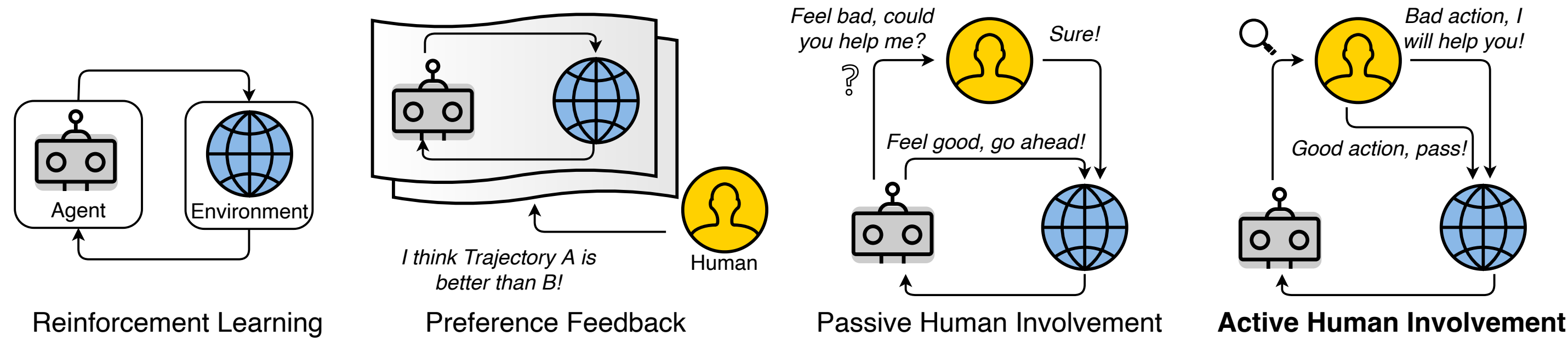




## Problem Formulation & Motivation

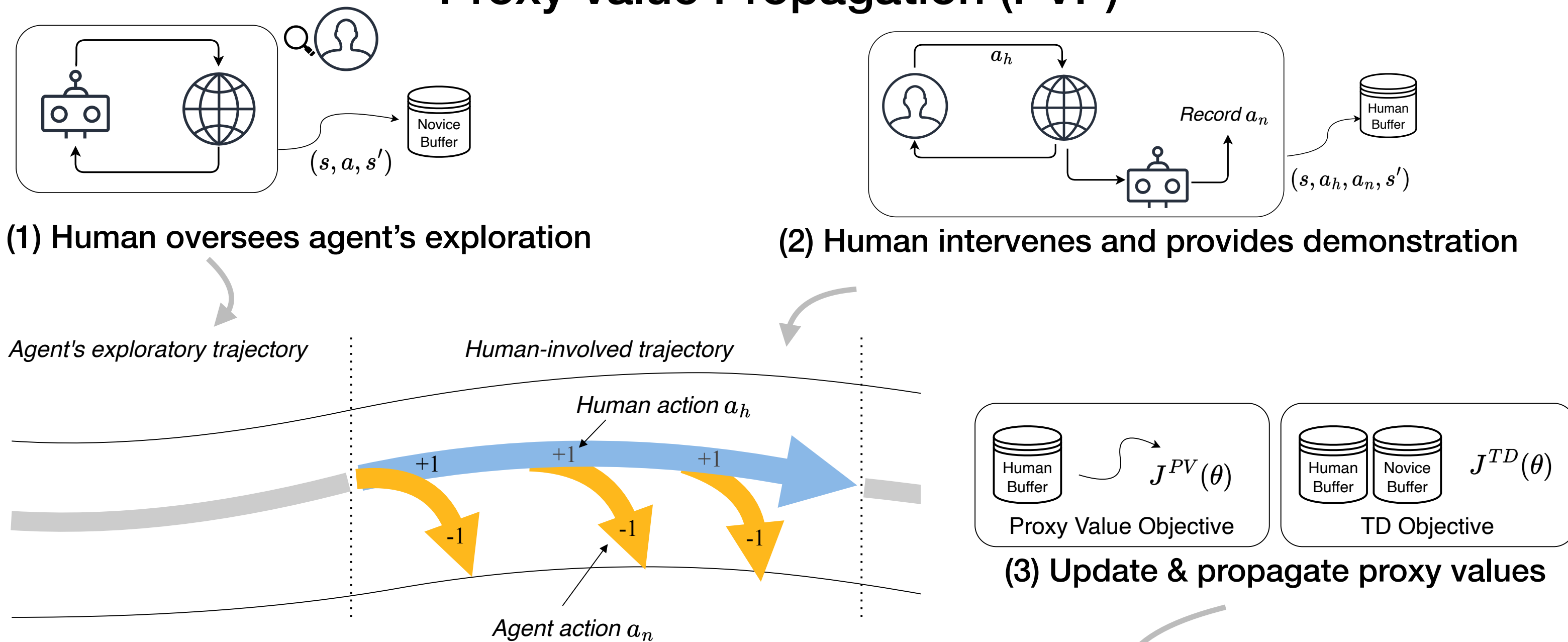


- Reward engineering is hard to encapsulate human intentions.
- **Human-in-the-loop** methods are promising to achieve alignment.
- To ensure safety, **active human involvement** enhances training-time safety.

Can we transform a value-based RL algorithm to a reward-free policy trainer learning from **online intervention & demonstration** from human expert?

**Policy learning in 10 minutes w/o reward via human-in-the-loop!**

## Proxy Value Propagation (PVP)



**Proxy Value Objective:** Apply loss to human-involved states

$$J^{PV}(\theta) = \mathbb{E}_{(s, a_n, a_h)} [ |Q_\theta(s, a_h) - 1|^2 + |Q_\theta(s, a_n) + 1|^2 ]$$

**Temporal Difference Objective:** Drop the reward term!

$$J^{TD}(\theta) = \mathbb{E}_{(s, a, s')} |Q_\theta(s, a) - \gamma \max_{a'} Q_{\hat{\theta}}(s', a')|^2$$

