

# Learning Realistic Traffic Agents in Closed-loop Supplementary Materials

Anonymous Author(s)

Affiliation

Address

email

## A Additional Results

**Metrics:** In order to measure the realism of our traffic models, we use a set of metrics which evaluate both the traffic models' ability to match human demonstration data in the nominal scenarios and avoid infractions in both nominal and simulated long-tail scenarios.

- *Reconstruction:* In nominal scenarios where expert demonstrations exist, we consider a set of metrics which evaluate how close a traffic model's simulation is to the real world conditioned on the same initial condition. We measure the final displacement error (**FDE**) [1], defined as the L2 distance between an agent's position in a simulated scenario vs the ground truth scenario after 5s. We also measure the along-track error (**ATE**) and cross-track error (**CTE**) of an agent's simulated position projected onto the ground truth trajectory. This decomposition disentangles speed variability and lateral deviations respectively.

- *Distributional:* While reconstruction metrics compare pairs of real and simulated logs, we can compute distributional similarity metrics as an additional method to gauge realism. We compute the Jensen-Shannon Divergence (**JSD**) [2] between histograms of scenario features to compute their distributional similarity. Features include agent kinematics like acceleration and speed, pair-wise agent interactions like distance to lead vehicle, and map interactions like lateral deviation from lane centerline.

- *Infraction Rate:* Finally, we measure the rate of traffic infractions made by agents controlled by a traffic model. Similar to prior work [3], we measure percentage of agents that end up in **collision** or drive **off-road**. As this metric does not require ground truth scenarios for pairing or computing statistics, it can be used in simulated long-tail scenarios that do not have ground truth.

**Comparison to state-of-the-art:** In our main paper, we presented select results from our comparison to state-of-the-art traffic models on both nominal and long-tail scenarios. Here, we include additional tradeoff plots for all metrics in Figure 1. We also include a table of detailed metrics for all methods in Table 1. Building on our observations in the main paper, we see that RTR outperforms and expands the existing Pareto frontier on all metrics and scenario sets. IL methods achieve strong reconstruction/distributional realism metrics but suffer from high infraction rates, while RL methods attain the opposite. RTR achieves the best of both worlds—a testament to its ability to learn human-like driving while avoiding unrealistic traffic infractions.

**Long-Tail Scenarios:** In our main paper, we evaluated our approach of using procedurally generated long-tail scenarios against the alternative of mining hard scenarios from data. Here, we include additional tradeoff plots for all metrics in Figure 2, with the detailed metrics in Table 2. We see that training on both nominal and long-tail scenarios outperforms the alternatives in most cases.

**Distributional Realism:** In Figure 3, we include additional plots showing the histograms used to compute JSD distributional realism metrics on the nominal scenario set. We can see that RL

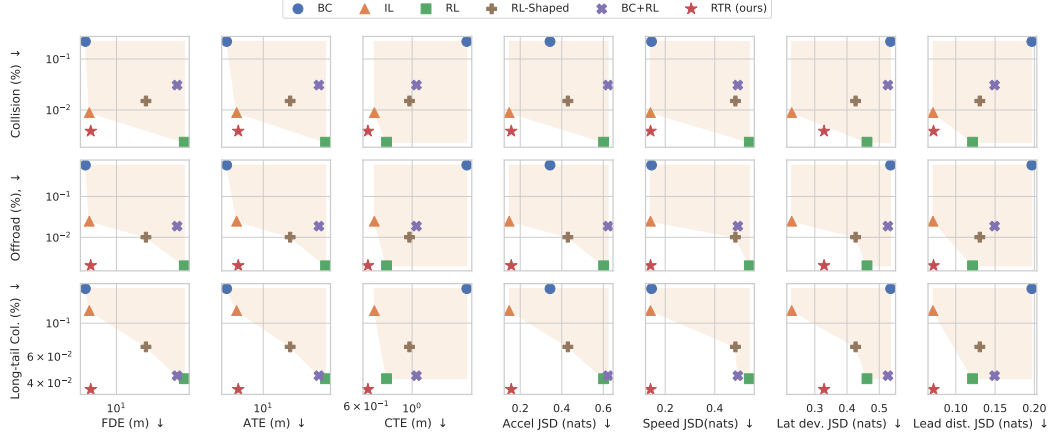


Figure 1: Additional plots comparing infraction / realism tradeoff of RTR compared to baseline models. We see that RTR outperforms and expands the existing Pareto frontier for all metrics.

