\log \sigma \left( \log \frac{\sum \checkmark}{\sum X} \right) \\
 &= \underbrace{-\log \sigma \left( \log \frac{\pi_{\theta}(y_l | x) \pi_{\theta}(y_w | x) + (1 - \pi_{\theta}(y_l | x))}{\pi_{\theta}(y_l | x) (1 - \pi_{\theta}(y_w | x))} \right)}_{\text{novel loss}}
 \end{aligned}$$

# How does the logic work? compilation

| $M(x, y_w)$ | $M(x, y_l)$ | CPO | ORPO | unCPO | <div> <math>P</math><br/> <math>\text{Implies}(</math><br/> <math>M(x, y_l), M(x, y_w)</math><br/> <math>)</math> </div> |
|-------------|-------------|-----|------|-------|--|
| T           | T           | ✓ X |      | ✓     |  |
| T           | F           | ✓   | ✓    | ✓     |  |
| F           | T           | X   | X    | X     |  |
| F           | F           |     |      | ✓     |  |

Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too

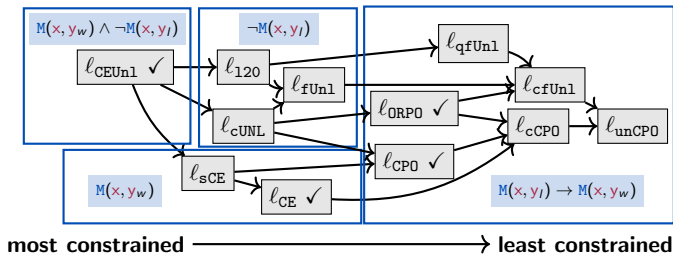
$$S(i) := \prod w(y) \cdot \prod 1 - w(y)$$

**note:**  $M(x, y_l) \rightarrow M(x, y_w) \equiv \neg M(x, y_l) \vee M(x, y_w)$

- Formula probability computed as a weighted count  $\sum \checkmark$  (Chavira and Darwiche, 2008), loss is  $-\log$  (Xu et al., 2018); generalizing:

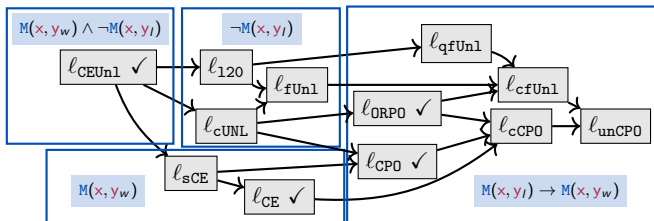
$$\begin{aligned}
 -\log P_{\theta}(P_{\text{unCPO}}) &:= -\log \sigma \left( \log \frac{\sum \checkmark}{\sum X} \right) \\
 &= -\log \sigma \left( \log \underbrace{\frac{\pi_{\theta}(y_l | x) \pi_{\theta}(y_w | x) + (1 - \pi_{\theta}(y_l | x))}{\pi_{\theta}(y_l | x) (1 - \pi_{\theta}(y_w | x))}}_{\text{novel loss}} \right)
 \end{aligned}$$

# The no reference loss landscape

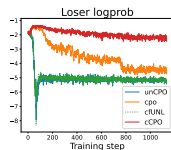
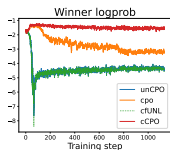




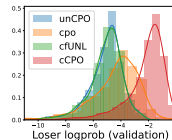
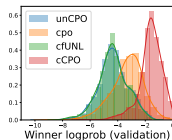
# The no reference loss landscape



most constrained  $\longrightarrow$  least constrained

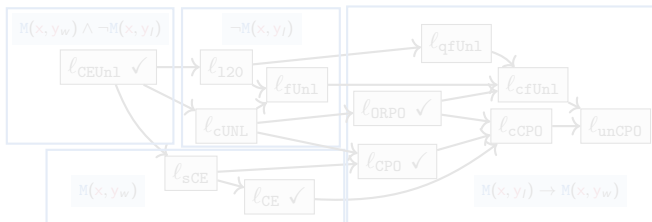


Training dynamics



Inference

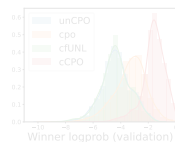
# The no reference loss landscape



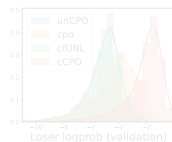
**question:** Are any of these losses good?



