# Learning Realistic Traffic Agents in Closed-loop

Chris Zhang, James Tu, Lunjun Zhang, Kelvin Wong, Simon Suo, Raquel Urtasun

UNIVERSITY OF TORONTO

## Traffic Simulation for Self-driving

**Motivation:** Developing self-driving in simulation is safer and more scalable than driving purely in the real world.

**Goal:** Learn models of how humans drive in order to use them as actor models in simulation.

**Task:** Given environmental information (e.g. high definition map, current actor positions and velocity), control how each actor should behave subsequently.

## Challenges and Existing Work

**Realistic** actor models must:
1. Capture nuances of **human driving**
2. **Avoid infractions** like collisions or driving off-road

Existing approaches have shortcomings which can result in a **trade-off** between the two.

**Imitation Learning:**
✔ Leverages offline data for realism
✘ No explicit knowledge of infractions

**Reinforcement Learning:**
✔ Explicit reward signal
✘ Manual reward design lacks realism

## Learning Objective

We model the problem with an MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, R, P, \gamma)$

A trajectory $\tau_{0:T} = (s_0, a_0, \ldots, s_{T-1}, a_{T-1}, s_T)$ is a state action sequence for all agents in the scene.

We aim to recover the expert distribution while satisfying an infraction-based constraint:
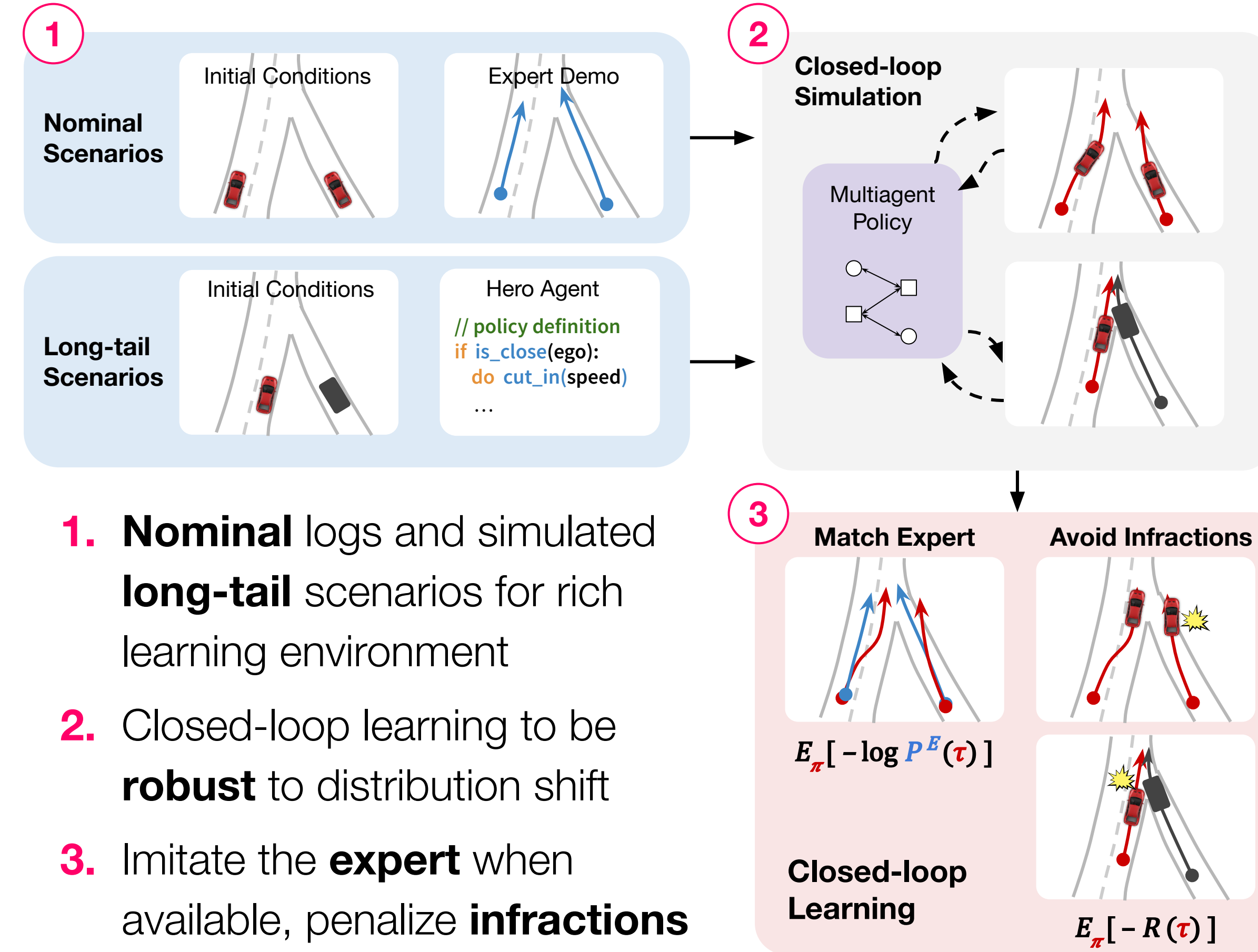
$$\arg\min_{\pi} \quad D_{\text{KL}}\left(P^{\pi}(\tau) \parallel P^{E}(\tau)\right)$$
$$\text{s.t.} \quad \mathbb{E}_{P^{\pi}}[R(\tau)] \geq 0 \qquad R^{(i)}(s, a^{(i)}) = \begin{cases} -1 & \text{if infraction} \\ 0 & \text{otherwise,} \end{cases}$$

