

# One-Shot Safety Alignment for Large Language Models via Optimal Dualization

Xinmeng Huang\*<sup>†</sup>

Osbert Bastani

Shuo Li\*<sup>†</sup>

Hamed Hassani

Edgar Dobriban

Dongsheng Ding<sup>†</sup>



## MOTIVATION

### Safety requirements for language models (LM)

1. **MUST NOT** contain offensive or discriminatory content
2. **MUST NOT** fabricate content and spread misinformation
- ...

### Constraining LMs w/ safety requirements

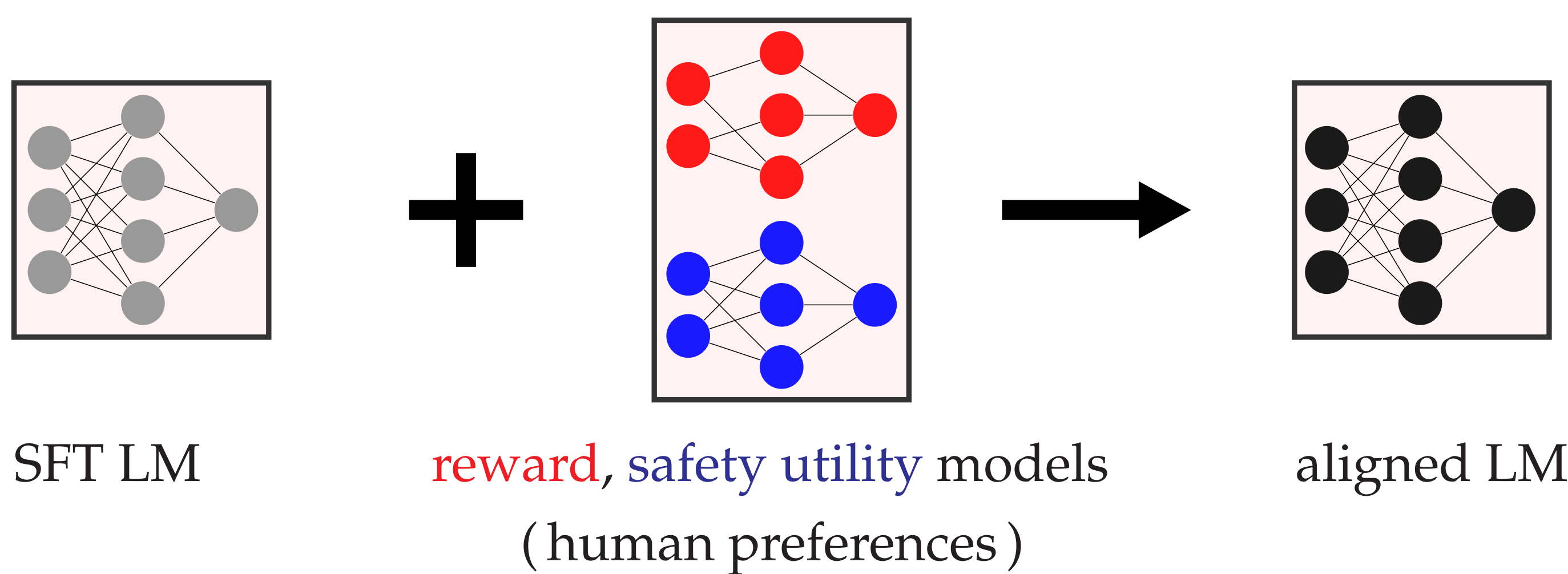
- Safe RLHF (ICLR 2024)
- Constrained RLHF (ICLR 2024)
- Constrained DPO (arXiv 2023)
- SACPO (arXiv 2024)

**Hurdle** • **instability** of iterative training • **no optimality** certificate

Can we align LMs w/ safety constraints in a **one-shot** way?

## PROBLEM FORMULATION

### LM alignment via human feedback



- $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$  – (prompt, response)
- $\pi_{\text{ref}}(\cdot | \mathbf{x}) \in \Delta(\mathcal{Y})$  – SFT LM
- $r(\mathbf{x}, \mathbf{y}), g_j(\mathbf{x}, \mathbf{y}), j = 1, \dots, m$  – reward, safety utility models

### Constrained alignment problem

$$\begin{aligned} \max_{\pi \in \Pi} \quad & \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})] - \beta \text{KL}(\pi(\cdot | \mathbf{x}) || \pi_{\text{ref}}(\cdot | \mathbf{x})) \right] \\ \text{s.t.} \quad & \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [g_j(\mathbf{x}, \mathbf{y})] - \mathbb{E}_{\mathbf{y} \sim \pi_{\text{ref}}(\cdot | \mathbf{x})} [g_j(\mathbf{x}, \mathbf{y})] \right] \geq b_j \\ & \bullet \Pi - \text{LM policy set} \quad \bullet b_j - \text{safety margin} \quad j = 1, \dots, m \end{aligned}$$

