

# Implicit Safety Alignment from Crowd Preferences

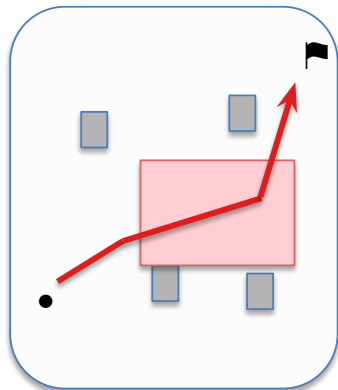
Qian Lin, Daniel S. Brown  
Kahlert School of Computing, University of Utah

ICML 2026

# Motivation / Problem Setting

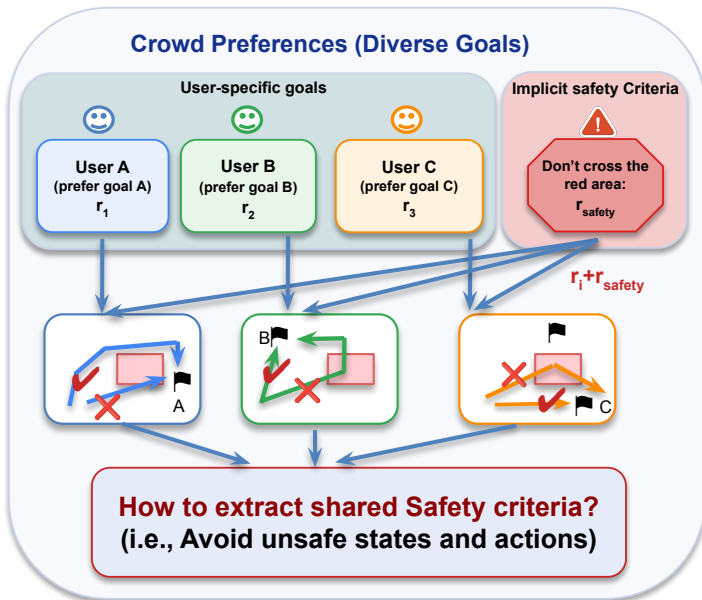
From crowd preferences, we extract shared safety and transfer to new tasks

**Problem:**  
Using safety-agnostic task reward leads to unsafe behavior

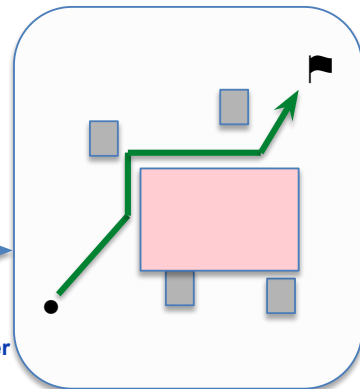


Policy optimized by task reward only  
→ High task reward, but unsafe

**Solution:** Crowd preferences from diverse users with shared safety



**Goal:** Combining shared safety with task reward yields safe behavior

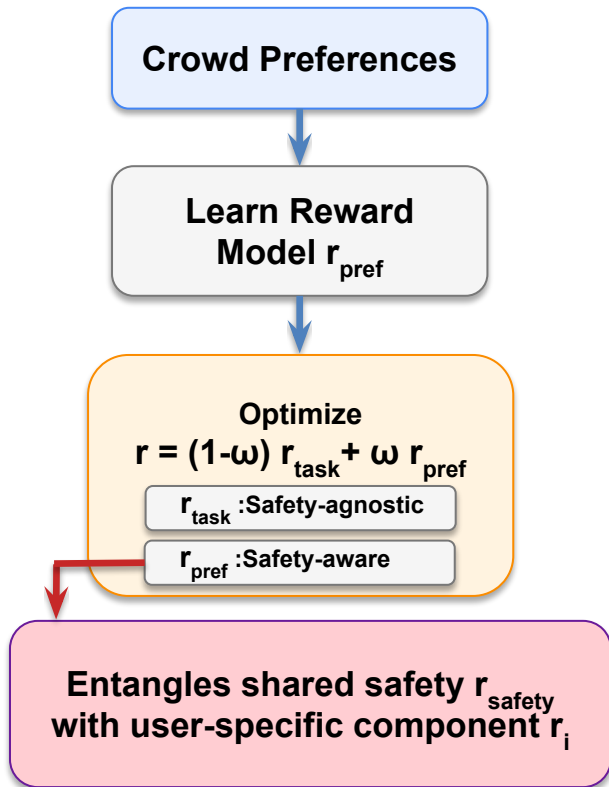


Policy with task reward + shared safety criteria  
→ High task reward and safe

● Start    🚩 Goal    ■ Obstacle    □ Danger Zone    --- Unsafe Trajectory    --- Safe Trajectory

