

# Understanding the Logic of Generative AI through Logic and Programming

**Kyle Richardson**

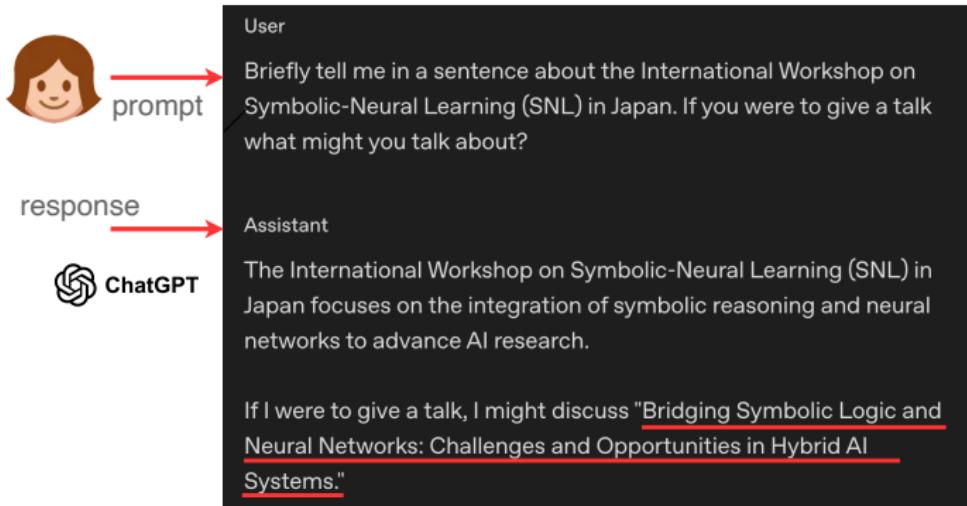
Allen Institute for AI (AI2)

October 2025

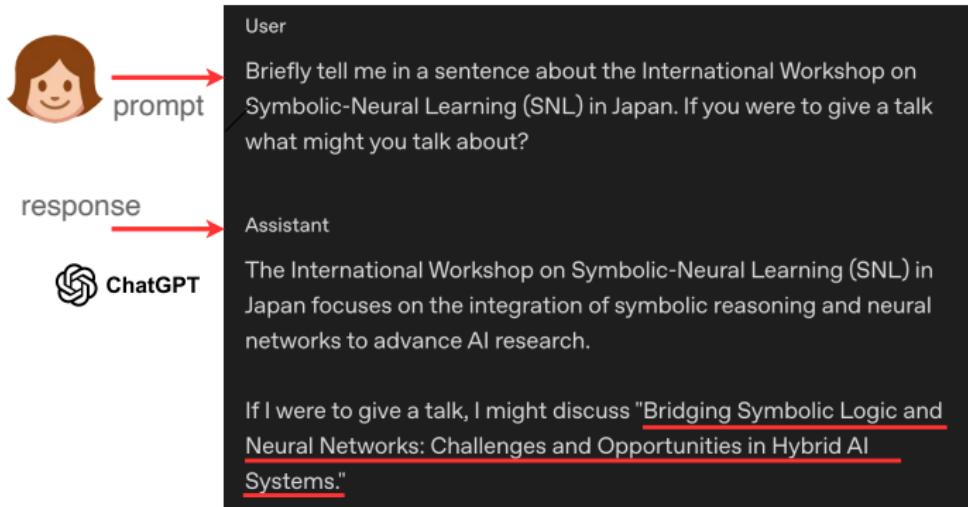
**Collaborators:** Ashish Sabharwal (AI2), Vivek Srikumar (University of Utah)



# General purpose large language models (LLMs)



# General purpose large language models (LLMs)



- ▶ **General purpose models:** trained at massive scales, used *as-is* and directly for a wide range of problems.

# General purpose large language models (LLMs)



- General purpose models: trained at massive scales, used *as-is* and directly for a wide range of problems.

# Language models as agent simulators

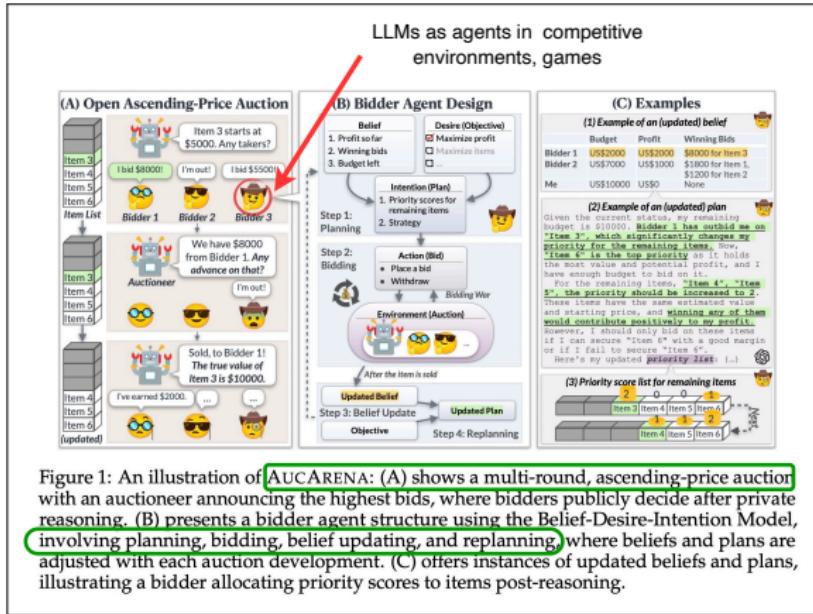
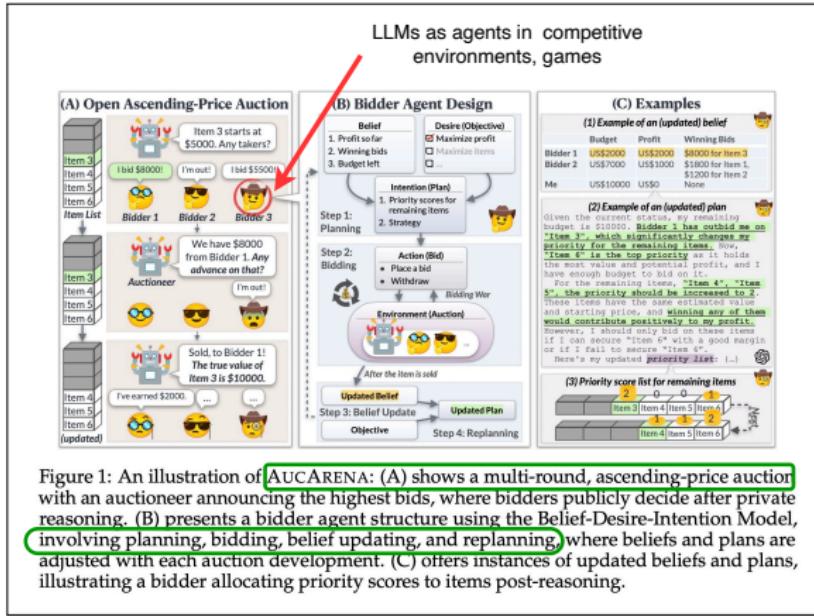


Figure 1: An illustration of AUCARENA: (A) shows a multi-round, ascending-price auction with an auctioneer announcing the highest bids, where bidders publicly decide after private reasoning. (B) presents a bidder agent structure using the Belief-Desire-Intention Model, involving planning, bidding, belief updating, and replanning, where beliefs and plans are adjusted with each auction development. (C) offers instances of updated beliefs and plans, illustrating a bidder allocating priority scores to items post-reasoning.

- ▶ Can we use LMs to simulate complex social dynamics? (Chen et al., 2023; Zhang et al., 2024; Yang et al., 2025)

# Language models as agent simulators



Valuable tool for running social science experiments, testing theories of language interaction, complex reasoning, adversarial language experts.

# Language models as part of complex systems

The screenshot shows the ML Experiment Execution Engine interface. On the left, there are three sections: "Machine learning experiment" (with a person icon), "Model generated code" (with a ChatGPT icon), and "Experiment automation" (with a wrench icon). The main area displays a query: "Query: Implement an encoder-only model using transformers that can do multiple-choice QA". Below it, "Step 1" shows a "System thought" block: "To implement an encoder-only model for multiple-choice question answering (QA) using the transformers library from Hugging Face, I'll first outline the steps in Python. This will involve loading a pre-trained encoder model, perhaps something like BERT or RoBERTa, and then adapting it to handle multiple-choice QA tasks. The model will receive a question and several possible answers as input and will have to select the most likely answer. I will use the transformers library for this purpose." An "Action (execute)" block contains Python code for loading a BERT model and tokenizer:

```
from transformers import AutoModelForMultipleChoice, AutoTokenizer  
  
model_name = 'bert-base-uncased'  
  
# Load model and tokenizer.  
tokenizer = AutoTokenizer.from_pretrained(model_name)  
model = AutoModelForMultipleChoice.from_pretrained(model_name)  
  
print("Model and tokenizer loaded successfully.")
```

The "Observation" section shows the output: "Model and tokenizer loaded successfully." Below it, "Step 2" shows a "Running step 2" status.

On the right, a Jupyter notebook titled "automated\_ml\_notebook Last Checkpoint: 23 minutes ago" shows the same code being run. A red arrow points to the progress bar for "Model and tokenizer loaded successfully", which is at 100% completion. Another red arrow points to the output cell, which shows the message "Model and tokenizer loaded successfully".

## Experiment automation

- SUPER ([Bogin et al., 2024](#)), benchmark for setting up and executing research code repositories, agent benchmarking ([Bragg et al., 2025](#)).

# Language models as part of complex systems, agents

## Language Modeling by Language Models

Junyan Cheng<sup>a,\*</sup> Peter Clark<sup>b</sup> Kyle Richardson<sup>b</sup>  
Allen Institute for AI<sup>a</sup> Dartmouth College<sup>b</sup>  
[jc.tbh@dartmouth.edu](mailto:jc.tbh@dartmouth.edu) [kyler@allenai.org](mailto:kyler@allenai.org)  
<https://github.com/allenai/genesys>

### Abstract

*Can we leverage LLMs to model the process of discovering novel language model architectures?* Inspired by real research, we propose a multi-agent LLM approach that simulates the conventional stages of research, from ideation and literature search (proposal stage) to design implementation (code generation), generative pre-training, and downstream evaluation (verification). Using ideas from scaling laws, our system *Genesys* employs a *Ladder of Scales* approach; new designs are proposed, adversarially reviewed, implemented, and selectively verified at increasingly larger model scales (14M–~250M parameters) with a narrowing budget (the number of models we can train at each scale). To help make discovery efficient and factorizable, *Genesys* uses a novel genetic programming backbone, which we show has empirical advantages over commonly used direct prompt generation workflows (e.g., ~86% percentage point improvement in successful design generation, a key bottleneck). We report experiments involving 1,162 newly discovered designs (1,062 fully verified through pre-training) and find the best designs to be highly competitive with known architectures (e.g., outperform GPT2, Manta2, etc., on 6/9 common benchmarks). We couple these results with comprehensive system-level ablations and formal results, which give broader insights into the design of effective autonomous LLM-driven discovery systems.

Cheng et al. (2025)

## TINYSCIENTIST: An Interactive, Extensible, and Controllable Framework for Building Research Agents

Haofei Yu<sup>1\*</sup> Keyang Xuan<sup>1,\*</sup> Fenghai Li<sup>1\*</sup>  
Kunlun Zhu<sup>1</sup> Zijie Lei<sup>1</sup> Jiaxun Zhang<sup>1</sup> Ziheng Qi<sup>1</sup>  
Kyle Richardson<sup>2</sup> Jiaxuan You<sup>1</sup>  
<sup>1</sup>University of Illinois Urbana-Champaign,  
<sup>2</sup>Allen Institute for Artificial Intelligence

### Abstract

Automatic research with Large Language Models (LLMs) is rapidly gaining importance, driving the development of increasingly complex workflows involving code generation, planning, tool usage, code execution, and human-agent interaction to accelerate research processes. However, as more researchers and developers begin to use and build upon these tools and platforms, the complexity and difficulty of extending and maintaining such agentic workflows have become a significant challenge, particularly as algorithms and architectures continue to advance. To address this growing complexity, TINYSCIENTIST identifies the essential components of the automatic research workflow and provides a modular, extensible, and controllable framework that easily adapts to new tools and supports iterative growth. We provide an open-source codebase<sup>2</sup>, an interactive web demonstration<sup>3</sup>, and a PyPI Python package<sup>4</sup> to make state-of-the-art auto-research pipelines broadly accessible to every researcher and developer.

en research pipelines (Jansen et al., 2023; Lu et al., 2024; Yamada et al., 2025; Li et al., 2024b; Cheng et al., 2025). Recent advances in this area leverage methods including multi-agent collaboration (Schmidgall et al., 2025), tool using (Skarinskis et al., 2024), and tree-based search (Yamada et al., 2025) to augment its performance.

In spite of this success, however, existing automatic research systems often design and use agentic frameworks that are overly complex and difficult to use and extend without significant technical expertise. These challenges stem from three key issues: (1) **lack of interactivity**: human researchers struggle to engage with the specific agent's research progress due to the complexity of research intents and unclear communication interfaces (Zou et al., 2025; Liu et al., 2025b), making feedback incorporation challenging. (2) **limited extensibility**: existing representative frameworks rely on rigid, tool-specific designs (Zhang et al., 2025), making it hard to integrate new tools or adapt to different research domains. (3) **insufficient controllability**: many

Yu et al. (2025)

A tool for scientific discovery, automated experiment execution, helping non-experts engage in research.

A tool for scientific discovery, automated experiment execution, helping non-experts engage in research.

Lots of optimism, hubris, Nobel prizes....

A tool for scientific discovery, automated experiment execution, helping non-experts engage in research.

Missing **semantic** and **algorithmic** foundations.

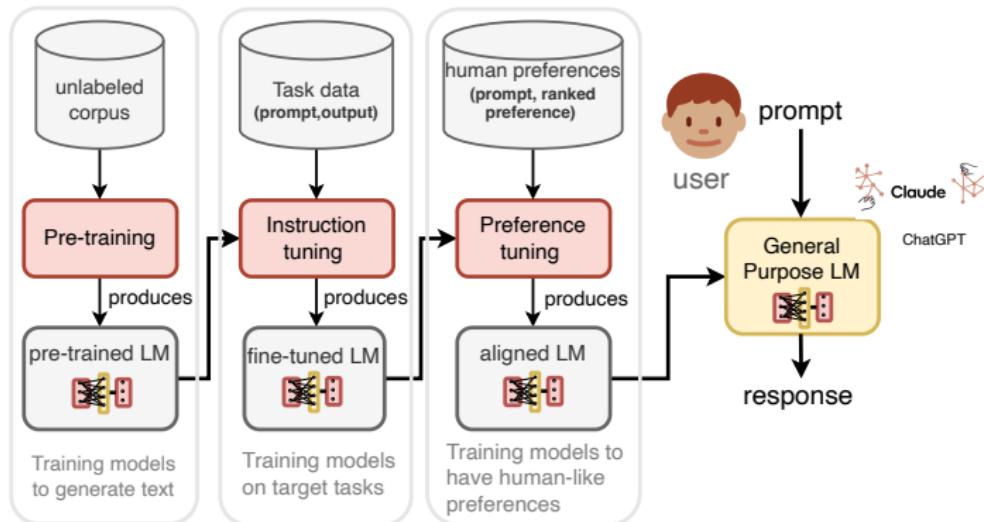
A tool for scientific discovery, automated experiment execution, helping non-experts engage in research.

Missing **semantic** and **algorithmic** foundations.

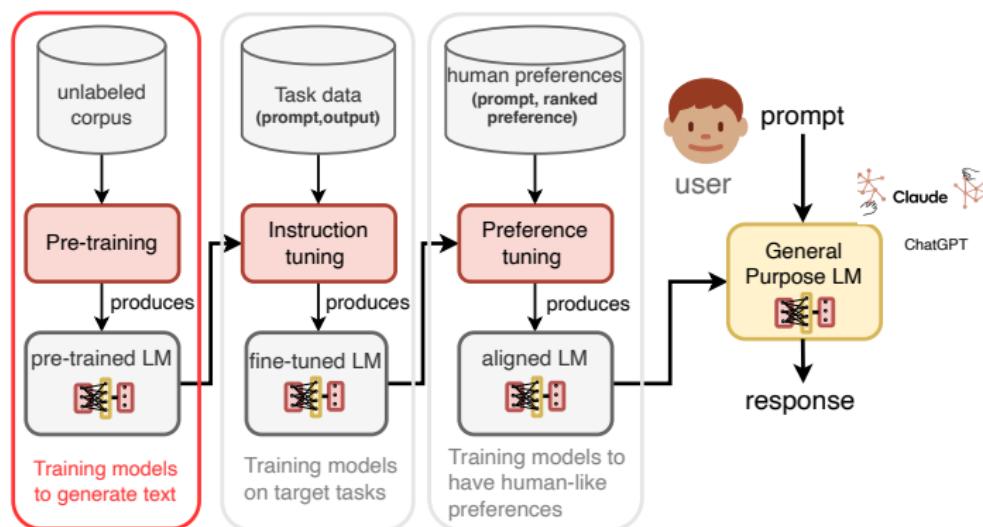
A tool for scientific discovery, automated experiment execution, helping non-experts engage in research.

Can symbolic techniques help?

