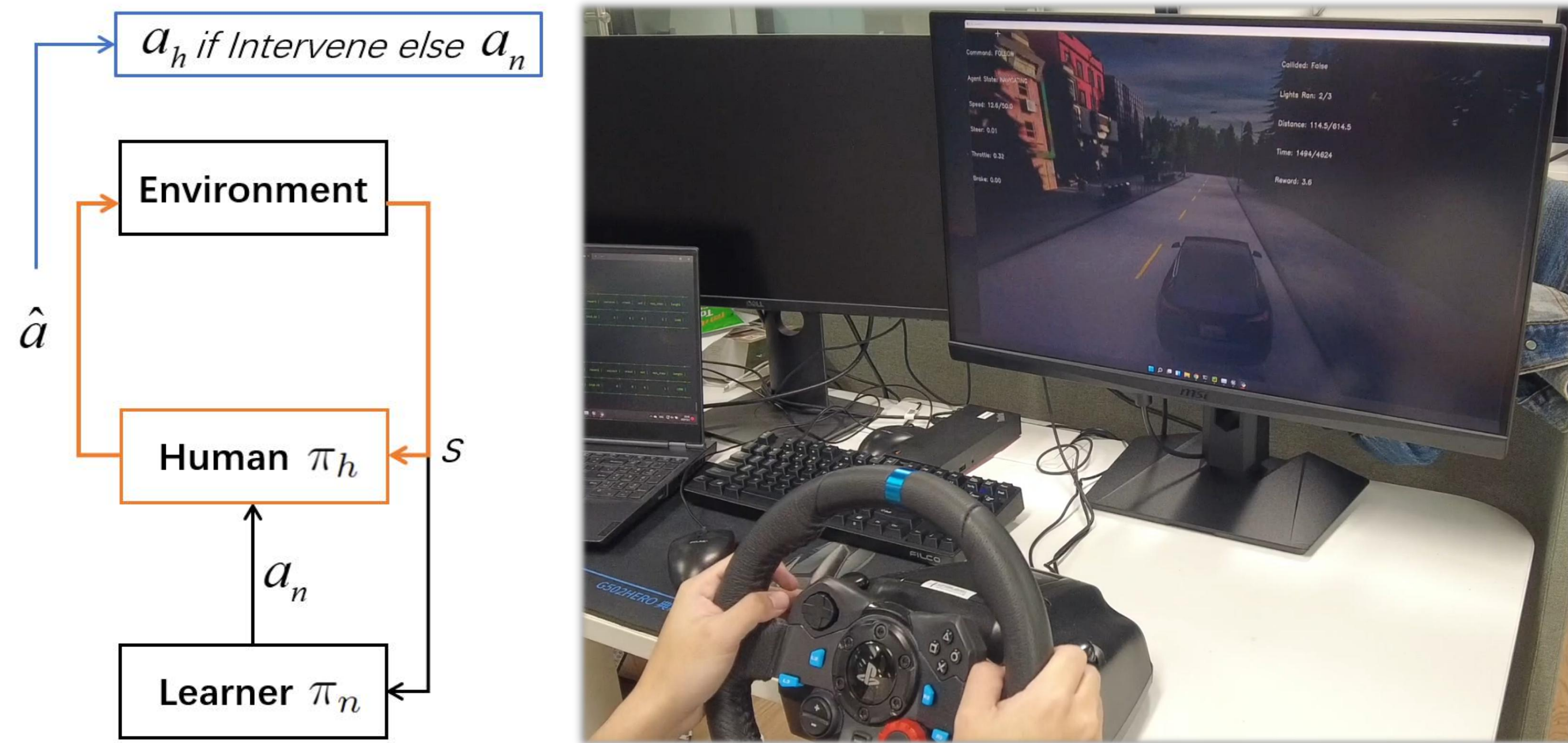
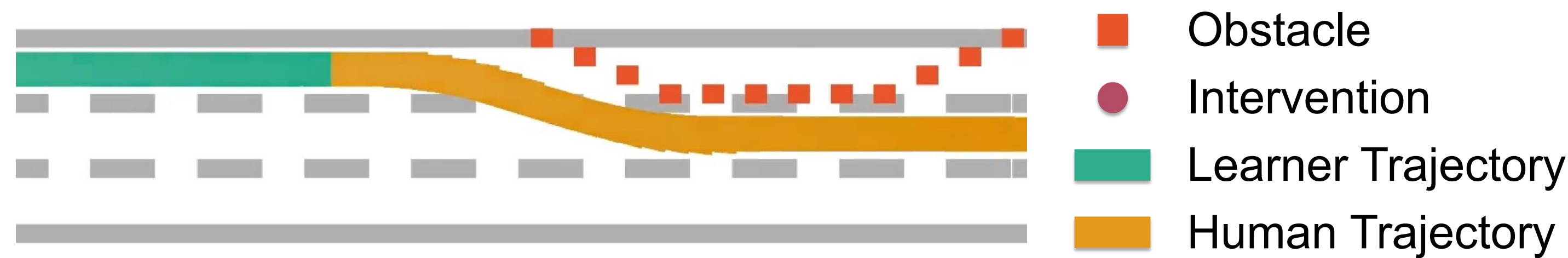


## Human-in-the-loop RL



$a_h \sim \pi_h$  Human policy     $a_n \sim \pi_n$  Learner policy     $\hat{a}$  Behavior action

We authorize human expert to take over or intervene, when Reinforcement Learning (RL) agents are in training. This paradigm is referred as **Human-in-the loop RL**.