## OPTIMAL DUALIZATION

### Dual problem

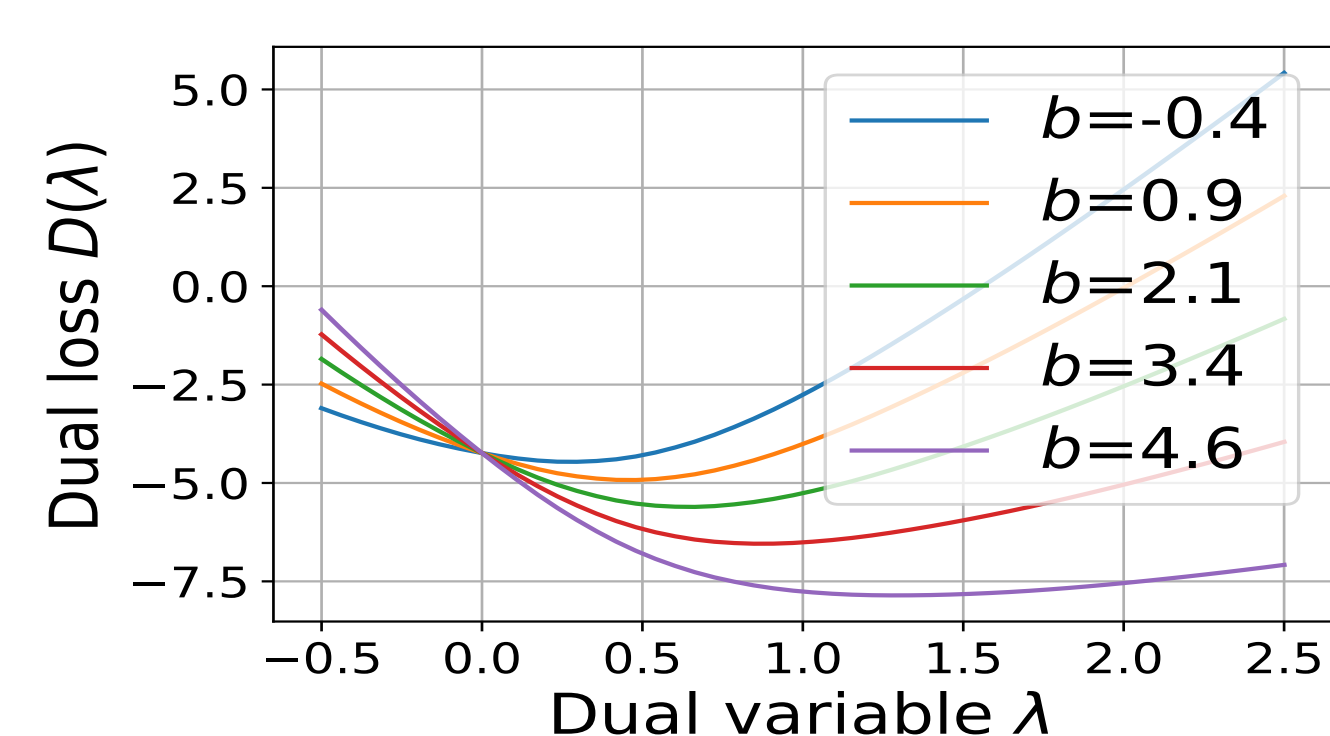
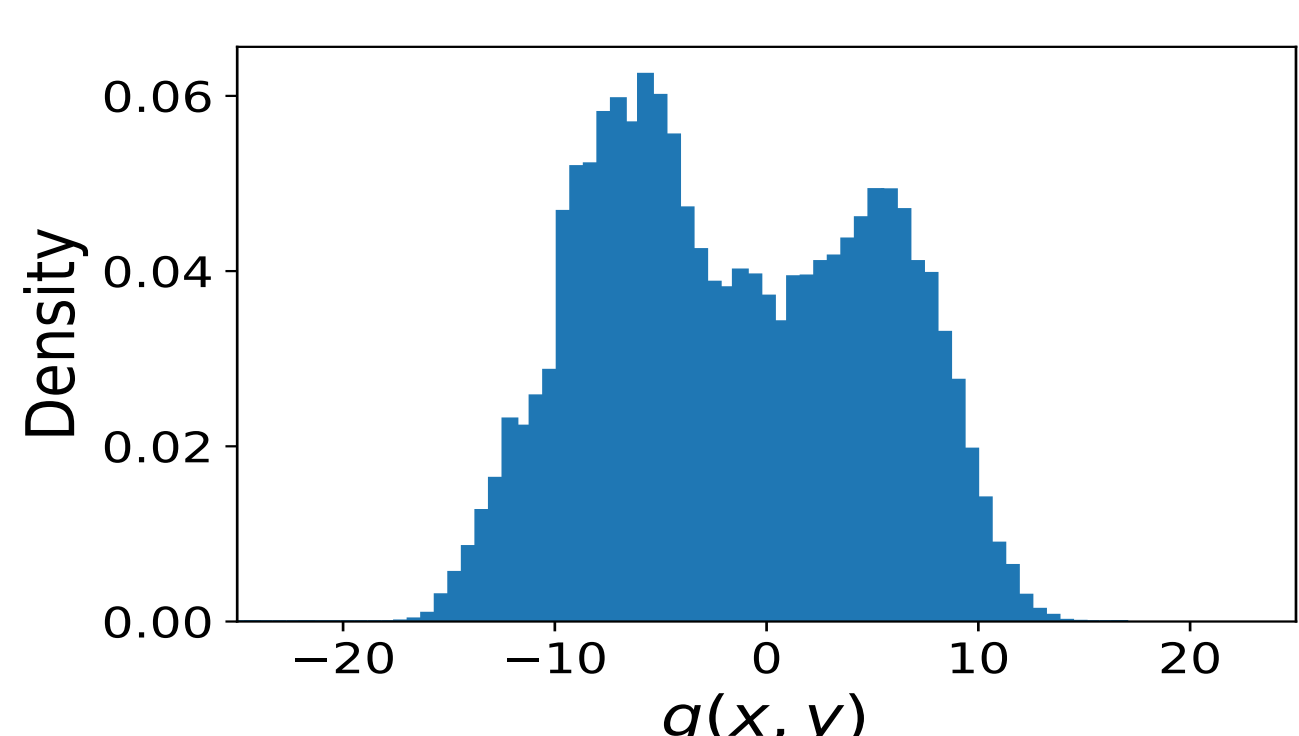
$$\text{minimize}_{\lambda \geq 0} D(\lambda) := \text{maximize}_{\pi \in \Pi} L(\pi, \lambda)$$