Taking the Lagrangian decomposes the objective into a combination of imitation and reinforcement learning:

$$\mathcal{L} = \mathbb{E}_{P^{\pi}}\left[\underbrace{-\log P^{E}(\tau)}_{\text{IL}} - \lambda \underbrace{R(\tau)}_{\text{RL}}\right] - H(\pi)$$
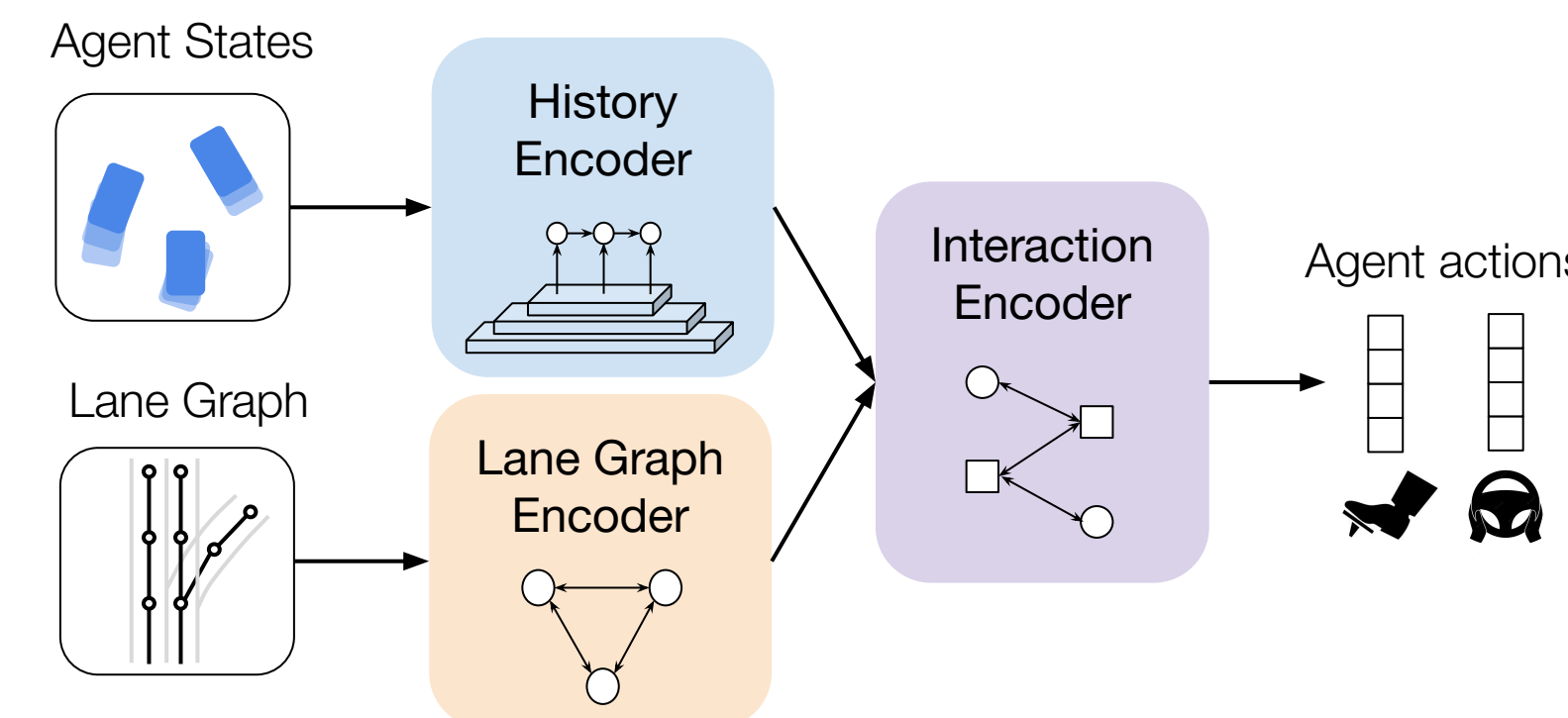
## Reinforcing Traffic Rules (RTR)