Method	Infraction (%)		Reconstruction (m)			JSD (nats)				LT-Inf. (%)
	Col.	Off Rd.	FDE	ATE	CTE	Acc.	Speed	Lat. D.	Ld. D.	Col.
BC	22.13	58.68	<b>4.50</b>	<b>3.60</b>	1.84	0.34	0.54	<b>0.14</b>	0.20	17.00
IL	0.89	2.48	4.98	4.75	0.66	<b>0.15</b>	<b>0.23</b>	<b>0.14</b>	<b>0.07</b>	12.13
RL	<b>0.23</b>	<b>0.20</b>	56.92	56.91	0.75	0.60	0.46	0.54	0.12	4.26
RL-Shaped	1.50	1.01	21.29	21.17	0.97	0.43	0.43	0.48	0.13	6.95
BC+RL	3.08	1.88	47.30	47.26	1.05	0.62	0.53	0.49	0.15	4.46
RTR	0.38	<b>0.20</b>	5.16	4.97	<b>0.61</b>	0.16	0.33	<b>0.14</b>	<b>0.07</b>	<b>3.61</b>

Table 1: Detailed breakdown of metrics. Metrics on the left (resp. right) are computed on nominal scenarios (resp. long-tail scenarios). IL methods achieve strong reconstruction/distributional realism metrics but suffer from high infraction rates, while RL methods attain the opposite. RTR achieves the best of both worlds, with high reconstruction/distributional realism and low infraction rates.

methods (RL, RL-Shaped, and BC + RL) struggle to capture human-like driving, particularly in speed and acceleration JSD where the RL methods tend to brake more often than humans. BC exhibits slightly better results overall, but it has worse map interaction reasoning due to distribution shift from compounding errors. In contrast, RTR captures human-like driving significantly better, closely matching IL in distributional realism while also improving on its infraction rate as seen in other results.

**Qualitative Results:** We include qualitative results comparing RTR against the baselines Figures 4, 5, 6, and 7. Across fork, merge, and long-tail scenarios, we see that RTR exhibits the greatest realism of the competing methods.

## B Learning

### B.1 Loss Derivation

In this section, we will provide more details on the loss derivation using the Lagrangian. Recall that we begin with the following optimization problem

$$\begin{aligned}
 & \arg \min_{\pi} D_{\text{KL}}(P^{\pi}(\tau) \parallel P^E(\tau)) \\
 & \text{s.t. } \mathbb{E}_{P^{\pi}}[R(\tau)] \geq 0
 \end{aligned} \tag{1}$$

