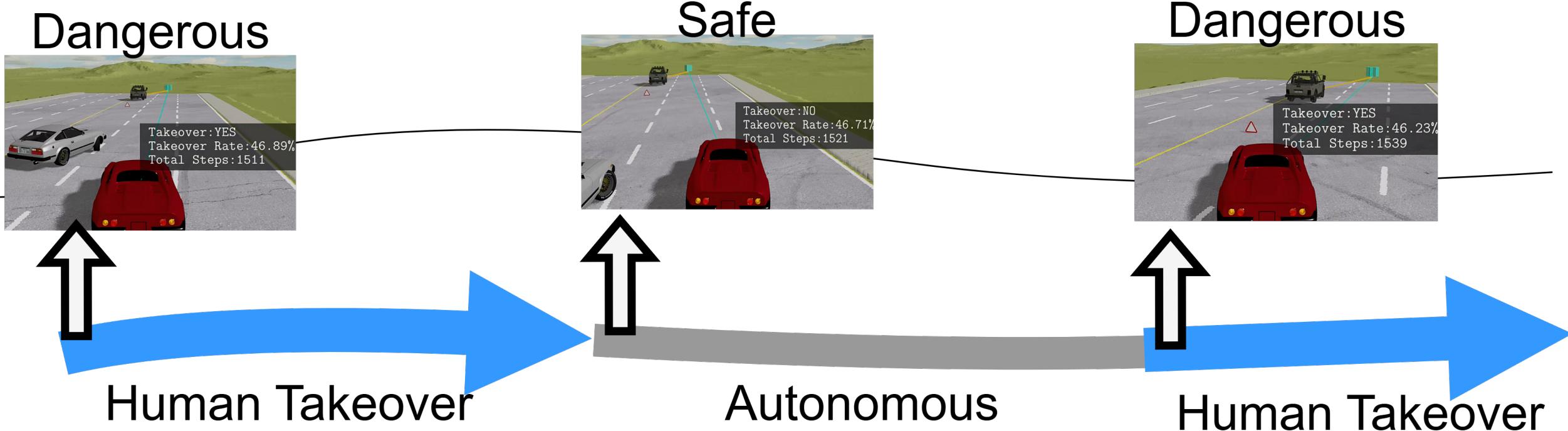




### Interactive Imitation Learning: Formulation & Motivation



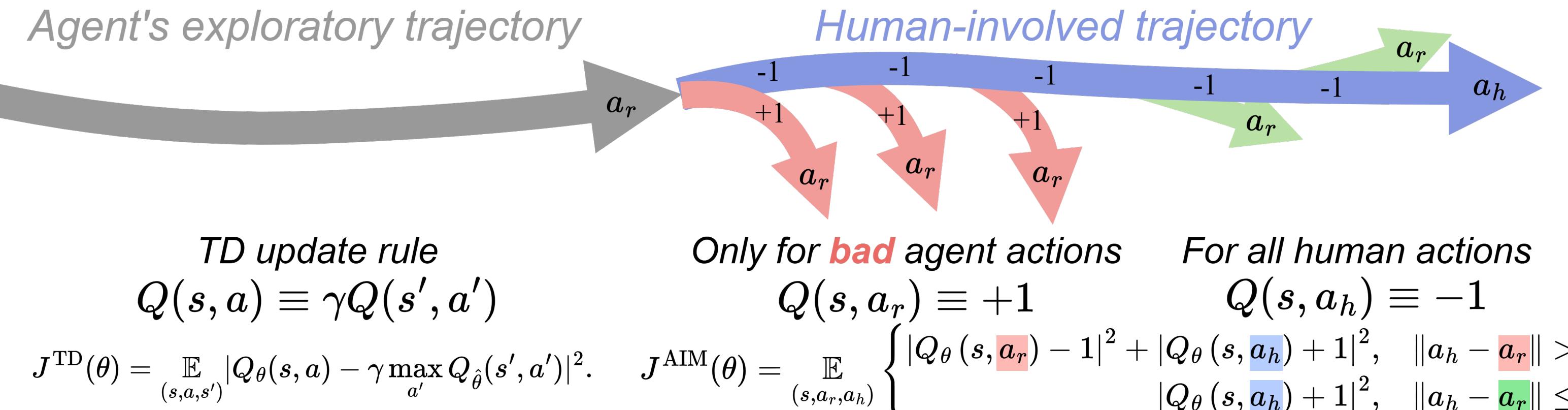
#### Human-Gated Intervention:



