

## 844 **Broader Impacts**

845 Safety-critical multi-agent cooperative policy learning has the potential to significantly benefit  
846 society by enhancing safety, efficiency, and coordination in various domains. In autonomous driving,  
847 robust multi-agent policies can prevent collisions, optimize traffic flow, and ensure pedestrian safety,  
848 contributing to safer urban mobility. In multi-drone systems, such as UAV swarms for search-and-  
849 rescue operations or environmental monitoring, agents can collaboratively navigate and perform tasks  
850 without risking collisions or environmental damage. Similarly, in industrial automation, multi-robot  
851 systems can coordinate safely in shared spaces, increasing productivity while reducing the risk of  
852 accidents.

853 However, the deployment of safety-critical multi-agent systems also introduces potential negative  
854 societal impacts. Poorly designed safety mechanisms can result in unsafe behaviors, leading to  
855 accidents or property damage. Our work addresses this issue by proposing a safety-guaranteed  
856 approach that ensures safety during both training and deployment phases, leveraging Control Barrier  
857 Functions (CBFs) to enforce strict safety constraints on agent behavior.

## 858 **A Appendix / supplemental material**

### 859 **A.1 Environment Details**

860 The five environments from the METADRIIVE simulator are illustrated in the figure 4. As mentioned  
861 these environments involve regions where agents can cross path with other agents, Hence, it is  
862 necessary for the agents to cooperate but at the safe time stay safe with each other to achieve the task  
i.e. reach destination successfully and timely.

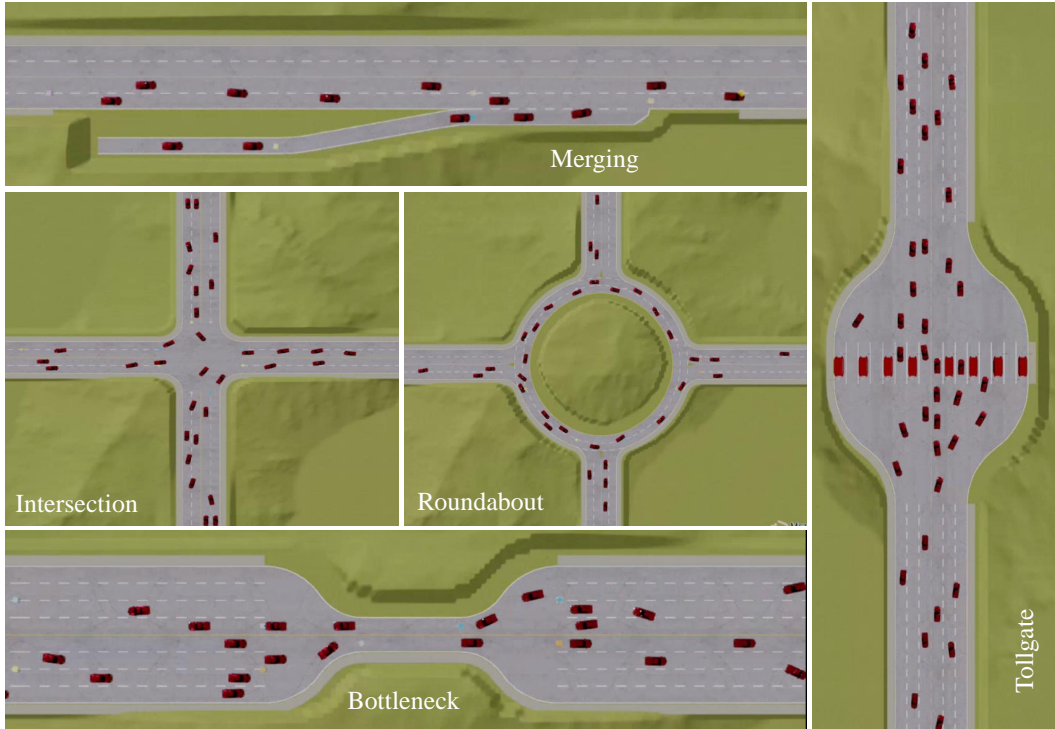


Figure 4: Illustration of the various METADRIIVE environments.

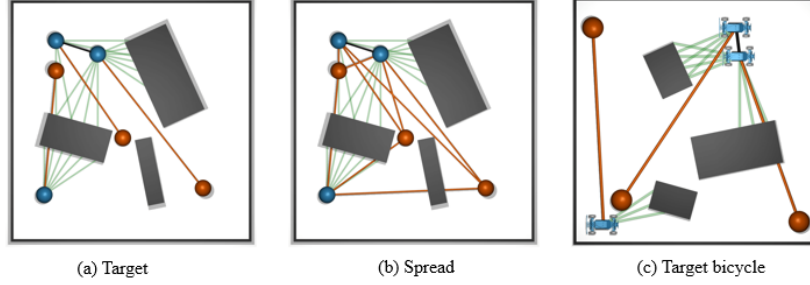


Figure 5: Illustration of the lidar based environments which include: (a) target environment, (b) spread environment and (c) target bicycle environment.

## A.2 Experimental Details and Additional Results

We begin this section by presenting the details of the training, evaluation and ablation experiments. Besides that, we include results for the generalization experiments used to validate our proposed method.

In addition, we evaluate our method on a suite of lidar-based environments introduced in [51], as illustrated in Figure 5. All environments are partially observable, with agents equipped with lidar sensors. In the *Target* environments, each agent is assigned a fixed goal and must reach it while avoiding collisions with other agents and obstacles. The primary distinction between the *Target* and *Target-Bicycle* environments lies in the underlying agent dynamics; further details are provided in Figure 5. In the *Spread* environment, agents must cooperate to cover all goal locations while avoiding both obstacles and one another.

### A.2.1 Training Setup and Results.

In all five environments used in our experiments, the vehicles are initialized at random spawn points, and assigned destinations to each of them randomly. The agents are terminated in three situations: reaching the destinations (deemed as successful), crashing with others or driving out of the road (deemed as failure). Once an existing vehicle is terminated new vehicles are spawned immediately to ensure continuous traffic flow. The terminated vehicles remains static for 10 steps (2s) in the original positions, creating impermeable obstacles. If a vehicle crashes into other terminated vehicles before it is removed from the environment, it is also considered as failure. The *success rate* is the fraction of the successful agents to the total number of spawned agents. Thus, the environment allows for the total number of vehicles to vary in time and exceed the the initial number of agents which mimics real-world challenging traffic scenarios.

We trained with 15 vehicles for merging, 30 vehicles for intersection and roundabout, 25 vehicles for bottleneck and 40 vehicles for tollgate respectively. As mentioned previously, the extrinsic reward corresponding to (9) was set to the environment reward, besides the regularization terms. The intrinsic reward is described in subsection A.5. We present the training results for the remaining three environments from figure 3 i.e. merging, roundabout and intersection in figure 6.