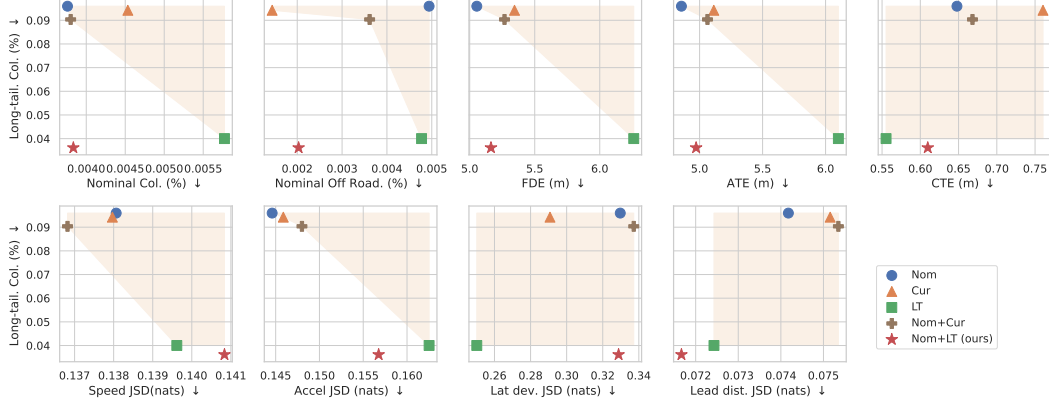


Figure 2: Additional plots showing the tradeoff between infraction rate on the long-tail set and other realism metrics on the nominal set, for models trained on different scenario sets. We see that for most metrics, training on both nominal and long-tail scenarios obtain the best tradeoff.

Training Scenarios	Infraction (%)		Reconstruction (m)			JSD (nats)				LT-Inf. (%)
	Col.	Off Rd.	FDE	ATE	CTE	Acc.	Speed	Lat. D	Ld. D.	Col.
Nominal	0.38	0.49	5.05	4.86	0.65	0.14	0.14	0.33	0.07	9.60
Curated	0.45	0.14	5.34	5.11	0.76	0.15	0.14	0.29	0.08	9.42
Long-tail	0.58	0.48	6.26	6.10	0.56	0.16	0.14	0.25	0.07	4.00
Nom. + Cur	0.38	0.30	5.27	5.06	0.67	0.15	0.14	0.34	0.08	9.04
Nom. + LT (ours)	0.38	0.20	5.16	4.97	0.61	0.16	0.14	0.33	0.07	3.61

Table 2: Detailed breakdown of realism and infraction metrics for training on different scenario sets.

49 We form the Lagrangian of the optimization problem

$$\mathcal{L}(\pi, \lambda) = D_{\text{KL}}(P^\pi(\tau) \parallel P^E(\tau)) + \lambda \mathbb{E}_{P^\pi}[R(\tau)] \quad (2)$$

$$= \mathbb{E}_{P^\pi} \left[ \log \frac{P^\pi(\tau)}{P^E(\tau)} - \lambda R(\tau) \right] \quad (3)$$

$$= \mathbb{E}_{P^\pi} [-\log P^E(\tau) - \lambda R(\tau)] - H(\pi) \quad (4)$$

50 where  $\lambda$  is a Lagrangian multiplier and

$$H(\pi) = -\mathbb{E}_{P^\pi} [\log P^\pi(\tau)] \quad (5)$$

$$= -\mathbb{E}_{P^\pi} \left[ \rho_0(s_0) \sum_{t=0}^{T-1} \log \pi(a^t | s^t) \right] \quad (6)$$

51 under deterministic dynamics is the causal entropy [4]. Using the Lagrangian, the optimization problem is converted to an unconstrained problem

$$\pi^* = \arg \min_{\pi} \max_{\lambda} \mathcal{L}(\pi, \lambda) \quad (7)$$

53 Equation 7 can be optimized in a number of ways, such as iteratively solving the inner maximization  
 54 over  $\lambda$  and outer minimization over  $\pi$ . We take a simplified approximate approach where we simply  
 55 set  $\lambda_{\text{fixed}} \geq 0$  as a hyperparameter, leading to what is ultimately a relaxed constraint or penalty  
 56 method.

$$\pi^* \approx \arg \min_{\pi} \mathbb{E}_{P^\pi} [-\log P^E(\tau) - \lambda_{\text{fixed}} R(\tau)] - H(\pi) \quad (8)$$

57 The causal entropy term is included as an entropy regularization term in some learning algorithms  
 58 such as PPO [5]. In practice, we found that it was not necessary to include.