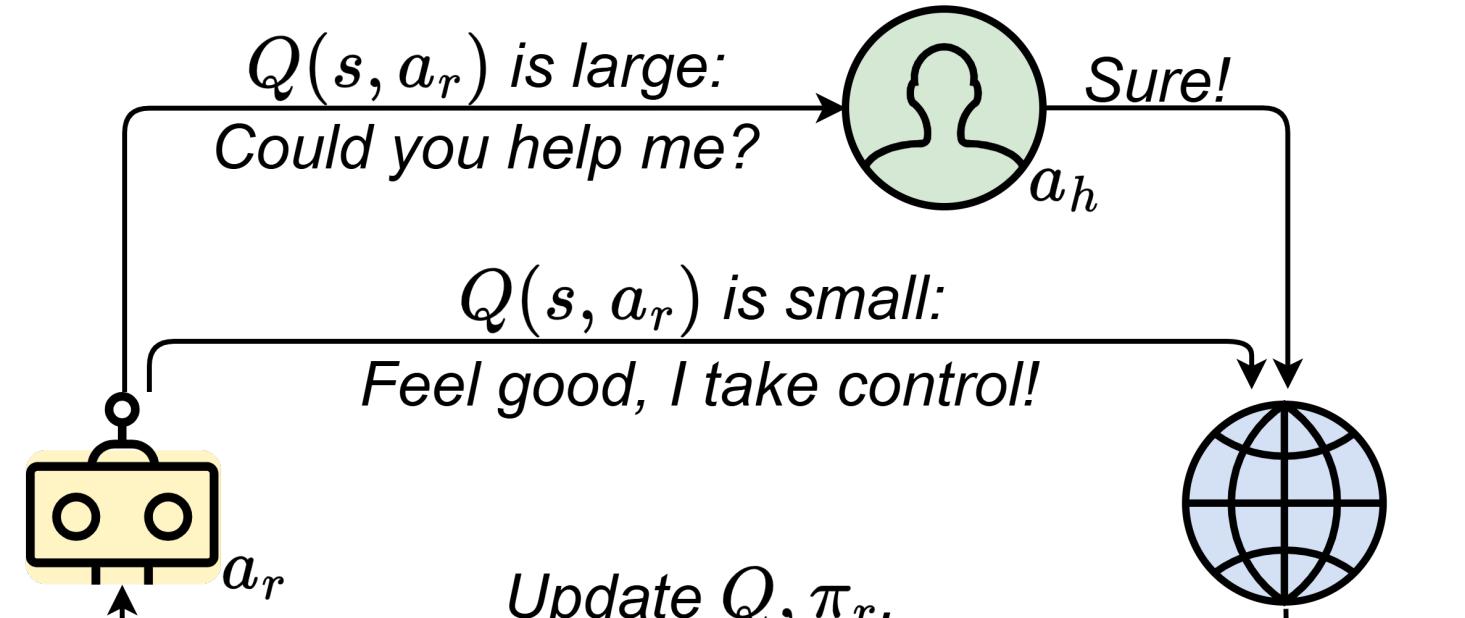
- Human-Gated**: human watch every step → **heavy cognitive load**.
- Robot-Gated**: agent **ask help when needed** → saves human effort.