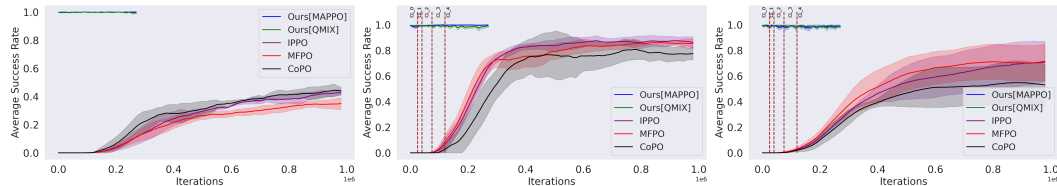


Figure 6: Plot of the success rate of the agents with our algorithm vs the baselines in merging, roundabout, and intersection environments respectively. CL refers to curriculum learning, four different curriculum were used with varying traffic densities. The benchmark methods use 2000+ hours of individual driving data prior to the training.

As can be noticed, our algorithm achieves a near perfect safety rate with both on-policy and off-policy algorithms during training, and, keeping consistent with the tollgate and bottleneck environments are trained for 300k iterations. However, in order to encourage exploration in roundabout and intersection environments for learning the cooperative behavior, we use curriculum learning by varying the arrival rates, while keeping the traffic density constant. The curriculum transitions happened at 30k, 50k, 80k and 125k iterations during training.

## A.2.2 Evaluation Setup and Results

**METADRIIVE Simulator.** As mentioned previously we used *success rate weighted average time and energy* for evaluating the performance of our algorithms. In order to compare against the baseline methods, for each environment and algorithm (e.g. on-policy, on-policy w/Grouping, off-policy and so on) we train using five seeds. Then we run five episodes for each seed and generate the evaluation results by computing the mean and standard deviation over the 25 runs. In addition to the presented results in table 2, we run the baseline algorithms by changing the penalty factors for the safety related terms (both inter-agent and environment safety) in the environment reward for the bottleneck environment. The results are presented in table 3. The primary takeaway is that the added emphasis on safety does not render better results for the baseline algorithms compared to our method. Besides that, different algorithms perform differently with increase in penalty factor, thus not exhibiting any trend between penalty and success rate.

Table 3: Summary of baseline comparison using various penalty factors for Bottleneck environment.  $\uparrow$ : higher is better;  $\downarrow$ : lower is better. **Bold**: the best performance for each metric.

Penalty Factors	Metrics	IPPO	CL	MFPO	CoPo	HMARL-CBF (on-policy)
1	Success $\uparrow$	0.09 $\pm$ 0.03	0.34 $\pm$ 0.15	0.48 $\pm$ 0.16	0.75 $\pm$ 0.07	<b>0.99<math>\pm</math>0.00</b>
	Time $\downarrow$	2644.44 $\pm$ 353.44	939.35 $\pm$ 147.21	679.83 $\pm$ 158.33	651.54 $\pm$ 63.52	<b>241.16<math>\pm</math>8.33</b>
	Energy $\downarrow$	69.66 $\pm$ 9.77	11.61 $\pm$ 2.02	9.95 $\pm$ 1.39	6.93 $\pm$ 0.61	<b>7.37<math>\pm</math>0.00</b>
5	Success $\uparrow$	0.58 $\pm$ 0.11	0.395 $\pm$ 0.01	0.615 $\pm$ 0.08	0.716 $\pm$ 0.02	<b>0.99<math>\pm</math>0.00</b>
	Time $\downarrow$	568 $\pm$ 56	646.5 $\pm$ 43.16	713 $\pm$ 273.7	529.13 $\pm$ 79	<b>241.16<math>\pm</math>8.33</b>
	Energy $\downarrow$	30 $\pm$ 41.5	10.8 $\pm$ 0.24	8.64 $\pm$ 0.65	7.69 $\pm$ 0.02	<b>7.37<math>\pm</math>0.00</b>
10	Success $\uparrow$	0.18 $\pm$ 0.08	0.55 $\pm$ 0.11	0.42 $\pm$ 0.23	0.655 $\pm$ 0.04	<b>0.99<math>\pm</math>0.00</b>
	Time $\downarrow$	1292.17 $\pm$ 647	606.49 $\pm$ 204.89	1414 $\pm$ 973	514 $\pm$ 126.92	<b>241.16<math>\pm</math>8.33</b>
	Energy $\downarrow$	40.97 $\pm$ 27	9.14 $\pm$ 0.93	22.12 $\pm$ 22.08	8.37 $\pm$ 0.04	<b>7.37<math>\pm</math>0.00</b>

## Lidar Environment Suite.

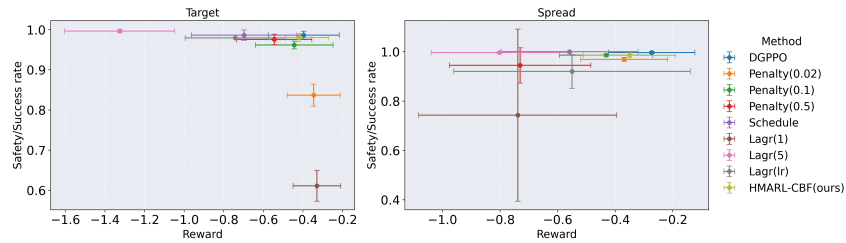


Figure 7: Success/safety rate vs reward for Target and Spread environment.

**Baselines.** We compare our method with DGPPPO [51], InforMARL [71] and MAPPO-Lagrangian [72], three state of the art methods for multi-agent safe control. For InforMARL, we evaluate fixed and scheduled penalty variants. Further details about the baselines can be found in [51]. For MAPPO-Lagrangian, we test two settings for the Lagrange multiplier and one learning rate. Performance is better if the plot is located in top right corner. As can be seen, our method achieves a high success rate while also achieving high reward compared to the baselines in [51]. Finally, we present the results for the bicycle target environment in 8 which shows that our method scales with number of agents.

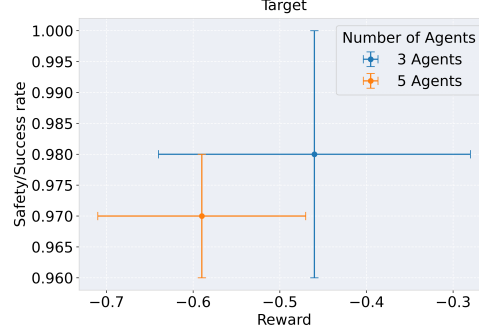


Figure 8: Success/safety rate vs reward for Target and Spread environment.

### A.2.3 Generalization Results

Besides the evaluation experiments, we conducted a series of experiments to assess the generalizability of our proposed approach. Firstly we systematically vary the lane widths in the environments. The corresponding success rates are reported in Table 4. Our analysis indicates that reducing lane width has a minimal impact on success rates in the roundabout and bottleneck environments. However, a more comparatively greater decline is observed in the intersection environment. This performance degradation is primarily attributed to the default intersection configuration (Figure 1), where two lanes are available in each direction. Agents are trained to execute lane changes from both lanes depending on their assigned routes. When lane width is reduced, agents attempting to perform lane changing from the inside lane at times gets stuck (and, eventually get eliminated from the episode due to timeout), or, hit road boundary which causes the drop in success rate. Conversely, increasing lane width generally maintains or enhances success rates across all environments.