## 60 B.2 Imitation Learning Loss

61 Recall that the imitation learning component of the loss is given as

$$\mathcal{L}^{\text{IL}} = \mathbb{E}_{\tau^E \sim D} \left[ \mathbb{E}_{\tau \sim P^\pi(\cdot | \mathbf{s}_0^E)} [D(\tau^E, \tau)] \right] \quad (9)$$

$$= \mathbb{E}_{(\mathbf{s}_0^E, \dots, \mathbf{s}_T^E) \sim D} \left[ \sum_{t=1}^T d(\mathbf{s}_t^E, \tilde{\mathbf{s}}_t) \right] \quad (10)$$

62 where

$$\tilde{\mathbf{a}}_t \sim \pi(\mathbf{a} | \tilde{\mathbf{s}}_t) \quad (11)$$

$$\tilde{\mathbf{s}}_{t+1} = \tilde{\mathbf{s}}_t + f(\tilde{\mathbf{s}}_t, \tilde{\mathbf{a}}_t) dt \quad (12)$$

63 Because the dynamics function  $f$  as described in Section B.7 is differentiable, Equation 10 com-  
 64 pletely differentiable using the reparameterization trick [6] when sampling from the policy. To  
 65 compute the inner expectation in Equation 9, we simply sample a single rollout. In practice, we  
 66 found that directly using the mean without sampling is also sufficient.

## 67 B.3 Reward Function

68 **Sparse reward:** Recall that we use the following reward function

$$R^{(i)}(\mathbf{s}, a^{(i)}) = \begin{cases} -1 & \text{if an infraction occurs} \\ 0 & \text{otherwise.} \end{cases} \quad (13)$$

69 In our experiments, we consider collisions events and driving off-road as infractions. Collisions are  
 70 computed by checking for overlap between the bounding boxes of agents. Off-road is computed  
 71 by checking if an agent’s bounding box still intersects with the road polygon. Note that when  
 72 optimizing the reward, we apply early termination of the scenario in the event of an infraction.

73 **Shaped reward:** For the RL-Shaped baseline, use the same reward in Equation 13 with an addi-  
 74 tional term which encourages driving at the speed limit.

$$R_{\text{shaped}}^{(i)}(\mathbf{s}, a^{(i)}) = R^{(i)}(\mathbf{s}, a^{(i)}) + 0.5(C - \delta)/C \quad (14)$$

75 where  $\delta = \text{abs}(\text{velocity} - \text{speed limit})$  and  $C = 30$ . We terminate the episode of  $\delta \geq C$ .

## 76 B.4 Reinforcement Learning Loss

77 We describe our factorized approach to multiagent PPO [5] in more detail. Starting off we compute  
 78 a per-agent probability ratio.

$$r^{(i)} = \frac{\pi(a^{(i)} | \mathbf{s})}{\pi_{\text{old}}(a^{(i)} | \mathbf{s})}. \quad (15)$$

79 Our centralized value-function uses the same architecture as our policy, and computes per-agent  
 80 value estimates  $\hat{V}^{(i)}(\mathbf{s})$ . Details of the architecture are found in Section B.6. The value model  
 81 is trained using per-agent value targets, which are computed with per-agent rewards  $R_t^{(i)} =$   
 82  $R^{(i)}(\mathbf{s}_t, a_t^{(i)})$

$$\mathcal{L}^{\text{value}} = \sum_i^N (\hat{V}^{(i)} - V^{(i)})^2 \quad (16)$$

$$V^{(i)} = \sum_{t=0}^T \gamma^t R_t^{(i)} \quad (17)$$

83 We can obtain a per-agent GAE using the value model as well,

$$A^{(i)} = \text{GAE}(R_0^{(i)}, \dots, R_{T-1}^{(i)}, \hat{V}^{(i)}(\mathbf{s}_T)) \quad (18)$$



