

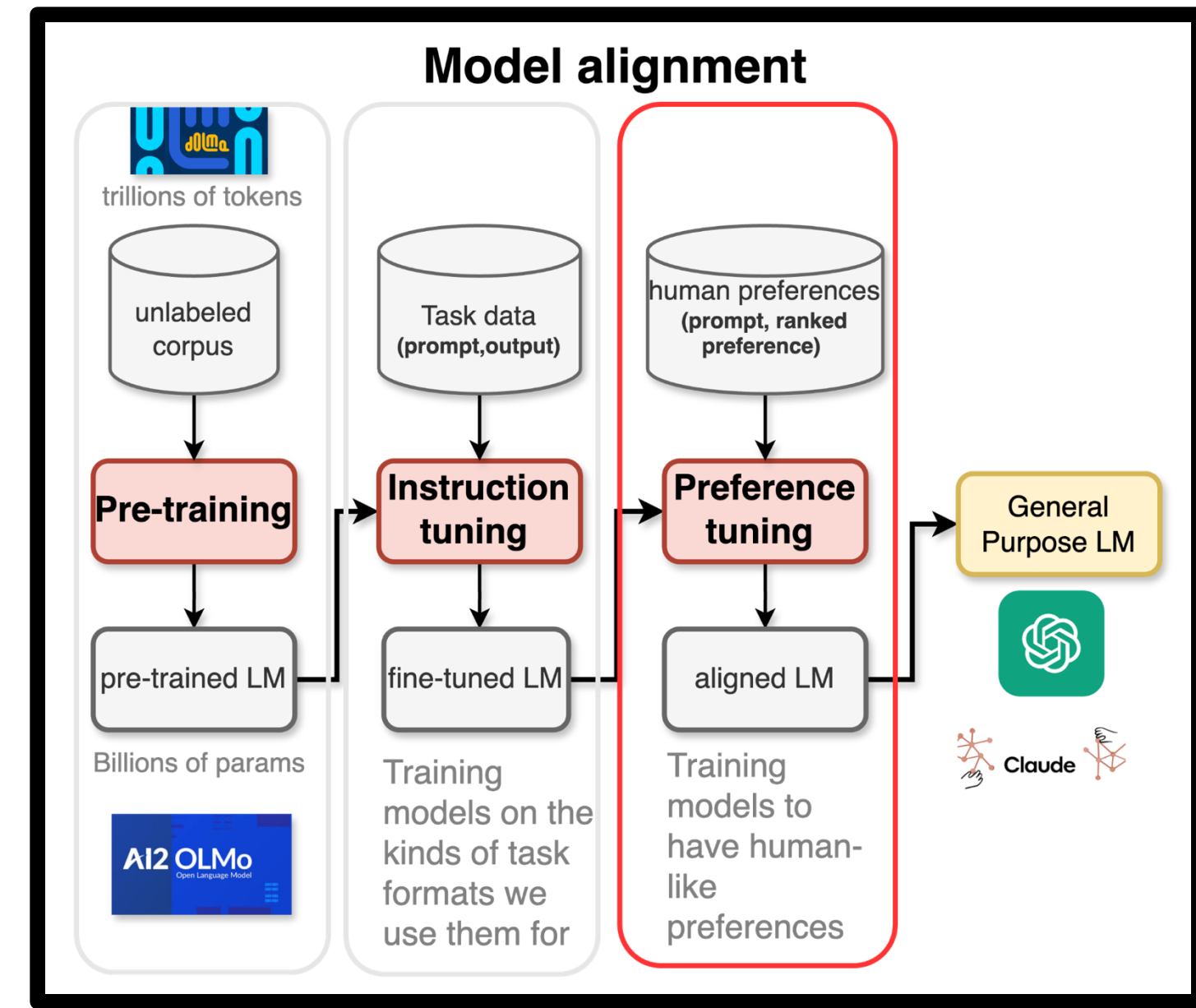
Understanding the Logic of Direct Preference Alignment through Logic



Kyle Richardson, Vivek Srikumar, Ashish Sabharwal
Allen Institute for Artificial Intelligence, University of Utah



Preference alignment for large language models (LLMs)