Secondly, we investigated skill transferability by training low-level skills in an intersection environment and subsequently using them to train a high-level policy across multiple environments. Namely, we investigate the transferability in the merging and roundabout environments. This is because the observation is different and distinct for the Bottleneck and Tollgate environments compared to the other environments. The results are presented in table 5. This analysis reveals that skills learned in one environment can generalize effectively to others producing comparatively similar performance. In the roundabout environment, we observe a minor decline in the success rate and increase in the weighted average energy and time; yet our approach outperforms the baselines presented in table 2.

Finally, we test the generalizability of our method for varying traffic density/ number of agents. The results for this test are presented in tables 6 and 7 respectively. For every environment we trained for the default density as shown in the table and tested for lower and higher values. From the results, we can conclude that the success rate is unaffected by traffic density. The average time of the agents increases with increasing traffic density as expected. Finally, the energy values show a decreasing trend with increase in traffic density across most environments. This was anticipated because energy has a negative correlation with time.

### A.2.4 Computer Details.

We conducted our experiments on a desktop equipped with a single NVIDIA RTX A5000 GPU and 32 CPU cores. Given that our environments involve a high number of agents (up to 45 agents in the tollgate environment), the primary computational cost is attributed to data collection. Training and convergence are achieved efficiently, as demonstrated in the training and evaluation results. Furthermore, we utilize shared neural network architectures across agents, which significantly reduces the overall computational cost and training time.

## A.3 Background

**Definition A.1** (Relative degree). The relative degree of a (sufficiently many times) differentiable function  $b : \mathbb{R}^n \rightarrow \mathbb{R}$  with respect to system (3) is the number of times it needs to be differentiated along its dynamics until the control  $a$  explicitly appears in the corresponding derivative.

Table 4: Success rate across different environments for varying lane widths.

Lane Width (m)	Success Rate				
	Merging	Intersection	Roundabout	Bottleneck	Tollgate
3.25	0.99 $\pm$ 0.01	0.90 $\pm$ 0.03	0.98 $\pm$ 0.02	0.97 $\pm$ 0.02	0.99 $\pm$ 0.001
3.5 (default)	0.99 $\pm$ 0.01	0.97 $\pm$ 0.04	0.99 $\pm$ 0.01	0.99 $\pm$ 0.003	0.99 $\pm$ 0.02
4	0.99 $\pm$ 0.01	0.97 $\pm$ 0.015	0.99 $\pm$ 0.01	0.99 $\pm$ 0.001	0.99 $\pm$ 0.02

Table 5: Summary of skill transfer results (\* indicates runs done by transferring skills).

Experiments	Success	Energy	Time
Merging	0.99 $\pm$ 0.01	14.24 $\pm$ 0.23	166.18 $\pm$ 2.64
Merging*	0.96 $\pm$ 0.016	14.49 $\pm$ 0.58	178.34 $\pm$ 5.90
Roundabout	0.99 $\pm$ 0.01	9.89 $\pm$ 0.23	403.18 $\pm$ 21.66
Roundabout*	0.96 $\pm$ 0.016	11.39 $\pm$ 0.08	404.16 $\pm$ 18.6

Table 6: Summary of results for varying number of agents in Merging, Intersection, and Roundabout environments.

Metric	Merging			Intersection			Roundabout		
	10	15 (default)	20	20	30 (default)	40	20	30 (default)	40
Success	0.997 $\pm$ 0.0	0.99 $\pm$ 0.01	0.98 $\pm$ 0.02	0.955 $\pm$ 0.02	0.97 $\pm$ 0.04	0.94 $\pm$ 0.04	0.99 $\pm$ 0.01	0.99 $\pm$ 0.0	0.98 $\pm$ 0.0
Energy	14.03 $\pm$ 0.36	14.24 $\pm$ 0.23	14.47 $\pm$ 0.18	5.85 $\pm$ 0.07	5.82 $\pm$ 0.04	5.94 $\pm$ 0.05	10.48	10.32	9.95
Time	165.18 $\pm$ 2.15	166.18 $\pm$ 2.64	169.55 $\pm$ 2.15	318.54 $\pm$ 19	332.3 $\pm$ 36	361.8 $\pm$ 20.46	389.34	404.38	415.03

Table 7: Summary of results for varying number of agents in the Bottleneck and Tollgate environments.

Metric	Bottleneck			Tollgate		
	20	25 (default)	30	30	40 (default)	45
Success	0.99 $\pm$ 0.01	0.99 $\pm$ 0.003	0.978 $\pm$ 0.02	0.99 $\pm$ 0.02	0.99 $\pm$ 0.02	0.99 $\pm$ 0.01
Energy	7.47 $\pm$ 0.03	7.44 $\pm$ 0.001	7.43 $\pm$ 0.02	8.51 $\pm$ 0.02	8.51 $\pm$ 0.048	8.50 $\pm$ 0.02
Time	227.22 $\pm$ 7.77	243.59 $\pm$ 8.33	244.62 $\pm$ 6.62	371.53 $\pm$ 6.78	371.05 $\pm$ 5.82	374.02 $\pm$ 8.48

955 For a constraint  $b(s) \geq 0$  with relative degree  $m$ ,  $b : \mathbb{R}^n \rightarrow \mathbb{R}$ , and  $\zeta_0(s) := b(s)$ , we define a  
 956 sequence of functions  $\zeta_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $i \in \{1, \dots, m\}$  :

$$\zeta_i(s) := \dot{\zeta}_{i-1}(s) + \alpha_i(\zeta_{i-1}(s)), \quad i \in \{1, \dots, m\}, \quad (16)$$

957 where  $\alpha_i(\cdot)$ ,  $i \in \{1, \dots, m\}$  denotes a  $(m-i)^{\text{th}}$  order differentiable class  $\mathcal{K}$  function. We further  
 958 define a sequence of sets  $C_i$ ,  $i \in \{1, \dots, m\}$  associated with (16) which take the following form,

$$C_i := \{s \in \mathbb{R}^n : \zeta_{i-1}(s) \geq 0\}, \quad i \in \{1, \dots, m\}. \quad (17)$$

959 **Definition A.2** (High Order CBF (HOCBF) [73, 74]). Let  $C_1, \dots, C_m$  be defined by (17) and  
 960  $\zeta_1(s), \dots, \zeta_m(s)$  be defined by (16). A function  $b : \mathbb{R}^n \rightarrow \mathbb{R}$  is a High Order Control Barrier  
 961 Function (HOCBF) of relative degree  $m$  for system (3) if there exists  $(m-i)^{\text{th}}$  order differentiable  
 962 class  $\mathcal{K}$  functions  $\alpha_i$ ,  $i \in \{1, \dots, m-1\}$  and a class  $\mathcal{K}$  function  $\alpha_m$  such that

$$\sup_{a \in \mathcal{U}} [L_f^m b(s) + L_g L_f^{m-1} b(s) a + \Omega(b(s)) + \alpha_m(\zeta_{m-1}(s))] \geq 0 \quad (18)$$

963 for all  $s \in \bigcap_{i=1}^m C_i$ . In (18),  $L_f^m$  and  $L_g$  denotes derivative along  $f$  and  $g$   $m$  times and one time  
 964 respectively, and  $S(\cdot)$  denotes the remaining Lie derivative along  $f$  with degree less than or equal to  
 965  $m-1$  (omitted for simplicity, see [73]).

