

Constrained Policy Optimization for Large Language Model Alignment

Dongsheng Ding

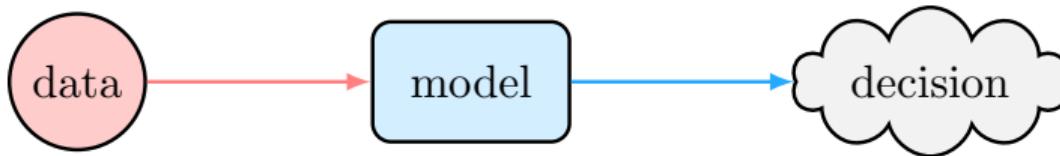
dongshed@utk.edu

EECS

Statistics and Data Science Seminar, UTK

November 13, 2025

The reality of machine learning



goal – loss / reward / likelihood

■ REQUIREMENTS

harmlessness



safety



robustness



fairness



Image sources: Stanford HAI, Waymo, WIRED, NBC

The risk of machine learning

AI chatbots might be sabotaging women by advising them to ask for lower salaries, study says

New York Post, JUL 29, 2025

NHTSA probes Waymo self-driving cars over school bus safety concerns

Reuters, OCT 21, 2025

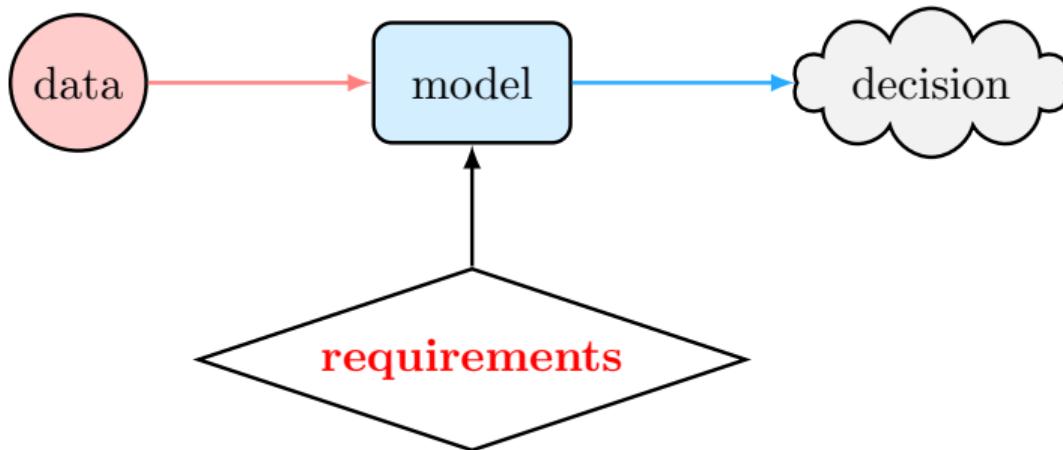
AI-Powered Robots Can Be Tricked Into Acts of Violence

WIRED, DEC 4, 2024

Study reveals why AI models that analyze medical images can be biased

MIT News, JUN 28, 2024

Requirement-driven machine learning



harmlessness



safety



robustness



fairness



Motivating application: Robotics

■ LLM-CONTROLLED ROBOTS



Figure A1

maximize **dos**
LLM policy
subject to **don'ts** \geq threshold

Motivating application: Healthcare

■ AI THERAPY CHATBOTS



Stanford HAI

maximize **helpfulness**
LLM policy
subject to **harmlessness \geq threshold**

REAL-WORLD CHALLENGE

Constraint satisfaction

OBJECTIVE

Find a **Large Language Model (LLM)** that
maximizes a performance metric
subject to a constraint on
another performance metric

Outline

■ CONSTRAINED LLM ALIGNMENT

- ★ constrained policy optimization

■ ALIGNMENT METHOD & THEORY

- ★ non-iterative & iterative methods
- ★ duality gap & optimality gap

■ EMPIRICAL STUDY

- ★ safety-alignment task

■ SUMMARY & OUTLOOK

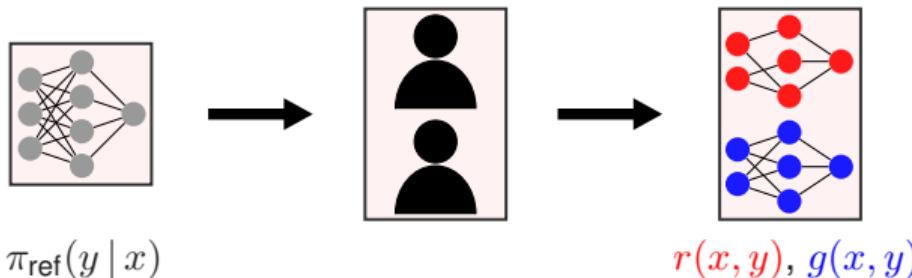
CONSTRAINED LLM ALIGNMENT

constrained policy optimization

Alignment framework

■ REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

- ★ reward modeling



$\pi_{\text{ref}}: \mathcal{X}$ (prompts) $\rightarrow \mathcal{Y}$ (responses) – reference LLM policy