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

LLM input winner loser

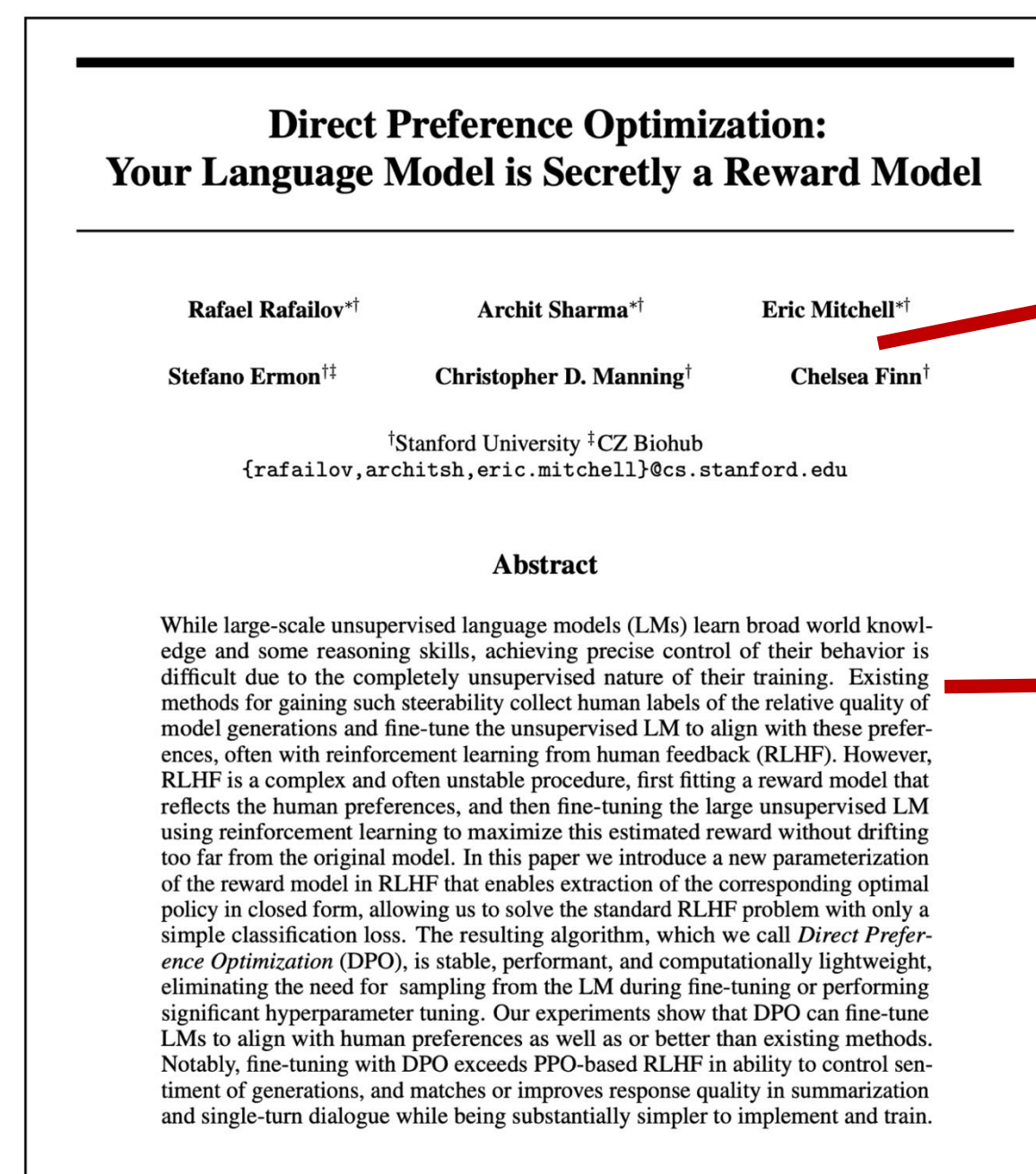
Safety example (Dai et al., 2024; Ji et al., 2024)

- 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,

- Important stage in LLM development (*post-training*), tuning from pairwise preferences

Direct preference alignment (DPA) approaches

- Recent approaches, such as DPO, take the form of **closed-form loss functions**, directly tune models to *offline* preference data (no RL). **Many variations**.
- Problem**: hard to interpret, understand relationships between variants, devise new approaches.



Original DPO loss

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

$y \sim \pi_{\theta}(\cdot | x)$ policy model
 $y \sim \pi_{\text{ref}}(\cdot | x)$ reference

Variations of DPO

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

From Meng et al. NeurIPS 2024

Understanding the DPA loss space

- Goals**: formal framework for characterizing the semantics of DPA losses, deriving new losses.
- Approach**: *decompiling losses to symbolic programs, discrete reasoning problems*

Preference loss ℓ (DPO)