84 The PPO policy loss is simply the sum of per-agent PPO loss,

$$\mathcal{L}^{\text{policy}} = \sum_{i=1}^N \min(r^{(i)} A^{(i)}, \text{clip}(r^{(i)}, 1 - \epsilon, 1 + \epsilon) A^{(i)}) \quad (19)$$

85 Finally, the overall loss is the sum of the policy and value learning loss.

$$\mathcal{L}^{\text{RL}} = \mathcal{L}^{\text{policy}} + \mathcal{L}^{\text{value}} \quad (20)$$

## 86 B.5 Input Parameterization

87 **Agent history:** Following [7], we adopt an viewpoint invariant representation of an agent’s past  
88 trajectory. We encode the past trajectory as a sequence of pair-wise relative positional encodings  
89 between the past waypoints and the current pose. Each relative positional encoding consists of the  
90 sine and cosine of distance and heading difference of a pair of poses. See [7] for details.

91 **Lane graph:** To construct our lane graph representation  $G = (V, E)$ , We first obtain the lane  
92 graph nodes by discretizing centerlines in the high-definition (HD) map into lane segments every  
93 10m. We use length, width, curvature, speed limit, and lane boundary type (e.g., solid, dashed)  
94 as node features. Following [8], we then connect nodes with 4 different relationships: successors,  
95 predecessors, left and right neighbors.

## 96 B.6 Model Architecture

97 Briefly, the RTR model architecture is composed of three main building blocks: (1) context encoders  
98 for embedding lane graph and agent history inputs; (2) interaction module for capturing scene-level  
99 interaction; and (3a) action decoder for parameterizing the per-agent policy and (3b) value decoder  
100 for the value model. Note that the policy model and the value model use the same architecture, but  
101 are trained completely separately and do not share any parameters.

102 **History encoder:** The *history encoder* consists of a 1D residual neural network (ResNet) followed  
103 by a gated recurrent unit (GRU) that extracts agent features  $h_a^{(i)} = f(s^{(i)})$  from a sliding window  
104 of past agent states  $s$ . Intuitively, the 1D CNN captures local temporal patterns, and the GRU  
105 aggregates them into a global feature.

106 **Lane graph encoder:** The *lane graph encoder* is a graph convolutional network (GCN) [8] that  
107 extracts map features  $h_m = g(\mathbf{m})$  from a given lane-graph  $G$  of map  $\mathbf{m}$ . We use hidden channel  
108 dimensions of [128, 128, 128, 128], layer normalization (LN), and max pooling aggregation.

109 **Interaction module:** To model scene-level interaction (i.e., agent-to-agent, agent-to-map, and  
110 map-to-map), we build a heterogeneous spatial graph  $G'$  by adding agent nodes to the original lane  
111 graph  $G$ . Besides the original lane graph edges, we connect agent nodes to their closest lane graph  
112 nodes. All agent nodes are also fully connected to each other. We use a *scene encoder* parameter-  
113 ized by a heterogeneous graph neural network (HeteroGNN) [7] to process map features and agent  
114 features into fused features,

$$\{h^{(1)}, \dots, h^{(N)}\} = \text{HeteroGNN}(\{h_a^{(1)}, \dots, h_a^{(N)}\}, h_m). \quad (21)$$

115 These fused features are then provided as input to the decoder.

116 **Action decoder:** Finally, we pass the fused features into a 4-layer MLP with hidden dimensions  
117 [128, 128, 128] to predict agent’s acceleration and steering angle distributions (parameterized as  
118 Normals).