$r(x, y), g(x, y)$ – reward/utility models

- ★ reward/utility models

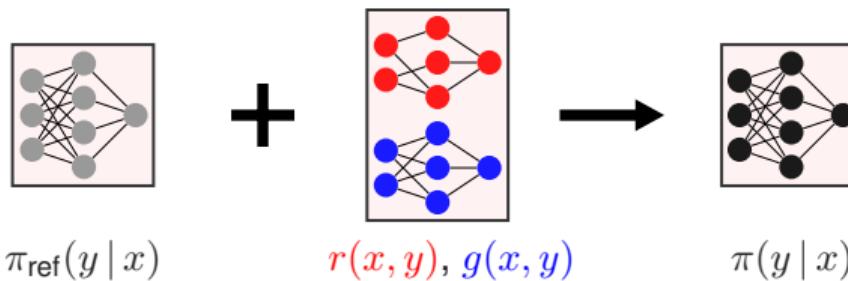
e.g., SafeRLHF: helpfulness and harmlessness

Dai et al., ICLR '24

Alignment framework

■ REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

- ★ policy optimization



$\pi_{\text{ref}}: \mathcal{X} \text{ (prompts)} \rightarrow \mathcal{Y} \text{ (responses)} - \text{reference LLM policy}$

$\pi: \mathcal{X} \text{ (prompts)} \rightarrow \mathcal{Y} \text{ (responses)} - \text{aligned LLM policy}$

e.g., direct preference optimization (preference-based)

Rafailov et al., NeurIPS '23

Response space

■ RESPONSE SPACE SIZE

$$|\mathcal{V}|$$

= (# total tokens)^{# tokens per sentence}

e.g., ChatGPT-3.5 / Llama 3: $4000^{30} \approx 10^{108} \gg 10^{80}$

exponentially large decision space

Constrained alignment problem

$$\underset{\pi}{\text{maximize}} \quad \mathbb{E}_x [\mathbb{E}_{y \sim \pi} [r(x, y)] - \beta D_{\text{KL}}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))]$$

$$\text{subject to} \quad \mathbb{E}_x [\mathbb{E}_{y \sim \pi} [g(x, y)]] \geq 0$$



KL-regularized objective

policy constraint

- * limit the policy space to **an inequality constraint**

e.g., harmless policy, safe policy

Constrained policy optimization

$$\underset{\pi}{\text{maximize}} \quad \mathbb{E}_x [\mathbb{E}_{y \sim \pi} [r(x, y)] - \beta D_{\text{KL}}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))]$$

$$\text{subject to} \quad \mathbb{E}_x [\mathbb{E}_{y \sim \pi} [g(x, y)]] \geq 0$$



KL-regularized objective



policy constraint

- ★ no transition dynamics
- ★ **concave** KL-regularized objective and **linear** constraint

Convex constrained policy optimization → **Strong duality**

Lagrangian relaxation

■ LAGRANGIAN

$$L(\pi, \lambda) = \mathbb{E}_x [\mathbb{E}_{y \sim \pi} [r(x, y) + \lambda g(x, y)] - \beta D_{\text{KL}}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))]$$

- ★ penalize violation via dual variable $\lambda \geq 0$

■ LAGRANGIAN MAXIMIZATION

$$\underset{\pi}{\text{maximize}} \quad L(\pi, \lambda)$$

convex conjugate

- ★ exponentially tilted distribution $\pi^*(\cdot | x; \lambda)$

$$\pi^*(\cdot | x; \lambda) \propto \pi_{\text{ref}}(\cdot | x) e^{(r(x, \cdot) + \lambda g(x, \cdot)) / \beta}$$

Existence of an optimal dual variable λ^*

Lagrangian dual function

■ UPPER ENVELOPE FUNCTION

$$D(\lambda) := \underset{\pi}{\text{maximize}} \ L(\pi, \lambda)$$

$$= \beta \mathbb{E}_x \left[\log \mathbb{E}_{y \sim \pi_{\text{ref}}} \left[e^{(r(x,y) + \lambda g(x,y)) / \beta} \right] \right]$$

cumulant-generating function

- ★ **convex**, and **smooth** function
- ★ **strictly convex**, and **locally strongly convex** function

$$\nabla^2 D(\lambda) \simeq \mathbb{E}_x \left[\text{Var}_{y \sim \pi^*(\cdot | x; \lambda)} [g(x, y)] \right]$$