966 Note that the HOCBF in (18) is a general form of the degree one CBF [54] ( $m=1$ ) and exponential  
 967 CBF in [75]. The following theorem on HOCBFs implies the forward invariance property of the  
 968 CBFs and the original safety set. The proof is omitted (see [73] for the proof).

**Definition A.3** (Control Lyapunov function (CLF)[54]). A continuously differentiable function  $V : \mathbb{R}^n \rightarrow \mathbb{R}$  is a globally and exponentially stabilizing CLF for (3) if there exists constants  $c_i \in \mathbb{R}_{>0}$ ,  $i = 1, 2$ , such that  $c_1 \|s\|^2 \leq V(s) \leq c_2 \|s\|^2$ , and the following inequality holds

$$\inf_{a \in \mathcal{A}} [L_f V(s) + L_g V(s) a + \eta(s)] \leq e, \quad (19)$$

where  $e$  makes this a soft constraint. If the dynamics are of the form:  $\dot{s} = f(s, a)$ , we can express it in the form:

$$\dot{s}' = f'(s', a) + g'(s') a' \quad (20)$$

where  $s' = \begin{bmatrix} s \\ a \end{bmatrix}$  and  $a'$  are the augmented system states and new primitive actions of the system respectively,  $f' = \begin{bmatrix} f(s, a) \\ 0_{\dim(a)} \end{bmatrix}$ , and  $g' = \begin{bmatrix} 0_{\dim(s) \times \dim(a')} \\ B \end{bmatrix}$ , and  $B \in \mathbb{R}^{\dim(a) \times \dim(a')}$ . Generally,  $B$  contains ones in its diagonal entries. The redefined dynamics can be used subsequently for defining CBFs.

#### A.4 Safe Skill learning

We provide the details about the computation of policy gradients necessary to learn the skill policy parameters in this subsection. We follow [76] by applying the Karush-Kuhn-Tucker condition on the Lagrangian to derive the gradient of the state feedback policy with respect to its parameters. Specifically, let  $\lambda$  denote the dual variables on the HOCBF and CLF constraints, let  $D(\cdot)$  create a diagonal matrix from a vector, and let  $\mu^*, \lambda^*$  denote the optimal solutions of  $\mu$  and  $\lambda$ , respectively. We can then write  $d_\mu$  and  $d_\lambda$  in the form:

$$\begin{bmatrix} d_\mu \\ d_\lambda \end{bmatrix} = \begin{bmatrix} H & G^T D(\lambda^*) \\ G & D(Ga^* - h) \end{bmatrix}^{-1} \begin{bmatrix} 1_{\dim(\mathcal{A}_i)} \\ 0 \end{bmatrix} \quad (21)$$

where  $G, h$  are concatenated by  $G_b^j, G_v^j, h_b^j, h_v^j$ ,

$$\begin{aligned} G_b^j &= -L_g L_f^{m-1} b_i^j(s^i, s^j); \quad j \in \mathcal{N}^i(s) \\ G_v^j &= -L_g V_i^j(s^i, s^j; \phi_v^{i,j}); \quad j = 1, 2, \dots \\ h_b^j &= L_f^m b_i^j(s^i, s^j) + \Omega(b_i^j(s^i, s^j); \phi_b^{i,j}) + \alpha_m(\zeta_{m-1}(s^i, s^j; \phi_b^{i,j})); j \in \mathcal{N}^i(s) \\ h_v^j &= L_f V_i^j(s^i, s^j) + \eta(s^i, s^j; \phi_v^{i,j}); \quad j = 1, 2, \dots \end{aligned} \quad (22)$$

As the control bounds in (14) are not trainable, they are not included in  $G$  and  $h$ . Then, the relevant gradients with respect to all the parameters are given by

$$\nabla_H \mu = \frac{1}{2} (d_u u^T + u d_u^T), \nabla_F \mu = d_u \quad (23)$$

$$\nabla_G \mu = D(\lambda^*) (d_\lambda u^T + \lambda d_u^T), \nabla_h \mu = -D(\lambda^*) d_\lambda \quad (24)$$

With the hierarchical composition of the policies ( $\pi_H \circ \pi_z$ ), we define the augmented state  $\hat{s} = [s^i, s^{-i}, z, s_{z,init}]$  where  $s^{-i}$  is the stage of the agents other than agent  $i$  and  $s_{z,init}$  is the joint state of the agents at the time of the initiation of the joint skill  $z$ . The initiation state is included as the joint termination of the skill depends on the initiation state. The state and action value functions for agent  $i$  corresponding to reward  $r_L^i$ , skill policy  $\pi_{z^i}$  (given the high level policy  $\pi_H$  and the skill policy other agents  $\pi_{z^{-i}}$ ) are defined as follows:

$$\begin{aligned} V^{\pi_{z^i}}(\hat{s}) &= \mathbb{E}[r_L^i(\hat{s}_0, a_0^i | \hat{s}_0 = \hat{s}) + \mathbb{E}_{\hat{s}_1}[Q(\hat{s}_1, a_1^i)]] \\ &= \mathbb{E}[r_L^i(s_0^i, a_0^i | s_0^i = s^i, z^i) + \mathbb{E}_{\hat{s}_1}[Q(\hat{s}_1, a_1^i)]] \end{aligned} \quad (25)$$

The state-action value function for agent  $i$ , given the hierarchical policy composition  $\pi_H \circ \pi_z$ , is defined as:

$$Q^{\pi_{z^i}}(\hat{s}, a^i) = \mathbb{E}[r_L^i(\hat{s}, a^i) + \gamma \mathbb{E}_{\hat{s}'}[V^{\pi_{z^i}}(\hat{s}')] ] \quad (26)$$

The corresponding advantage function is given by:

$$A^{\pi_{z^i}}(\hat{s}, a^i) = Q^{\pi_{z^i}}(\hat{s}, a^i) - V^{\pi_{z^i}}(\hat{s}) \quad (27)$$

997 Based on the above definition of value functions and advantage function, as well the gradient of QP  
 998 controller output with respect to its parameters, standard on-policy and off-policy algorithms can be  
 999 used to derive the policy gradient to learn the controller parameters.