**Can we build a robot-gated intervention mechanism that reliably detects danger and requests help before a crash?**

#### Adaptive Intervention Mechanism (AIM)



**Value-Based Intervention Mechanism:** Request human help when  $Q_\theta(s, a_r) > (1 - \delta)$ -th quantile



Total objective for Q network

$$J(\theta) = J^{\text{AIM}}(\theta) + J^{\text{TD}}(\theta).$$

BC loss for policy update

$$J(\pi_r) = - \mathbb{E}_{(s, a_h)} \|\pi_r(s) - a_h\|^2.$$

$Q_\theta$  - proxy Q network

$a_r$  - agent action

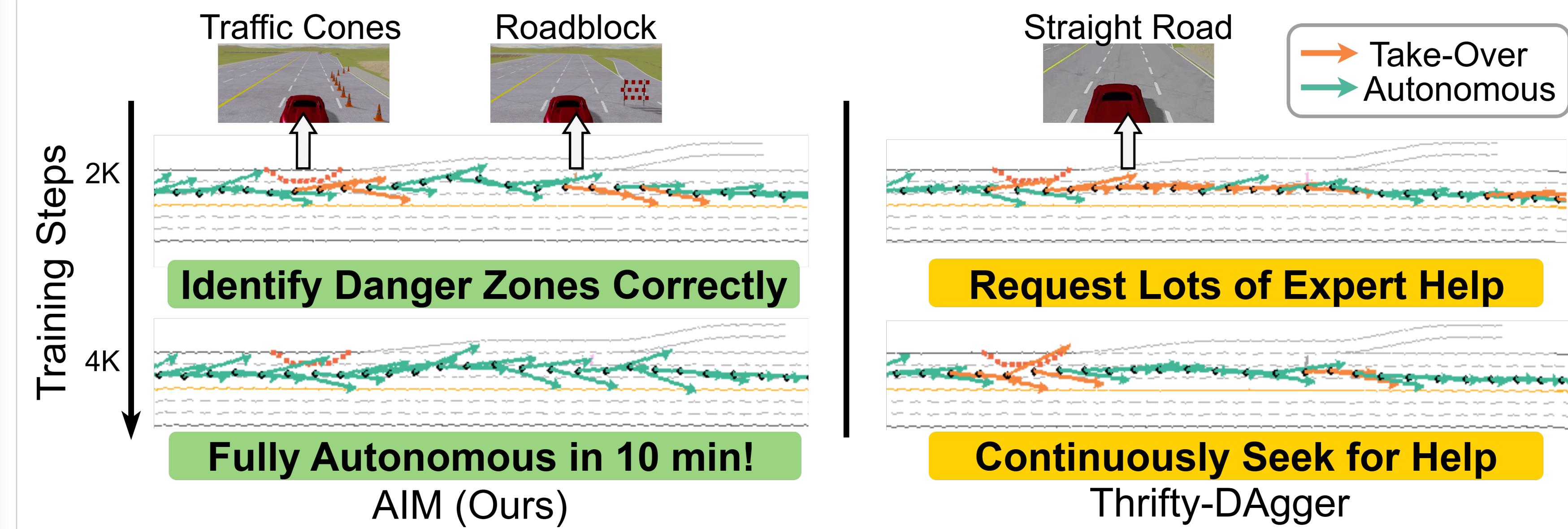
$a_h$  - human action

$s$  - state

$\pi_r$  - agent policy

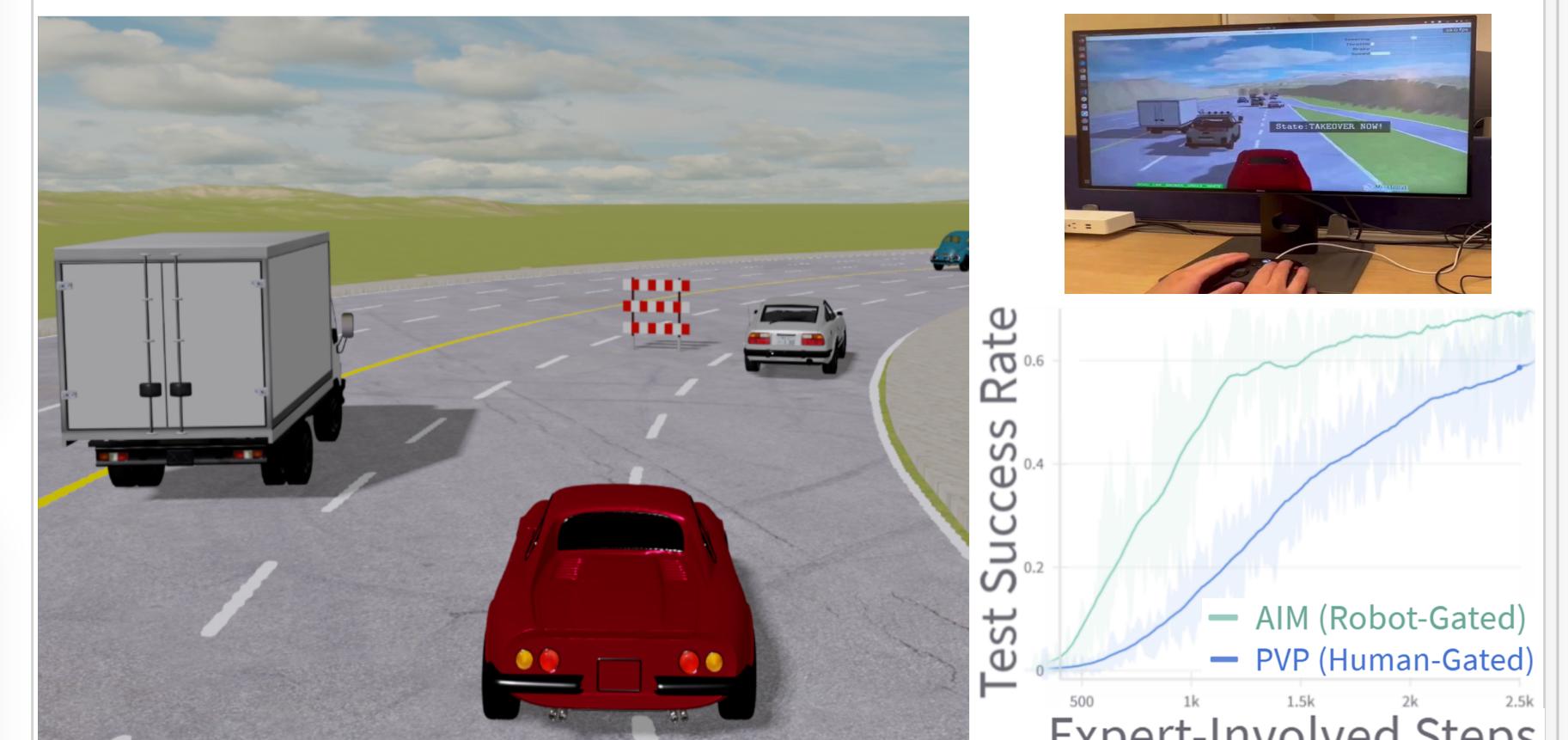
$\delta$  - switching threshold

### Request Less Expert Help Than Robot-Gated Baselines



- AIM precisely pinpoints where expert intervention is most necessary.
- AIM adaptively reduces help requests as the agent becomes proficient.