Copilot trajectory

When training with the learner policy, human's duties are:

1. **Protecting** the agent as a guardian
2. **Teaching** the agent by providing demonstration

## Human-AI-Copilot Optimization (HACO)

### 1. Learning from Demonstration

HACO learns from human-provided demonstrations by applying CQL loss to train proxy value function:

$$\min_{\phi} \mathbb{E}_{\mathcal{B}} [I(s, a_n) (Q(s, a_n; \phi) - Q(s, a_h; \phi))]$$

which is updated through the TD-target.

### 2. Intervention Minimization

To minimize intervention, HACO additionally learns a intervention cost value function to estimate expected accumulative intervention cost:

$$Q^I(s, a_n) = C(s, a_n) + \gamma \mathbb{E}_{a' \sim \pi_n(\cdot|s')} [Q^I(s', a')]$$

The intervention cost is calculated by cosine similarity:

$$C(s, a_n) = 1 - \frac{a_n^T a_h}{\|a_n\| \|a_h\|}, \quad a_h \sim \pi_h(\cdot|s)$$

### 3. Policy Optimization

The policy optimization goal is to maximize proxy value and minimize the intervention cost:

$$\max_{\theta} \mathbb{E}_{\mathcal{B}} [Q(s, a_n) - Q^I(s, a_n)], \quad a_n \sim \pi_n(\cdot|s; \theta)$$

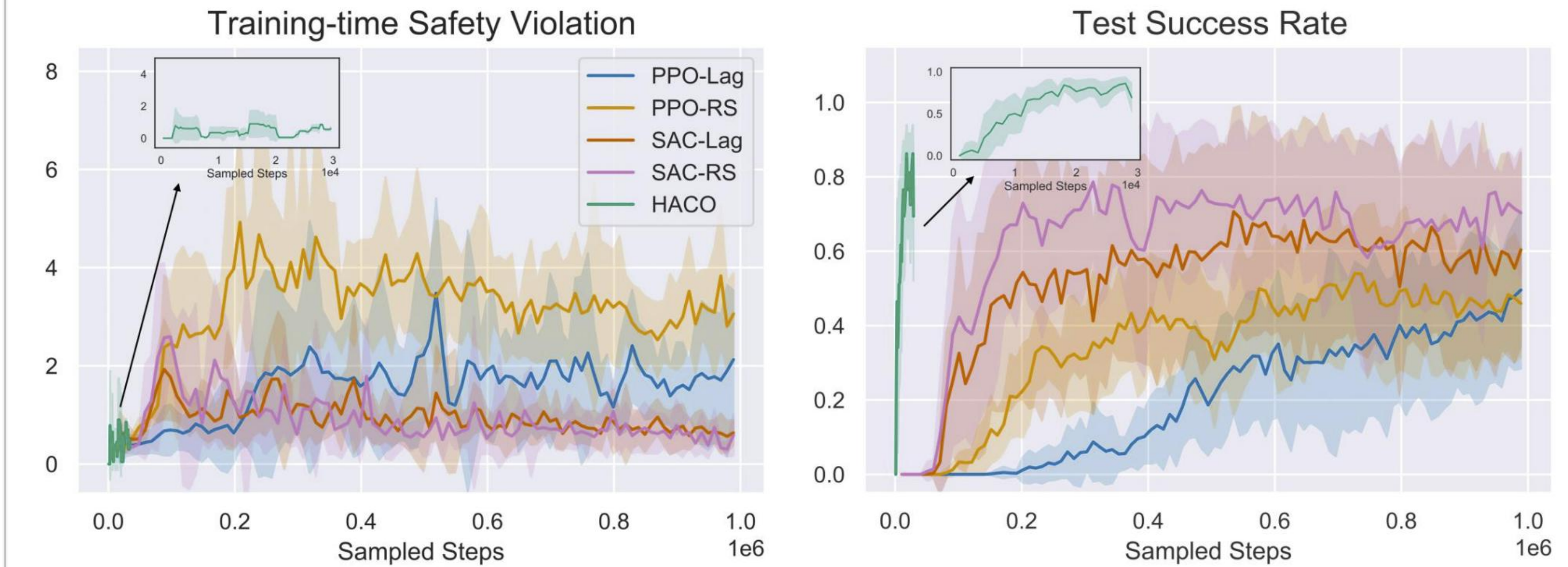
## Experiment Results

### MetaDrive Safe RL Environment

We evaluate HACO on safe RL suite of MetaDrive:



### Comparison with RL baselines



Category	Method	Total Training Safety Violation	Training Data Usage	Test Return	Test Safety Violation	Test Success Rate
RL	SAC-RS	2.76K $\pm$ 0.95K	1M	<b>386.77</b> $\pm$ 35.1	0.73 $\pm$ 1.18	0.82 $\pm$ 0.18
	PPO-RS	24.34K $\pm$ 3.56K	1M	335.39 $\pm$ 12.41	3.41 $\pm$ 1.11	0.69 $\pm$ 0.08
Safe RL	SAC-Lag	1.84K $\pm$ 0.49K	1M	351.96 $\pm$ 101.88	<b>0.72</b> $\pm$ 0.49	0.73 $\pm$ 0.29
	PPO-Lag	11.64K $\pm$ 4.16K	1M	299.99 $\pm$ 49.46	1.18 $\pm$ 0.83	0.51 $\pm$ 0.17
	CPO	4.36K $\pm$ 2.22K	1M	194.06 $\pm$ 108.86	1.71 $\pm$ 1.02	0.21 $\pm$ 0.29
Ours	<b>HACO</b>	<b>30.14</b> $\pm$ 11.36	<b>30K*</b>	349.25 $\pm$ 11.45	0.79 $\pm$ 0.31	<b>0.83</b> $\pm$ 0.04

Comparison with IL/Offline RL, CARLA Experiments and videos is available at: <https://decisionforce.github.io/HACO/>