

Stepwise Alignment for Constrained Language Model Policy Optimization



Akifumi Wachi^{*1}, Thien Q. Tran^{*1}, Rei Sato¹, Takumi Tanabe¹, Youhei Akimoto^{2,3}

¹LY Corporation (*Equal contribution)

²University of Tsukuba

³RIKEN AIP

Correspondence: akifumi.wachi@lycorp.co.jp

Paper

Code & Models

Overview

Alignment of language models (LMs) is inherently multifaceted:

Helpfulness vs. *Harmlessness*

We consider the following constrained LM alignment problem:

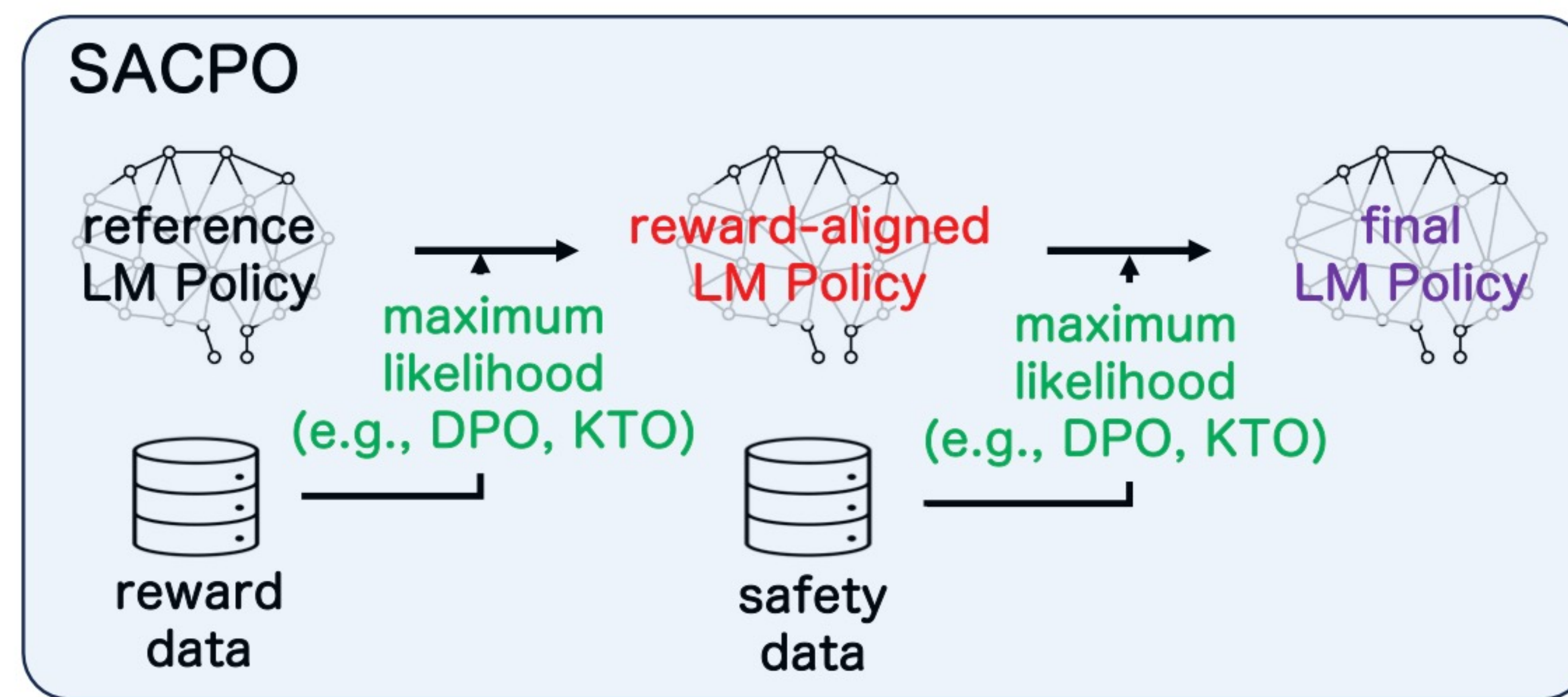
$$\max_{\pi} \mathbb{E}_{\rho, \pi} [r^*(x, y)] - \beta \mathbb{D}_{\text{KL}}[\pi(y | x) \| \pi_{\text{ref}}(y | x)] \quad \text{subject to} \quad \mathbb{E}_{\rho, \pi} [g^*(x, y)] \geq b$$

Typical objective of RLHF or DPO
Safety constraint

SACPO: Stepwise Alignment for Constrained Policy Optimization

Key Idea: Reward alignment → Safety Alignment (or vice versa)

We can use RL-free alignment algorithms (e.g., DPO, KTO) for each alignment



SACPO's stepwise approach is theoretically justified!