### MetaDrive Safety Benchmark



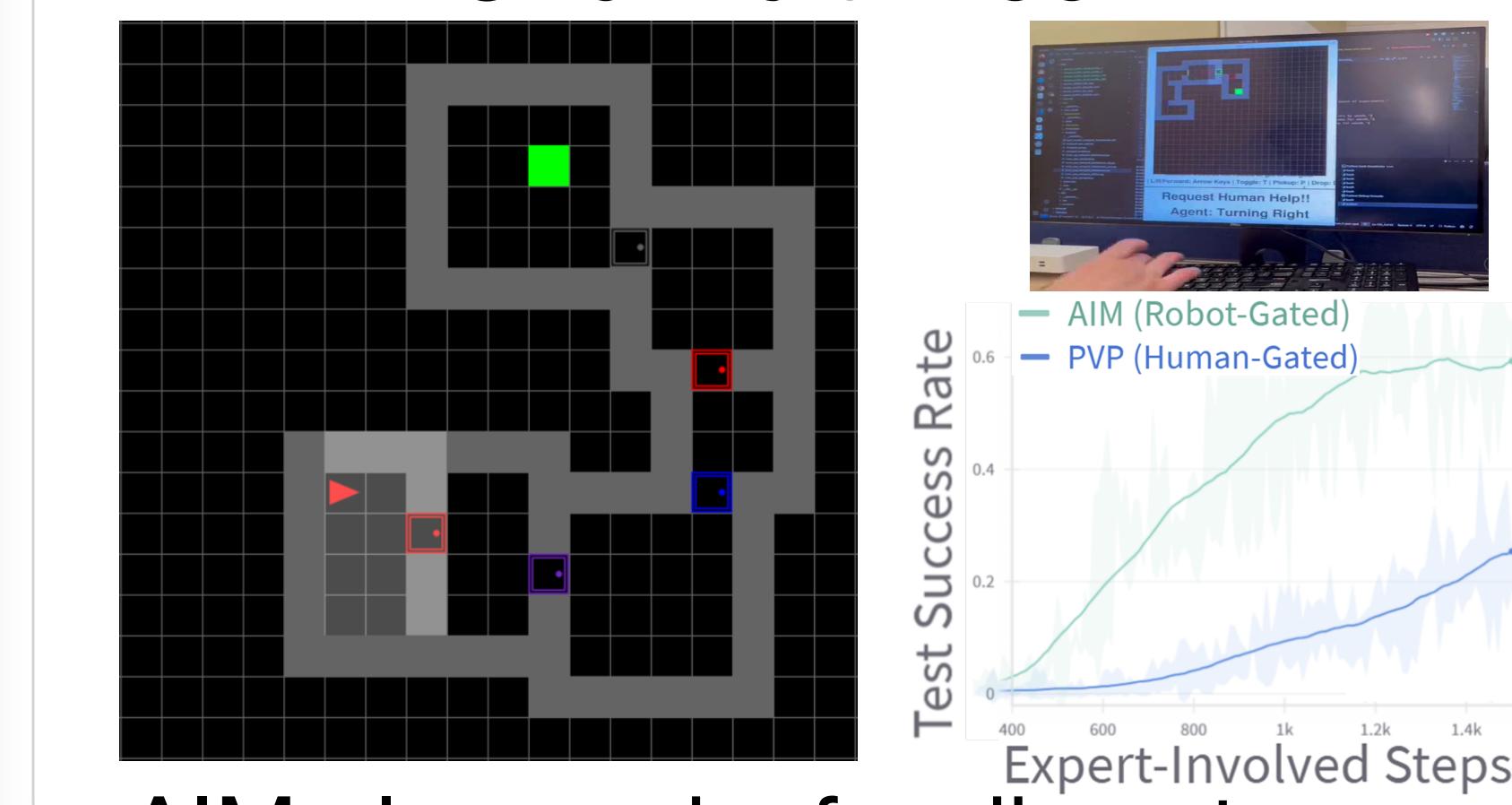
AIM saves expert cognitive effort.

### Efficiency & Performance

MetaDrive Env	Training		Testing	
	Expert-Involved Steps	Take-Over Rate	Episodic Return	Success Rate
Neural Expert	-	-	336.5	0.84
BC	2K	-	243.0	0.33
HG-Dagger	2K	0.45	310.8	0.61
PVP	2K	0.19	270.4	0.62
Ensemble-Dagger	2K	0.55	267.4	0.60
Thrifty-Dagger	2K	0.21	250.0	0.58
AIM (Ours)	1.9K	0.24	328.4	0.82

MiniGrid Env	Training		Testing	
	Expert-Involved Steps	Take-Over Rate	Total Steps	Success Rate
Neural Expert	-	-	-	0.78
BC	2K	-	-	0.01
HG-Dagger	2K	0.12	-	0.20
PVP	2K	0.12	-	0.34
Ensemble-Dagger	2K	0.36	5.6K	0.38
Thrifty-Dagger	2K	0.27	7.4K	0.42
AIM (Ours)	0.4K	0.09	4.4K	0.63

### MiniGrid Multi-Room Env



AIM also works for discrete env.

- AIM learns **most performance** agent with **least expert involvement**.
- AIM outperforms both human-gated and robot-gated imitation baselines.

Code is available at [metadrive.github.io/aim](https://metadrive.github.io/aim)