- $L(\pi, \lambda)$  – Lagrangian
- $$= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})] - \beta \text{KL}(\pi(\cdot | \mathbf{x}) || \pi_{\text{ref}}(\cdot | \mathbf{x}))$$
 objective
- $$+ \sum_{j=1}^m \lambda_j \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [h_j(\mathbf{x}, \mathbf{y})]$$
 safety violation
- $h_j(\mathbf{x}, \mathbf{y}) := g_j(\mathbf{x}, \mathbf{y}) - \mathbb{E}_{\pi_{\text{ref}}}[g_j(\mathbf{x}, \mathbf{y})] - b_j$  – shifted safety utility
- $L(\pi^*, 0) = D(\lambda^*)$  for an optimal pair  $(\pi^*, \lambda^*)$  – **strong duality**

### Explicit dual function

$$D(\lambda) = \beta \mathbb{E}_{\mathbf{x}} \left[ \ln \mathbb{E}_{\mathbf{y} \sim \pi_{\text{ref}}(\cdot | \mathbf{x})} \left[ \exp \left( \frac{r(\mathbf{x}, \mathbf{y}) + \lambda^{\top} \mathbf{h}(\mathbf{x}, \mathbf{y})}{\beta} \right) \right] \right]$$

- $(\pi^*, \lambda^*)$  – **uniqueness** of optimal primal-dual pair
- $D$  – **smooth** & **strongly convex** at the unique  $\lambda^*$