# Why Vanilla RLHF Fails

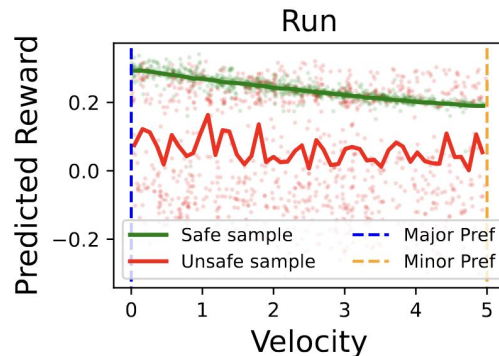
## Natural solution: reward combination



## Two major limitations

### 1. Preference imbalance bias

- Safe trajectories get higher reward
- Majority preferences get over-rewarded



### 2. Sensitive to $\omega$

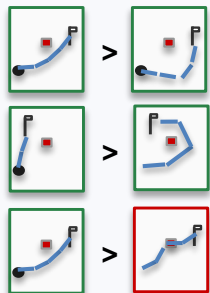
- Small  $\omega$   $\rightarrow$  unsafe
- Large  $\omega$   $\rightarrow$  poor task performance

# Our Method: Composing Safe Skills

Instead of combining rewards, we compose safe behaviors (skills)

## Stage 1: Skill Discovery

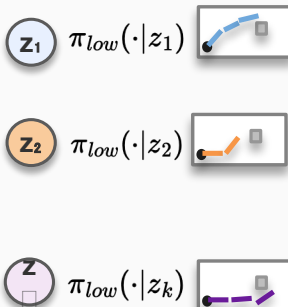
Crowd Preference Data



(A) Safe VPL

VAE Reward Model  $r(\cdot, z)^{[1]}$

Preference-Aligned (Safe) Skill Set



(B) Safe CPL: Directly learn VAE policy via CPL<sup>[2]</sup>

Directly learn safe, preference-aligned skills from crowd preferences, instead of learning a reward

transfer safety info

## Stage 2: Skill Combination

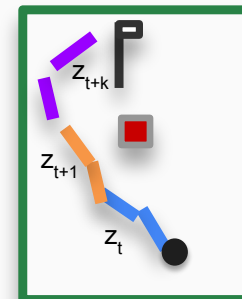
Optimize Task Reward  $r_{\text{task}}(s, a)$

Freeze  $\pi_{\text{low}}(\cdot | z)$

Train High-level Policy  $\pi_{\text{high}}(z | s)$

Regularize  $z$  toward the prior to avoid OOD unsafe skills

Resulting Behavior (Safe & Task-Aligned)



A high-level policy selects among safe skills to maximize task reward.

key idea: the high-level policy solves the task, and the low-level skills provide safety.

1. Poddar, Sriyash, et al. "Personalizing reinforcement learning from human feedback with variational preference learning." NeurIPS 2024.

2. Hejna, Joey, et al. "Contrastive preference learning: Learning from human feedback without reinforcement learning." ICLR 2024.

# Experimental Results

Safe skill composition lowers safety cost while keeping task performance comparable to oracle.

## Evaluation Setup



### Six safe RL environments

Reach, Run, Circle, Ant-vel,  
Swimmer-vel, HalfCheetah-vel



### Proof-of-concept LLM-style evaluation

Avoid harmful responses



### Preference data

Different goals, shared safety  
(e.g., users prefer different navigation  
targets but consistently avoid unsafe areas)



### Downstream Task

A safety-agnostic reward  
specifies an unseen goal



### Metrics

- Normalized reward  $\uparrow$
- Normalized safety cost  $\downarrow$



### Baselines

Task-only, SOPL, Oracle,  
Reward Combination (RC),

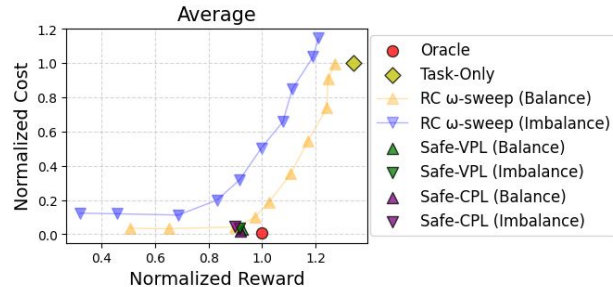
**Ours: Safe-VPL / Safe-CPL**

## Main Results (Average over 6 Environments)

Stats	Oracle	Task-Only	SOPL	RC( $\omega = 0.5$ )	Safe-VPL	Safe-CPL
Norm Rew	1.00	1.46	1.04	0.82	0.93	0.92
Norm Cost	0.01	1.00	0.01	0.05	0.03	0.02

- ✘ Task-Only  $\rightarrow$  severe safety violations
- ✘ Oracle / SOPL  $\rightarrow$  strong, but require explicit safety reward or preference
- ✘ RC  $\rightarrow$  struggles to balance safety and task performance
- ✓ Ours  $\rightarrow$  comparable reward to Oracle / SOPL with  $\sim 97\%$  lower safety cost than Task-Only

## Sensitivity to preference imbalance



- ✘ RC is sensitive to the trade-off weight  $\omega$  and degrades under preference imbalance.
- ✓ Under both balanced and imbalanced preferences, our methods stay on the low-cost frontier and are robust to imbalance.