**Total Objective for Q network:**

$$J(\theta) = J^{PV}(\theta) + J^{TD}(\theta)$$

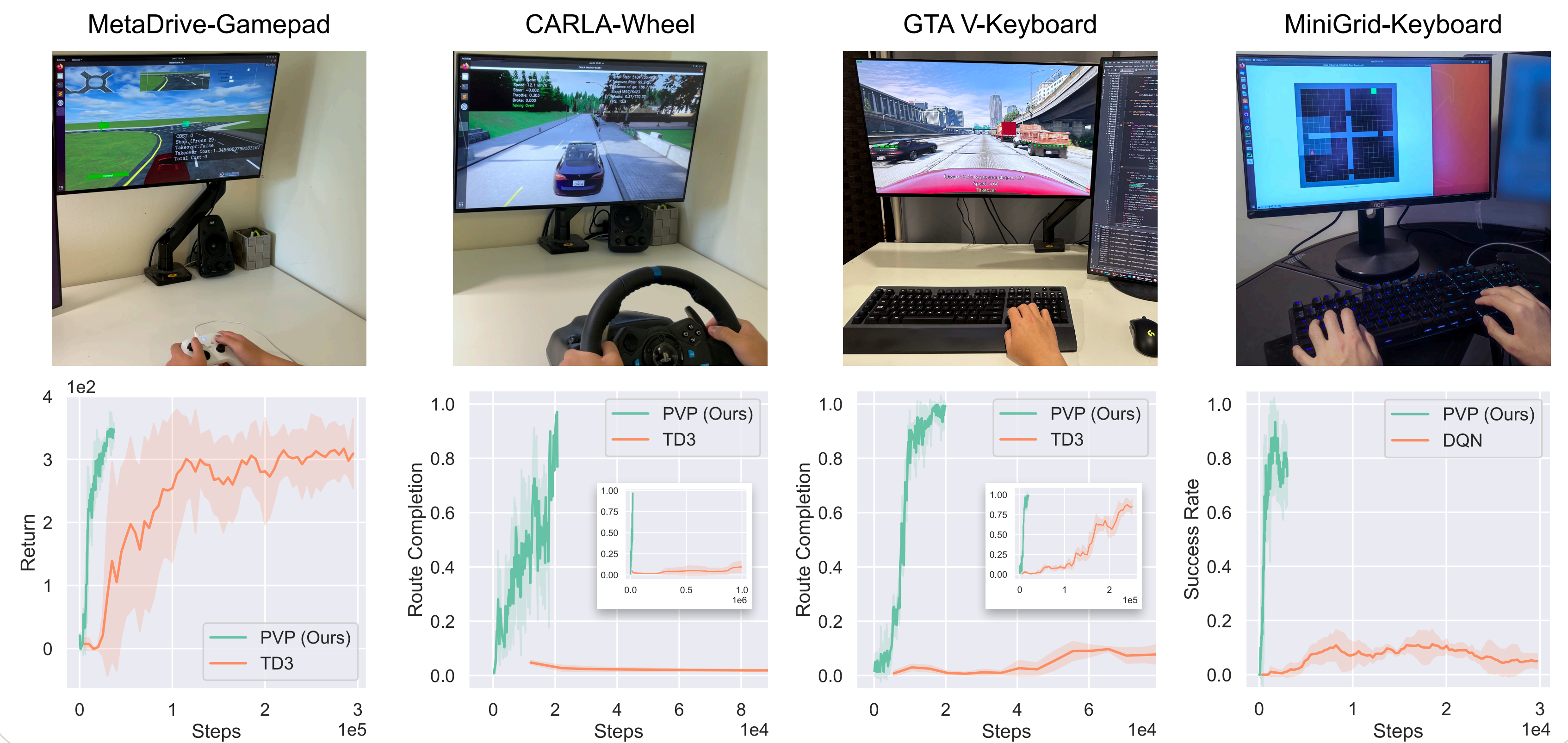
**Objective for Policy:**

$$J(\phi) = -\mathbb{E}_s Q_\theta(s, \pi_\phi(s))$$

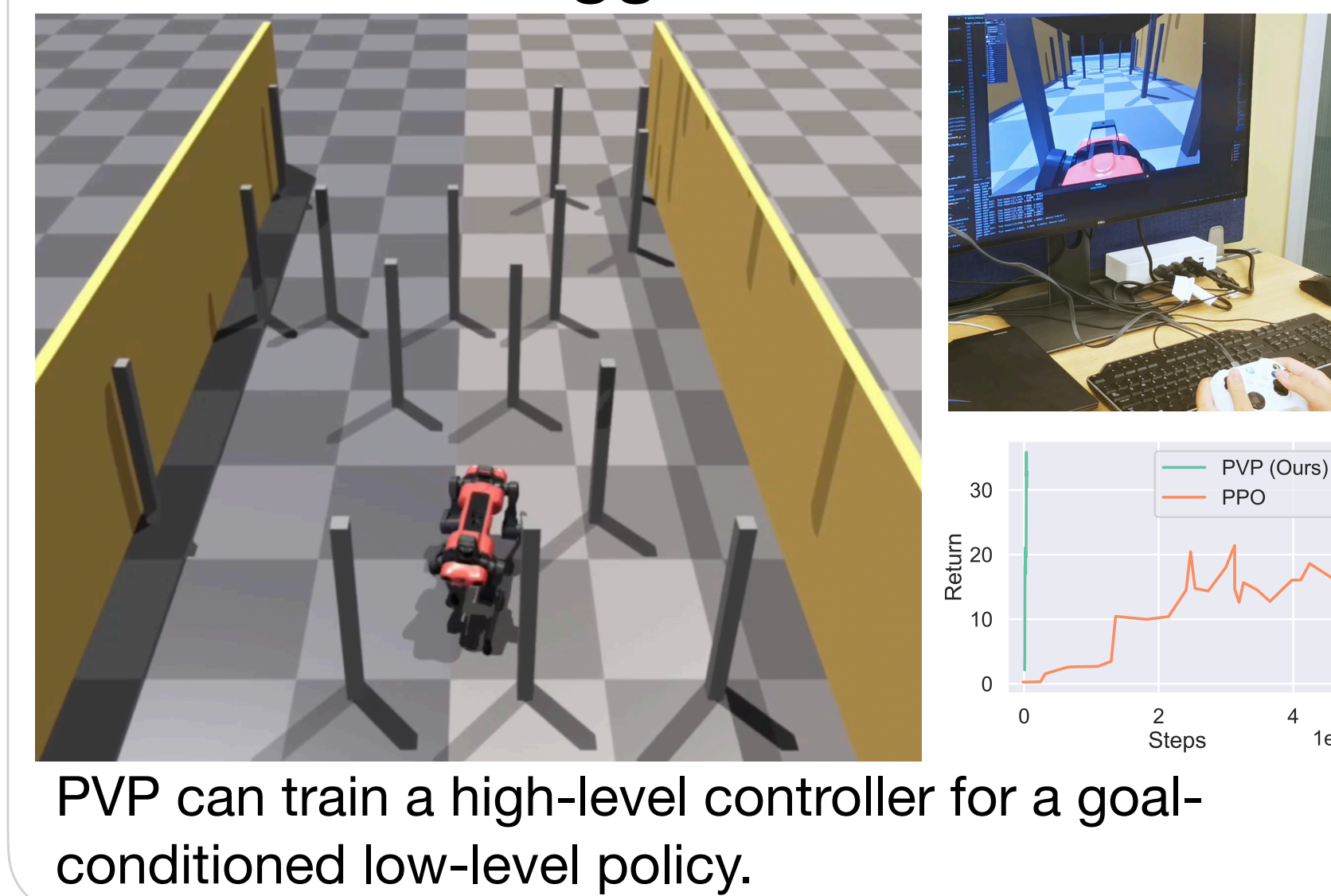
$Q_\theta$  - Q network  
 $a_n$  - agent's action  
 $a_h$  - human's action  
 $s$  - state  
 $\pi_\phi$  - agent's policy

Code is available at [metadriverse.github.io/pvp](https://metadriverse.github.io/pvp)

## 👏 10x Faster than RL Baseline



## 🐕 Legged Robot



PVP can train a high-level controller for a goal-conditioned low-level policy.

## 😊 Human Friendly

User Study	HGDagger	IWR	HACO	PVP
Compliance	3.0	4.0	3.0	<b>4.8</b>
Performance	2.2	3.7	3.3	<b>4.8</b>
Stress	3.2	4.5	2.3	<b>4.7</b>