Lagrangian dual problem

■ LAGRANGIAN DUAL MINIMIZATION

$$\underset{\lambda \geq 0}{\text{minimize}} \quad D(\lambda)$$

convex and smooth optimization

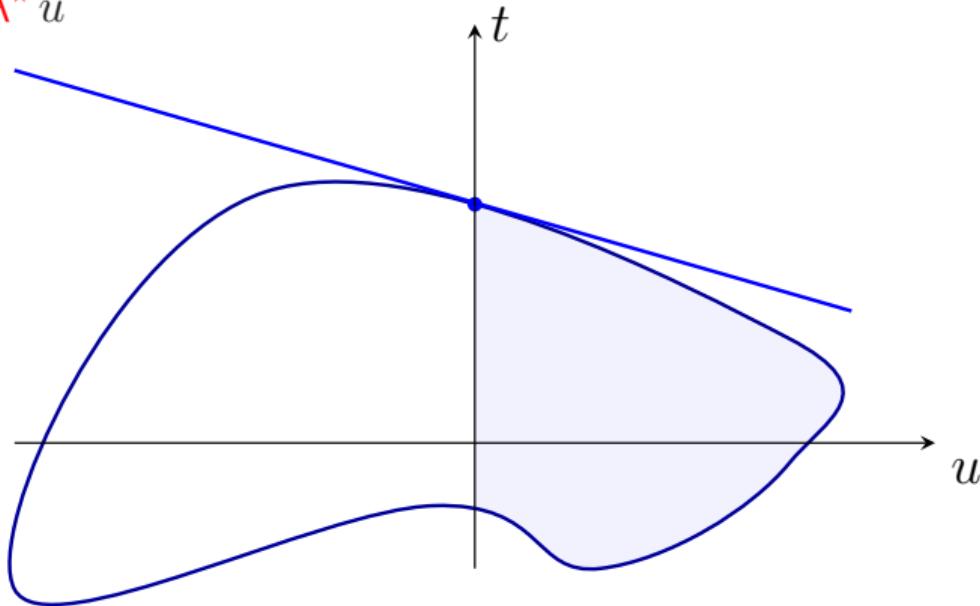
- * gradient descent finds an optimal dual variable λ^*
- * uniqueness of an optimal dual variable λ^*

Recovery of an optimal constrained policy π^*

$$\pi^*(\cdot | x) = \pi^*(\cdot | x; \lambda^*) \propto \pi_{\text{ref}}(\cdot | x) e^{(r(x, \cdot) + \lambda^* g(x, \cdot)) / \beta}$$

■ GEOMETRIC INTERPRETATION OF STRONG DUALITY

$$D^* = t + \lambda^* u$$



$$\text{image } \mathcal{E} = \{ (g(\pi), r_{\text{KL}}(\pi)) \mid \pi \}$$

Optimal hyperplane touches \mathcal{E} at an optimal policy: $D^* = r_{\text{KL}}(\pi^*) := P^*$

Overview of our results

ZLHB \mathbf{D}^\dagger R, NeurIPS '25

$\mathbf{H}^\dagger \mathbf{L}^\dagger \mathbf{D} \mathbf{B} \mathbf{L} \mathbf{H} \mathbf{D}^\dagger$, NeurIPS '24

■ DUALIZATION-BASED ALIGNMENT METHODS

- ★ non-iterative & iterative methods
- ★ duality gap
- ★ optimality gap objective and constraint

Dual methods find **an optimal constrained LLM policy**,
up to **a parametrization gap**

ALIGNMENT METHOD & THEORY

non-iterative method

Dualization-based alignment

$H^\dagger L^\dagger DBLH D^\dagger$, NeurIPS '24

■ STAGE #1: FIND AN OPTIMAL DUAL VARIABLE

$$\lambda^* = \underset{\lambda \geq 0}{\operatorname{argmin}} D(\lambda)$$

convex and smooth optimization

■ STAGE #2: SEARCH FOR AN LLM POLICY

$$\pi^* = \underset{\pi}{\operatorname{argmax}} L(\pi, \lambda^*)$$

unconstrained alignment

Computational efficiency

Search for an optimal dual variable

Lagrangian maximizer: $\pi^*(\lambda) = \operatorname{argmax}_{\pi} L(\pi, \lambda)$

Gradient: $\nabla D(\lambda) = \nabla_{\lambda} L(\pi, \lambda) |_{\pi = \pi^*(\lambda)}$

■ PROJECTED GRADIENT DESCENT

$$\lambda^+ \leftarrow [\lambda - \eta \nabla D(\lambda)]$$

* π -independent gradient $\nabla D(\lambda)$

$$\nabla D(\lambda) = \frac{\mathbb{E}_{y \sim \pi_{\text{ref}}} \left[e^{(r(x,y) + \lambda g(x,y))/\beta} g(x,y) \right]}{\mathbb{E}_{y' \sim \pi_{\text{ref}}} \left[e^{(r(x,y') + \lambda g(x,y'))/\beta} \right]}$$

Offline biased estimate

Search for an LLM policy

Lagrangian maximizer: $\pi^*(\lambda^*) \in \operatorname*{argmax}_{\pi} L(\pi, \lambda^*)$

■ POLICY PARAMETRIZATION

$$\pi(y | x) \leftarrow \pi_{\theta}(y | x)$$

model parameter θ

■ PARAMETRIZED LAGRANGIAN MAXIMIZER

$$\pi_{\theta^*(\lambda^*)} \in \operatorname*{argmax}_{\theta} L(\pi_{\theta}, \lambda^*)$$

QUESTION: Optimality of λ^* -recovered model $\pi_{\theta^*(\lambda^*)} := \pi_p^*(\lambda^*)$?