#### 1000 **Proof of Theorem 5.1**

1001 *Proof.* Given, the constraints  $b_i^j$ 's are satisfied at the initial time, the satisfaction of the CBF con-  
 1002 straints makes the corresponding safety sets forward invariant per Theorem (3.4), thus guaranteeing  
 1003 their satisfaction at all future times. This makes the policies corresponding to the skills of the agents  
 1004 safe, thus, guaranteeing the safety of our proposed HMARL approach.  $\square$

1005 **High-level task specific reward to low level learning.** The low-level policies should be such that  
 1006 they optimize the entire trajectory given by the problem in (5). Hence, in addition to introducing  
 1007 intrinsic motivation through the low-level reward, inspired by [23], we introduce an advantage term  
 1008 to the low-level reward resulting in the following revised low-level reward:

$$\tilde{r}_L^i(\mathbf{s}_t^i, \mathbf{a}_t^i | z_{t'}^i) = \lambda \frac{A^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'})}{\mathcal{N}} + (1 - \lambda) r_L^i(\mathbf{s}_t^i, \mathbf{a}_t^i | z_{t'}^i) \quad (28)$$

1009 where  $\lambda \in [0, 1]$  and  $t' \leq t$  is the time when the joint skill at time  $t$  was initiated. The high-level  
 1010 advantage  $A^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'})$  is defined as below:

$$A^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'}) = Q^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'}) - V^{\pi_H}(\mathbf{s}_{t'}) \quad (29)$$

1011 where  $Q^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'})$  and  $V^{\pi_H}(\mathbf{s}_{t'})$  are as defined in (2) and (1) respectively. The parameter  $\lambda$  can be  
 1012 used to balance between skill learning and optimization of the joint behavior of the agents. From [23]  
 1013 we have that for high level policies the following holds:

$$J(\tilde{\pi}_H) = J(\pi_H) + \mathbb{E}_{\mathbf{s}_{t'}, \mathbf{z}_{t'} \sim \pi_H} \left[ \sum_{t'=0} \gamma^{t'} A^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'}) \right] \quad (30)$$

1014 **Proposition A.4.** Given the low level reward in (28), the revised low level learning objective  $\tilde{J}^i$  for  
 1015 agent  $i$  can be expressed as:

$$\tilde{J}^i(\pi_H, \pi_{\mathbf{z}^{-i}}, \pi_{\mathbf{z}^i}) \approx \frac{\lambda}{\mathcal{N}(1 - \gamma)} \mathbb{E} \left[ \sum_{t'=0} \gamma^{t'} A^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'}) \right] + (1 - \lambda) J^i(\pi_H, \pi_{\mathbf{z}^{-i}}, \pi_{\mathbf{z}^i}) \quad (31)$$

1016 *Proof.* Based on 28, we can write our low level policy's learning objective for agent  $i$  as the following:

$$\begin{aligned} \tilde{J}^i(\pi_H, \pi_{\mathbf{z}^{-i}}, \pi_{\mathbf{z}^i}) &= \mathbb{E} \left[ \sum_{t'=0} \mathbb{E}_{\mathbf{s}_t^i, \mathbf{a}_t^i \sim \pi_{\mathbf{z}^i}^i} \left[ \sum_{t=0}^{k_{t'}-1} \gamma^{t+k_{t'}} \tilde{r}_L^i(\mathbf{s}_t^i, \mathbf{a}_t^i | z_{t'}^i) \right] \right] \\ &= \mathbb{E} \left[ \sum_{t'=0} \mathbb{E}_{\mathbf{s}_t^i, \mathbf{a}_t^i \sim \pi_{\mathbf{z}^i}^i} \left[ \sum_{t=0}^{k_{t'}-1} \gamma^{t+k_{t'}} \left( \lambda \frac{A^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'})}{\mathcal{N}} + (1 - \lambda) r_L^i(\mathbf{s}_t^i, \mathbf{a}_t^i | z_{t'}^i) \right) \right] \right] \\ &= \frac{\lambda}{\mathcal{N}} \mathbb{E} \left[ \sum_{t'=0} \frac{\gamma^{t'} (1 - \gamma^{k_{t'}})}{1 - \gamma} A^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'}) \right] + (1 - \lambda) \mathbb{E} \left[ \sum_{t'=0} \mathbb{E}_{\mathbf{s}_t^i, \mathbf{a}_t^i \sim \pi_{\mathbf{z}^i}^i} \left[ \sum_{t=0}^{k_{t'}-1} \gamma^{t+k_{t'}} r_{L_A}^i(\mathbf{s}_t^i, \mathbf{a}_t^i | z_{t'}^i) \right] \right] \\ &\approx \frac{\lambda}{\mathcal{N}(1 - \gamma)} \mathbb{E} \left[ \sum_{t'=0} \gamma^{t'} A^{\pi_H}(\mathbf{s}_{t'}, \mathbf{z}_{t'}) \right] + (1 - \lambda) J^i(\pi_H, \pi_{\mathbf{z}^{-i}}, \pi_{\mathbf{z}^i}) \end{aligned}$$

1018  $\square$

1019 Similarly, the low level reward for every agent can be revised as above to align skill learning process  
 1020 with the task learning process, as the first term is related to the high level policy return.

## 1021 A.5 Algorithm Implementation

### 1022 A.5.1 Algorithm details

1023 As mentioned previously, our method can be implemented using both existing on-policy and off-  
 1024 policy RL algorithms. The bilevel RL problems presented in (10) and (15), respectively, should be  
 1025 solved sequentially. However, as previous works illustrated [77], jointly optimizing both policies  
 1026 by updating their corresponding parameters alternatively renders similar performances, which is  
 1027 consistent with our observation from the experiments. Based on this observation, the algorithm is  
 1028 presented in (1). The high-level policy  $\pi_{H_\theta}$  is decomposed into agent-level parameterized policies  
 $\pi_{H_\theta^i}^i$  to learn decentralized policies.

---

#### Algorithm 1 Hierarchical Multi-Agent Reinforcement Learning with CBFs (On-Policy)

---

```

1: Initialize randomly the high-level policy parameters  $\theta$  and parameters  $\nu$  of the low-level policy
   parameters  $\phi$  //network parameter sharing done between agents.
2: Set Grouping flag.
3: for each episode do
4:    $z_0^i \sim \pi_{H_\theta}(z^i | \mathbf{o}_0^i), \forall i \in \mathcal{N}$  // sample initial skill
5:   Set data buffer  $\mathcal{D} = \{\}$ 
6:   for  $t \leftarrow 0, \dots, T$  do
7:     Get action for agent  $i$  based on the selected skill  $a_t^i = \mu^i(\mathbf{o}_t^i | z_t^i; \phi_\nu) \forall i \in \mathcal{N}$  // sample
       primitive action
8:     Get  $s_{t+1}$  and  $r_H$  from environment and  $r_L^i \forall i \in \mathcal{N}$ 
9:     Sample skill for agent  $i, z_{t+1}^i \sim \pi_{H_\theta}^i(z^i | \mathbf{o}_t^i), \forall i \in \mathcal{N}_z(t)$  //  $\mathcal{N}_z(t)$  is the list of agents who
       have to update their skills based on termination scheme used.
10:    if  $\sim \text{Grouping}$  then
11:       $\mathcal{D} = \mathcal{D} \cup \{s_t, \{\mathbf{o}_t^i, z_t^i, \mathbf{a}_t^i, \phi_t^i, r_L^i, \mathbf{o}_{t+1}^i\}_{i \in \mathcal{N}}, r_H, s_{t+1}\}$ 
12:    else
13:       $\mathcal{D} = \mathcal{D} \cup \{s_t^{G^j}, \{\mathbf{o}_t^i, z_t^i, \mathbf{a}_t^i, \phi_t^i, r_L^i, \mathbf{o}_{t+1}^i\}_{i \in G^j}, r_H^{G^j}, s_{t+1}^{G^j}\}_{j \in \mathcal{N}_G}$ 
14:    end if
15:  end for
16:  Update high level policy network  $\theta$  according to (11), or, (12), using sampled data.
17:  Update low level policy network parameters  $\nu$  based on (23) and (24), using sampled data
18: end for