$$(\mu^{(i)}, \sigma^{(i)}) = \text{MLP}(h^{(i)}) \quad (22)$$

$$\pi(a^{(i)} | \mathbf{s}) = \mathcal{N}(\mu^{(i)}, \sigma^{(i)}) \quad (23)$$

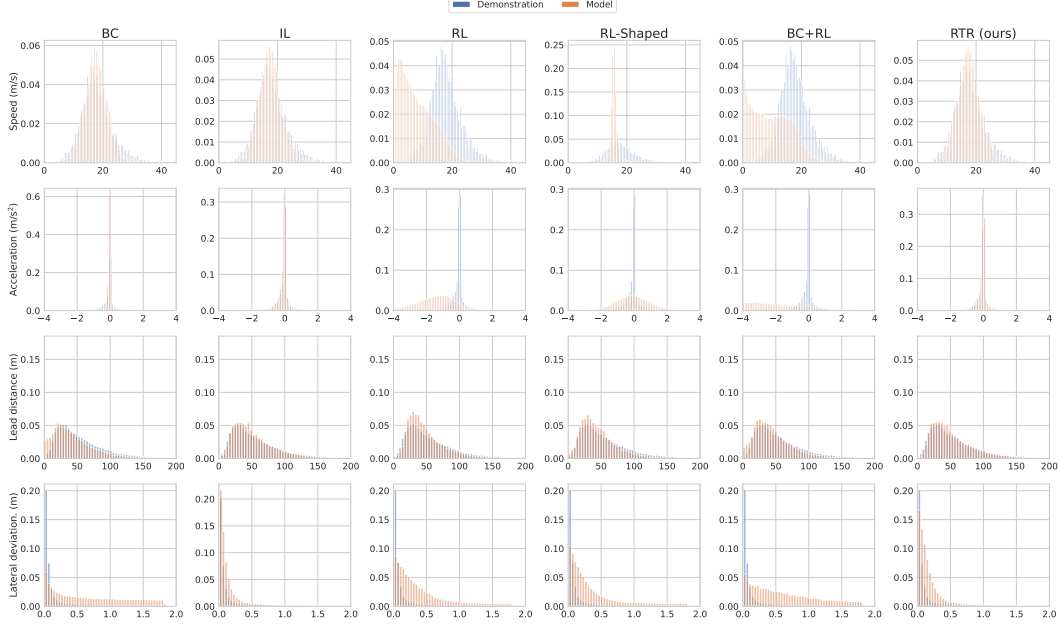


Figure 3: Histograms of scenario features for all methods used to compute JSD distributional realism metrics. We see that BC and RL methods often struggle with capturing the data distribution compared to IL and RTR. Notably, RTR closely matches IL performance in distributional realism, while greatly improving infraction rate as seen in other results.

119 **Value decoder:** For the value model, a 4-layer MLP instead regresses a single scalar value repre-  
 120 senting the value

$$\hat{V}^{(i)} = \text{MLP}_{\text{value}} \left( h_{\text{value}}^{(i)} \right) \quad (24)$$

## 121 B.7 Kinematic Bicycle Model

122 We use a kinematic bicycle model [9] for our environment dynamics. The bicycle model state is  
 123 given as

$$s = (x, y, \theta, v) \quad (25)$$

124 where  $x, y$  is the position of the center of the rear axel,  $\theta$  is the yaw, and  $v$  is the velocity. The bicycle  
 125 model actions are

$$a = (u, \phi) \quad (26)$$

126 where  $u$  is the acceleration, and  $\phi$  is the steering angle. The dynamics function  $\dot{s} = f(s, a)$  is then  
 127 defined as

$$\dot{x} = v \cos(\theta) \quad (27)$$

$$\dot{y} = v \sin(\theta) \quad (28)$$

$$\dot{\theta} = \frac{v}{L} \tan(\phi) \quad (29)$$

$$\dot{v} = u \quad (30)$$

128 where  $L$  is wheelbase length, i.e. the distance between the rear and front axel. We can use a simple  
 129 finite difference approach to computing the next state

$$s_{t+1} = s_t + f(s_t, a_t)dt \quad (31)$$

130 where  $dt$  is chosen to be 0.5 seconds in practice. We can apply the bicycle model to each agent  
 131 individually to obtain the joint state dynamics function.