Training dynamics



Inference



## Adding a reference model

```
P:= Implies(  
  And(M(x,yl), Ref(x,yw)),  
  And(M(x,yw), Ref(x,yl))  
)
```

Whenever the model being tuned deems the loser to be a valid generation and the reference model deems the winner to be valid, the tuned model should deem the winner to be valid too, and *the reference should deem the loser to be valid.*

# Adding a reference model

```
P:= Implies(  
  And(M(x, yl), Ref(x, yw)),  
  And(M(x, yw), Ref(x, yl))  
)
```

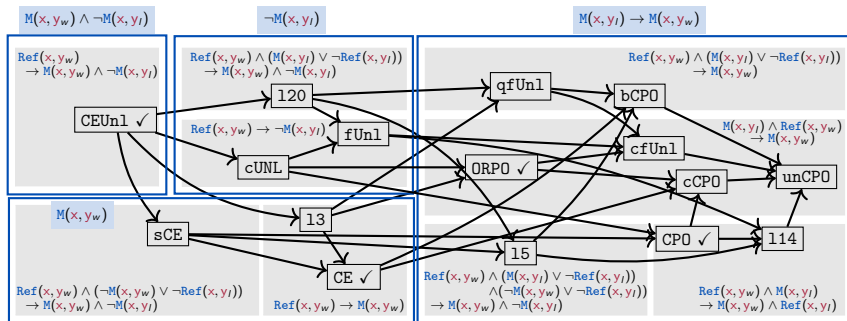
Whenever the model being tuned deems the loser to be a valid generation and the reference model deems the winner to be valid, the tuned model should deem the winner to be valid too, and *the reference should deem the loser to be valid.*

- **Peculiar semantics**, but the logic makes sense, e.g., we want to maximize

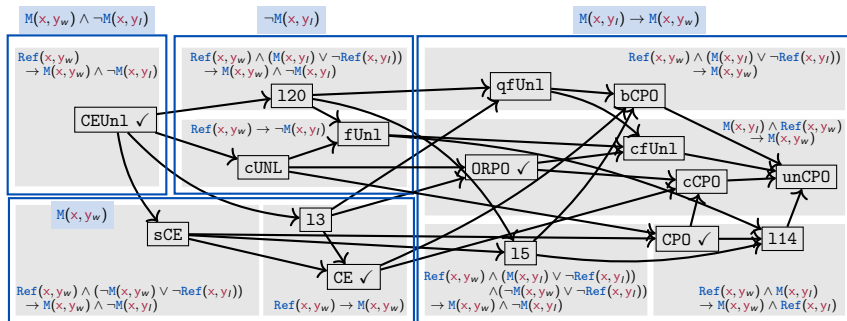
$$\sigma \left( \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\theta}(y_l | x)} - \frac{\pi_{\text{ref}}(y_w | x)}{\pi_{\text{ref}}(y_l | x)} \right)$$

negating left side of implication (i.e., making  $M(x, y_l)$  and  $\text{Ref}(x, y_w)$  false) and making the right side true is logical.

# The full landscape, reference approaches



# The full landscape, reference approaches



- Many new losses to explore and experiment with!

# Conclusions

- ▶ New ideas about formalizing preference loss functions using symbolic techniques, developed new technical tools for this.

# Conclusions

- ▶ New ideas about formalizing preference loss functions using symbolic techniques, developed new technical tools for this.
  1. Understanding the full space of loss functions (finding: it's a huge space, many novel variations yet to be explored)
  2. Understanding the structure of the space and relationships between different losses (finding: tied to the semantics of the losses).