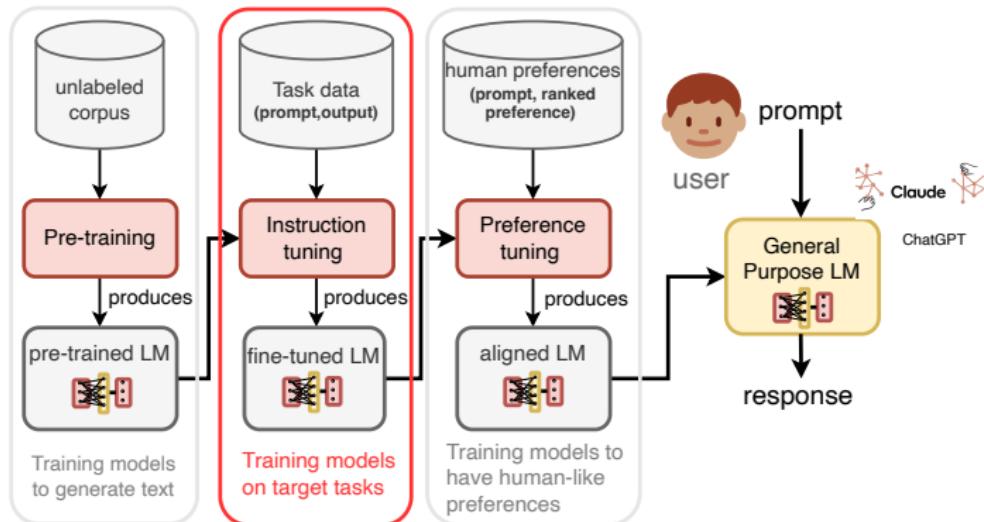
# How do we get to general purpose LLMs? recipe



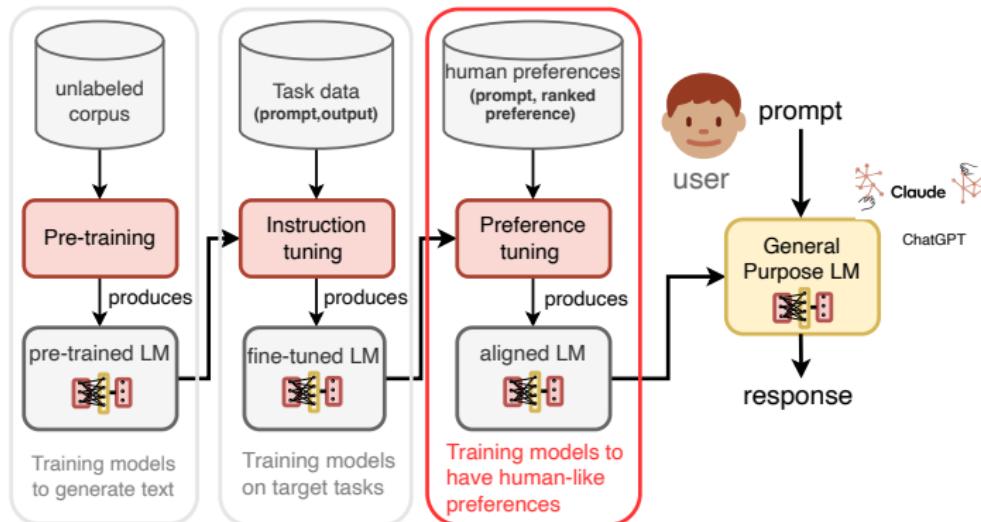
# How do we get to general purpose LLMs? recipe



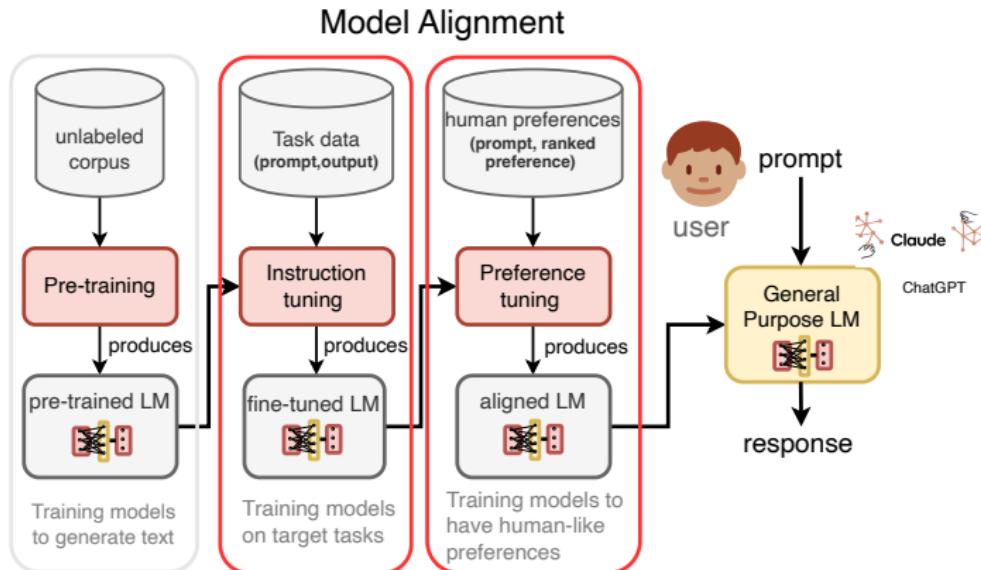
# How do we get to general purpose LLMs? recipe



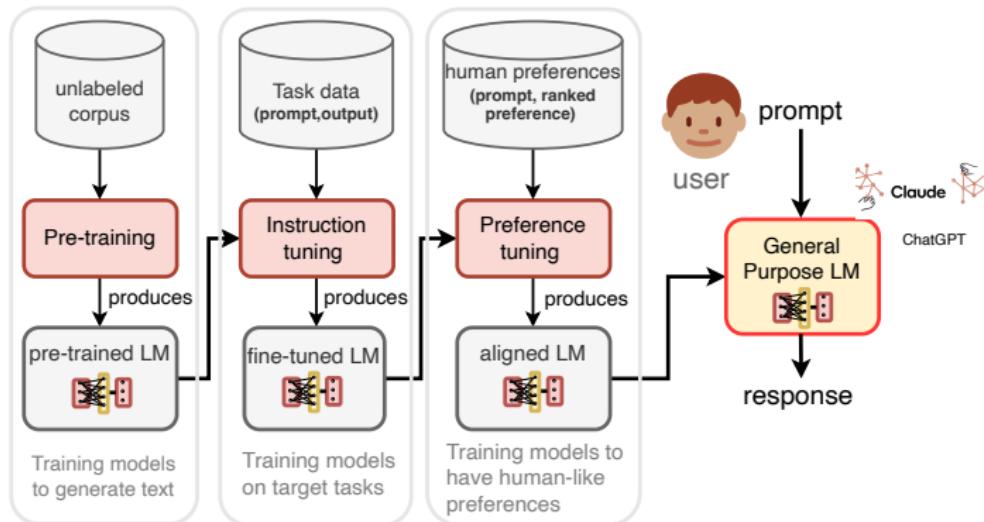
# How do we get to general purpose LLMs? recipe



# How do we get to general purpose LLMs? recipe



# How do we get to general purpose LLMs? recipe



- ▶ Rough approximation of the kinds of general purpose models we use.

# OLMo: fully open-source general purpose LMs

huggingface.co/allenai/OLMo-2-1124-13B-Instruct

allenai/OLMo-2-1124-13B-Instruct · like 46 · Following · A2 · 3.6k

Text Generation · Transformers · Safetensors · allenai/RUV-MATH · English · olmo2 · conversational · anvl/2501.09696 · anvl/2411.15124 · License: apache-2.0

Model card · Files and versions · Community · Settings

All code and artifacts

NOTE: 1/3/2025 UPDATE:

Upon the initial release of OLMo-2 models, we realized the post-trained models did not share the pre-tokenization logic that the base models use. As a result, we have trained new post-trained models. The new models are available under the same names as the original models, but we have made the old models available with a postfix "-preview". See [OLMo-2 Preview Post-trained Models](#) for the collection of the legacy models.

Release Documentation

OLMo 2 13B Instruct November 2024 is post-trained variant of the [OLMo-2 13B November 2024](#) model, which has undergone supervised finetuning on an OLMo-specific variant of the [Tulu 3 dataset](#) (<https://huggingface.co/datasets/allenai/tulu-3-sh-olmo-2-mixture>) and further DPO training on [this dataset](#), and finally RUVR training using [this data](#). Tulu 3 is designed for state-of-the-art performance on a diversity of tasks in addition to chat, such as MATH, GSM8K, and IFFEval. Check out the [OLMo 2 paper](#) or [Tulu 3 paper](#) for more details!

Edit model card

Downloads last month: 10,023 · View full history

Safetensors

Model size: 13.7B params · Tensor type: BF16 · Chat template · Files info

Inference Providers

Text Generation

This model isn't deployed by any Inference Provider. Ask for provider support

Model tree for allenai/OLMo-2-1124-13B-Instruct

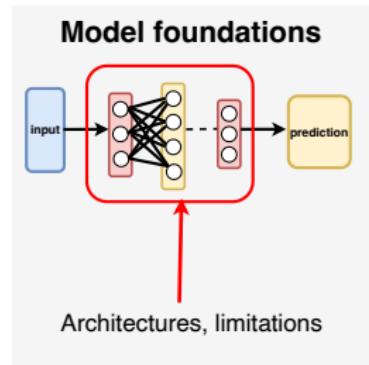
- Base model
  - Finetuned
    - Finetuned
      - Finetuned
        - Finetuned (1)
          - Adapters
          - Finetunes
          - Quantizations
  - allenai/OLMo-2-1124-7B
  - allenai/OLMo-2-1124-7B-SFT
  - allenai/OLMo-2-1124-7B-DPO
  - allenai/OLMo-2-1124-13B-Instruct-RUV1
  - allenai/OLMo-2-1124-13B-Instruct-RUV2

this model

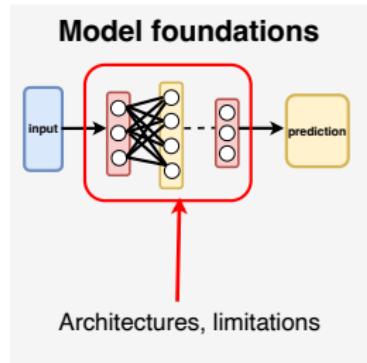
  - 4 models
  - 2 models
  - 29 models

<https://allenai.org/olmo>

# The landscape of Generative AI research

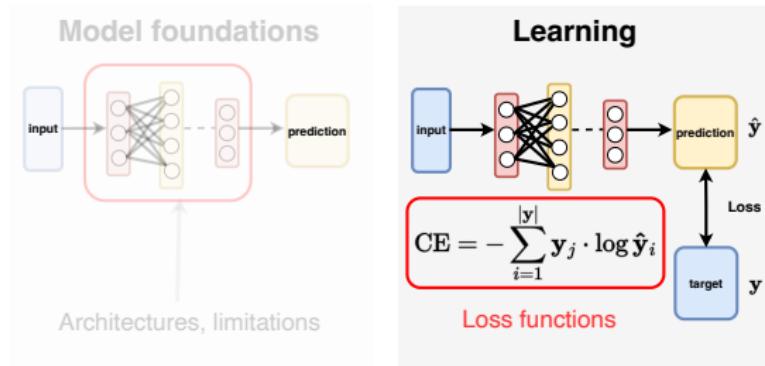


# The landscape of Generative AI research



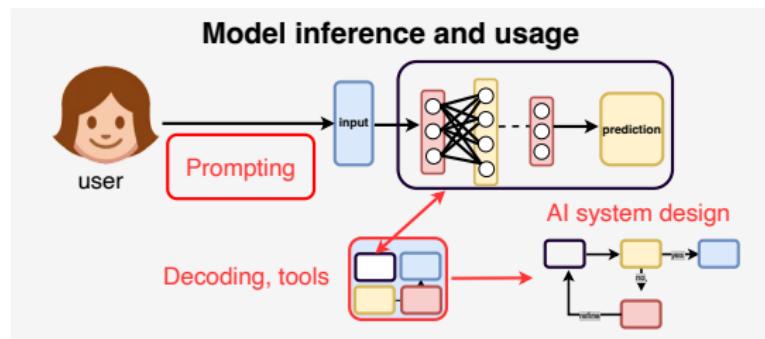
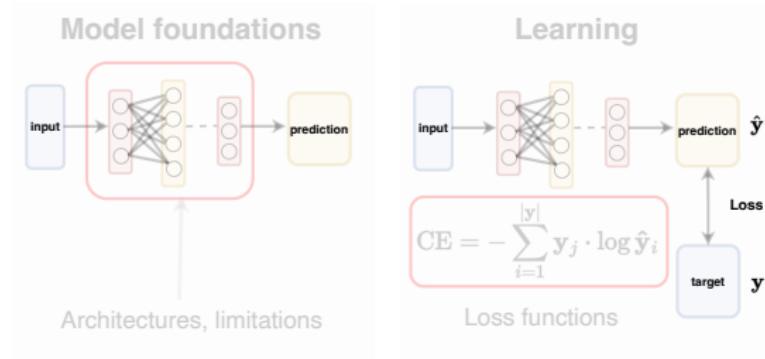
- ▶ What model to use? What kinds of computational problems can models solve? Limitations

# The landscape of Generative AI research

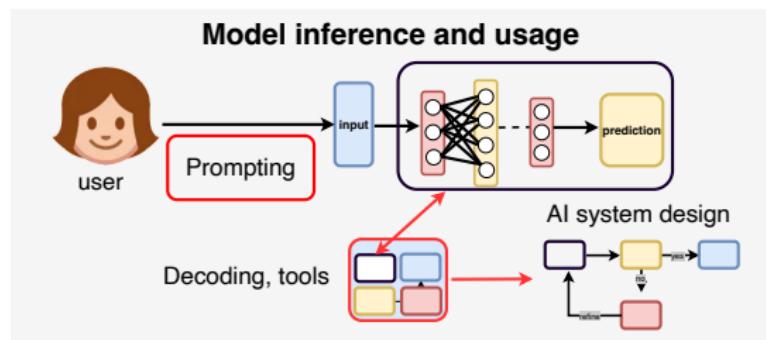
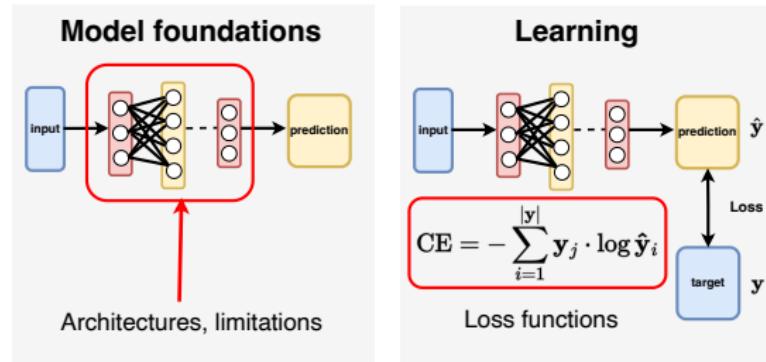


- ▶ How do we train and tune models? Loss function design, optimization algorithms.

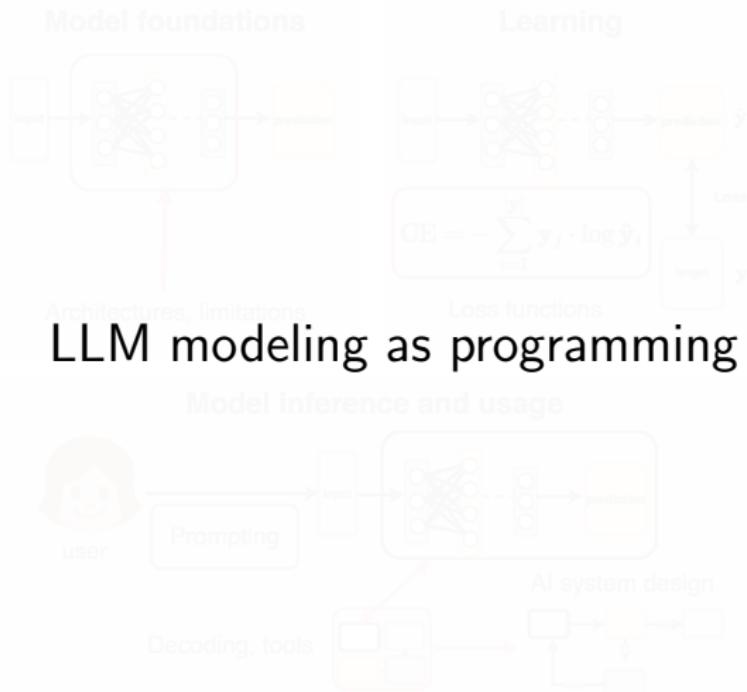
# The landscape of Generative AI research



# The landscape of Generative AI research



# The landscape of Generative AI research



# Programming techniques for model development

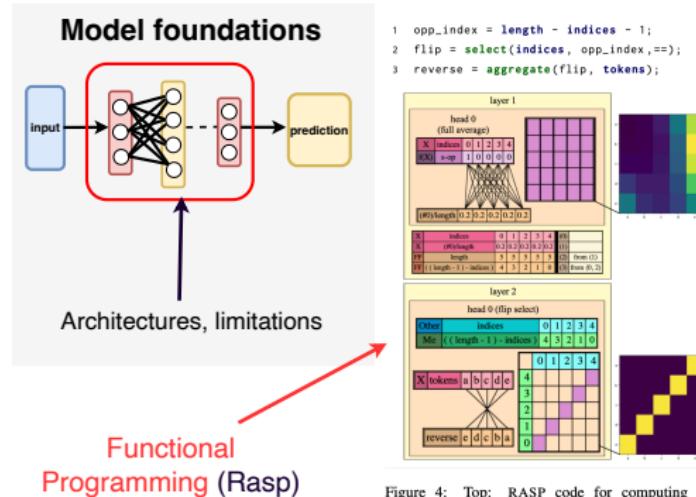
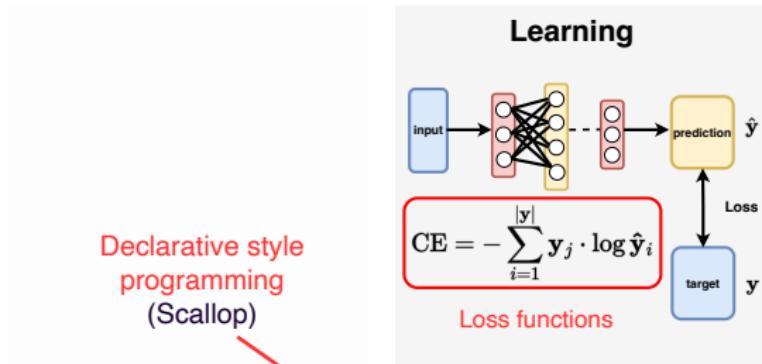


Figure 4: Top: RASP code for computing reverse (e.g., reverse("abc")="cba"). Below, its compilation to a transformer architecture (left, obtained through draw(reverse, "abcde") in the RASP REPL), and the attention heatmaps of a transformer trained on the same task (right), both visualised on the same input. Visually, the attention

- ▶ Programming languages for expressing transformer computation (Weiss et al., 2021; Yang et al., 2024; Yang and Chiang, 2024).