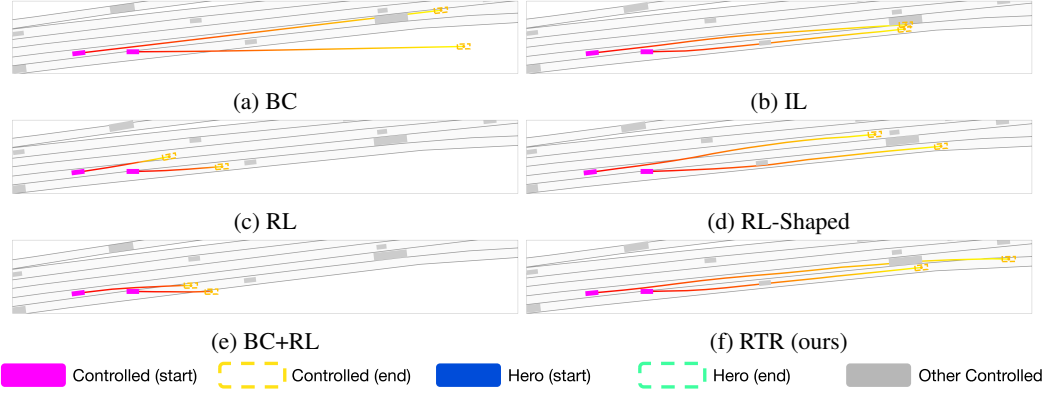


Figure 4: Qualitative results on a fork scenario. BC drives off the road, IL results in a collision while RL and BC+RL slow down. RL-Shaped drives straight and loses the interesting lane change behavior.

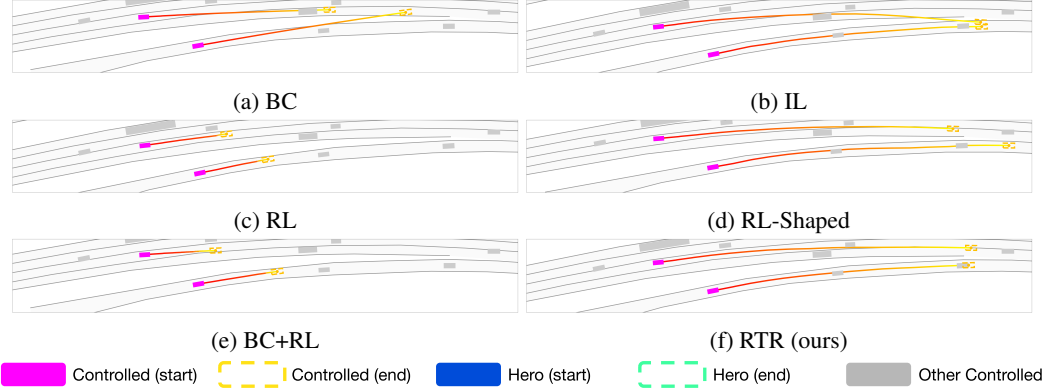


Figure 5: Qualitative results on a merge scenario. We see that RL methods slow down unrealistically. IL results in a collision while RTR maintains realism.