Constrained parameter optimization

$$\underset{\theta}{\text{maximize}} \quad \mathbb{E}_x [\mathbb{E}_{y \sim \pi_{\theta}} [r(x, y)] - \beta D_{\text{KL}}(\pi_{\theta}(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))]$$

$$\text{subject to} \quad \mathbb{E}_x [\mathbb{E}_{y \sim \pi_{\theta}} [g(x, y)]] \geq 0$$



KL-regularized objective



policy constraint

- * decision of model parameter θ

CHALLENGE

Nonconvex constrained optimization → Lack of strong duality

METHOD & THEORY

iterative method

Parametrized Lagrangian dual problem

■ LAGRANGIAN DUAL FUNCTION

$$D_p(\lambda) := \underset{\theta}{\text{maximize}} \ L(\pi_\theta, \lambda)$$

- * convex, and nondifferentiable function

■ LAGRANGIAN DUAL MINIMIZATION

$$\underset{\lambda \geq 0}{\text{minimize}} \ D_p(\lambda)$$

Existence of an optimal parametrized dual variable λ_p^*

Search for an optimal parametrized dual variable

Lagrangian maximizer: $\theta^*(\lambda) \in \operatorname{argmax}_{\theta} L(\pi_\theta, \lambda)$

Subgradient: $u(\lambda) = \nabla_\lambda L(\pi_\theta, \lambda) |_{\theta = \theta^*(\lambda)}$

■ PROJECTED SUBGRADIENT DESCENT

$$\lambda^+ \leftarrow [\lambda - \eta u(\lambda)]_+$$

* explicit subgradient $u(\lambda) = \mathbb{E}_{y \sim \pi_{\theta^*(\lambda)}} [g(x, y)]$

Online unbiased estimate

Iterative dualization-based alignment

ZLHBDR, NeurIPS '25

■ ITERATION #1: COMPUTE A LAGRANGIAN MAXIMIZER

$$\theta^*(\lambda) \in \operatorname{argmax}_{\theta} L(\pi_{\theta}, \lambda)$$

■ ITERATION #2: PERFORM A SUBGRADIENT DESCENT STEP

$$\lambda^+ \leftarrow \left[\lambda - \eta \mathbb{E}_{y \sim \pi_{\theta^*(\lambda)}} [g(x, y)] \right]_+$$

QUESTION: Optimality of λ_p^* -recovered model $\pi_{\theta^*(\lambda_p^*)} := \pi_p^*(\lambda_p^*)$?

METHOD & THEORY

duality gap & optimality gap

Duality gap

Duality gap: $|P^* - D_p^*|$

Theorem (informal)

★ **Duality gap** is dominated by

ν

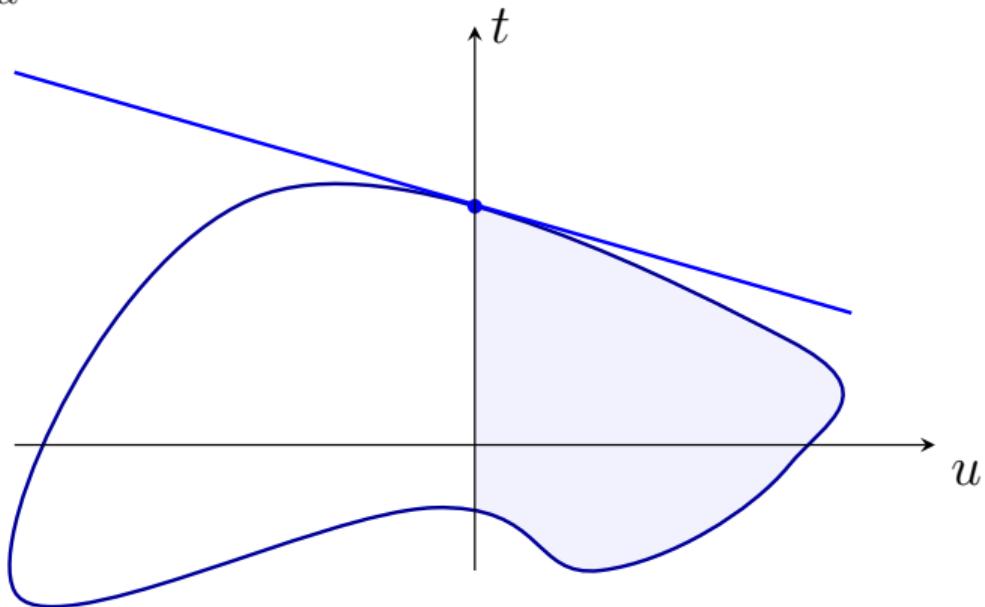
parametrization gap $\nu := \max_{\pi} \min_{\theta} \text{dist}_1(\pi, \pi_{\theta})$