# Programming techniques for model development



```
1 // File path_planner.scl
2 type actor(x: i32, y: i32), goal(x: i32, y: i32), enemy(x: i32, y: i32)
3
4 const UP = 0, DOWN = 1, RIGHT = 2, LEFT = 3
5 rel safe_cell(x, y) = range(0, 5, x), range(0, 5, y), not enemy(x, y)
6 rel edge(x, y, x, yp, UP) = safe_cell(x, y), safe_cell(x, yp), yp == y + 1
7 // Rules for DOWN, RIGHT, and LEFT edges are omitted...
8
9 rel next_pos(p, q, a) = actor(x, y), edge(x, y, p, q, a)
10 rel path(x, y, x, y) = next_pos(x, y, _)
11 rel path(x1, y1, x3, y3) = path(x1, y1, x2, y2), edge(x2, y2, x3, y3, _)
12 rel next_action(a) = next_pos(p, q, a), path(p, q, r, s), goal(r, s)
```

Fig. 3. The logic program of the PacMan-Maze application in Scallop.

- ▶ Loss design via logical and probabilistic programming, neuro-symbolic modeling (Li et al., 2023; Manhaeve et al., 2018).

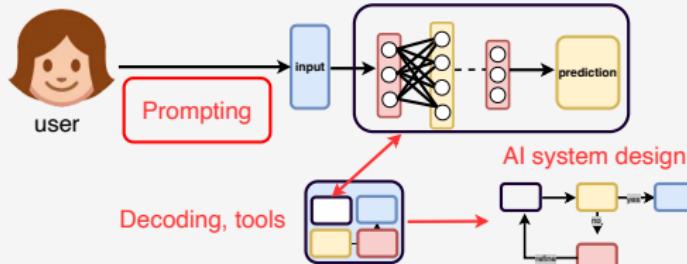
# Programming techniques for model development

## Structured imperative programming (LMQL)

```
gql.query
def meaning_of_life():
    # top-level strings are prompts
    "Q: What is the answer to life, the
     universe and everything?"
    # generation via (constrained) variables
    "# {ANSWER}" where {
        len(ANSWER) < 120 and STOP$_AT(ANSWER, ".")
    }
    # results are directly accessible
    print("LM returned", ANSWER)
    # use typed variables for guaranteed
    # output format
    # "The answer is {NUM: int}"
    # query programs are just functions
    return NUM

    # so from Python, you can just do this
meaning_of_life() # 42
```

## Model inference and usage



- ▶ Prompting as (imperative) programming (Beurer-Kellner et al., 2023).

# Programming techniques for model development

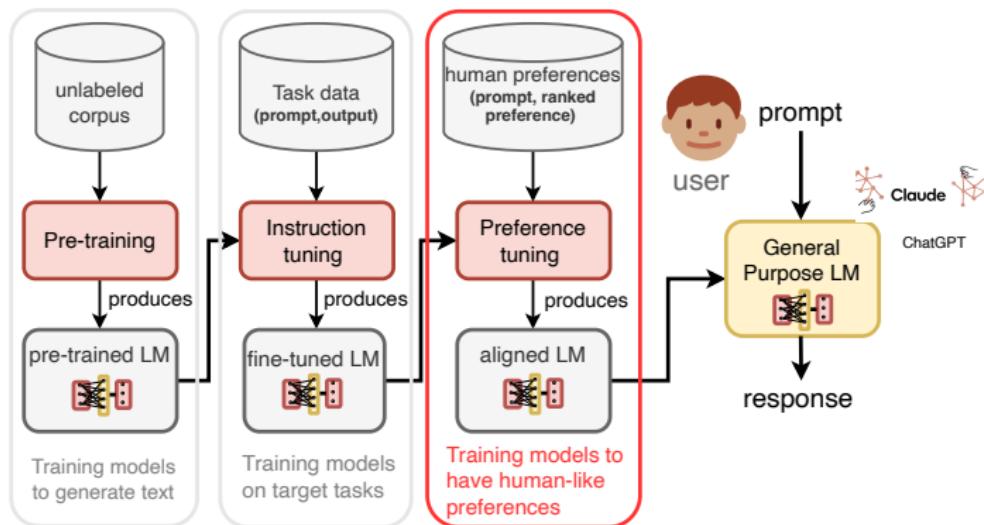


Declarative-style programming, loss function design.



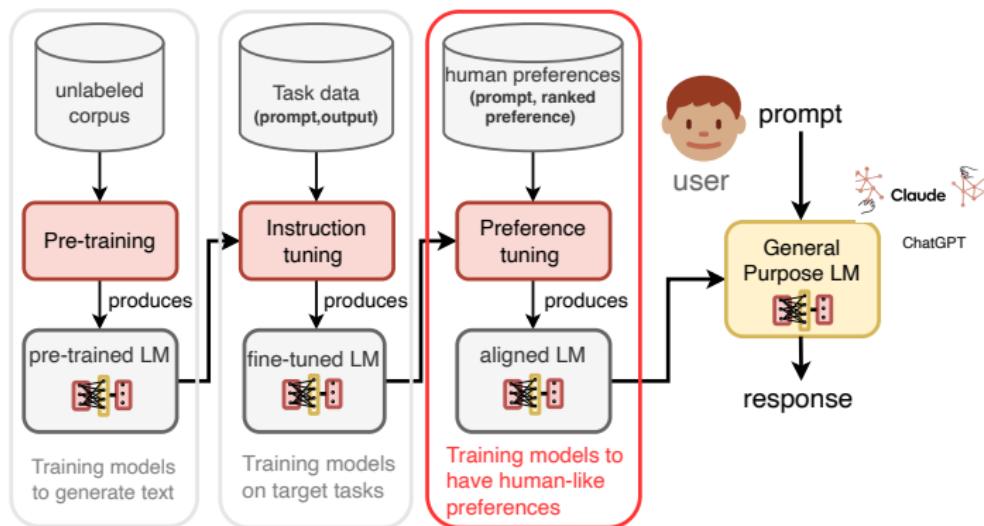
- ▶ Prompting as (imperative) programming (Beurer-Kellner et al., 2023)

# Declarative-style programming for preference modeling



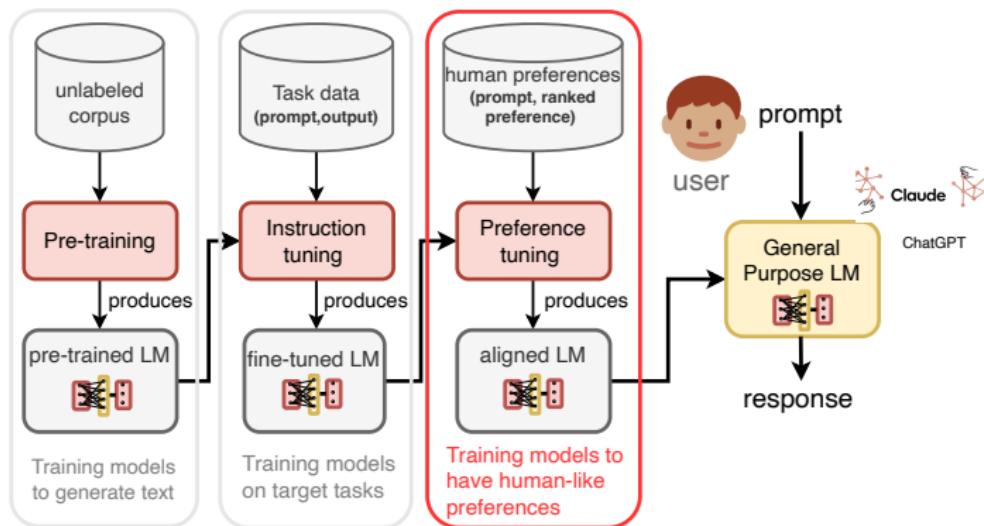
**Today:** Logical and probabilistic programming of preference losses, semantic characterizations.

# Declarative-style programming for preference modeling



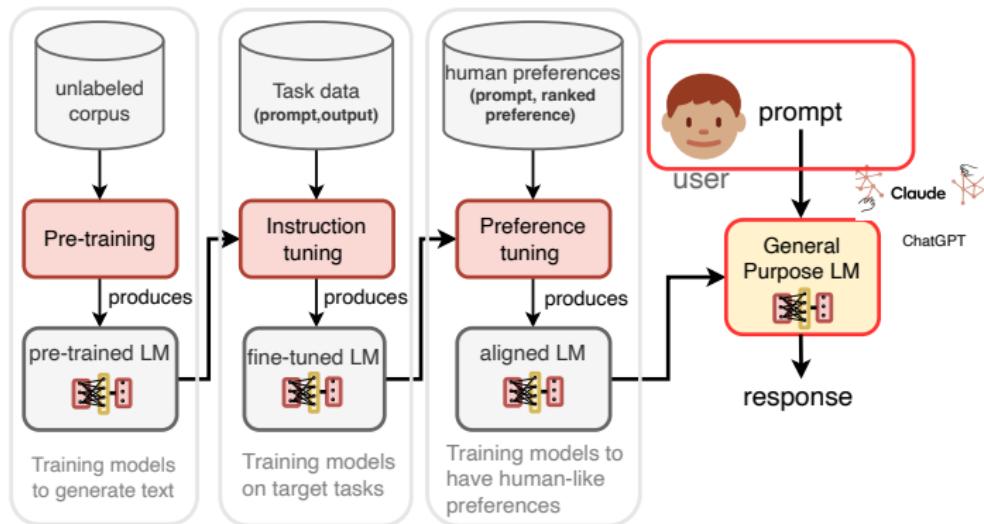
**Questions:** What do we do when we tune models to preferences? Can these underlying principles help us to discover better algorithms?

# Declarative-style programming for preference modeling



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

# ~~Probabilistic programming for other applications~~



**If time permits:** Probabilistic programming techniques and languages for prompting and chain-of-thought.

# Probabilistic programming for other applications



## 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, ....

# 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, ....

**Note:** What constitutes a *winner* or *loser* is fuzzy, datasets are very noisy and aggregate many different kinds of 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)

## 1. Unclear what is actually in our datasets

$y_l$  : 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, ....

**Note:** What constitutes a *winner* or *loser* is fuzzy, datasets are very noisy and aggregate many different kinds of preferences.

# 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  
`{rafailev, 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>, Michael A. Lewis<sup>†</sup>, Chelsea Finn<sup>†</sup>

## DPO loss function

<sup>\*</sup>Stanford University <sup>†</sup>CZ Biohub

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

model generates 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 measures the human preference, and then the reward is backpropagated to LM. This is far from the original model. In this paper we introduce a new parameterization of the reward model in RLHF that enables extraction of the components in control of the reward function, which we call *Direct Preference Optimization* (DPO), and 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,architeh,eric.mitchell)@cs.stanford.edu

## 2. These equations are not easy to understand

Abstract. Large-scale unsupervised language models (LMs) learn to do what is known: edge 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 sterility 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>, Michael D. Noll<sup>†</sup>, Chelsea Finn<sup>†</sup>

### DPO loss function

<sup>\*</sup>Stanford University <sup>†</sup>CZ Biohub

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

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 policies are trained on, and then using the learned reward to fine-tune the language model. This can lead to instability, especially if the reward model is 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 directly from the language model. This reduces the 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.

**Question:** What kind of discrete reasoning problems do these losses encode?

# The many varieties of DPO

## 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  
`{rafailev, architsb, 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.

## DPO loss

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



# The many varieties of DPO

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

Rafael Rafailev<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  
rafailev,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.

## DPO loss

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

## DPO variants

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

from Meng et al. (2024)

- No reference approaches (e.g., CPO, ORPO, *only involves a single model*) versus multi-model, reference approaches (DPO).

# The many varieties of DPO

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

Rafael Rafailev<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  
rafailev,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.

## DPO loss

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

## DPO variants

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

from Meng et al. (2024)

**Questions:** How are all these variations related to one another, nature of the space of losses?

# The many varieties of DPO

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

Rafael Radulov<sup>a†</sup>

Archit Sharma<sup>a\*</sup>

Eric Mitchell<sup>a†</sup>

Stefano Ermon<sup>b‡</sup>

Christopher D. Manning<sup>b</sup>

Chelsea Finn<sup>b</sup>

<sup>a</sup>Stanford University <sup>b</sup>CZ Labub  
(rafrafa@cs.stanford.edu, architsharma@cs.stanford.edu, eric.mitchell1@cs.stanford.edu)

### Abstract

While large-scale model training is a well-studied problem, aligning language models to human preferences is much more difficult due to the complex, subjective nature of human feedback. Most prior work requires humans to provide explicit methods for guiding such stochasticity: collect human labels of the relative quality of model generations and fine-tune the unsupervised LM to align with those preferences, often with reinforcement learning from human feedback (RLHF). However, RLHF is a complex and often costly procedure, first fitting a reward model that reflects the human preferences, and then fine-tuning the large unsupervised LM using reinforcement learning to maximize the 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 single classification loss. The resulting algorithm, which we call *Direct Preference Optimization* (DPO), is safe, 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, the training with DPO exhibits PPO-based RLHF in safety, in control sentiment of generations, and in lack of improvements to quality in summarization and single-line generation while being subsumed by the global alignment and train-

## Formalization of preference losses

### DPO loss

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

### DPO variants

Method	Objective
DPO [6]	$-\log \sigma \left( \beta \log \frac{\pi_\theta(y_\theta x)}{\pi_{ref}(y_{ref} x)} - \beta \log \pi_\theta(y_l x) \right)$
IPO [6]	$\left( \log \frac{\pi_\theta(y_\theta x)}{\pi_{ref}(y_{ref} x)} - \log \pi_\theta(y_l x) - \frac{1}{\alpha} \right)^2$
CPO [88]	$-\log \sigma \left( \beta \log \pi_\theta(y_\theta x) - \beta \log \pi_\theta(y_l x) \right) - \lambda \log \pi_\theta(y_l x)$
KTO [29]	$-\lambda_{opt} \left( \beta \log \frac{\pi_\theta(y_\theta x)}{\pi_{ref}(y_{ref} x)} - \beta_{opt} \right) + \lambda_{opt} \left( \beta_{opt} - \beta \log \frac{\pi_\theta(y_\theta x)}{\pi_{ref}(y_{ref} x)} \right)$ , where $\beta_{opt} = \mathbb{E}_{(x,y_\theta,y_l)\sim D} [\text{ReLU}(\pi_\theta(y_\theta x)/\pi_\theta(y_l x))]$
ORPO [42]	$-\log \pi_\theta(y_\theta x) - \lambda \log \sigma \left( \log \frac{\pi_\theta(y_\theta x)}{\pi_{ref}(y_{ref} x)} - \log \frac{\pi_\theta(y_l x)}{\pi_{ref}(y_l x)} \right)$ , where $\pi_\theta(y x) = \exp \left( \frac{1}{\alpha} \log \pi_\theta(y x) \right)$
R-DPO [64]	$-\log \sigma \left( \beta \log \frac{\pi_\theta(y_\theta x)}{\pi_{ref}(y_{ref} x)} - \beta \log \frac{\pi_\theta(y_l x)}{\pi_{ref}(y_l x)} + (\alpha y_\theta - \alpha y_l) \right)$
SimPO	$-\log \sigma \left( \frac{\beta}{\alpha} \log \pi_\theta(y_\theta x) - \frac{\beta}{\alpha} \log \pi_\theta(y_l x) - \gamma \right)$

from Meng et al. (2024)

# The many varieties of DPO

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

Rafael Radulov<sup>a†</sup>

Archit Sharma<sup>a\*</sup>

Eric Mitchell<sup>a†</sup>

Stefano Ermon<sup>b‡</sup>

Christopher D. Manning<sup>b</sup>

Chelsea Finn<sup>b</sup>

<sup>a</sup>Stanford University <sup>b</sup>CZB Biobab  
(rafrafa@cs.stanford.edu, architsharma, eric.mitchell1)@cs.stanford.edu

### Abstract

While edge cases are well-known, it is less clear how to train LMs to perform difficult tasks. One approach is to collect human annotations, then train learning methods for gaining such annotations collect human labels of the relative quality of model generations and fine-tune the untrained LM to align with these preferences, often with reinforcement learning from human feedback (RLHF). However, RLHF is a complex and often costly procedure, first fitting a reward model that reflects the human preferences, and then fine-tuning the large untrained LM using reinforcement learning to maximize the 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 single classification loss. The resulting algorithm, which we call *Direct Preference Optimization* (DPO), is safe, performant, and computationally lightweight, eliminating the need for collecting data 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, the training with DPO requires PPO-based RLHF to only re-control sentiment of generations, and none of the human feedback quality summarization and single-line annotations being subsumed by the global alignment and train-

## DPO loss

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

## DPO variants

Method	Objective
RLHF [1]	$-\log \sigma \left( \beta \log \frac{\pi_\theta(y_u x)}{\pi_{ref}(y_u x)} - \beta \log \frac{\pi_\theta(y_l x)}{\pi_{ref}(y_l x)} \right) - \lambda \log \pi_\theta(y_u x)$
RLHF [2]	$-\log \sigma \left( \beta \log \frac{\pi_\theta(y_u x)}{\pi_{ref}(y_u x)} - \beta \log \pi_\theta(y_l x) \right) + \lambda \log \pi_\theta(y_u x)$
DPO [6]	$-\log \sigma \left( \beta \log \frac{\pi_\theta(y_u x)}{\pi_{ref}(y_u x)} - \beta \log \pi_\theta(y_l x) \right)$
IPO [6]	$\left( \log \frac{\pi_\theta(y_u x)}{\pi_{ref}(y_u x)} - \log \frac{\pi_\theta(y_l x)}{\pi_{ref}(y_l x)} - \frac{1}{\alpha} \right)^2$
CPO [88]	$-\log \sigma \left( \beta \log \pi_\theta(y_u x) - \beta \log \pi_\theta(y_l x) \right) - \lambda \log \pi_\theta(y_u x)$
KTO [29]	$-\lambda_{opt} \sigma \left( \beta \log \frac{\pi_\theta(y_u x)}{\pi_{ref}(y_u x)} - \beta \log \pi_\theta(y_l x) \right) + \lambda_{opt} \left( \beta \log \frac{\pi_\theta(y_u x)}{\pi_{ref}(y_u x)} \right)$ , where $\lambda_{opt} = Z_{\text{KL}}(\pi_\theta(y_u x)\  \pi_{ref}(y_u x))$
ORPO [42]	$-\log \pi_\theta(y_u x) - \lambda \log \sigma \left( \log \frac{\pi_\theta(y_u x)}{\pi_{ref}(y_u x)} - \log \frac{\pi_\theta(y_l x)}{\pi_{ref}(y_l x)} \right)$ , where $\pi_\theta(y x) = \exp \left( \frac{1}{\alpha} \log \pi_\theta(y x) \right)$
R-DPO [64]	$-\log \sigma \left( \beta \log \frac{\pi_\theta(y_u x)}{\pi_{ref}(y_u x)} - \beta \log \frac{\pi_\theta(y_l x)}{\pi_{ref}(y_l x)} + (\alpha y_u - \alpha y_l) \right)$
SimPO	$-\log \sigma \left( \frac{\beta}{\alpha} \log \pi_\theta(y_u x) - \frac{\beta}{\alpha} \log \pi_\theta(y_l x) - \eta \right)$

from Meng et al. (2024)