# Conclusions

- ▶ New ideas about formalizing preference loss functions using symbolic techniques, developed new technical tools for this.
  1. Understanding the full space of loss functions (finding: it's a huge space, many novel variations yet to be explored)
  2. Understanding the structure of the space and relationships between different losses (finding: tied to the semantics of the losses).

**The procedure:** write a (high-level) symbolic program, or modify an existing one, compile into a loss and experiment (then repeat)

Thank you.

# References I

- Amini, A., Vieira, T., and Cotterell, R. (2024). Direct preference optimization with an offset. *arXiv preprint arXiv:2402.10571*.
- Azar, M. G., Rowland, M., Piot, B., Guo, D., Calandriello, D., Valko, M., and Munos, R. (2023). A general theoretical paradigm to understand learning from human preferences. *arXiv preprint arXiv:2310.12036*.
- Chavira, M. and Darwiche, A. (2008). On probabilistic inference by weighted model counting. *Artificial Intelligence*, 172(6-7):772–799.
- 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*.
- Friedman, D., Wettig, A., and Chen, D. (2023). Learning transformer programs. *Advances in Neural Information Processing Systems*, 36:49044–49067.
- Gu, Y., Mishra, B. D., and Clark, P. (2023). Do language models have coherent mental models of everyday things? *Proceedings of ACL*.
- Hong, J., Lee, N., and Thorne, J. (2024). Reference-free monolithic preference optimization with odds ratio. *arXiv preprint arXiv:2403.07691*.
- Jeffrey, R. C. (1965). *The logic of decision*. University of Chicago press.
- 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.

# References II

- Kassner, N., Tafjord, O., Sabharwal, A., Richardson, K., Schuetze, H., and Clark, P. (2023). Language models with rationality. *Proceedings of EMNLP*.
- Meng, Y., Xia, M., and Chen, D. (2024). Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*.
- Merrill, W. and Sabharwal, A. (2023). A logic for expressing log-precision transformers. *Advances in neural information processing systems*, 36:52453–52463.
- Pal, A., Karkhanis, D., Dooley, S., Roberts, M., Naidu, S., and White, C. (2024). Smaug: Fixing failure modes of preference optimisation with dpo-positive. *arXiv preprint arXiv:2402.13228*.
- Park, R., Rafailov, R., Ermon, S., and Finn, C. (2024). Disentangling length from quality in direct preference optimization. *arXiv preprint arXiv:2403.19159*.
- Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C. D., and Finn, C. (2024). Direct preference optimization: Your language model is secretly a reward model. *Proceedings of Neurips*.
- Rescher, N. (1967). *The logic of decision and action*. University of Pittsburgh Pre.
- Tang, Y., Guo, Z. D., Zheng, Z., Calandriello, D., Munos, R., Rowland, M., Richemond, P. H., Valko, M., Pires, B. Á., and Piot, B. (2024). Generalized preference optimization: A unified approach to offline alignment. *arXiv preprint arXiv:2402.05749*.
- Weiss, G., Goldberg, Y., and Yahav, E. (2021). Thinking like transformers. In *International Conference on Machine Learning*, pages 11080–11090. PMLR.

# References III

- Xu, H., Sharaf, A., Chen, Y., Tan, W., Shen, L., Van Durme, B., Murray, K., and Kim, Y. J. (2024). Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. *arXiv preprint arXiv:2401.08417*.
- Xu, J., Zhang, Z., Friedman, T., Liang, Y., and Broeck, G. (2018). A Semantic Loss Function for Deep Learning with Symbolic Knowledge. In *International Conference on Machine Learning*, pages 5498–5507.
- Yang, A. and Chiang, D. (2024). Counting like transformers: Compiling temporal counting logic into softmax transformers. *arXiv preprint arXiv:2404.04393*.
- Yuan, H., Yuan, Z., Tan, C., Wang, W., Huang, S., and Huang, F. (2024). Rhhf: Rank responses to align language models with human feedback. *Advances in Neural Information Processing Systems*, 36.
- Zhao, Y., Khalman, M., Joshi, R., Narayan, S., Saleh, M., and Liu, P. J. (2022). Calibrating sequence likelihood improves conditional language generation. In *The eleventh international conference on learning representations*.