* ν -parametrization gap yields ν -duality gap

linear independence

■ GEOMETRIC INTERPRETATION OF DUALITY GAP

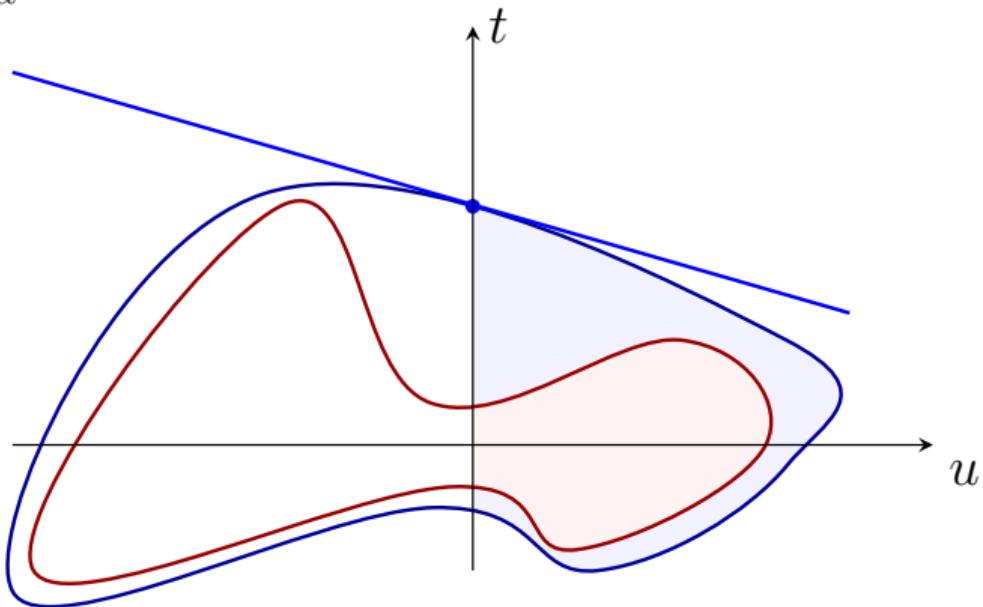
$$D^* = t + \lambda^* u$$



$$\text{image } \mathcal{E} = \{ (g(\pi), r_{\text{KL}}(\pi)) \mid \pi \}$$

■ GEOMETRIC INTERPRETATION OF DUALITY GAP

$$D^* = t + \lambda^* u$$



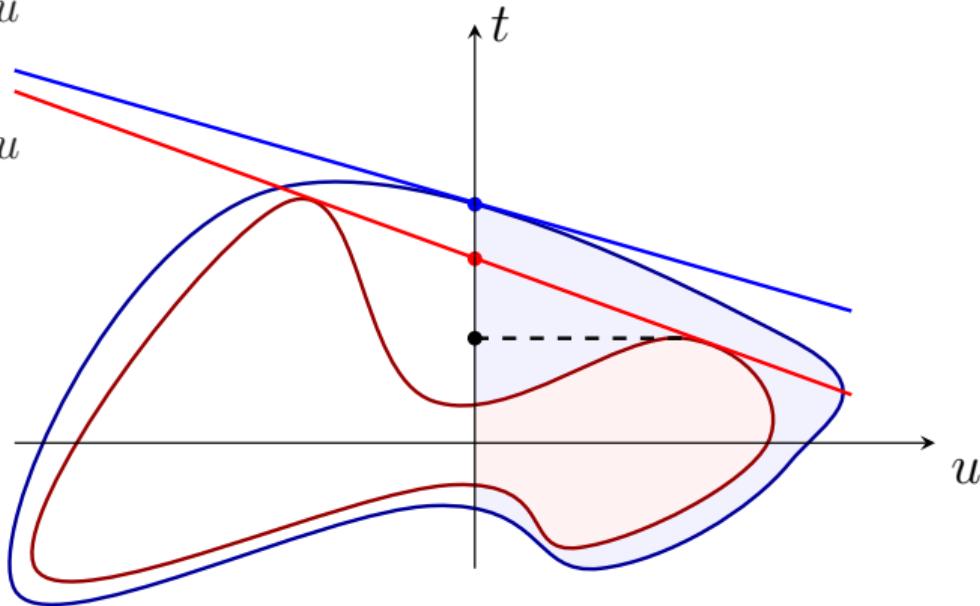
$$\text{image } \mathcal{E} = \{ (g(\pi), r_{\text{KL}}(\pi)) \mid \pi \}$$

$$\text{image } \mathcal{E}_p = \{ (g(\pi_\theta), r_{\text{KL}}(\pi_\theta)) \mid \theta \}$$