# Preference learning as a discrete reasoning problem

## Loss Function $\ell$

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

# Preference learning as a discrete reasoning problem

**Loss Function  $\ell$**

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

Two models, four predictions

# Preference learning as a discrete reasoning problem

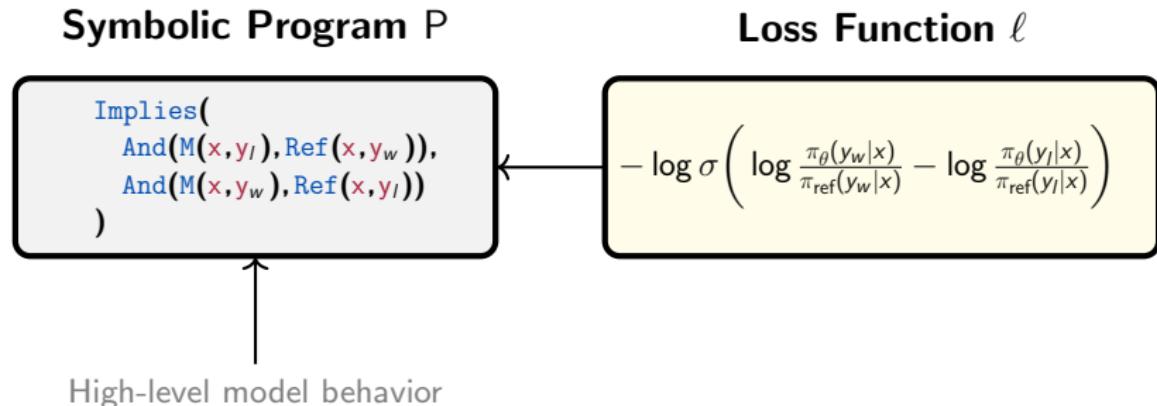
## Loss Function $\ell$

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

Two models, four predictions

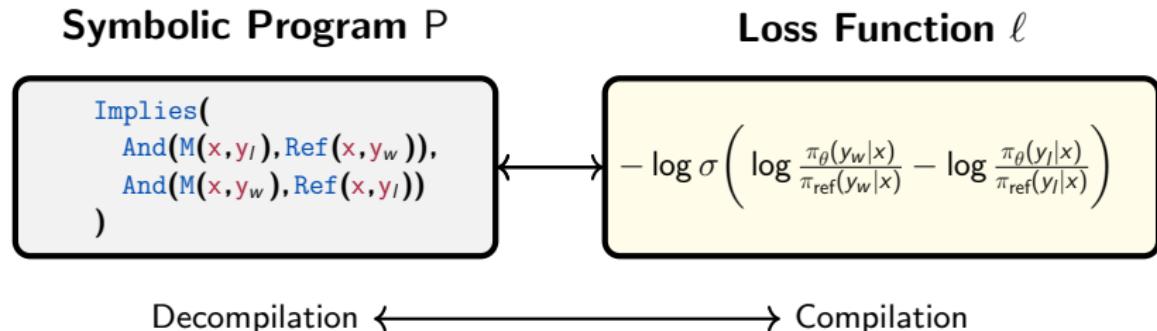
- ▶ **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?

# Preference learning as a discrete reasoning problem



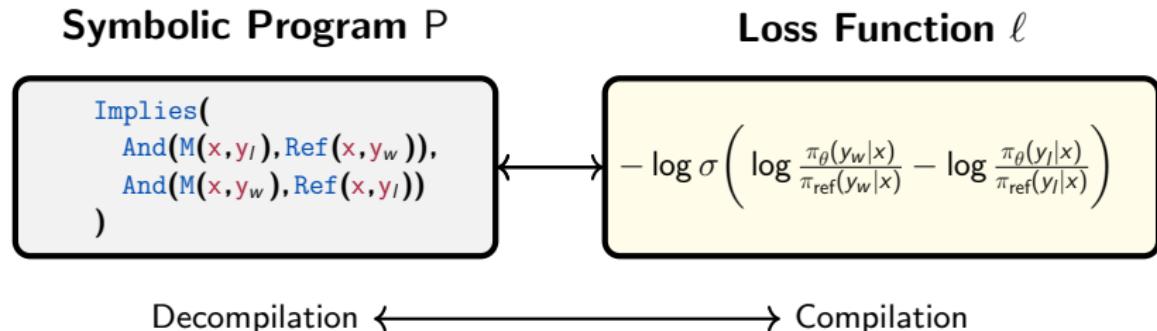
- ▶ **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?

# Preference learning as a discrete reasoning problem



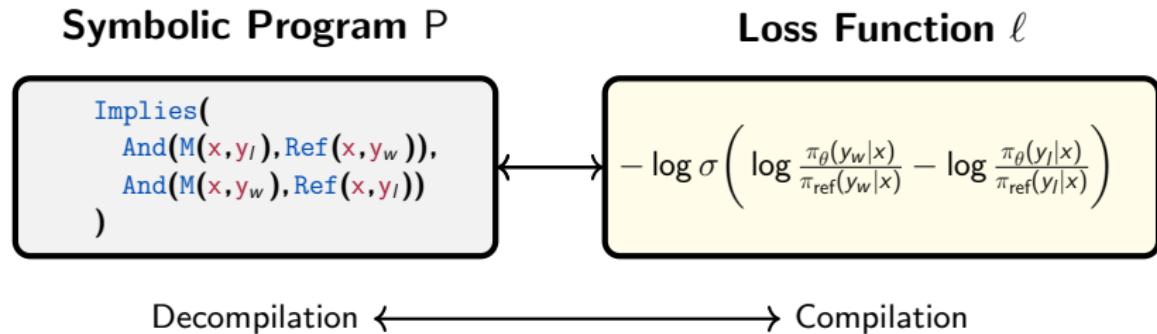
- ▶ **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?

# Preference learning as a discrete reasoning problem



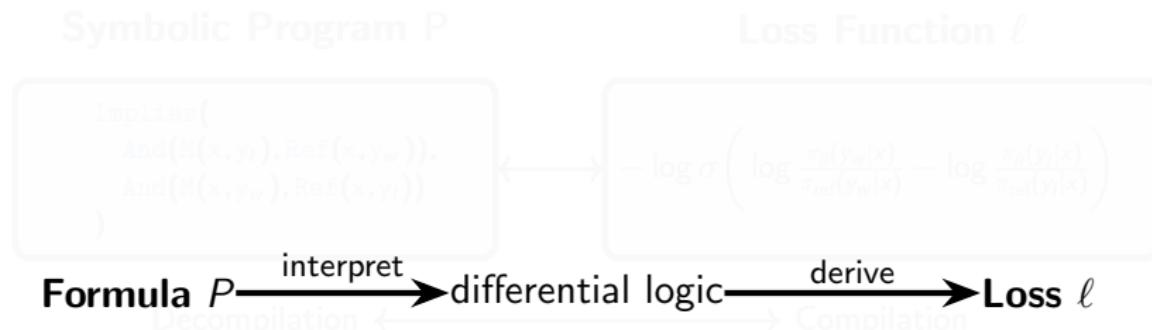
- ▶ **Problem:** Given some loss function, can we derive a symbolic program or expression that characterizes the semantics of that loss?

# Preference learning as a discrete reasoning problem



- ▶ **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.

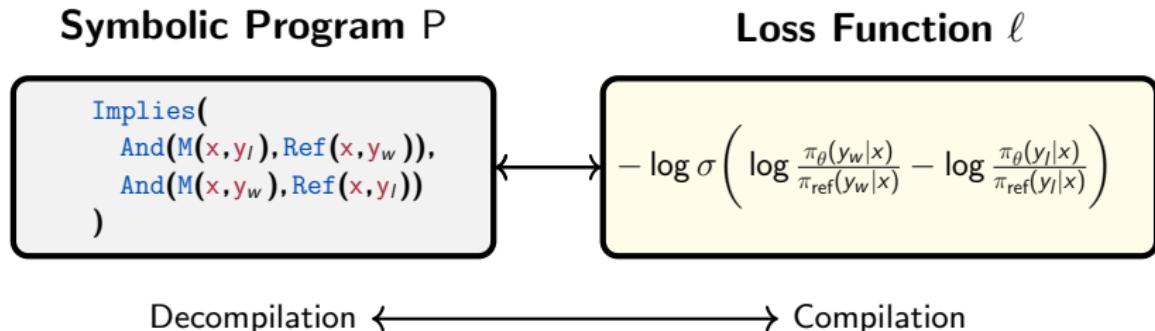
# Preference learning as a discrete reasoning problem



- Problem: Given a symbolic program or expression that characterizes the semantics of a loss
- Logic as loss, learning to satisfy* ([Marra et al., 2024](#))

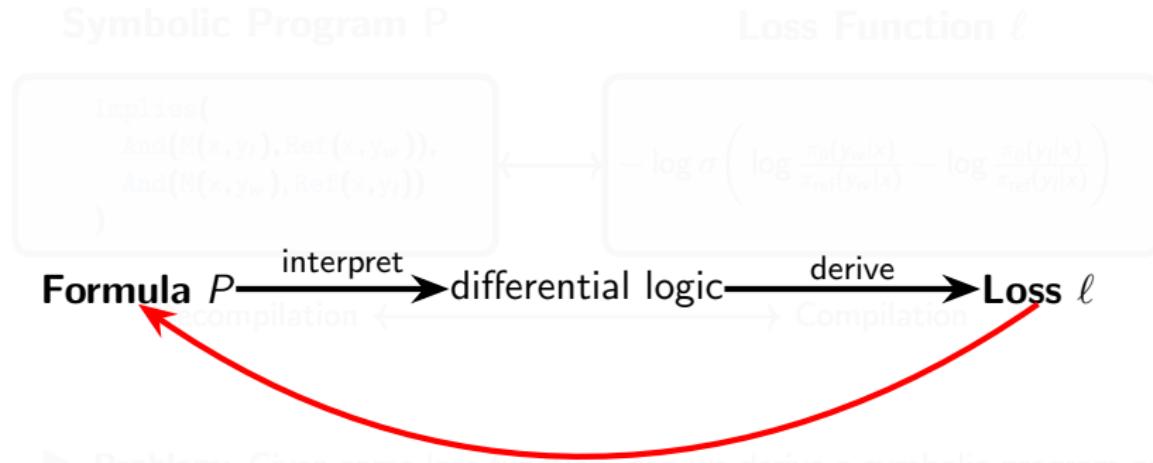
1. Compilation: Translating specifications into loss, well studied.

# Preference learning as a discrete reasoning problem



- ▶ **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.

# Preference learning as a discrete reasoning problem



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

**Loss as logic** (Richardson et al., 2025)

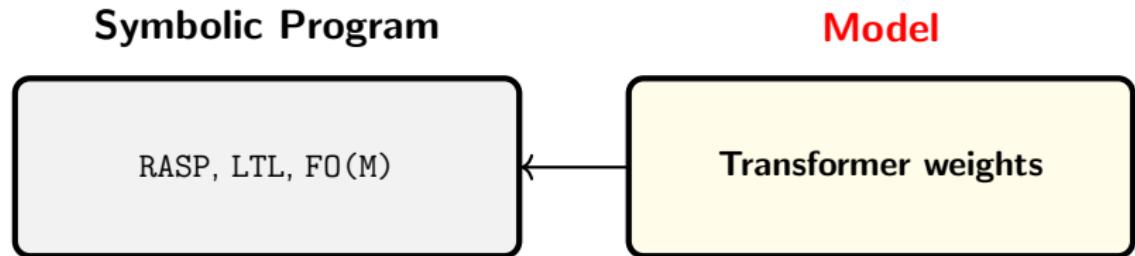
1. Compilation: Translating specifications into loss, well studied.
2. Decomposition: Losses to specifications (inverse), less explored.

# Formal analysis via decompilation in general

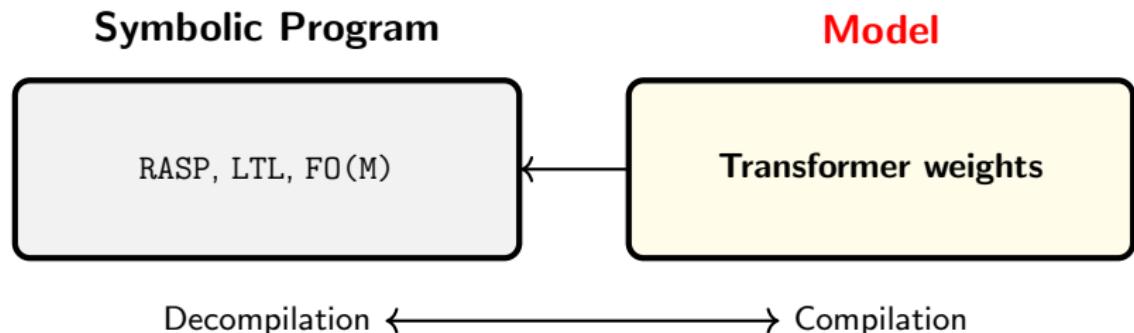
Model

Transformer weights

## Formal analysis via decompilation in general



# Formal analysis via decompilation in general



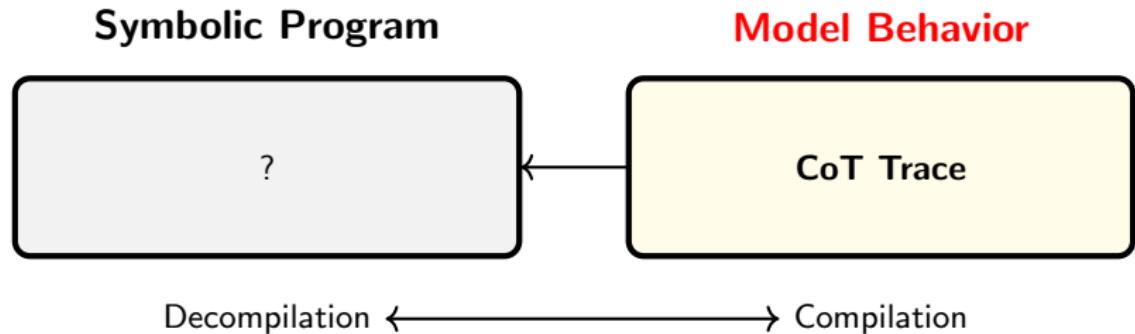
- ▶ We know what the *target languages* are ([Weiss et al., 2021](#); [Merrill and Sabharwal, 2023](#); [Yang and Chiang, 2024](#)), how to compile, decompile ([Friedman et al., 2023](#)).

Formal analysis via decompilation in general

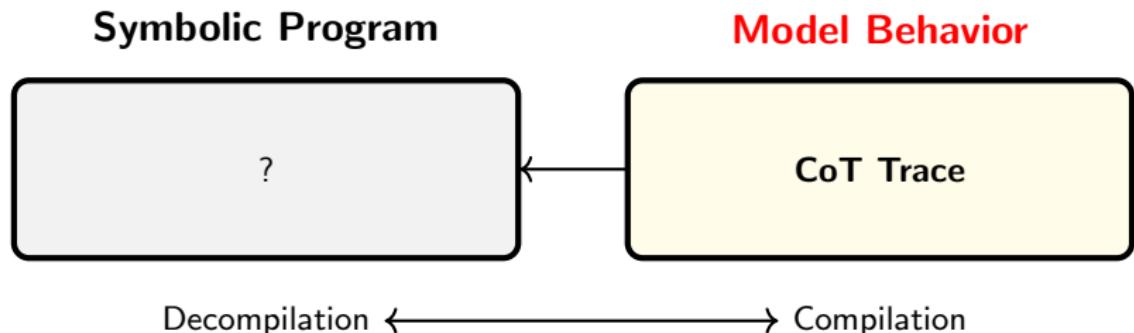
**Model Behavior**

**CoT Trace**

## Formal analysis via decompilation in general



# Formal analysis via decompilation in general



- ▶ Not always clear what the target programming language is or should be.

**Language model programming:** the languages and formal interfaces used for doing such analysis (Richardson and Wijnholds, 2025).

# Language model programming: ESSLLI 2025

## Lecturers

[Kyle Richardson](#) (Allen Institute for AI)

[Gijs Wijnholds](#) (Leiden Institute of Advanced Computer Science)

## Slides

[lecture 1](#): course overview, language modeling basics, [transformers](#), RASP.

[lecture 2](#): declarative approaches to model training and fine-tuning, the [semantic loss](#) and [weighted model counting](#), [other](#) approaches.

[lecture 3](#): high-level programming techniques for [direct preference alignment](#) and [LLM alignment](#), [formal characterizations](#) of known loss functions.

[lecture 4](#): declarative and probabilistic approaches to test-time inference, [LLM self-correction](#), [consistency](#), distilling LLMs to tractable models, [logic programming](#).