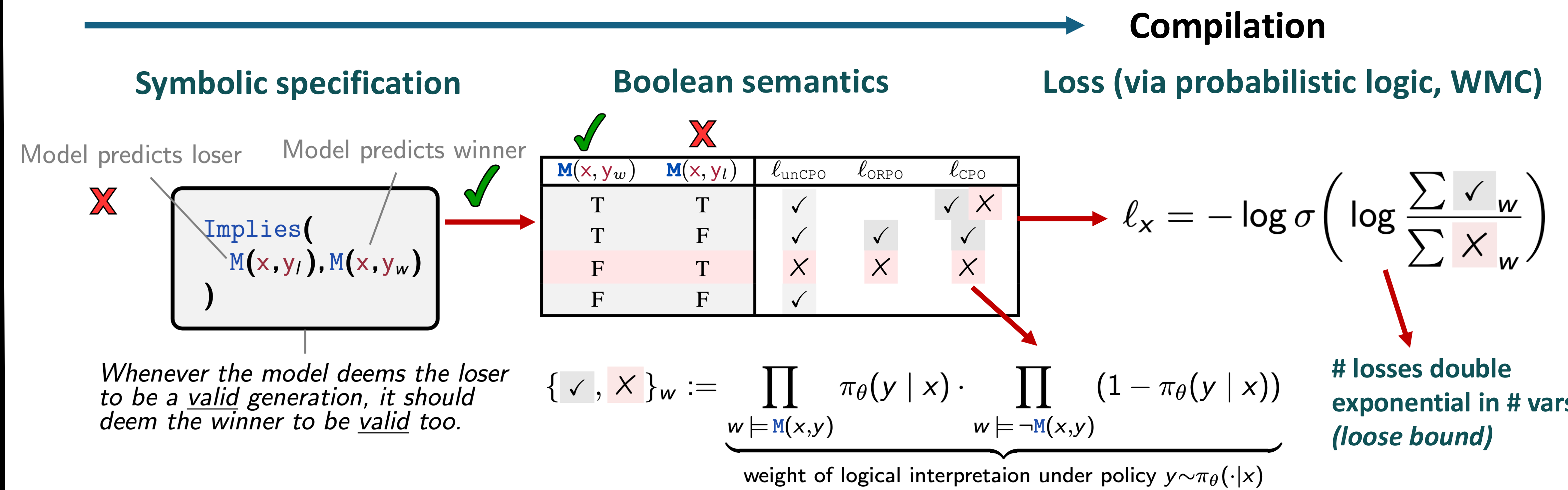
$$-\log \sigma \left(\log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right)$$

Symbolic program

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

Compilation ← → Decompile

From symbolic programs to losses (and back)



Decompilation

- How?** Devised novel logic, compositional translation from DPA losses into that logic.
- Preference structure**: Boolean encoding, express losses as symbolic programs + hard constraints.

| 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 | $P := \text{Implies}(M(x, y_l), M(x, y_w))$ $P_C := \text{Or}(M(x, y_l), M(x, y_w))$ |
| ORPO | $P := \text{Implies}(M(x, y_l), M(x, y_w))$ $P_C := \text{XOR}(M(x, y_l), M(x, y_w))$ |

Core semantics

Differing conditioning constraints

Deriving new losses from first principles

- Why is this useful?** high-level programming language for deriving new losses, modifying existing ones.

Symbolic Program

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

DPO Loss

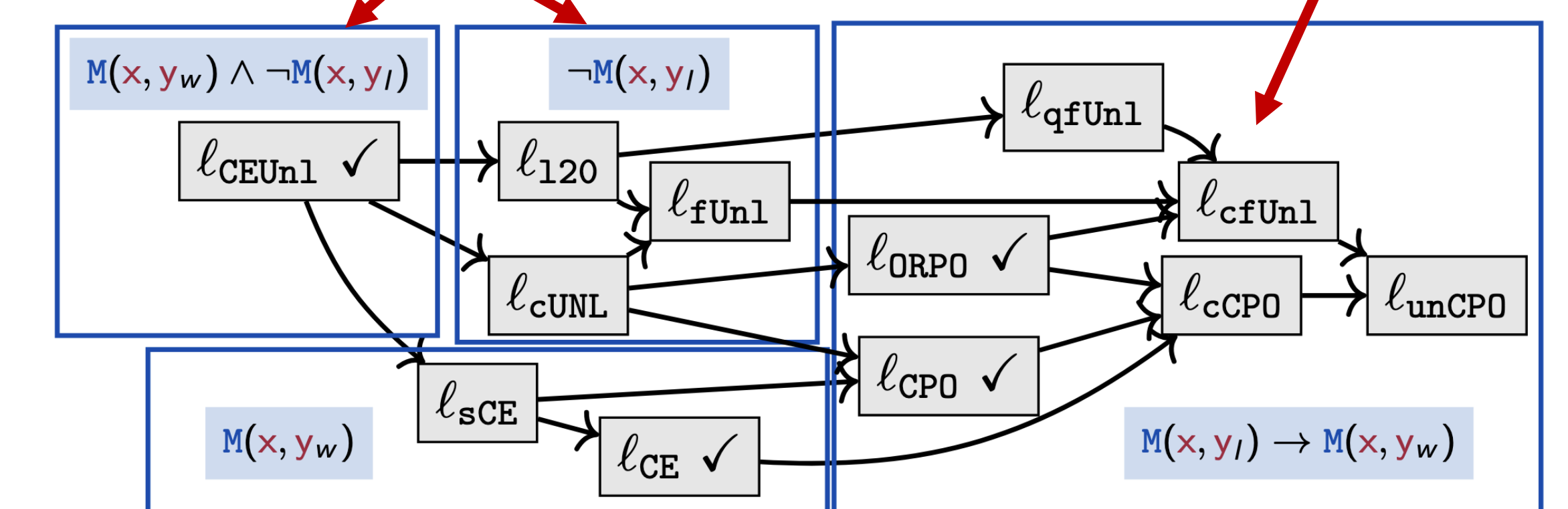
$$-\log \sigma \left(\log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right)$$

Novel loss

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

Semantic regions

New losses, entailment



most constrained → least constrained

- Loss lattice**: structured representation of loss space for exploration, small empirical **case study**.