## ONE-SHOT SAFETY ALIGNMENT

### Constrained Alignment via dualization (CAN)

**Stage 1** Optimal dual:  $\lambda^* = \arg \min_{\lambda \geq 0} D(\lambda)$

**Stage 2** Update LM:  $\pi^* = \arg \max_{\pi \in \Pi} L(\pi, \lambda^*)$

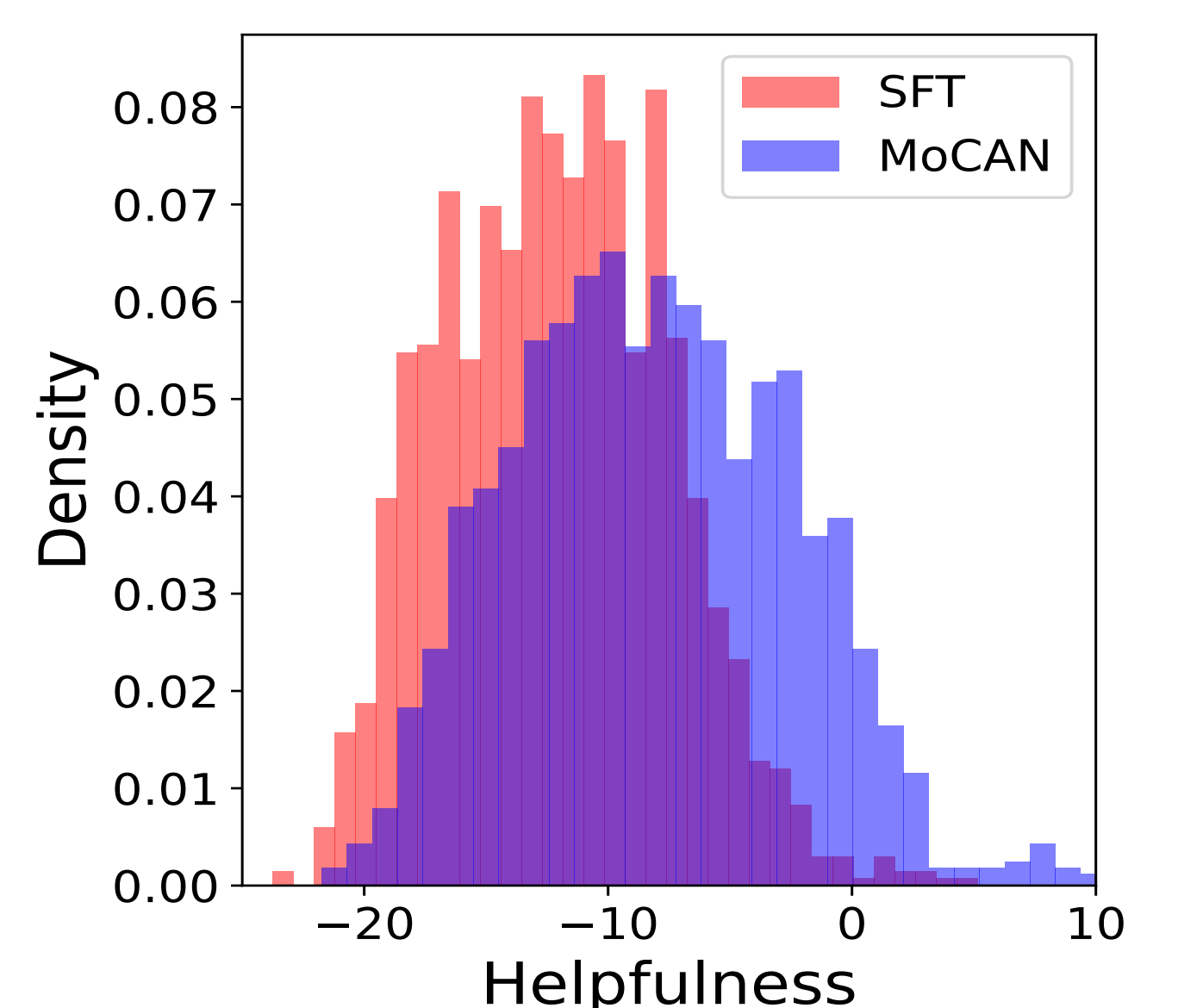
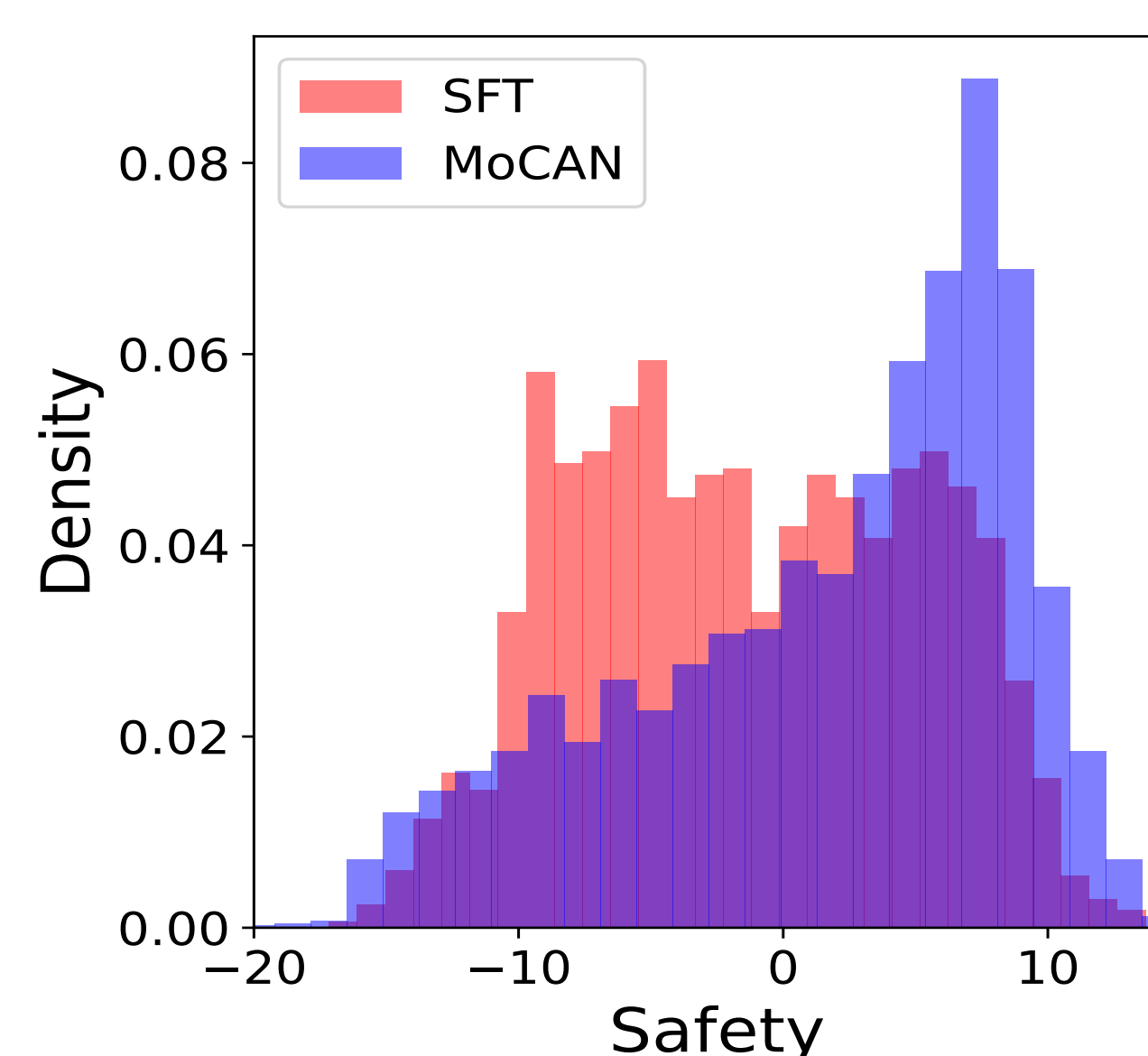