[lecture 5](#): Advanced prompting, [chain-of-thought](#), [imperative model programming](#), [\(discrete\) probabilistic programming](#).

background [logic notes](#), [extended notes on transformers](#)

extra lectures [Prompting as programming](#), [Grammar-constrained decoding](#)

## Helpful Resources

Below are some pointers to code resources:

- languages [scallop](#), [problog](#), [pyDatalog](#), [limgl](#), [rasp](#), [NumPy Rasp](#), [deepproblog](#)
- automated reasoning tools/circuits [Z3 solver](#), [python-sat](#), [pysdd](#), [cirkit](#)
- NLP and general ML [transformers](#), [PyTorch](#), [pylon-lib](#), [hf datasets](#), [hf hub](#)
- other useful utilities [sympy](#)

Useful tutorials: [Transformers from scratch](#) (some examples/ideas used in lecture 1), [Lectures on Probabilistic Programming](#), [Tractable Probabilistic Models](#)

<https://github.com/yakazimir/LMProgramming>

# Language model programming: ESSLLI 2025

## Lecturers

Hao Heinecke (Allen Institute for AI)

Günther Wärthen (Leiden Institute of Advanced Computer Science)

## Slides

lecture 1: course overview, language modeling basics, [transformers](#), [PyTorch](#)

lecture 2: declarative approaches to model training and fine-tuning, the [semantic tree](#) and [weighted model counting](#), other approaches.

lecture 3: high-level programming techniques for direct preference alignment and LLM alignment, formal characterizations of known loss functions.

lecture 4: direct model-to-model alignment, fast-time inference, LLM and preference consistency, defining LLMs as programs.

# What is the right programming language for preference?

background: [probabilistic programming](#), [probabilistic circuits](#)

extra lectures: [Prompting as programming](#), [Grammar-constrained decoding](#)

## Helpful Resources

Below are some pointers to code resources:

- languages: [Julia](#), [Python](#), [JavaScript](#), [TensorFlow](#), [PyTorch](#), [TensorFlow.js](#), [TensorFlow Lite](#)
- automated reasoning tools/circuits: [Z3](#), [CVC4](#), [EUSolver](#), [SMTInterpol](#), [SAT4J](#)
- NLP and general ML: [Transformers](#), [PyTorch](#), [TensorFlow](#), [TensorFlow.js](#), [TensorFlow Lite](#)
- other useful utilities: [Jupyter](#)

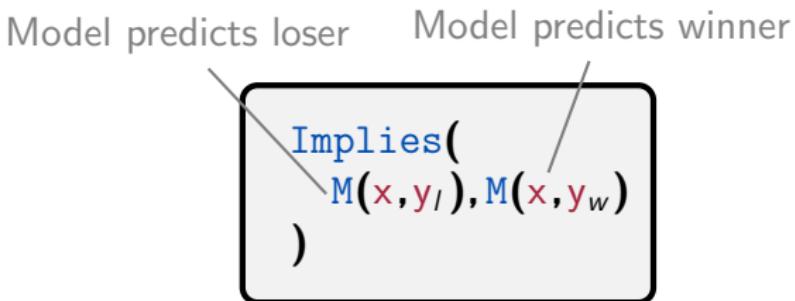
Useful tutorials: [Transformers from scratch](#) (some examples/code used in Lecture 1), [Lectures on Probabilistic Programming](#), [Probabilistic Models](#)

<https://github.com/yakazimir/LMProgramming>

## What do these programs tell us?

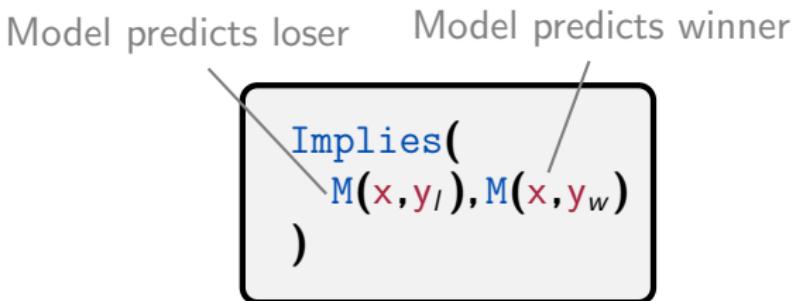
```
Implies(  
    M(x,yl), M(x,yw)  
)
```

# What do these programs tell us?



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

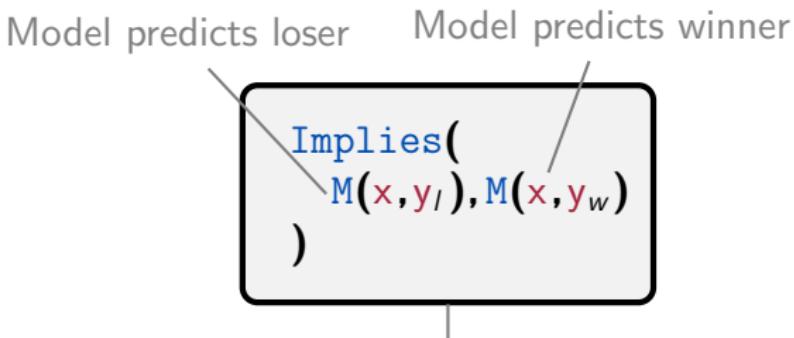
## What do these programs tell us?



$$w(M(x, y)) = \pi_M(y | x)$$

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

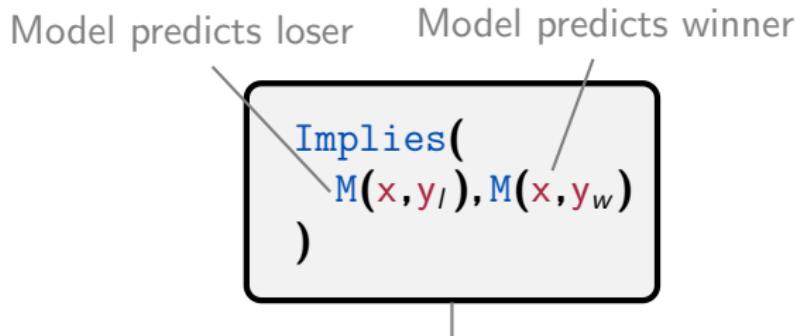
## What do these programs tell us?



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

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

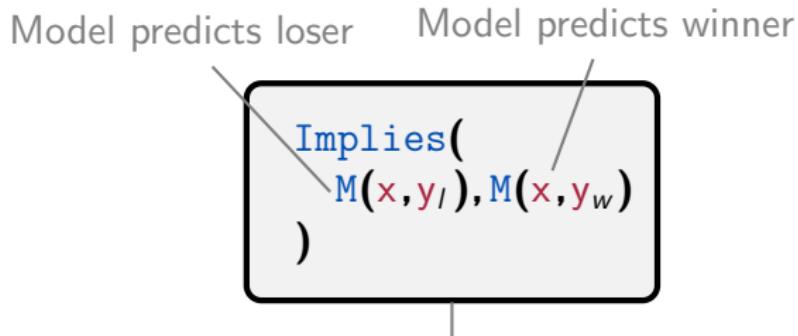
# Uncovering the natural logic of these algorithms



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.

## Uncovering the natural logic of these algorithms



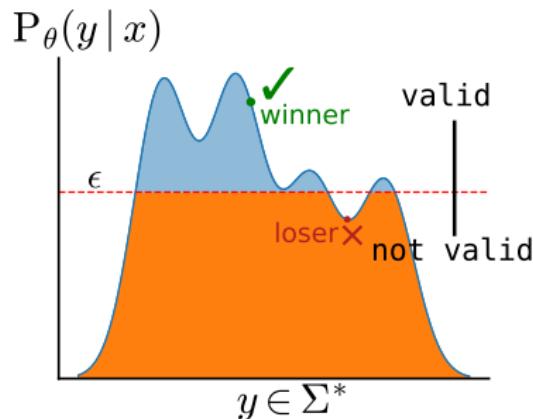
Whenever the model deems the loser to be a valid generation, it should deem the winner to be valid too.

**Running example:** This program and semantics is foundational to many DPO-style losses.

# Uncovering the natural logic of these algorithms

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.

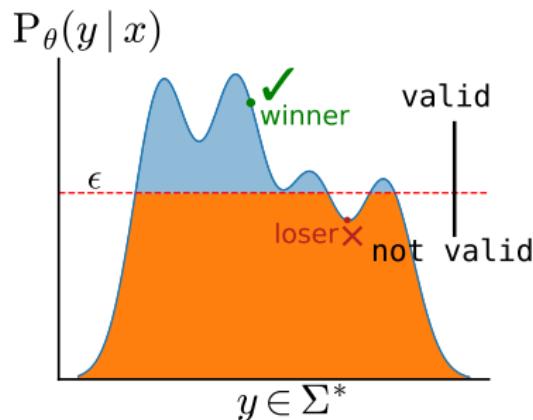


**Model behavior:** Programs tell us about the structure of the model's output distribution (right).

# Uncovering the natural logic of these algorithms

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

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

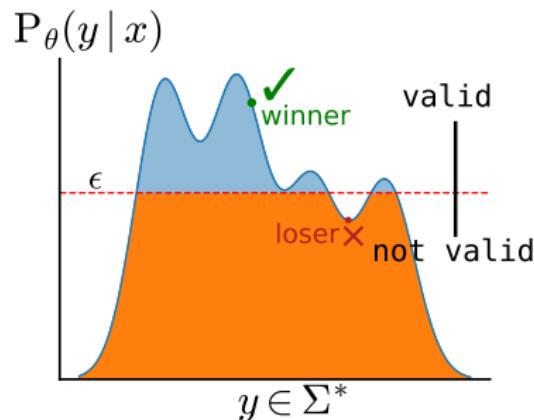


**Model behavior:** Programs tell us about the structure of the model's output distribution (right).

# Uncovering the natural logic of these algorithms

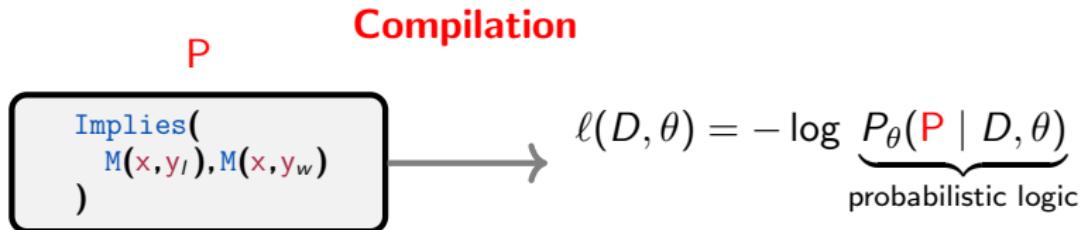
Implies(  
M( $x, y_l$ ), M( $x, y_w$ ))

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



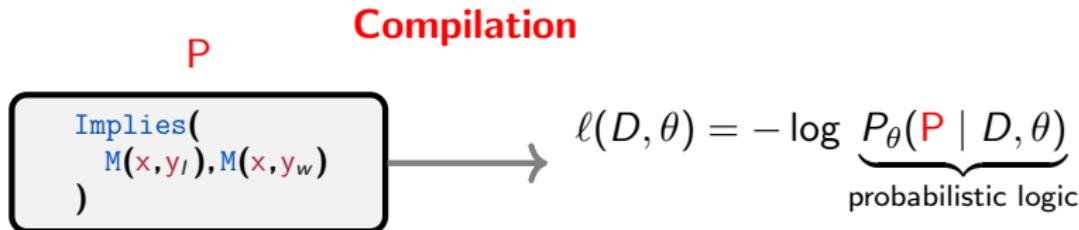
**Observation:** The second program is more strict than the first, involves semantic entailment.

# Compilation and decompilation again



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

# Compilation and decompilation again



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, interpret formulas in that logic to derive differentiable losses.

# Compilation and decompilation again

**Decompilation**

P

Implies(

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

)

$\ell_{\text{CP0}} = - \log \sigma \left( \log \frac{\pi_\theta(y_w|x)}{\pi_\theta(y_l|x)} \right)$

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

# Compilation and decompilation again

**Decompilation**

**P**

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

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

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

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

# Compilation and decompilation again

**Decompilation**

P

Implies(

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

)

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

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

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

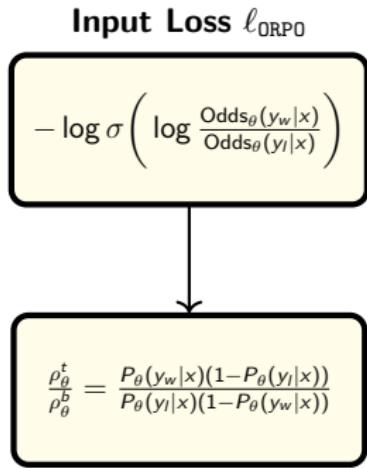
**The second thing we did:** Defined a mechanical procedure for decompilation, proved its correctness, invariance to choice of  $f$ .

## Illustration of approach and results (Richardson et al., 2025)

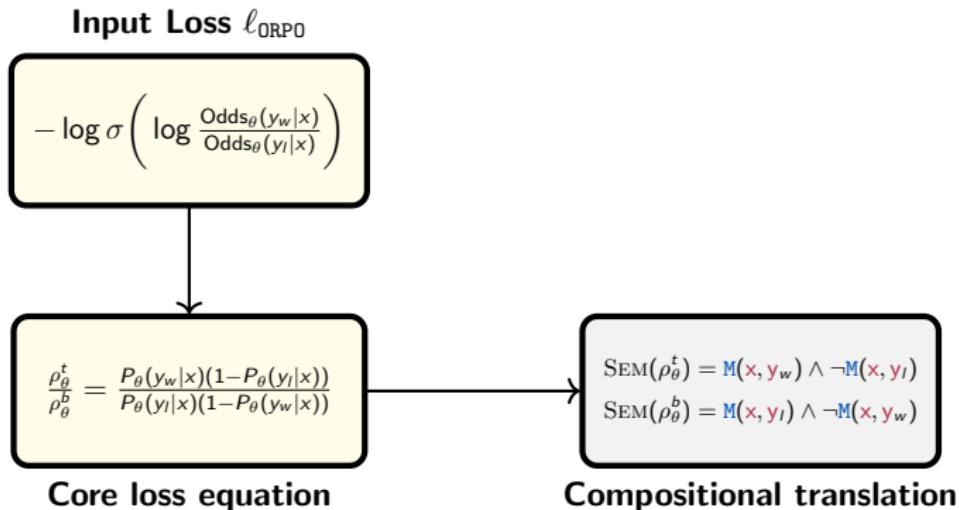
**Input Loss**  $\ell_{\text{ORPO}}$

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

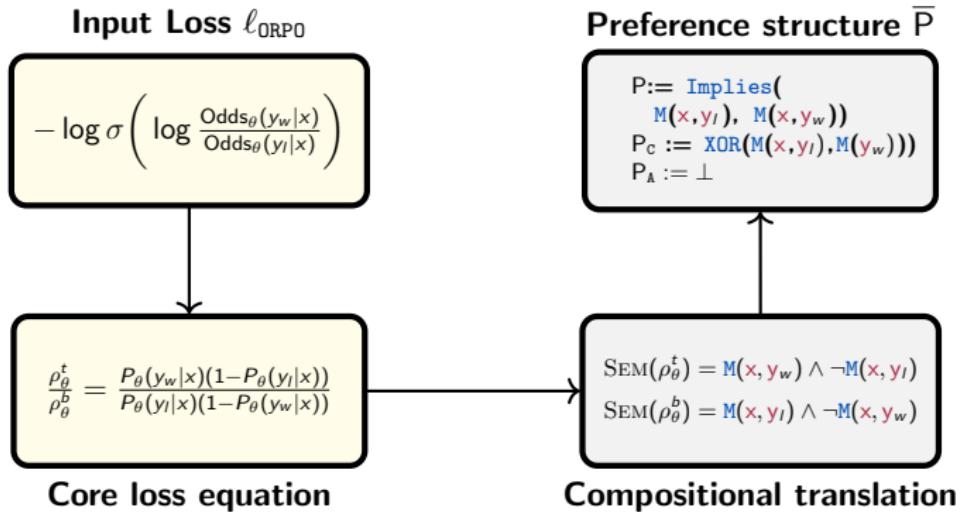
# Illustration of approach and results (Richardson et al., 2025)



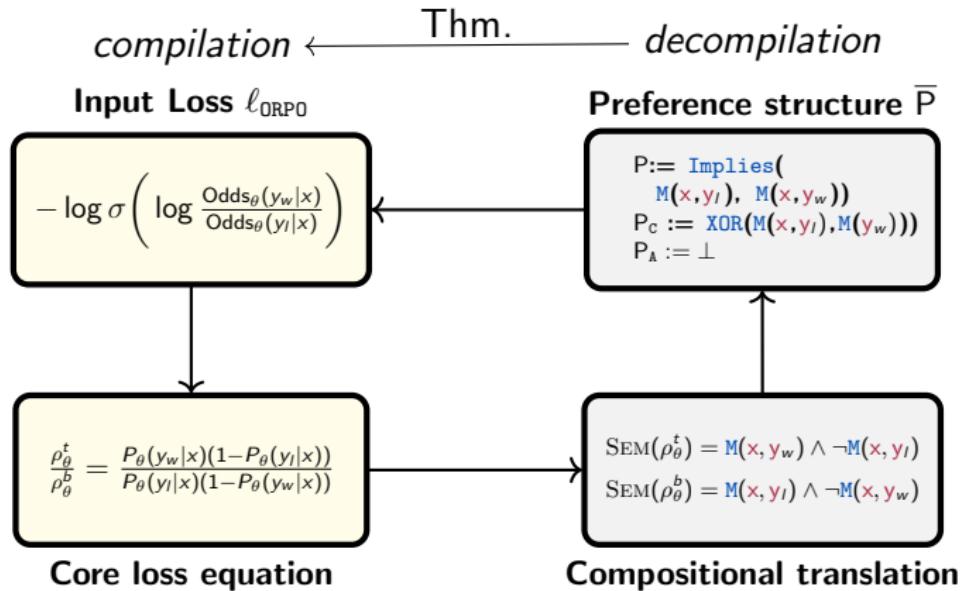
# Illustration of approach and results (Richardson et al., 2025)