■ GEOMETRIC INTERPRETATION OF DUALITY GAP

$$D^* = t + \lambda^* u$$

$$D_p^* = t + \lambda_p^* u$$



$$\text{image } \mathcal{E} = \{ (g(\pi), r_{\text{KL}}(\pi)) \mid \pi \}$$

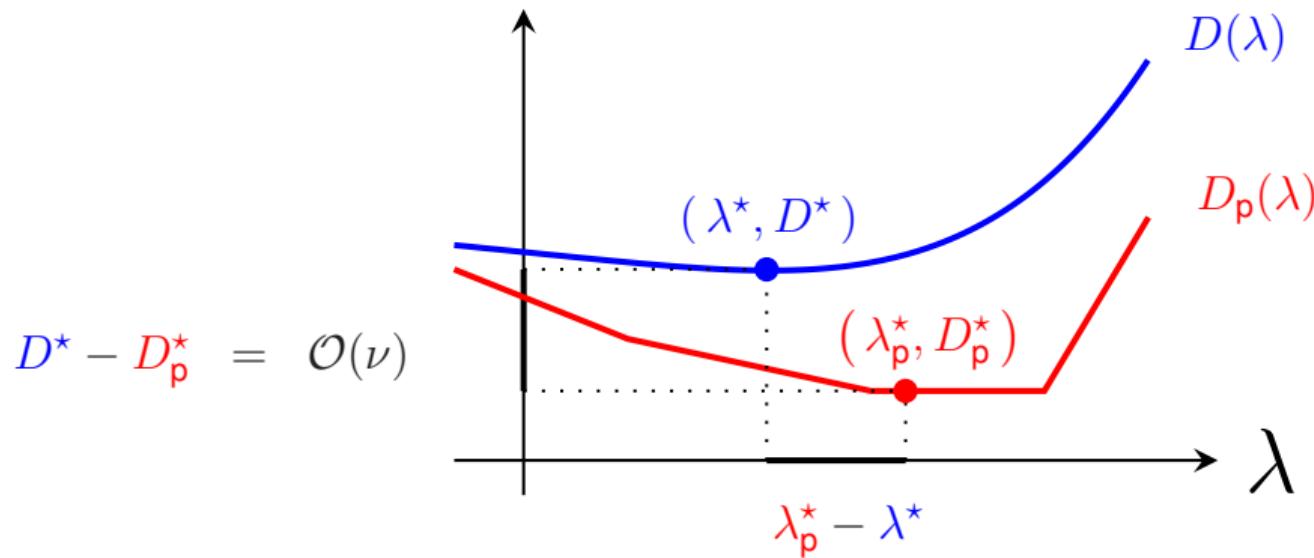
$$\text{image } \mathcal{E}_p = \{ (g(\pi_\theta), r_{\text{KL}}(\pi_\theta)) \mid \theta \}$$

Optimal hyperplane touches \mathcal{E}_p w/ t -intercept D_p^* :

$$D^* - D_p^* = \mathcal{O}(\nu)$$

Gap between optimal dual variables

■ GAP BETWEEN (UN)PARAMETRIZED DUAL FUNCTIONS



Optimal dual variables: λ^* , λ_p^* are close:

$$\|\lambda^* - \lambda_p^*\| = \mathcal{O}(\sqrt{\nu})$$

Optimality gap for iterative method

Objective optimality: $|r_{\text{KL}}(\pi_p^*(\lambda_p^*)) - r_{\text{KL}}(\pi^*)|$

Constraint feasibility: $|g(\pi_p^*(\lambda_p^*)) - g(\pi^*)|$

Implication (informal)

★ **Objective optimality & Constraint feasibility** are dominated by

$$\sqrt{\nu}$$

parametrization gap $\nu := \max_{\pi} \min_{\theta} \text{dist}_1(\pi, \pi_{\theta})$

Root-scaling of parametrization gap

Optimality gap for non-iterative method

Objective optimality: $|r_{\text{KL}}(\pi_p^*(\lambda^*)) - r_{\text{KL}}(\pi^*)|$

Constraint feasibility: $|g(\pi_p^*(\lambda^*)) - g(\pi^*)|$

Implication (informal)

★ **Objective optimality & Constraint feasibility** are dominated by

$$\sqrt{\nu}$$

parametrization gap $\nu := \max_{\pi} \min_{\theta} \text{dist}_1(\pi, \pi_{\theta})$