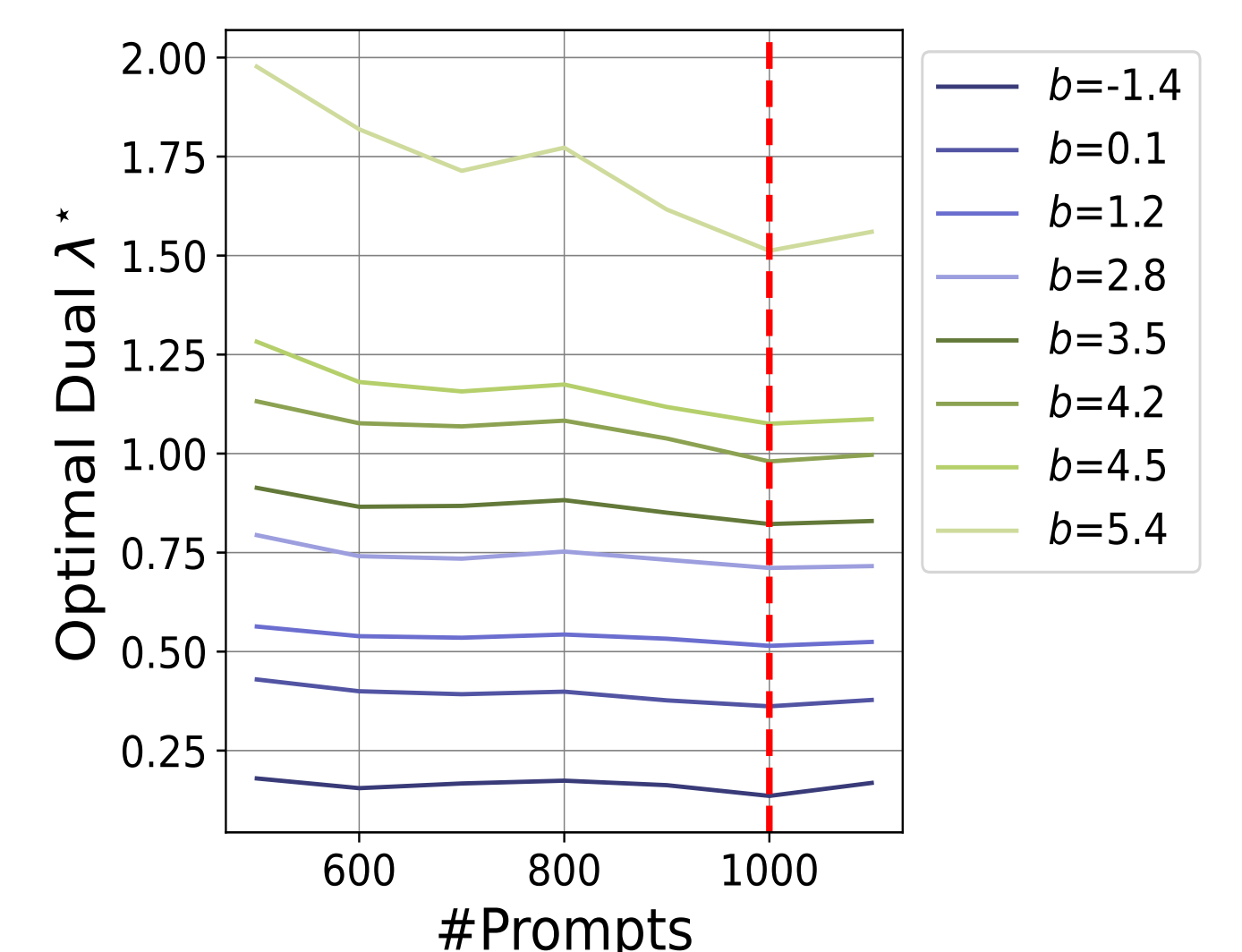
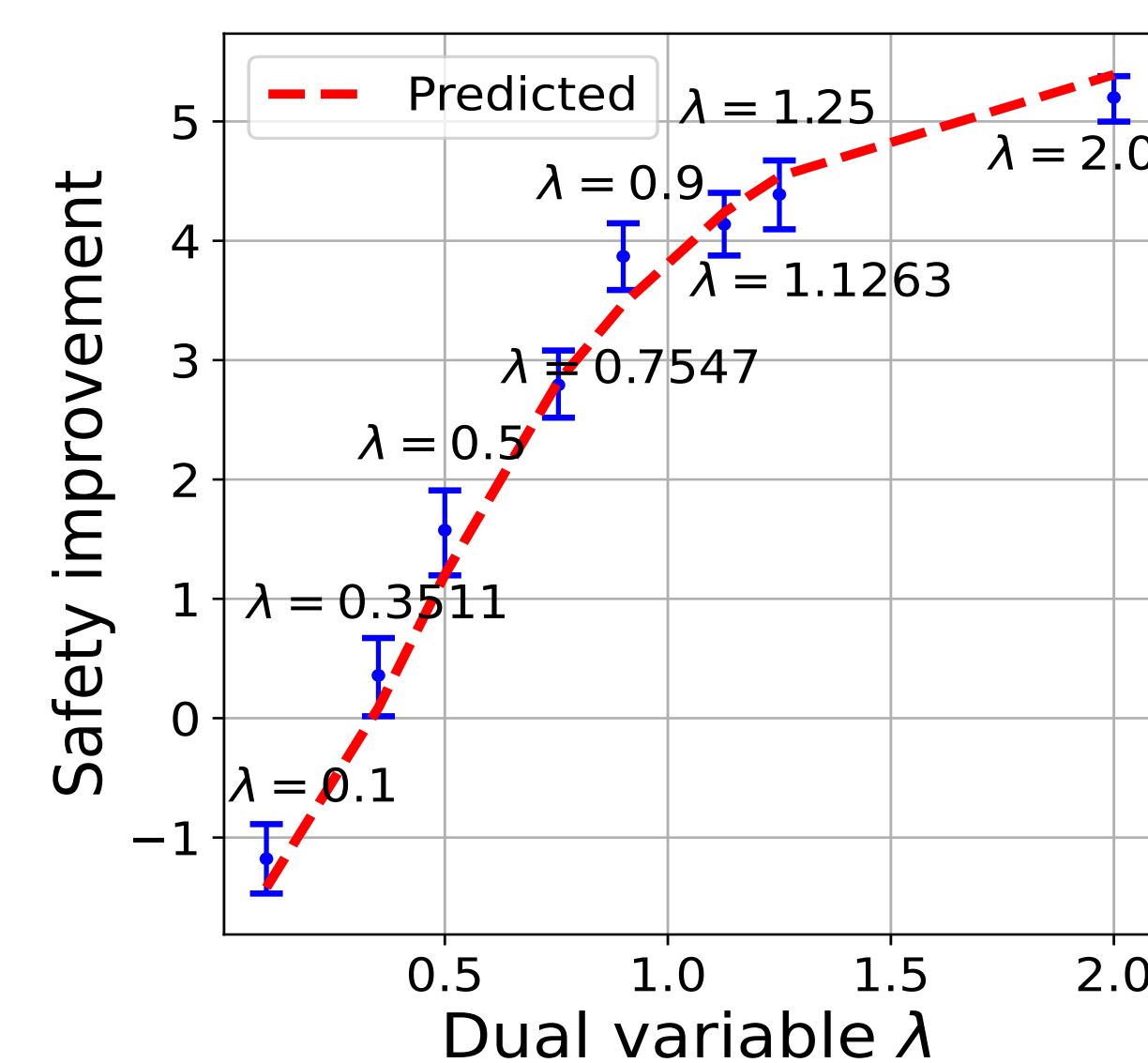
smooth convex optimization & unconstrained alignment