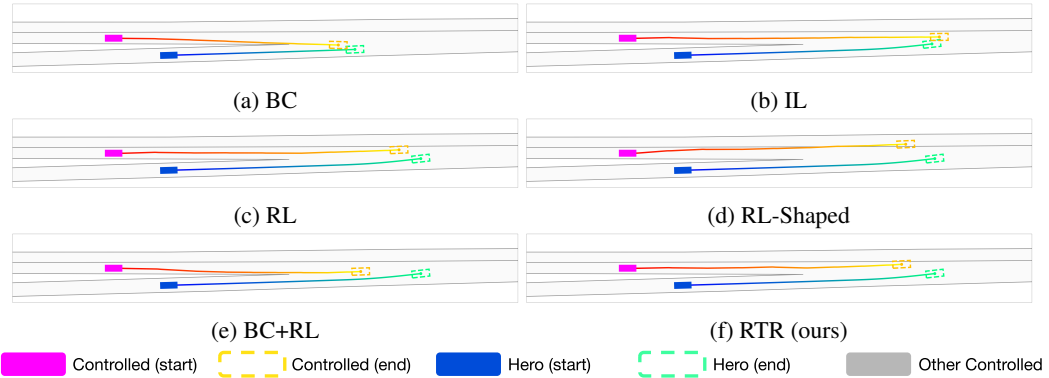


Figure 6: Qualitative results on procedurally generated merge scenario. IL and BC result in a collision. RTR maintains realism.

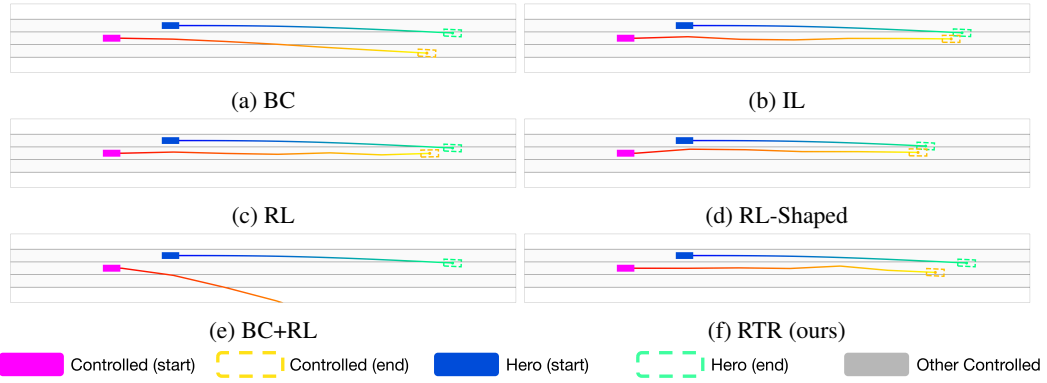


Figure 7: Qualitative results on a procedurally generated cut-in scenario. BC+RL drives off the road, while IL and RL-shaped result in a collision. RTR maintains realism.

## References

- [1] S. Casas, C. Gulino, S. Suo, K. Luo, R. Liao, and R. Urtasun. Implicit latent variable model for scene-consistent motion forecasting. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII 16*, pages 624–641. Springer, 2020.
- [2] M. Igl, D. Kim, A. Kuefler, P. Mougin, P. Shah, K. Shiarlis, D. Anguelov, M. Palatucci, B. White, and S. Whiteson. Symphony: Learning realistic and diverse agents for autonomous driving simulation, 2022. URL <https://arxiv.org/abs/2205.03195>.
- [3] S. Suo, S. Regalado, S. Casas, and R. Urtasun. Trafficsim: Learning to simulate realistic multi-agent behaviors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10400–10409, June 2021.
- [4] J. Ho and S. Ermon. Generative adversarial imitation learning. *Advances in neural information processing systems*, 29, 2016.
- [5] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [6] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [7] A. Cui, S. Casas, K. Wong, S. Suo, and R. Urtasun. Gorela: Go relative for viewpoint-invariant motion forecasting. *arXiv preprint arXiv:2211.02545*, 2022.
- [8] M. Liang, B. Yang, R. Hu, Y. Chen, R. Liao, S. Feng, and R. Urtasun. Learning lane graph representations for motion forecasting. In A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part II*, volume 12347 of *Lecture Notes in Computer Science*, pages 541–556. Springer, 2020. doi:10.1007/978-3-030-58536-5\_32. URL [https://doi.org/10.1007/978-3-030-58536-5\\_32](https://doi.org/10.1007/978-3-030-58536-5_32).
- [9] S. M. LaValle. *Planning algorithms*. Cambridge university press, 2006.