```

---

1029 We implement the algorithm using two different approaches inspired by existing works in the area of  
 1030 MARL. Notably, we implemented the algorithm in the following ways. It is important to note that  
 1031 under all these approaches, the hierarchical policy is learned agent-wise to allow for decentralized  
 1032 execution.  
 1033

- 1034 1. Joint policy learning: This approach involves solving (10) by optimizing the high-level policies  
 1035 of the agents jointly by learning the policy parameters  $\theta^i$ s.
- 1036 2. Group policy learning: This approach is inspired by the idea of learning the value function  
 1037 by decomposing the joint value function into the value function of groups of agents. In our  
 1038 fully cooperative setting, the total step reward can be composed as a sum of the rewards of the  
 1039 individual agents, i.e.  $r_H(s, \mathbf{a}) = \sum_{i=1}^N r_H^i(s^i, \mathbf{a}^i)$ . Let  $\{G^j\}_{j=1}^{N_G}$  represent the groups with each  
 1040  $G^j \subset \mathcal{N}$  such that,  $G^j \cap G^{j'} = \emptyset$  and  $\mathcal{N} = \cup_{i=1}^{N_G} G^i$ . We can define the group reward as the sum  
 1041 over the individual reward of the agents in the group. Then, under this approach we can express  
 1042 the objective in (10) as below (We only include the necessary symbols in the expectation to avoid



1043 symbol overload):

$$\begin{aligned}
J(\pi_H, \pi_z) &= \mathbb{E} \left[ \sum_{t'=0}^{\infty} \mathbb{E}_{\mathbf{a}_t \sim \pi_{z_{t'}}} \left[ \sum_{t=0}^{k_{t'}-1} \gamma^{t+k_{t'}} r_H(\mathbf{s}_t, \mathbf{a}_t) \right] \right] \\
&= \mathbb{E} \left[ \sum_{t'=0}^{\infty} \mathbb{E}_{\mathbf{a}_t \sim \pi_{z_{t'}}} \left[ \sum_{t=0}^{k_{t'}-1} \gamma^{t+k_{t'}} \underbrace{\sum_{j=1}^{N_G} \sum_{i=1}^N r_H^i(\mathbf{s}_t^i, \mathbf{a}_t^i) \mathbb{I}[i \in G^j]}_{=r_H^{G^j}(\mathbf{s}_t^{G^j}, \mathbf{a}_t^{G^j})} \right] \right] \\
&= \sum_{j=1}^{N_G} \mathbb{E} \left[ \sum_{t'=0}^{\infty} \mathbb{E}_{\mathbf{a}_t^{G^j} \sim \pi_{z_{t'}^{G^j}}} \left[ \sum_{t=0}^{k_{t'}-1} r_H^{G^j}(\mathbf{s}_t^{G^j}, \mathbf{a}_t^{G^j}) \right] \right] \tag{32}
\end{aligned}$$

1044 Here,  $\mathbf{s}^{G^j}$  and  $\mathbf{a}^{G^j}$  are the states and actions of the agents in group  $G^j$ . The agent skill policies  
1045 are learnt independently of each other. Hence, the joint skills policies can be written as  
1046  $\pi_z = \prod_{i=1}^N \pi_{z^i} = \prod_{j=1}^{N_G} \prod_{i \in G^j} \pi_{z^i} = \prod_{j=1}^{N_G} \pi_{z^{G^j}}$  where  $\pi_{z^{G^j}}$  is the joint skills policy of the agents  
1047 in the group  $G^j$ . Then, (32) can be written as:

$$\begin{aligned}
J(\pi_H, \pi_z) &= \sum_{j=1}^{N_G} \mathbb{E} \left[ \sum_{t'=0}^{\infty} \mathbb{E}_{\mathbf{a}_t^{G^j} \sim \pi_{z_{t'}^{G^j}}} \left[ \sum_{t=0}^{k_{t'}-1} r_H^{G^j}(\mathbf{s}_t^{G^j}, \mathbf{a}_t^{G^j}) \right] \right] \\
&= J^1(\pi_H^{G^1}) + \dots + J^{N_G}(\pi_H^{G^{N_G}})
\end{aligned}$$

1048 where,  $k_{t'}$  now becomes the time at which any agent in the group changes selects/switches to  
1049 a new skill. We drop the low-level policy corresponding to skills for brevity. This approach  
1050 enables better scalability as the overall problem is decomposed into one that learns policies  
1051 for the groups; however, this comes at the expense of performance as the original optimization  
1052 problem involves all the agents.

1053 Note that all these implementations retain the safety guarantee as it is offered by our safe skill-learning  
1054 approach.

### 1055 A.5.2 Implementation Details: METADRIVE Simulator

1056 Metadrive environments illustrated in figure 4 by default, include 3 basic sensors: lidar, sideDetector  
1057 and laneLineDetector which are used for detecting moving objects, sidewalks/solid lines, and broken/  
1058 solid lines respectively. The lidar state observation returns a state vector containing necessary  
1059 information for navigation tasks. The details of the observation can be found here <sup>2</sup>. The environ-  
1060 mental safety constraints include i. the road boundaries and ii. the booth for tollgate environment  
1061 only. And, each agent has to stay safe to all neighboring agents which are within the lidar sensing  
1062 radius. It is important to mention that we only use a dynamic bicycle model for the CBF constraints,  
1063 the environment uses pybullet physics engine for simulation. The kinematic bicycle model of the  
1064 vehicle is provided in (33).

$$\begin{aligned}
\dot{p} &= \frac{v \cos(\psi)}{1 - d\kappa(p)}, \\
\dot{d} &= v \sin(\psi), \\
\dot{\psi} &= \frac{v}{L} \tan(\delta) - \kappa(p)\dot{p}, \\
\dot{v} &= a.
\end{aligned} \tag{33}$$

1065 where  $p$  represents the position of the vehicle along the reference path,  $d$  is the lateral deviation of the  
1066 vehicle from the reference path,  $\psi$  denotes the yaw angle of the vehicle (heading angle relative to the  
1067 path), and  $v$  is the longitudinal velocity of the vehicle. The control inputs are  $\delta$ , which is the steering  
1068 angle, and  $a$ , which is the acceleration (or deceleration). The parameter  $L$  denotes the wheelbase of  
1069 the vehicle, and  $\kappa(p)$  is the curvature of the reference path at position  $s$ . It is important to add that,  
1070 these states are provided by the environment for each ego vehicle/agent in its observation.