Root-scaling of parametrization gap

EMPIRICAL STUDY

safety-alignment task

Practical implementation

■ NON-ITERATIVE DUALIZATION ALIGNMENT

model-based setting

offline model-based dual
pseudo-preference optimization

preference-based setting

offline preference-based dual
pseudo-preference optimization

■ ITERATIVE DUALIZATION ALIGNMENT

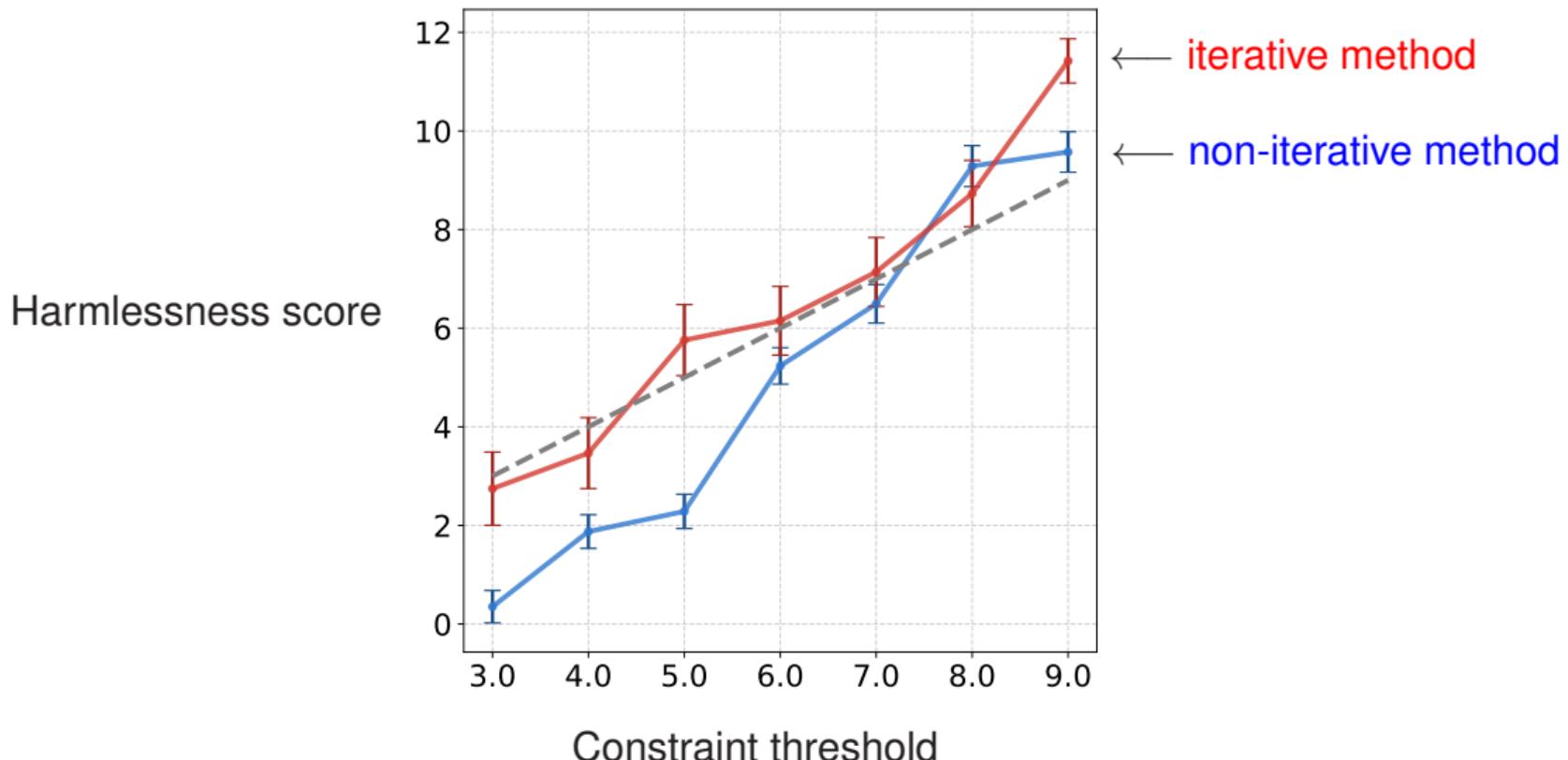
model-based setting

online model-based dual
pseudo-preference optimization

preference-based setting

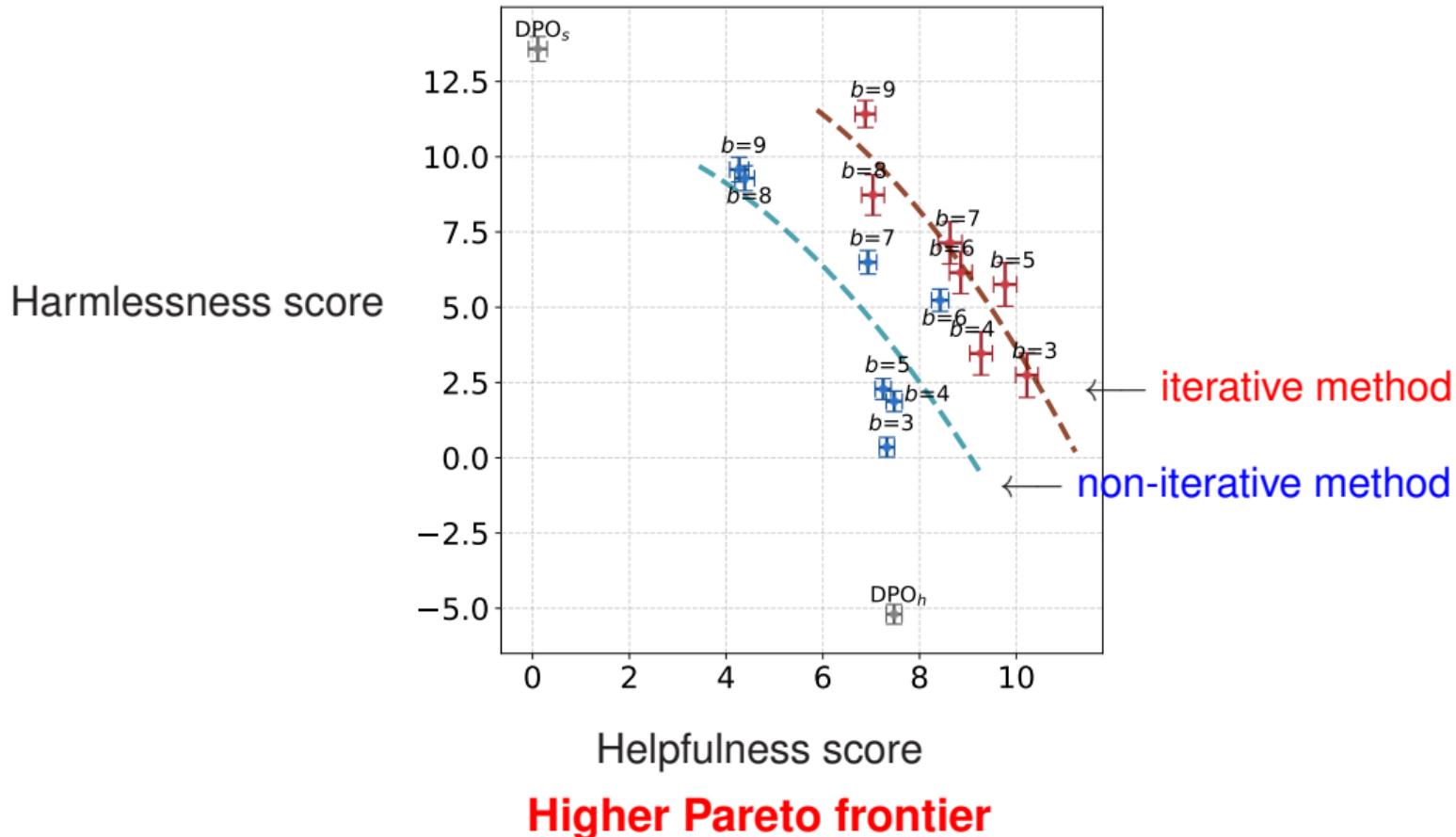
online preference-based dual
pseudo-preference optimization