We **combine RL and IL** to learn robust policies in closed-loop.



1. **Nominal** logs and simulated **long-tail** scenarios for rich learning environment
2. Closed-loop learning to be **robust** to distribution shift
3. Imitate the **expert** when available, penalize **infractions**

$$E_{\pi}[-\log P^{E}(\tau)]$$

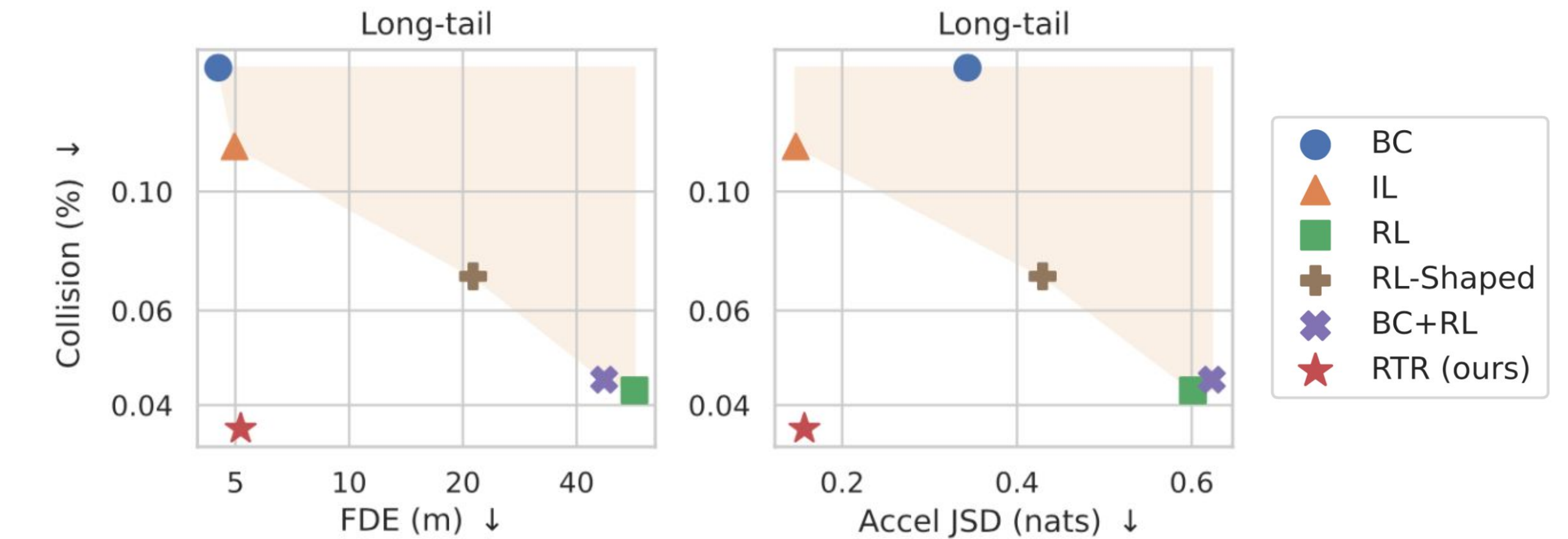$$E_{\pi}[-R(\tau)]$$

**Closed-loop Learning**

**Architecture:**

We use an efficient **multi-agent** architecture to extract features and jointly predict all agent actions.
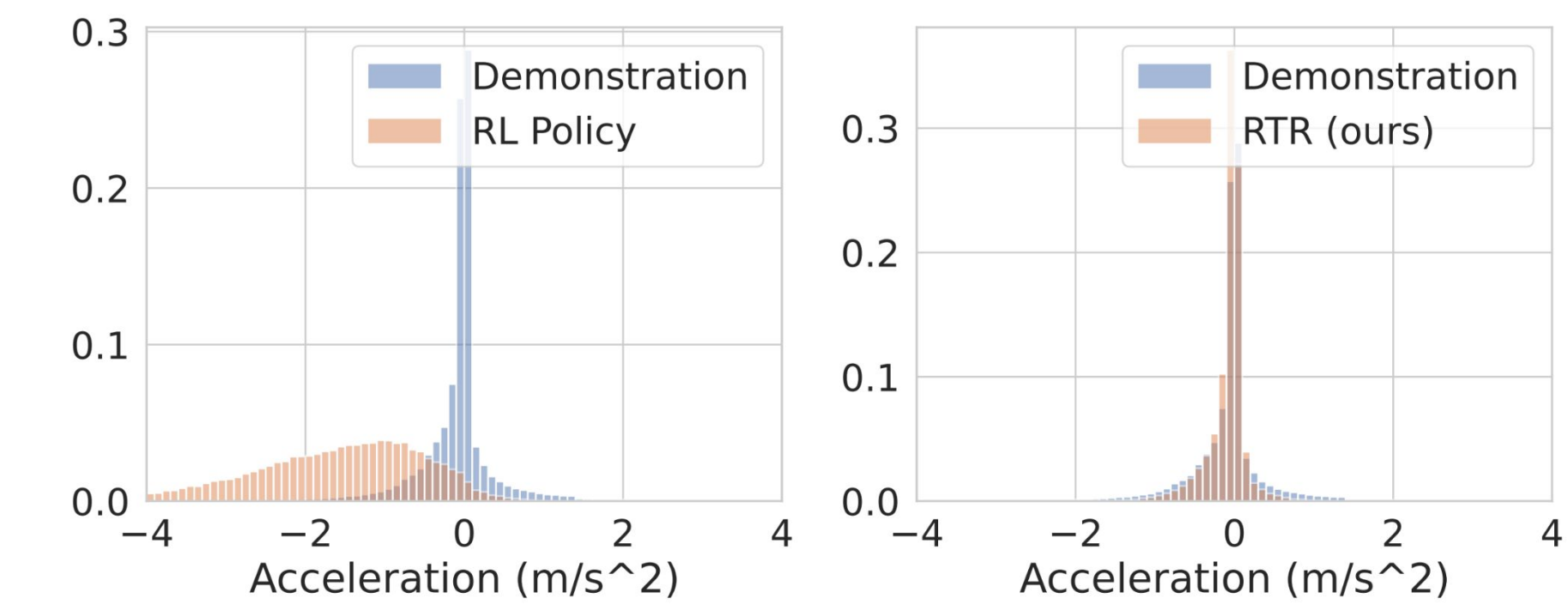
*Value network design is the same as policy network but regresses value targets instead.*