**Advantages** • **optimality** of LM • **stability** of safety training

## PRACTICAL IMPLEMENTATION & EVALUATION

### Model-based CAN (MoCAN)

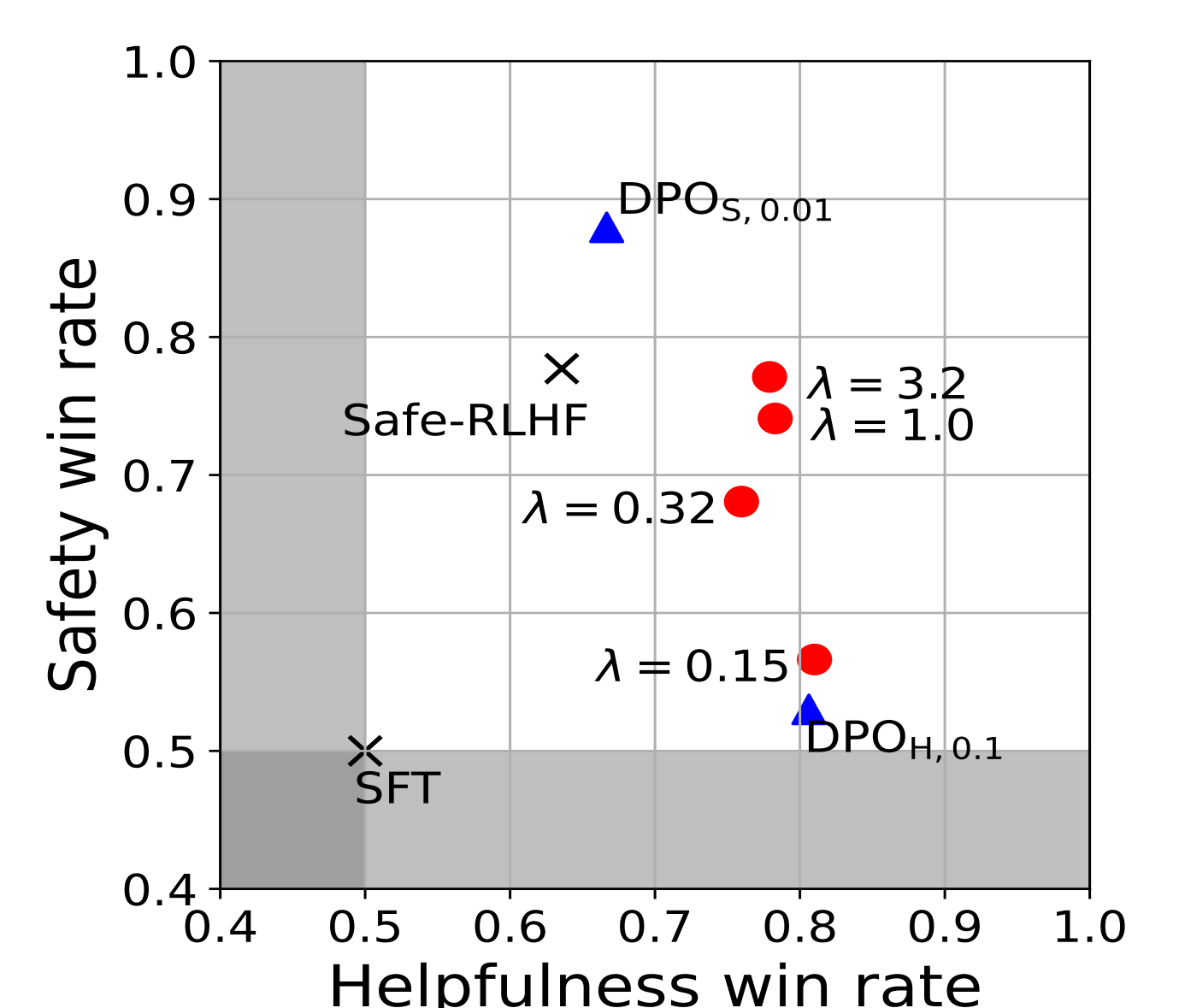
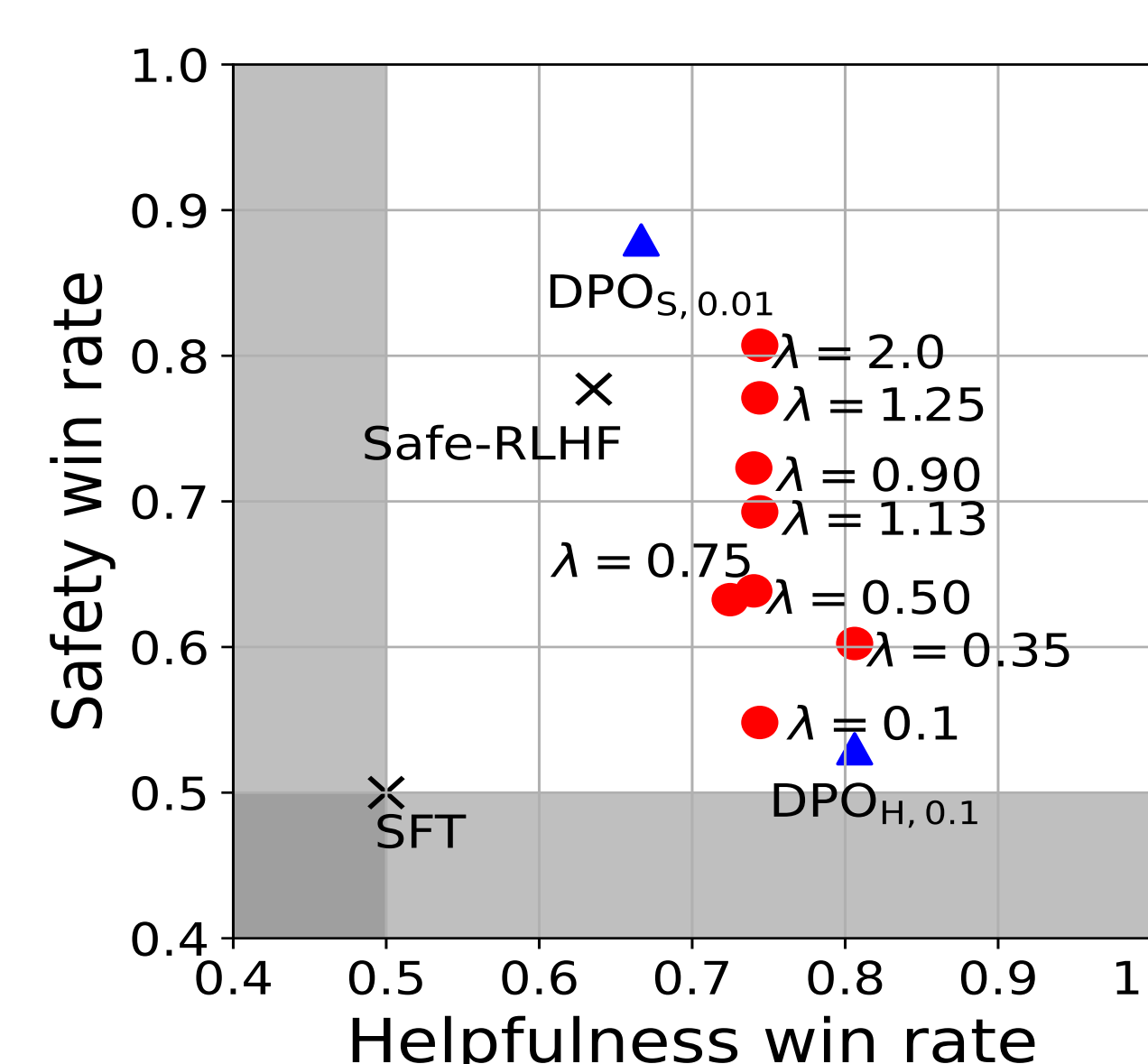
1. collect **offline data**  $(r(\mathbf{x}, \mathbf{y}), g(\mathbf{x}, \mathbf{y}))$ -pairs; estimate  $h(\mathbf{x}, \mathbf{y})$
2. find the optimal dual  $\lambda^*$  using **model-based**  $D(\lambda)$
3. update LM w/ **pseudo preference** from  $r + (\lambda^*)^{\top} g$



### Preference-based CAN (PeCAN)

1. obtain unconstrained **pre-aligned models**  $(\pi_{\theta_r}, \pi_{\theta_g})$
2. collect **offline**  $(\ln \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}), \ln \pi_{\theta_r}(\mathbf{y} | \mathbf{x}), \ln \pi_{\theta_g}(\mathbf{y} | \mathbf{x}))$ -tuples
3. estimate **KL terms**  $\text{KL}(\pi_{\text{ref}} || \pi_{\theta_{g_j}}), j = 1, \dots, m$
4. find the optimal dual  $\lambda^*$  using **preference-based**  $D(\lambda)$
5. update LM w/ **pseudo preference** from  $\beta \ln \frac{\pi_{\theta_r}}{\pi_{\text{ref}}} + \beta (\lambda^*)^{\top} \ln \frac{\pi_{\theta_g}}{\pi_{\text{ref}}}$

### Safety / helpfulness tradeoff (L: MoCAN, R: PeCAN)



### Key takeaways

- efficient dual optimization for safety improvement
- empirically optimal dual variable quickly stabilizes
- **MoCAN** finds the optimal safe LM in one-shot
- **PeCAN** performs similarly if pre-aligned models are accurate