Constraint satisfaction



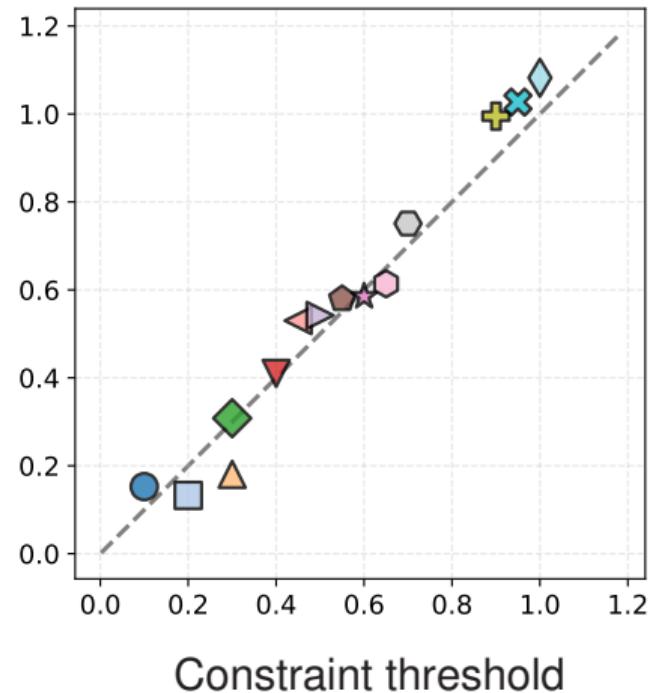
Better constraint satisfaction

Helpfulness and Harmlessness tradeoff

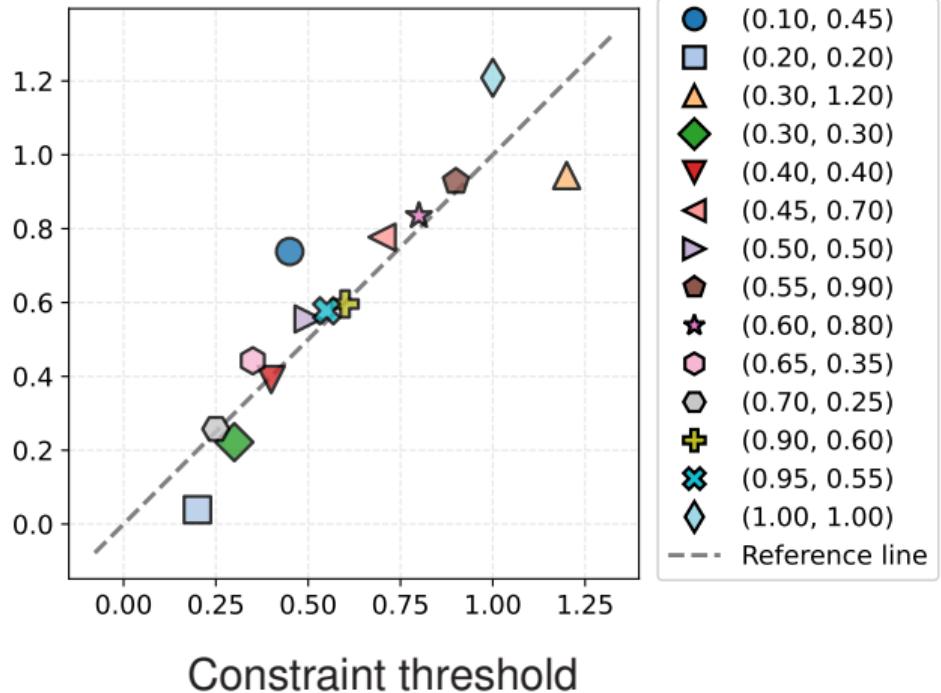


Harmlessness and Humor constraints

Harmlessness score



Humor score



Multi-constraint satisfaction

Summary & outlook

ZLHB \mathbf{D}^\dagger R, NeurIPS '25

$\mathbf{H}^\dagger \mathbf{L}^\dagger \mathbf{D} \mathbf{B} \mathbf{L} \mathbf{H} \mathbf{D}^\dagger$, NeurIPS '24

■ DUALIZATION-BASED ALIGNMENT METHODS

- ★ non-iterative & iterative methods
- ★ duality gap & optimality gap

■ OPEN CHALLENGES

- ★ optimal sample complexity
- ★ multi-turn alignment

Thank you for your attention.