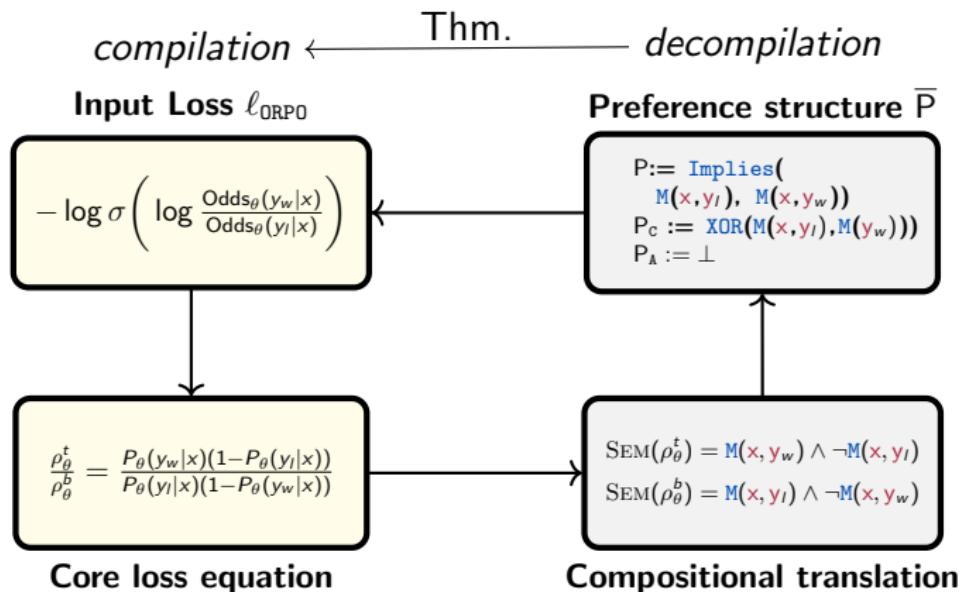
# Illustration of approach and results (Richardson et al., 2025)



# Illustration of approach and results (Richardson et al., 2025)

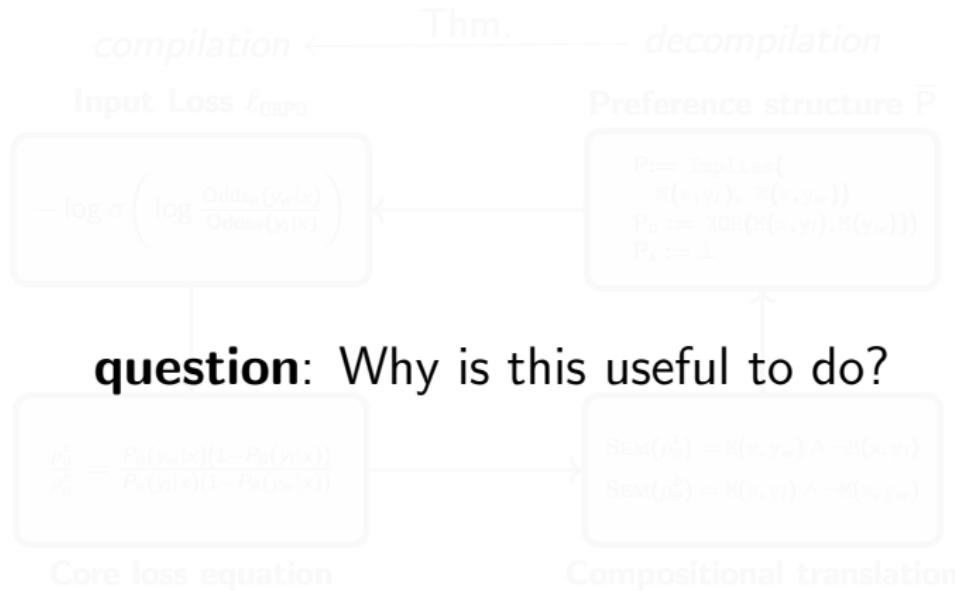


## Illustration of approach and results (Richardson et al., 2025)



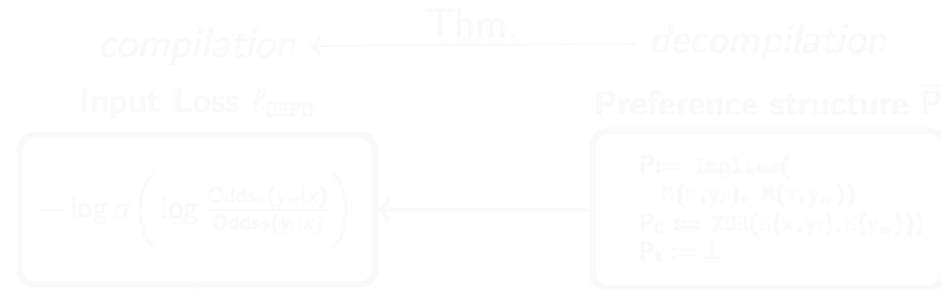
- ▶ **Preference structure**, a core construct in our logic, encoding for preference losses, has a natural Boolean interpretation.

# Illustration of approach and results (Richardson et al., 2025)



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

# Illustration of approach and results (Richardson et al., 2025)



How many preference loss functions are there?



Core loss equation

Compositional translation

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

# Why is this useful? understanding the space

$P^{(1)}$

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

$P^{(2)}$

And(  
     $M(x, y_w),$   
    Not( $M(x, y_I)$ )))

Boolean functions, 2 variables

$M(x, y_w)$	$M(x, y_I)$	$P^{(1)}$	$P^{(2)}$
T	T	✓	X
T	F	✓	✓
F	T	X	X
F	F	✓	X

# Why is this useful? understanding the space

Boolean functions, 2 variables

$M(x, y_w)$	$M(x, y_I)$	$P^{(1)}$	$P^{(2)}$
T	T	✓	X
T	F	✓	✓
F	T	X	X
F	F	✓	X

$P^{(1)}$

$P^{(2)}$

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

$\text{And}(\text{M}(x, y_w), \text{Not}(\text{M}(x, y_I)))$

- 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

Boolean functions, 2 variables

$P^{(1)}$	$M(x, y_w)$	$M(x, y_I)$	$P^{(1)}$	$P^{(2)}$
T	T	T	✓	X
T	F	F	✓	✓
F	T	F	X	X
F	F	F	✓	X

$P^{(2)}$

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

$\text{And}(\text{M}(x, y_w), \text{Not}(\text{M}(x, y_I)))$

Loss creation will end up being equivalent to drawing different sets of ✓ s and X (or blank marks) in a truth table.

# Why is this useful? understanding the space

Boolean functions, 2 variables

$M(x, y_w)$	$M(x, y_I)$	$P^{(1)}$	$P^{(2)}$
T	T	✓	X
T	F	✓	✓
F	T	X	X
F	F	✓	X

$P^{(1)}$

$P^{(2)}$

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

$\text{And}(\text{M}(x, y_w), \text{Not}(\text{M}(x, y_I)))$

**no reference: 256 losses**

Loss creation will end up being equivalent to drawing different sets of ✓ s and X (or blank marks) in a truth table.

# Loss functions as truth tables

Implies(  
And( $M(x, y_I)$ ,  $Ref(x, y_w)$ ),  
And( $M(x, y_w)$ ,  $Ref(x, y_I)$ ))  
)

4 variables

$Ref(x, y_w)$	$M(x, y_I)$	$Ref(x, y_I)$	$M(x, y_w)$
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( $M(x, y_I)$ , Ref( $x, y_w$ )),  
And( $M(x, y_w)$ , Ref( $x, y_I$ )))  
)

4 variables

**answer:** loads.

Ref( $x, y_w$ )	$M(x, y_I)$	Ref( $x, y_I$ )	$M(x, y_w)$
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
F	T	T	T
T	F	F	F
T	F	T	F
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(If( $x, y_I$ ), Ref( $x, y_R$ )),  
And( $x, y_R$ ))

And

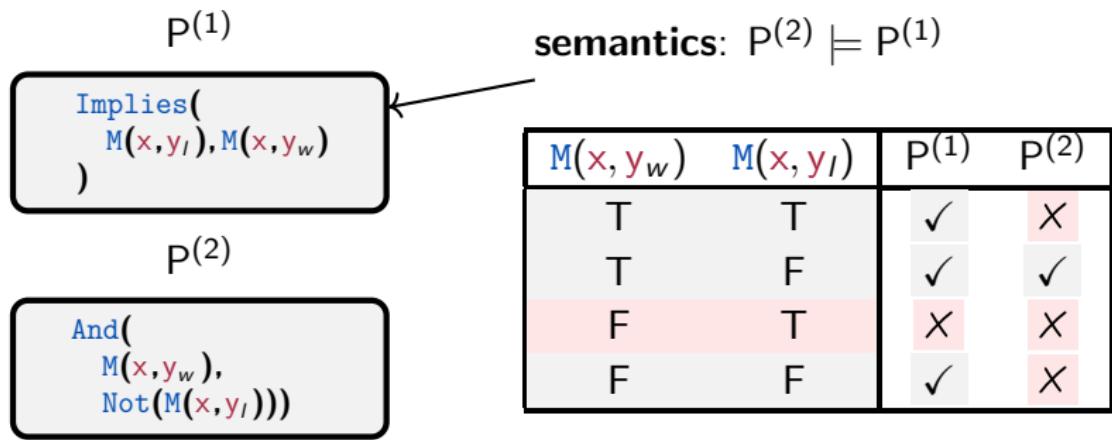
**question:** How are losses related to one another?

4 variables

Ref( $x, y_R$ )	If( $x, y_I$ )	Ref( $x, y_I$ )	If( $x, y_R$ )
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	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

$P^{(1)}$                                     semantics:  $P^{(2)} \models P^{(1)}$

		$M(x, y_w)$	$M(x, y_I)$	$P^{(1)}$	$P^{(2)}$
		T	T	✓	X
		T	F	✓	✓
		F		X	X
		F		✓	X

$P^{(1)}$

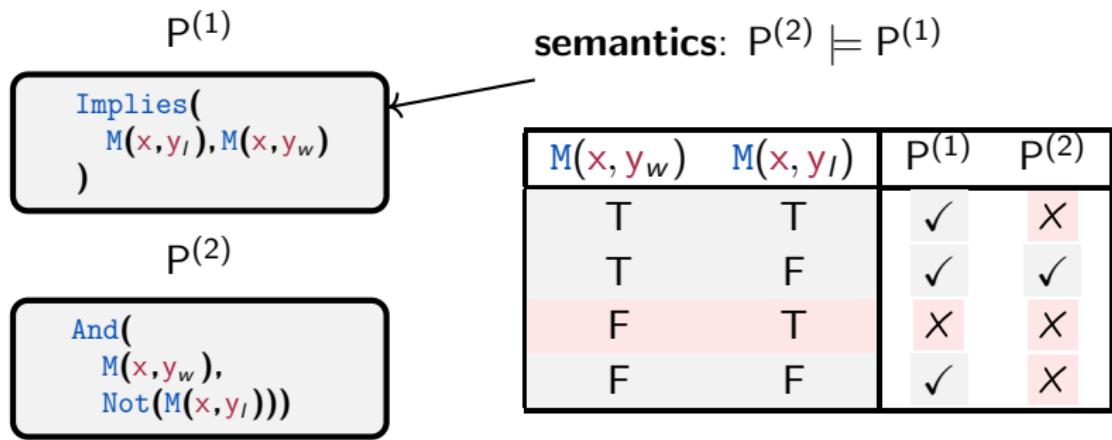
Implies(  
     $M(x, y_I), M(x, y_w)$   
)

$P^{(2)}$

And(  
     $M(x, y_w),$   
    Not( $M(x, y_I)$ ))

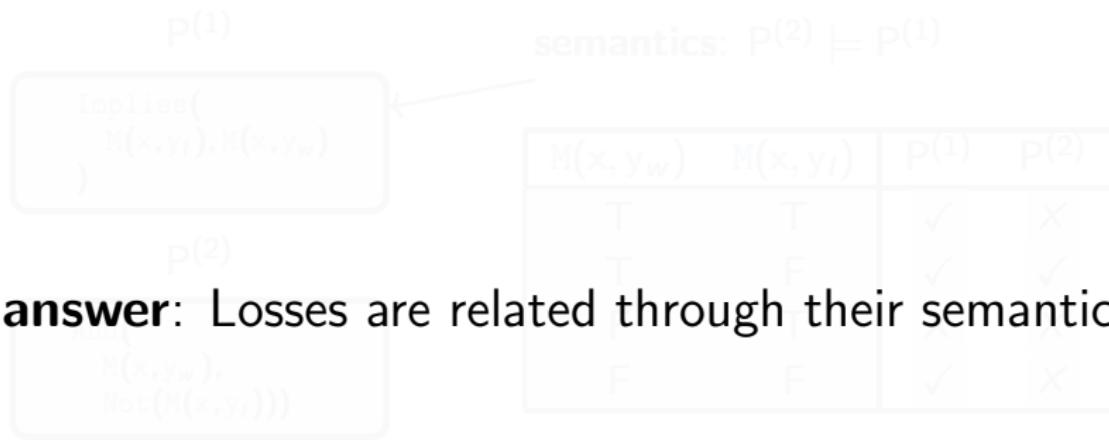
**Proposition** (Xu et al., 2018): Loss is equivalent under semantic equivalence: If  $P^{(2)} \equiv P^{(1)}$  then  $\ell(D, \theta, P^{(2)}) = \ell(D, \theta, P^{(1)})$ .

# Why is this useful? understanding the structure



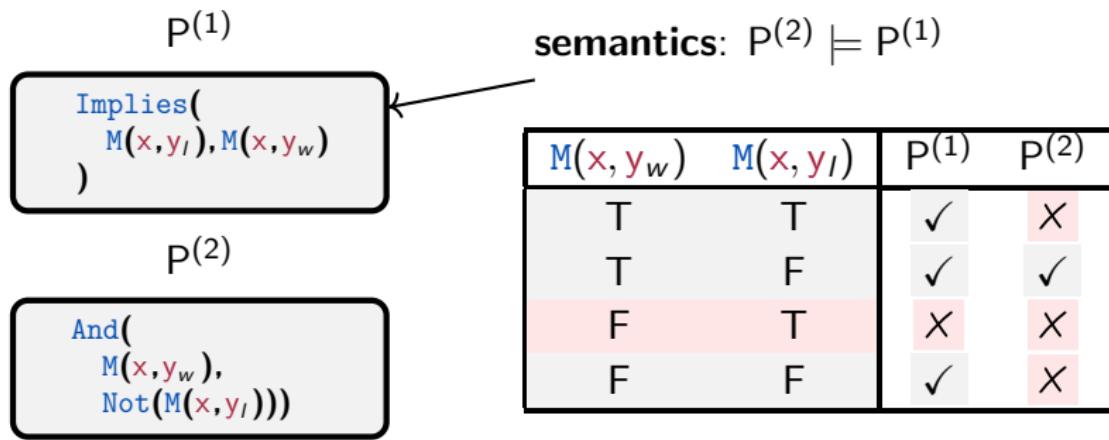
**Theorem:**  $\ell(D, \theta, P^{(2)}) > \ell(D, \theta, P^{(1)})$  (the loss of  $P^{(1)}$  is contained in the loss of  $P^{(2)}$ ).

## 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.

# 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))  
)
```

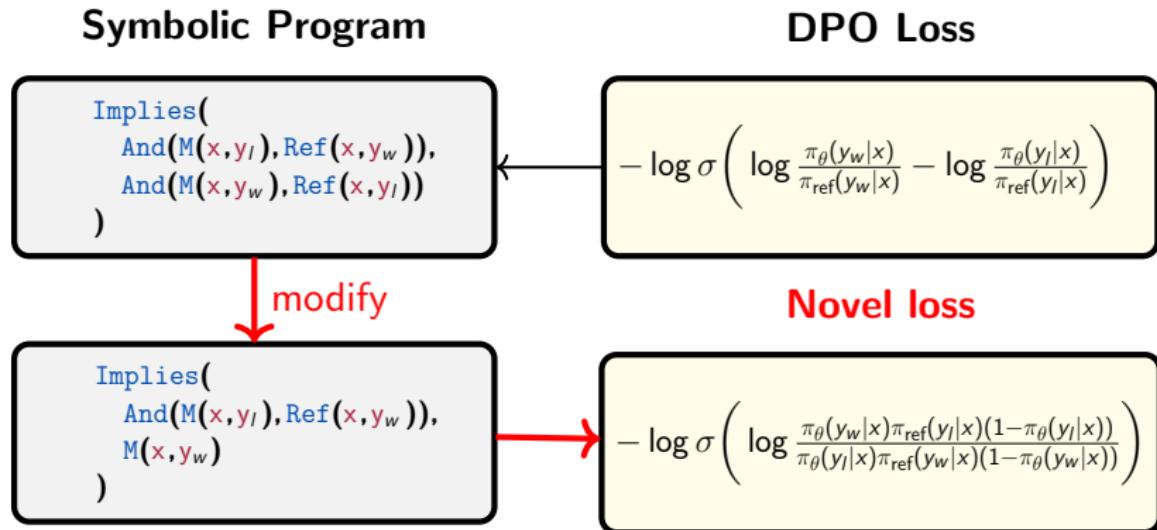
## DPO Loss

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

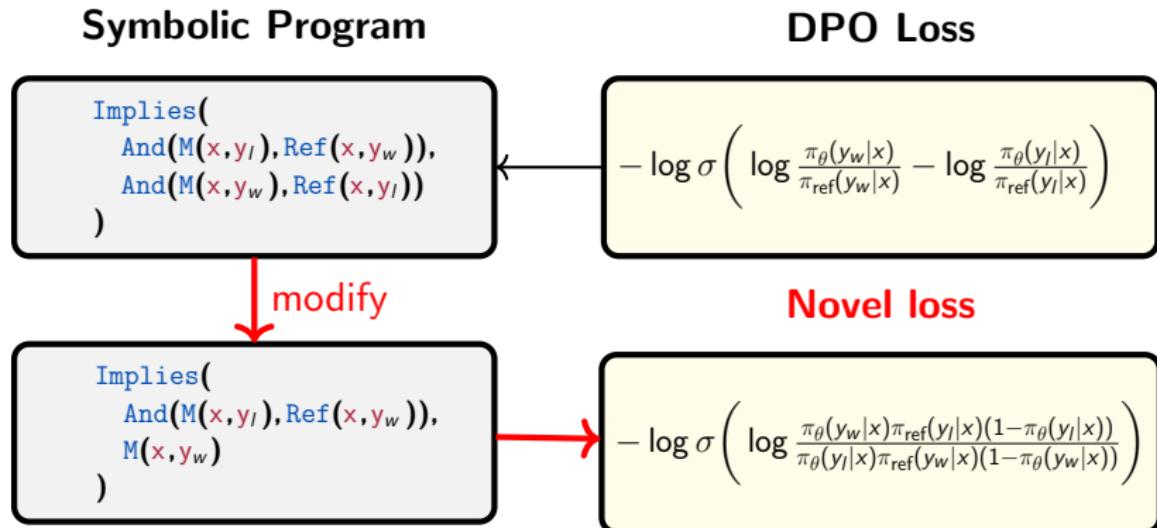
modify  


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

# Deriving new losses symbolically, from first principles



# Deriving new losses symbolically, from first principles



- ▶ High-level programming language for defining new losses.

# Deriving new losses symbolically, from first principles

## Symbolic Program