<sup>2</sup><https://metadrive-simulator.readthedocs.io/en/latest/obs.html>

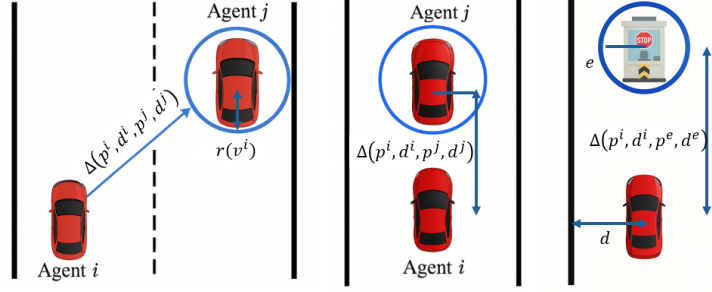


Figure 9: From left to right: illustration of i. inter agent safety constraints on another lane, ii. inter agent safety constraints on same lane and ii. environmental constraints.

In order to incorporate safety with respect to other agents and environmental obstacles, we define simple circles-based constraints as illustrated in figure 9. And, for the road boundaries we use the lateral deviation from the center lane to define safety functions. Let  $p^j$  and  $d^j$  denote the distance along the path and lateral deviation of another agent  $j$  with respect to agent  $i$ , located within within the sensing radius of agent  $i$ . This information is provided in the observation of agent  $i$  i.e.  $(d^i, p^j, d^j) \in \mathcal{O}^i$ . Note that, for an ego vehicle  $p^i$  is always 0 in its body frame. let  $\Delta(p^i, d^i, p^j, d^j)$  denote the euclidean distance between agent  $i$  and  $j$  respectively. Similarly, let  $p^e$  and  $d^e$  denote the longitudinal and lateral distance of an environmental obstacle, like a tollgate booth as illustrated in figure 9 from agent  $i$ , and  $\Delta(p^i, d^i, p^e, d^e)$  be the euclidean distance. Then, the safety constraints corresponding to other agent  $j$ , environmental constraints for an agent  $i$  are given below:

$$\Delta(p^i, d^i, p^j, d^j)^2 \geq r(v^i), \Delta(p^i, d^i, p^e, d^e)^2 \geq e; , \text{ and } |d^i| \leq c \quad (34)$$

where  $r(v^i)$  is a monotonically increasing function of  $v^i$ ,  $e \in \mathbb{R}_{>0}$  and  $c \in \mathbb{R}_{>0}$ .

Consider the CBF for the constraint for inter-agent safety defined as  $h(\mathcal{O}^i) = \Delta(p^i, d^i, p^j, d^j) - r(v^i)$  where  $\Delta(p^i, d^i, p^j, d^j) = \sqrt{(p^i - p^j)^2 + (d^i - d^j)^2}$ , then the corresponding CBF constraint is given by:

$$\underbrace{2(p^i - p^j) \left( \frac{v^i \cos(\psi^i)}{1 - d^i \kappa(p^i)} - \dot{p}^j \right) + 2(d^i - d^j) (v^i \sin(\psi^i) - \dot{d}^j) + 2r(v^i) \frac{dr(v^i)}{dv^i} a^i}_{L_f h(\mathcal{O}^i)} + \underbrace{\alpha \left( (p^i - p^j)^2 + (d^i - d^j)^2 - r(v^i)^2 \right)}_{h(\mathcal{O}^i)} \geq 0 \quad (35)$$

Additionally, we define CLFs for state reference tracking namely velocity and heading reference tracking. As we define termination state set corresponding to each skill, we define CLFs to drive the system to the termination state to be able to successfully execute the skill. However, it is important to note that the CLFs are not defined specific to skills and, hence, are generalized to the task. The CLFs corresponding to velocity and heading reference are given by the equation below. where  $v_t^i$  and  $v_{t_f}^i$  are the velocities at We parameterize the corresponding CLF constraint which in case of the velocity and heading constraints becomes the following:

$$V(v^i, v_{des}^i) = (v^i - v_{des}^i)^2; \quad (36)$$

$$V(\psi^i, \psi_{des}^i) = (\psi^i - \psi_{des}^i)^2 \quad (37)$$

where  $v_{des}$  is the desired velocity and  $\psi_{des}$  is the desired heading of the vehicle.

**Skills description.** We define five skills for each agent for the METADRIIVE environments which are: i. cruising, ii. accelerating (speed up), iii. yielding (slowing down), iv. left lane changing, and v. right lane changing.

**Cruising.** As the name suggests, this skill causes the vehicle to move at a constant speed for a fixed distance. The velocity CLF constraint in (36) is incorporated in the skill QP based policy by setting  $v_{des} = 0$ .

1099 *Accelerating/speeding up* This skill corresponds to an agent increasing its longitudinal velocity. Let  
 1100 the initial velocity at the time of the initiation of the skill be  $v_0$ . The skill terminates when the velocity  
 1101 reaches  $v_T = v_0 + dv$ , where  $dv \in \mathbb{R}_{>0}$  denotes the desired velocity increment and  $v_T$  is the velocity  
 1102 of the agent at the time of termination. Hence, the  $v_{des}$  is set to  $v_0 + dv$  for this skill in the low level  
 1103 QP based skill policy. This skill facilitates timely task completion and helps avoid congestion by  
 1104 preventing agents from unnecessarily slowing down the traffic flow.

1105 *Yielding/slowng down. Yielding.* This skill enables an agent to reduce its longitudinal velocity. Let  
 1106  $v_0$  be the velocity at skill initiation. The skill terminates when the velocity decreases to  $v_T = v_0 - dv$ ,  
 1107 where  $dv \in \mathbb{R}_{>0}$  is the desired reduction and  $v_T$  is the velocity at the time of termination of the skill.  
 1108 Thus, the velocity CLF in (36) is incorporated to the QP based skill policy by setting  $v_{des} = v_0 - dv$   
 1109 Yielding allows agents to safely navigate dense traffic and cooperate by yielding to agents approaching  
 1110 from other directions, thereby facilitating efficient joint task completion.

1111 *Left Lane Changing.* This skill enables an agent to switch to the adjacent left lane. Let  $d_{left,0}$  denote  
 1112 the distance to the left road boundary and  $d_0$  be the lateral deviation from the center lane at the  
 1113 initiation of the skill. Let  $d_{lane}$  denote the lane width, provided by the environment. The target  
 1114 lateral position is computed as  $d_{left,T} = d_{left,0} - d_{lane} + d_0$ . We define the desired heading angle as  
 1115 a function of  $e$  i.e.  $\phi_{des} = f(e)$ , where  $e = d_{left,T} - d_{left,0}$ , and incorporate it into the CLF constraint  
 1116 on heading in (37) via a QP-based low-level skill control policy. This skill is essential for scenarios  
 1117 such as tollgates, bottlenecks, and merging zones, where adaptive lateral movement is required to  
 1118 maintain flow and avoid congestion.