PVP agents make human feels better (*compliance*), stronger (*performance*) and less stressful (*stress*) in shared control. It also makes human takes over less (right table).

## 🎉 Safety, Efficiency, Performance

MetaDrive-Keyboard Env	Training			Testing	
	Human Data	Total Data	Total SafetyCost	Episodic SafetyCost	Success Rate
SAC	—	1M	2.76K	0.73	0.82
PPO	—	1M	24.34K	3.41	0.69
TD3	—	1M	1.74K	<b>0.47</b>	0.70
SAC-Lag	—	1M	1.84K	0.72	0.73
PPO-Lag	—	1M	11.64K	1.18	0.51
CPO	—	1M	4.36K	1.71	0.21
<b>HumanDemo</b>	<b>30K</b>	-	<b>39</b>	<b>0.39</b>	<b>0.97</b>
BC	30K	-	—	2.17	0.07
GAIL	30K	2M	25.90K	1.30	0.0
HGDagger	39.0K	51K	56	1.97	0.04
IWR	35.8K	45K	<b>52</b>	1.45	0.46
HACO	19.2K	40K	130	1.64	0.13
<b>PVP (Ours)</b>	<b>14.6K</b>	<b>40K</b>	<b>76.8</b>	<b>0.89</b>	<b>0.85</b>

- Human-in-the-loop makes training much safer, even than Safe RL baselines.
- PVP learns most performant agent.