```
implies(  
    And(Not(x,y_l), Ref(x,y_w)),  
    And(Not(x,y_w), Ref(x,y_l)))  
)
```

## DPO Loss

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

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

```
implies(  
    And(Not(x,y_l), Ref(x,y_w)),  
    Not(x,y_w))  
)
```

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

## How does the logic work?

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

$M(x, y_w)$	$M(x, y_I)$	CPO	ORPO	unCPO	P
T	T	✓ X		✓	
T	F	✓	✓	✓	
F	T	X	X	X	
F	F			✓	

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

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

$$\{ \text{✓}, \text{X} \}_w := \prod_{w \models M(x,y)} \pi_\theta(y | x) \cdot \prod_{w \models \neg M(x,y)} 1 - \pi_\theta(y | x)$$

# How does the logic work?

$M(x, y_w)$	$M(x, y_I)$	CPO	ORPO	unCPO
T	T	✓ X		✓
T	F	✓	✓	✓
F	T	X	X	X
F	F			✓

P  
Implies(  
 $M(x, y_I), M(x, y_w)$   
)

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

$$\{\checkmark, X\}_w := \prod_{w \models M(x,y)} \pi_\theta(y | x) \cdot \prod_{w \models \neg M(x,y)} 1 - \pi_\theta(y | x)$$

- ▶ Formula probability P computed as a weighted count  $\sum \checkmark_w$  (Chavira and Darwiche, 2008), loss is  $-\log$ , *semantic loss* (Xu et al., 2018).

# How does the logic work?

$M(x, y_w)$	$M(x, y_I)$	CPO	ORPO	unCPO
T	T	✓	X	
T	F			✓ ↛

P

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



Available online at [www.sciencedirect.com](http://www.sciencedirect.com)



Artificial Intelligence 172 (2008) 772–799

Artificial  
Intelligence

[www.elsevier.com/locate/artint](http://www.elsevier.com/locate/artint)

## On probabilistic inference by weighted model counting $\star$

Mark Chavira <sup>\*</sup>, Adnan Darwiche algorithmic foundation

*Computer Science Department, UCLA, Los Angeles, CA 90095, USA*

Received 25 August 2006; received in revised form 22 July 2007; accepted 5 November 2007

## A Semantic Loss Function for Deep Learning with Symbolic Knowledge

### Losses computed from weighted model counts

Jingyi Xu<sup>1</sup> Zilu Zhang<sup>2</sup> Tal Friedman<sup>1</sup> Yitao Liang<sup>1</sup> Guy Van den Broeck<sup>1</sup>

#### Abstract

This paper develops a novel methodology for using symbolic knowledge in deep learning. From first principles, we derive a semantic loss function that bridges between neural output vectors and logical constraints. This loss function captures how close the neural network is to satisfying the constraints on its output. An experimental evaluation shows that it effectively guides the learner to achieve (near-)state-of-the-art results on semi-supervised multi-class classification. Moreover, it significantly increases the ability of the neural network to predict structured objects, such as rankings and paths. These discrete concepts are tremendously difficult to learn, and benefit from a tight integration of deep learning and symbolic reasoning methods.

This paper considers learning in domains where we have symbolic knowledge connecting the different outputs of a neural network. This knowledge takes the form of a constraint (or sentence) in Boolean logic. It can be as simple as an exactly-one constraint for one-hot output encodings, or as complex as a structured output prediction constraint for intricate combinatorial objects such as rankings, subgraphs, or paths. Our goal is to augment neural networks with the ability to learn how to make predictions subject to these constraints, and use the symbolic knowledge to improve the learning performance.

Most neuro-symbolic approaches aim to simulate or learn symbolic reasoning in an end-to-end deep neural network, or capture symbolic knowledge in a vector-space embedding. This choice is partly motivated by the need for smooth *differentiable* models; adding symbolic reasoning code (e.g., SAT solvers) to a deep learning pipeline de-

# How does the logic work?

$M(x, y_w)$	$M(x, y_I)$	CPO	ORPO	unCPO
T	T	✓ X		✓
T	F	✓	✓	✓
F	T	X	X	X
F	F			✓

P  
 Implies(  
 $M(x, y_I), M(x, y_w)$   
 )

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

$$\{ \text{✓}, \text{X} \}_w := \prod_{w \models M(x,y)} \pi_\theta(y | x) \cdot \prod_{w \models \neg M(x,y)} 1 - \pi_\theta(y | x)$$

$$\underbrace{\ell_x}_{\text{column}} := \underbrace{-\log \sigma \left( \log \frac{\sum \text{✓}_w}{\sum \text{X}_w} \right)}_{\text{log ratio of } \text{✓}_w \text{ and } \text{X}_w \text{ counts}}$$

# How does the logic work?

$M(x, y_w)$	$M(x, y_l)$	CPO	ORPO	unCPO	P
T	T	✓ X		✓	Implies( $M(x, y_l), M(x, y_w)$ )
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.

$$\{ \text{✓}, \text{X} \}_w := \prod_{w \models M(x,y)} \pi_\theta(y | x) \cdot \prod_{w \models \neg M(x,y)} 1 - \pi_\theta(y | x)$$

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

# How does the logic work?

$M(x, y_w)$	$M(x, y_l)$	CPO	ORPO	unCPO	P
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.

$$\{ \text{✓}, \text{X} \}_w := \prod_{w \models M(x,y)} \pi_\theta(y | x) \cdot \prod_{w \models \neg M(x,y)} 1 - \pi_\theta(y | x)$$

$$\underbrace{\ell_x}_{\text{column}} := -\log \sigma \left( \log \frac{\sum_w \text{✓}_w}{\sum_w \text{X}_w} \right)$$

$$= \underbrace{-\log \sigma \left( \log \frac{\pi_\theta(y_w | x)}{\pi_\theta(y_l | x)} \right)}_{\ell_{\text{CPO}}, \sim P_\theta(P | \text{one true})}$$

# How does the logic work?

$H(x, y_w)$	$H(x, y_I)$	CPO	ORPO	unCPO	P
T	T	✓ X		✓	Implies( $H(x, y_I), H(x, y_w)$ )
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.

{ ✓ **observation:** losses differ in hard constraints

$$\ell_{\text{CPO}} = -\log \sigma(\pi_\theta(y_w | x))$$

$$\ell_x := -\log \sigma \left( \log \frac{\sum_y \pi_\theta(y_w)}{\sum_y \pi_\theta(y_I)} \right)$$

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

$\ell_{\text{CPO}}, \sim P_\theta(\text{P lone true})$

# How does the logic work?

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

$$\{ \text{✓}, \text{X} \}_w := \prod_{w \models M(x,y)} \pi_\theta(y | x) \cdot \prod_{w \models \neg M(x,y)} 1 - \pi_\theta(y | x)$$

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

# How does the logic work?

$M(x, y_w)$	$M(x, y_I)$	CPO	ORPO	unCPO
T	T	✓ X		✓
T	F	✓	✓	✓
F	T	X	X	X
F	F			✓

P  
Implies(  
 $M(x, y_I)$ ,  $M(x, y_w)$ )

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

$$\{ \text{✓}, \text{X} \}_w := \prod_{w \models M(x,y)} \pi_\theta(y | x) \cdot \prod_{w \models \neg M(x,y)} 1 - \pi_\theta(y | x)$$

- ▶ **Preference structure:** equivalent way of expressing truth table representations (Richardson et al., 2025),

$$\bar{P} := \left( \underbrace{P}_{\text{core}}, \underbrace{P_c, P_A}_{\text{constraints}} \right)$$

# How does the logic work?

$M(x, y_w)$	$M(x, y_I)$	CPO	ORPO	unCPO	P
T	T	✓ X		✓	Implies( $M(x, y_I), M(x, y_w)$ )
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.

$$\{ \text{✓}, \text{X} \}_w := \prod_{w \models M(x,y)} \pi_\theta(y | x) \cdot \prod_{w \models \neg M(x,y)} 1 - \pi_\theta(y | x)$$

$$\begin{aligned}
 \underbrace{\ell_x}_{\text{column}} &:= -\log \sigma \left( \log \frac{\sum_w \text{✓}_w}{\sum_w \text{X}_w} \right) \\
 &= -\log \sigma \left( \log \underbrace{\frac{\pi_\theta(y_I | x) \pi_\theta(y_w | x) + (1 - \pi_\theta(y_I | x))}{\pi_\theta(y_I | x) (1 - \pi_\theta(y_w | x))}}_{\text{novel loss without constraints, } P_\theta(P | \top)} \right)
 \end{aligned}$$

# How does the logic work?

$H(x, y_w)$	$H(x, y_l)$	CPO	ORPO	unCPO	P
T	T	✓	X	✓	Implies( $H(x, y_l), H(x, y_w)$ )
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.

**observation:** real losses are highly constrained

$$\begin{aligned}
 \ell_x &:= -\log \sigma \left( \log \frac{\sum y_w}{\sum X_w} \right) \\
 &= -\log \sigma \left( \log \underbrace{\frac{\pi_\theta(y_l | x)\pi_\theta(y_w | x) + (1 - \pi_\theta(y_l | x))}{\pi_\theta(y_l | x)(1 - \pi_\theta(y_w | x))}}_{\text{novel loss without constraints, } P_\theta(P|T)} \right)
 \end{aligned}$$

# How does the logic work?

$M(x, y_w)$	$M(x, y_I)$	CPO	ORPO	unCPO	P
T	T	✓ X			
T	F	✓	✓		
F	T	X	X		
F	F			✓ ✓ ✓ X ✓	Implies( $M(x, y_I), M(x, y_w)$ )

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

**note:**  $\underbrace{M(x, y_I) \rightarrow M(x, y_w)}_{\text{goal of optimization}} \equiv \neg M(x, y_I) \vee M(x, y_w)$

$$\ell_x := -\log \sigma \left( \log \frac{\sum \checkmark_w}{\sum X_w} \right)$$

column

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

novel loss without constraints,  $P_\theta(P|T)$

# How does the logic work?

$M(x, y_w)$	$M(x, y_I)$	CPO	ORPO	unCPO	P
T	T	✓ X			Implies( $M(x, y_I), M(x, y_w)$ )
T	F	✓	✓		
F	T	X	X		
F	F			✓ ✓ X ✓	Whenever the model deems the loser to be a valid generation, it should deem the

**note:**  $M(x, y_I) \rightarrow M(x, y_w) \equiv \underbrace{\neg M(x, y_I) \vee M(x, y_w)}_{\text{drive probability to zero}}$

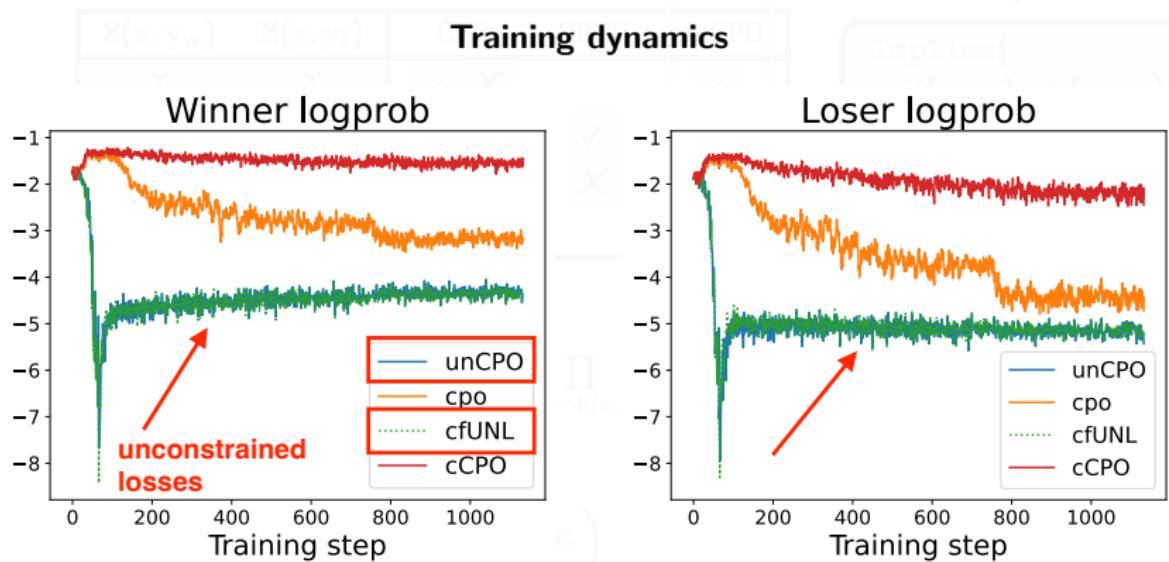
$$\{\checkmark, X\}_w := \prod_{w \vdash \mathbb{M}(x, y)} \pi_\theta(y | x) \cdot \prod_{w \not\models \mathbb{M}(x, y)} 1 - \pi_\theta(y | x)$$

$$\ell_x := \underbrace{-\log \sigma \left( \log \frac{\sum \checkmark_w}{\sum X_w} \right)}_{\text{column}}$$

$$= -\log \sigma \left( \log \underbrace{\frac{\pi_\theta(y_I | x) \pi_\theta(y_w | x) + (1 - \pi_\theta(y_I | x))}{\pi_\theta(y_I | x) (1 - \pi_\theta(y_w | x))}}_{\text{novel loss without constraints, } P_\theta(\mathbb{P} | \mathbb{T})} \right)$$

novel loss without constraints,  $P_\theta(\mathbb{P} | \mathbb{T})$

# How does the logic work?



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

novel loss without constraints,  $P_\theta(P|T)$

# How does the logic work?

$H(x, y_w)$	$H(x, y_l)$	CPO	ORPO	unCPO	P
T	T	✓	X	✓	Implies( $H(x, y_l), H(x, y_w)$ )
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.

{✓, X} Constrainedness is an important property

$$\pi_\theta(y_l | x) = \pi_\theta(y_w | x)$$

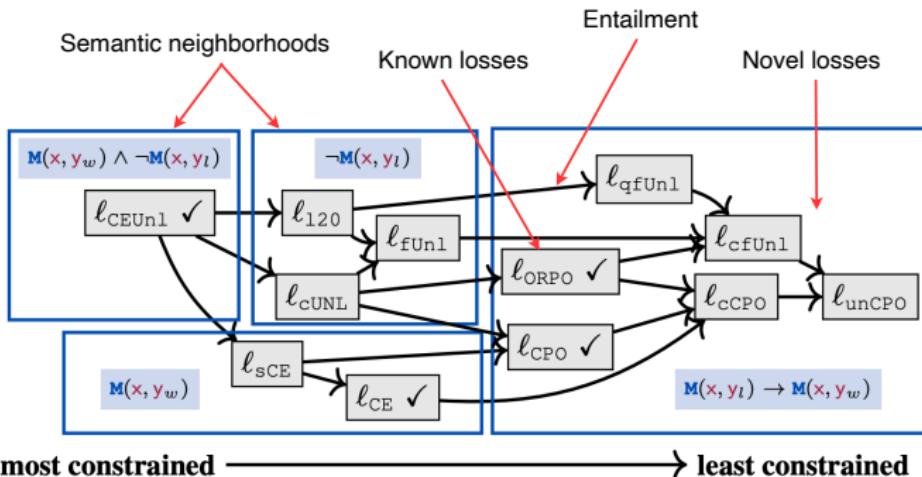
$$\ell_x := -\log \sigma \left( \log \frac{\sum y_w}{\sum X_w} \right)$$

column

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

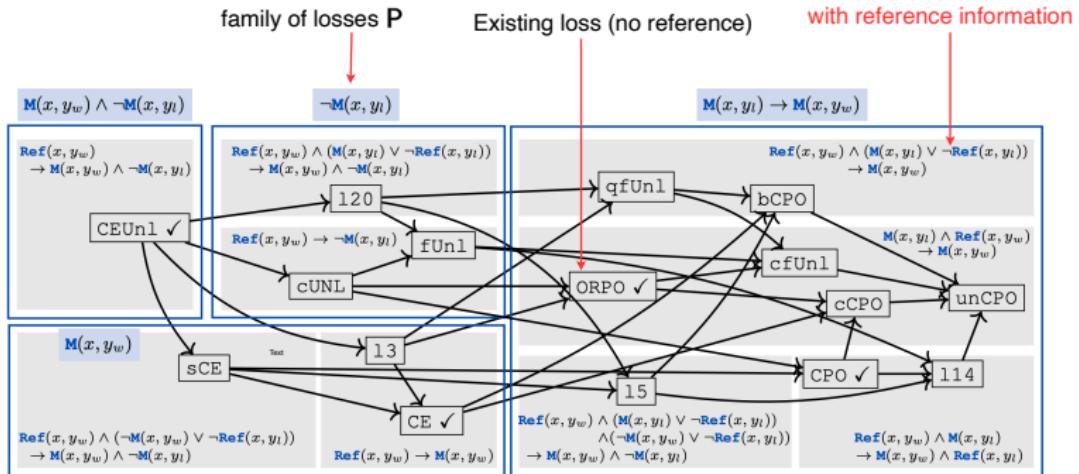
novel loss without constraints,  $P_\theta(P|T)$