SACPO is computationally efficient and stable!

Detailed Steps

Step 1: Reward Alignment

- Align an LM reference policy using reward data via an RL-free alignment algorithm (e.g., DPO, KTO)
- This step is same as typical alignment by DPO or KTO. For example,

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

Step 2: Safety Realignment

- Realign the reward-aligned LM policy using safety data using DPO or KTO
- Note : λ^* is the optimal Lagrangian multiplier

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}, \pi_{r^*}, \beta/\lambda^*) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}_g} \left[\log \sigma \left(\frac{\beta}{\lambda^*} \log \frac{\pi_{\theta}(y_w | x)}{\pi_{r^*}(y_w | x)} - \frac{\beta}{\lambda^*} \log \frac{\pi_{\theta}(y_l | x)}{\pi_{r^*}(y_l | x)} \right) \right]$$

	Step1 (reward alignment)	Step 2 (safety realignment)
Loss function	$\mathcal{L}_{\text{DPO}}(\pi_{\theta}, \pi_{\text{ref}}, \beta)$	$\mathcal{L}_{\text{DPO}}(\pi_{\theta}, \pi_{r^*}, \beta/\lambda^*)$
LM policy to be aligned	π_{ref} (typically SFT models)	π_{r^*} (reward-aligned LM)
KL penalty parameter	β	β/λ^*

Why is SACPO Theoretically Justified?

The optimal policy of our constrained LM alignment problem satisfies:

$$\pi^*(y | x) \propto \pi_{\text{ref}}(y | x) \exp \left(\frac{1}{\beta} (r^*(x, y) + \lambda^* g^*(x, y)) \right)$$

$$\propto \pi_{r^*}^*(y | x) \exp \left(\frac{\lambda^*}{\beta} g^*(x, y) \right)$$

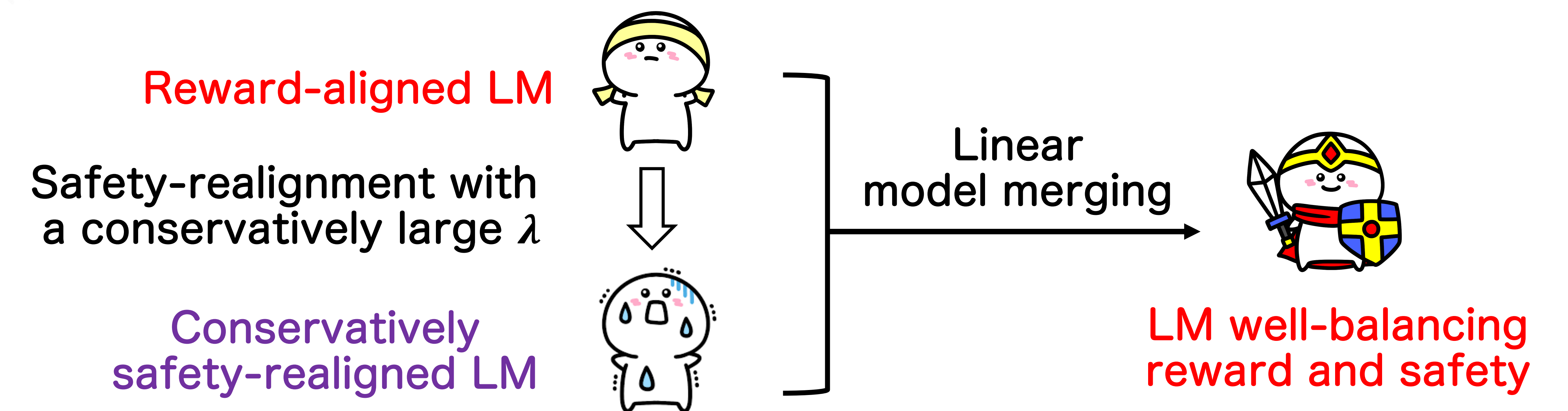
Reward-aligned LM policy
Safety function

The optimal policy can be obtained by realigning reward-aligned LM policy regarding the safety function with a KL penalty parameter β/λ^* .