## Realism and Infraction Avoidance



RTR achieves the **best tradeoff**, outperforming the Pareto frontier of baselines which vary between IL, RL and IL + RL

RTR learns to **avoid infractions** while still capturing **human**-like driving.
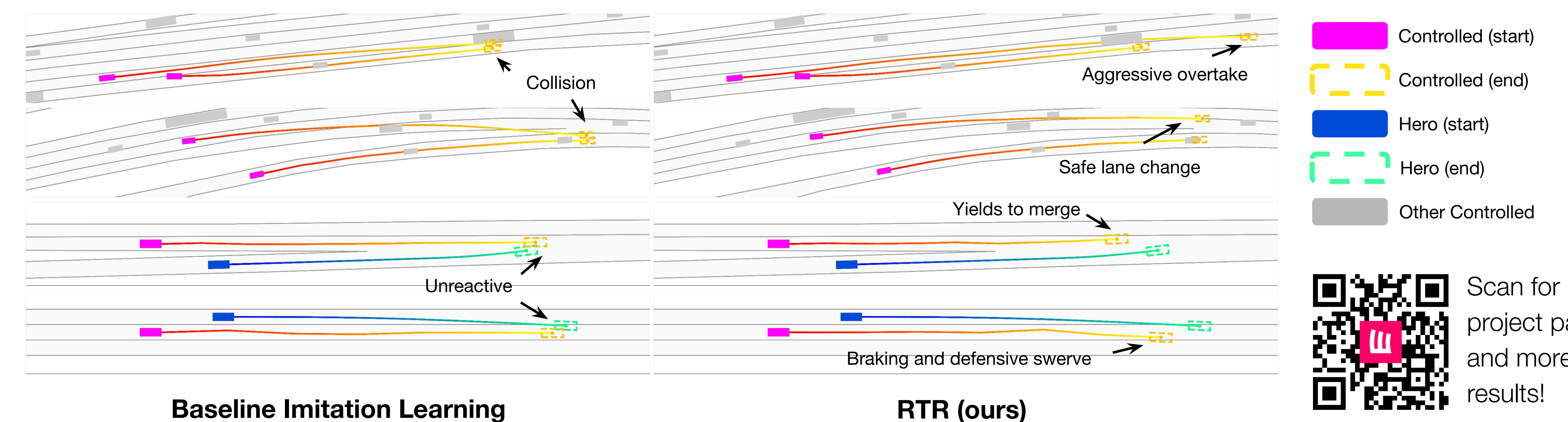


## Downstream evaluation

We train a prediction model on **actor-simulated** data and evaluate them on **real data**.

RTR simulations have **lower domain gap** vs. baselines

| Method | FDE (m) | CTE (m) |
|--------|---------|---------|
| BC | $2.44 \pm 0.05$ | $0.90 \pm 0.04$ |
| IL | $1.75 \pm 0.06$ | $0.28 \pm 0.01$ |
| RL | $15.42 \pm 1.21$ | $0.32 \pm 0.02$ |
| RL-Shp | $6.66 \pm 0.26$ | $0.33 \pm 0.01$ |
| BC+RL | $9.06 \pm 0.50$ | $0.42 \pm 0.03$ |
| RTR | $\mathbf{1.58 \pm 0.05}$ | $\mathbf{0.27 \pm 0.03}$ |

## Qualitative Results



Collision

Unreactive

**Baseline Imitation Learning**

Aggressive overtake

Safe lane change

Yields to merge

Braking and defensive swerve

**RTR (ours)**

- Controlled (start)
- Controlled (end)
- Hero (start)
- Hero (end)
- Other Controlled

Scan for project page and more results!