# The no reference loss landscape



- ▶ **Loss lattice:** semantic structure of loss space, ordering.

# The reference loss landscape



- ▶ The semantics of DPO-style reference losses can be straightforwardly computed from no reference approaches, much less explored.

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

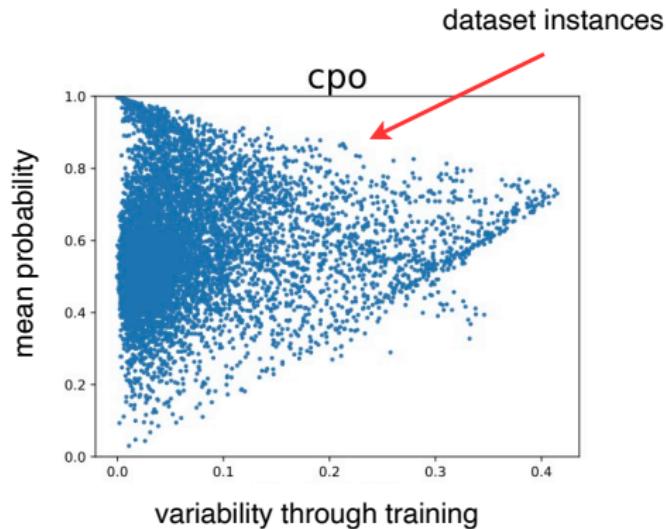
# Characterizing dataset semantics?

loss	WR% ( $\ell_{\text{cpo}}$ )	evol	false-qa	flan	sharegpt	ultrachat
$\ell_{\text{CFUNL}}$	46.1 ( $\pm 0.4$ )	46.1 ( $\pm 2.2$ )	51.6 ( $\pm 2.9$ )	46.4 ( $\pm 1.7$ )	46.2 ( $\pm 1.2$ )	44.1 ( $\pm 1.0$ )
$\ell_{\text{QFUNL}}$	48.9 ( $\pm 0.8$ )	45.3 ( $\pm 1.9$ )	34.7 ( $\pm 6.3$ )	57.9 ( $\pm 1.2$ )	46.8 ( $\pm 2.4$ )	41.3 ( $\pm 1.4$ )
$\ell_{\text{CCPO}}$	52.0 ( $\pm 0.6$ )	50.7 ( $\pm 0.5$ )	50.2 ( $\pm 0.7$ )	57.2 ( $\pm 1.1$ )	47.2 ( $\pm 1.8$ )	53.1 ( $\pm 1.9$ )
$\ell_{\text{uncPO}}$	46.0 ( $\pm 0.2$ )	45.8 ( $\pm 0.3$ )	52.1 ( $\pm 3.0$ )	45.7 ( $\pm 0.6$ )	46.2 ( $\pm 2.1$ )	44.8 ( $\pm 2.1$ )

Table 5: Comparing performance of Qwen-0.5B tuned on new losses (rows) against  $\ell_{\text{CPO}}$  based on aggregate win-rate (WR % (std)) on ultrafeedback test (second column) and different test subsets (columns 2-6).

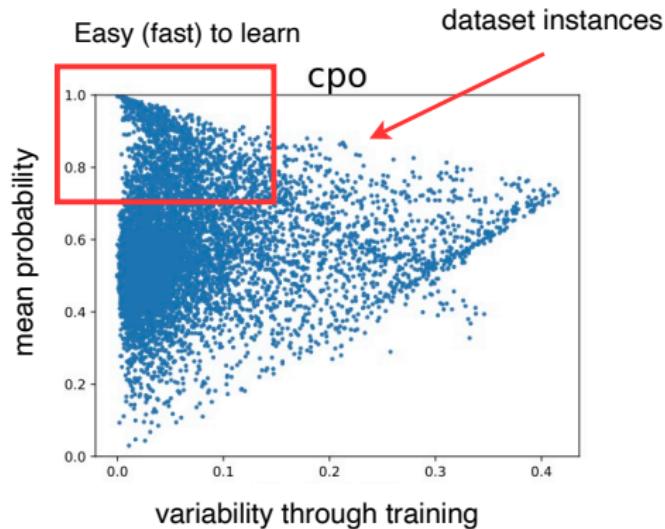
- ▶ **Finding:** Different losses perform better/worse on different subsets of data, reflecting the different semantics in preference data.

# Characterizing dataset semantics?



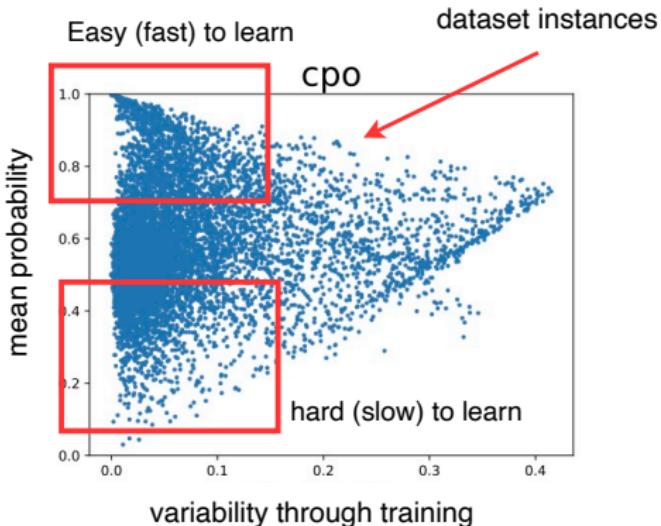
- ▶ Tracking instance-level training dynamics ([Swayamdipta et al., 2020](#)) and behavior across losses, use to reverse engineer data semantics.

# Characterizing dataset semantics?



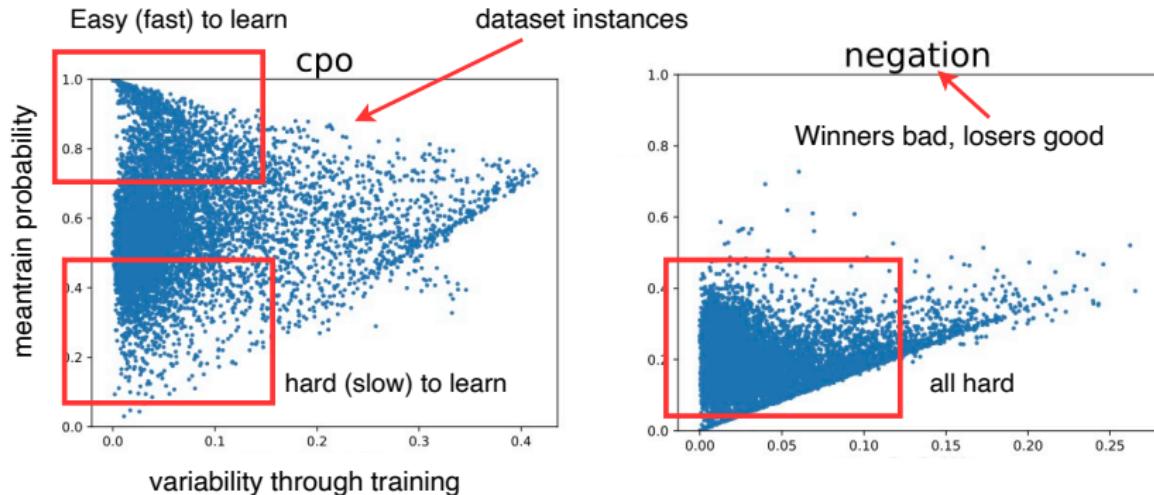
- ▶ **Intuition:** The speed/ease of training is a proxy for *goodness of semantic fit*, similar to small-loss criterion in noisy-label learning (Gui et al., 2021).

# Characterizing dataset semantics?



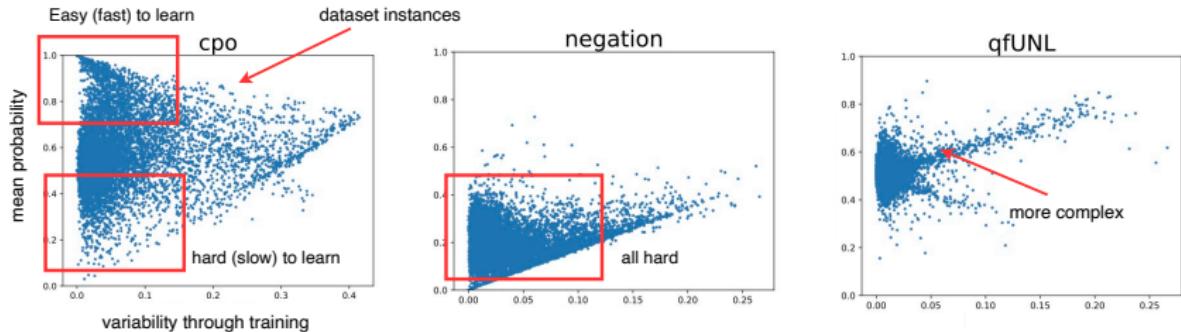
- ▶ **Intuition:** The speed/ease of training is a proxy for *goodness of semantic fit*, similar to small-loss criterion in noisy-label learning (Gui et al., 2021).

# Characterizing dataset semantics?



- ▶ **Intuition:** The speed/ease of training is a proxy for *goodness of semantic fit*, similar to small-loss criterion in noisy-label learning (Gui et al., 2021).

# Characterizing dataset semantics?



- ▶ **Intuition:** The speed/ease of training is a proxy for *goodness of semantic fit*, similar to small-loss criterion in noisy-label learning (Gui et al., 2021).

## Characterizing dataset semantics?

- ▶ Intuition: The speed/ease of training is a proxy for goodness of semantic

Blueprint for much empirical exploration of loss space

# Conclusions

- ▶ New ideas about using symbolic techniques to formally characterize the semantics of LLM algorithms, preference learning.

# Conclusions

- ▶ New ideas about using symbolic techniques to formally characterize the semantics of LLM algorithms, preference learning.
  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 using symbolic techniques to formally characterize the semantics of LLM algorithms, preference learning.
  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).

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

# Conclusions

- ▶ New ideas about using symbolic techniques to formally characterize the semantics of LLM algorithms, preference learning.
  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).

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

- ▶ **Decompiling models to symbolic programs:** *semantics of data, reinforcement learning, chain-of-thought, LLM agents ...*

Thank you.

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

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.**

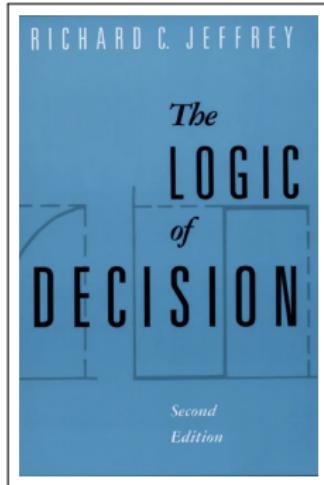
- ▶ **Peculiar semantics**, but the logic makes sense, e.g., we want to maximize

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

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

# Classical work on preference

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



(Jeffrey, 1965)

Preference Principle	THE STATUS OF VARIOUS PREFERENCE PRINCIPLES					
	Von Wright	Chisholm Sosa	Martin	$P^\#$	$P^*$	$P^w$
1. $\hat{P}Pq \rightarrow \sim(\hat{q}P\hat{p})$	✓	✓	✓	+	+	+
2. $(\hat{P}Pq \& \hat{q}Pr) \rightarrow \hat{P}Pr$	✓	✓	✓	+	+	+
3. $\hat{P}Pq \rightarrow \sim qP \sim \hat{p}$		x	✓	(+) <sup>1</sup>	+	+
4. $\sim qP \sim \hat{p} \rightarrow \hat{P}Pq$		x	✓	(+) <sup>1</sup>	+	+
5. $\hat{P}Pq \rightarrow (\hat{p} \& \sim \hat{q}) P(\sim \hat{p} \& \hat{q})$		✓	✓	+	+	+
6. $(\hat{p} \& \sim \hat{q}) P(\sim \hat{p} \& \hat{q}) \rightarrow \hat{P}Pq$		✓	x	+	+	+
7. $[\sim(\hat{p}P \sim \hat{p}) \& \sim(\sim \hat{p}P)] \& \sim(\hat{q}P \sim \hat{q}) \& \sim(\sim \hat{q}P \sim \hat{q}) \rightarrow [\sim(\hat{p}Pq) \& \sim(\hat{q}P\hat{p})]$		✓	✓	+	+	+
8. $[\sim(\hat{q}P \sim \hat{q}) \& \sim(\sim \hat{q}Pq) \& \hat{p}Pq] \rightarrow \hat{P}P \sim \hat{p}$		✓	✓	+	+	-
9. $[\sim(\hat{q}P \sim \hat{q}) \& \sim(\sim \hat{q}Pq) \& \hat{q}P \sim \hat{p}] \rightarrow \hat{P}P \sim \hat{p}$		✓	+	+	-	-
10. $\hat{P}Pq \rightarrow [(\hat{p} \& r) P(q \& r) \& (\hat{p} \& \sim r)]$		✓		-	-	+
11. $[(\hat{p} \& r) P(q \& r) \& (\hat{p} \& \sim r) P(q \& \sim r)] \rightarrow \hat{P}Pq$		✓		(+) <sup>2</sup>	(+) <sup>3</sup>	+
12. $[\sim(\hat{p}Pq) \& \sim(\hat{q}Pr)] \rightarrow \sim(\hat{p}Pr)$			✓	+	+	-
13. $(\hat{p}Pr \vee \hat{q}Pr) \rightarrow (\hat{p} \vee \hat{q}) Pr$			✓	-	-	-
14. $(\hat{p} \vee \hat{q}) Pr \rightarrow [\hat{p}Pr \& \hat{q}Pr]$	✓			-	-	-
15. $[\hat{p}Pr \& \hat{q}Pr] \rightarrow (\hat{p} \vee \hat{q}) Pr$	✓			-	-	-
16. $(\hat{p} \vee \hat{q}) Pr \rightarrow (\hat{p}Pr \vee \hat{q}Pr)$				-	-	-
17. $\hat{P}P(g \vee r) \rightarrow (\hat{P}Pq \& \hat{P}Pr)$		✓		-	-	-
18. $(\hat{p}Pq \& \hat{p}Pr) \rightarrow \hat{P}P(g \vee r)$				-	-	-
19. $(\hat{p}Pr \& \hat{q}Pr) \rightarrow (\hat{p} \& \hat{q}) Pr$				-	-	-

Semantic foundations for the logic of preference Rescher (1967)

# References |

- Beurer-Kellner, L., Fischer, M., and Vechev, M. (2023). Prompting is programming: A query language for large language models. *Proceedings of the ACM on Programming Languages*, 7(PLDI):1946–1969.
- Bogin, B., Yang, K., Gupta, S., Richardson, K., Bransom, E., Clark, P., Sabharwal, A., and Khot, T. (2024). Super: Evaluating agents on setting up and executing tasks from research repositories. *Proceedings of EMNLP*.
- Bragg, J., D'Arcy, M., Balepur, N., Bareket, D., Dalvi, B., Feldman, S., Haddad, D., Hwang, J. D., Jansen, P., Kishore, V., Majumder, B. P., Naik, A., Rahamimov, S., Richardson, K., Singh, A., Surana, H., Tiktinsky, A., Vasu, R., Wiener, G., Anastasiades, C., Candra, S., Dunkelberger, J., Emery, D., Evans, R., Hamada, M., Huff, R., Kinney, R., Latzke, M., Lochner, J., Lozano-Aguilera, R., Nguyen, C., Rao, S., Tanaka, A., Vlahos, B., Clark, P., Downey, D., Goldberg, Y., Sabharwal, A., and Weld, D. S. (2025). Astabench: Rigorous benchmarking of ai agents with a scientific research suite.
- Chavira, M. and Darwiche, A. (2008). On probabilistic inference by weighted model counting. *Artificial Intelligence*, 172(6-7):772–799.
- Chen, J., Yuan, S., Ye, R., Majumder, B. P., and Richardson, K. (2023). Put your money where your mouth is: Evaluating strategic planning and execution of llm agents in an auction arena. *arXiv preprint arXiv:2310.05746*.
- Cheng, J., Clark, P., and Richardson, K. (2025). Language modeling by language models. *Proceedings of Neurips*.

## References II

- 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.
- Gui, X.-J., Wang, W., and Tian, Z.-H. (2021). Towards understanding deep learning from noisy labels with small-loss criterion. *arXiv preprint arXiv:2106.09291*.
- 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.
- Li, Z., Huang, J., and Naik, M. (2023). Scallop: A language for neurosymbolic programming. *Proceedings of the ACM on Programming Languages*, 7(PLDI):1463–1487.
- Manhaeve, R., Dumancic, S., Kimmig, A., Demeester, T., and De Raedt, L. (2018). Deepproblog: Neural probabilistic logic programming. *Advances in neural information processing systems*, 31.
- Marra, G., Dumančić, S., Manhaeve, R., and De Raedt, L. (2024). From statistical relational to neurosymbolic artificial intelligence: A survey. *Artificial Intelligence*, page 104062.

## References III

- 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.
- Rescher, N. (1967). *The logic of decision and action*. University of Pittsburgh Pre.
- Richardson, K., Srikumar, V., and Sabharwal, A. (2025). Understanding the logic of direct preference alignment through logic. *Proceedings of ICML*.
- Richardson, K. and Wijnholds, G. (2025). Lectures on language model programming.
- Swayamdipta, S., Schwartz, R., Lourie, N., Wang, Y., Hajishirzi, H., Smith, N. A., and Choi, Y. (2020). Dataset cartography: Mapping and diagnosing datasets with training dynamics. *arXiv preprint arXiv:2009.10795*.
- Weiss, G., Goldberg, Y., and Yahav, E. (2021). Thinking like transformers. In *International Conference on Machine Learning*, pages 11080–11090. PMLR.
- 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*.
- Yang, A., Chiang, D., and Angluin, D. (2024). Masked hard-attention transformers recognize exactly the star-free languages. *Advances in Neural Information Processing Systems*, 37:10202–10235.

## References IV

- Yang, R., Chen, J., Zhang, Y., Yuan, S., Chen, A., Richardson, K., Xiao, Y., and Yang, D. (2025). Selfgoal: Your language agents already know how to achieve high-level goals. *Proceedings of NAACL*.
- Yu, H., Xuan, K., Li, F., Zhu, K., Lei, Z., Zhang, J., Qi, Z., Richardson, K., and You, J. (2025). Tinyscientist: An interactive, extensible, and controllable framework for building research agents. *Proceedings of EMNLP*.
- Zhang, Y., Yuan, S., Hu, C., Richardson, K., Xiao, Y., and Chen, J. (2024). Timearena: Shaping efficient multitasking language agents in a time-aware simulation. *Proceedings of ACL*.