Practical SACPO (P-SACPO)

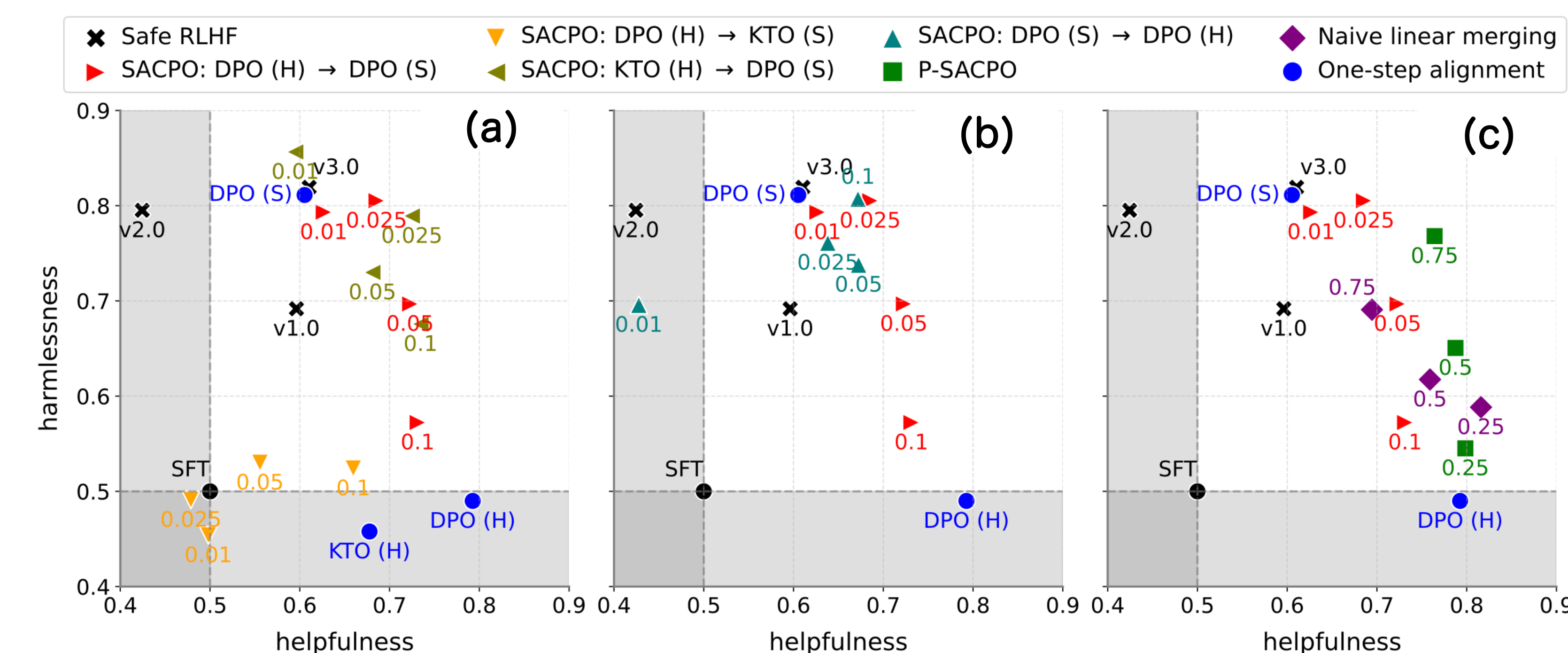
It is still costly to apply DPO (or KTO) w/ various β/λ^* in SACPO

Control the balance between reward and safety by tuning the merging ratio!



Experiments

- Compare the performance of (P-)SACPO with Safe RLHF
- Reward = Helpfulness (H), Safety = Harmlessness (S)
- SFT model = Alpaca-7b, Dataset = PKU-SafeRLHF-30K



Win-rate against SFT model. In (a) and (b), the numbers indicate β/λ . In (c), the numbers for the red triangles represent β/λ , while those for the green and purple squares represent the model merging ratio.

Results

- SACPO well-balances helpfulness and safety, which performs better than Safe RLHF!
- In SACPO, difference alignment algorithms can be used (this is consistent with theory)
- P-SACPO empirically performs well with reduced computational time and stable learning