1119 *Right Lane Changing.* This skill enables an agent to switch to the adjacent right lane. Let  $d_{right,0}$   
 1120 denote the distance to the right road boundary and  $d_0$  be the lateral deviation from the center lane at  
 1121 the initiation of the skill. Let  $d_{lane}$  denote the lane width, provided by the environment. The target  
 1122 lateral position is computed as  $d_{right,T} = d_{right,0} - d_{lane} + d_0$ . We define the desired heading angle  
 1123 as a function of  $e$ , i.e.,  $\phi_{des} = f(e)$ , where  $e = d_{right,T} - d_{right,0}$ , and incorporate it into the CLF  
 1124 constraint on heading in (37) via a QP-based low-level skill control policy. This skill is useful in  
 1125 scenarios such as avoiding obstacles, preparing for highway exits, or facilitating merging from the  
 1126 right, enabling smoother and safer lane coordination.

1127 As mentioned previously we include a timeout to ensure finite time termination of the skills for  
 1128 algorithmic sanity.

1129 **Intrinsic reward** The intrinsic reward serves as an internal shaping signal that guides  
 1130 each agent toward desirable driving behaviors—smoothness, lane-centering, and reference-  
 1131 tracking—independently of any external task reward. By combining several negative penalties,  
 1132 the agent is discouraged from abrupt maneuvers, drifting off its lane, or deviating from its target  
 1133 speed and heading.

$$r_1^i = -c_1 \|\mathbf{a}^i\|^2 - c_2 \left( \frac{v^i - v_{des}}{v_{des}} \right)^2 - c_3 \left( \frac{\psi^i - \psi_{des}}{\pi} \right)^2 - c_4 \|d^i\|^2 \quad (38)$$

1134 Where  $c_1, c_2, c_3, c_4, c_5$  are the coefficients of each terms. The intrinsic reward combines penalties on  
 1135 aggressive acceleration and steering, deviations from reference speed and heading, and lateral offset  
 1136 from the lane center, with a small forward-progress bonus. Together, these terms encourage smooth,  
 1137 centered, and efficient driving without abrupt maneuvers.

1138 **CBF-Based Safe Skill Implementation.** As described earlier, we parameterize the QP-based  
 1139 deterministic policy  $\mu(o \mid z, \phi)$  for each skill  $z$ , as detailed in Section 5.2. The QP formulation  
 1140 is parameterized through the objective, as well as the CBF and CLF constraints. In particular, we  
 1141 introduce parameters for the class- $\mathcal{K}$  function in the CBF constraint, the shape and size of the circular  
 1142 safety regions, and the convergence rate of the CLF constraint. To ensure validity, all parameters are  
 1143 constrained to be nonnegative using box constraints of the form  $(0, \phi_H]$ , where  $\phi_H$  is a fixed upper  
 1144 bound applied uniformly across all parameters.

### 1145 A.5.3 Implementation Details: Lidar Suite Simulator

1146 **Target.** Agents dynamics in this environment are governed by double integrator dynamics with  
 1147 continuous-time states  $\mathbf{x}_i = [p_x^i, p_y^i, v_x^i, v_y^i]^\top \in \mathbb{R}^4$ , where  $\mathbf{p}_i = [p_x^i, p_y^i]^\top \in \mathbb{R}^2$  denotes position and  
 1148  $\mathbf{v}_i = [v_x^i, v_y^i]^\top \in \mathbb{R}^2$  denotes velocity; control inputs  $\mathbf{u}_i = [a_x^i, a_y^i]^\top \in \mathbb{R}^2$  represent accelerations

1149 along the Cartesian axes, and the agent dynamics follow  $\dot{\mathbf{x}}_i = [v_x^i, v_y^i, a_x^i, a_y^i]^\top$ , with velocity and  
 1150 control constraints bounded in  $[-10, 10]$  and  $[-1, 1]$  respectively.

1151 **Spread.** The underlying dynamics, state representation, control parametrization, and bounds are  
 1152 identical to those in the Target environment, i.e., agents evolve under double integrator dynamics  
 1153 with four-dimensional state vectors comprising positions and velocities, and two-dimensional control  
 1154 inputs that modulate accelerations; the bounded constraints on velocities and controls likewise remain  
 1155  $[-10, 10]$  and  $[-1, 1]$ , respectively.

1156 **Target Bicycle** In the Target Bicycle environment, agents dynamics are modeled using a bicy-  
 1157 cle model where the state vector  $\mathbf{x}_i = [p_x^i, p_y^i, \cos \theta^i, \sin \theta^i, v^i]^\top \in \mathbb{R}^5$  encapsulates the agent's  
 1158 position, heading orientation via its sine and cosine, and scalar speed; the control input  $\mathbf{u}_i =$   
 1159  $[\delta^i, a^i]^\top$  includes the steering angle  $\delta^i$  and linear acceleration  $a^i$ , with dynamics specified as  
 1160  $\dot{\mathbf{x}}_i = [v \cos \theta, v \sin \theta, -v \sin \theta \tan \delta, v \cos \theta \tan \delta, a]^\top$ , where  $v \in [-10, 10]$  and  $\delta \in [-1.47, 1.47]$   
 1161 are the admissible bounds on speed and steering angle.

1162 Similar to the METADRIVE environments, we define CBFs based on distance/proximity as in  
 1163 (34). In our formulation, we do not employ a separate heading CLF; instead, we regulate agent  
 1164 behavior using only the velocity CLF by aligning the desired velocity vector with the desired heading  
 1165 direction. Specifically, given a desired heading direction  $\theta_{\text{des}}$ , we compute a target velocity vector  
 1166 as  $\mathbf{v}_{\text{des}} = v_{\text{mag}} [\cos(\theta_{\text{des}}), \sin(\theta_{\text{des}})]^\top$ , where  $v_{\text{mag}}$  is a task-specific magnitude (e.g., maximum safe  
 1167 speed or adaptive goal-directed speed). This  $\mathbf{v}_{\text{des}}$  is then used in the velocity CLF described below to  
 1168 guide both the speed and heading behavior of the agent.

1169 For the Target and Spread environments, where the agent's longitudinal and lateral velocities  $v_x^i, v_y^i$   
 1170 are explicit components of the state, we define the Control Lyapunov Function (CLF) to track a  
 1171 desired velocity vector  $\mathbf{v}_{\text{des}} = [v_{x,\text{des}}, v_{y,\text{des}}]^\top \in \mathbb{R}^2$  as:

$$V(v_x^i, v_y^i) = \frac{1}{2} \left( (v_x^i - v_{x,\text{des}})^2 + (v_y^i - v_{y,\text{des}})^2 \right).$$

1172 Given the double integrator dynamics, where the control inputs  $a_x^i$  and  $a_y^i$  directly affect the velocity  
 1173 via  $\dot{v}_x^i = a_x^i$  and  $\dot{v}_y^i = a_y^i$ , the time derivative of the CLF becomes:

$$\dot{V} = (v_x^i - v_{x,\text{des}})a_x^i + (v_y^i - v_{y,\text{des}})a_y^i.$$

1174 To enforce convergence, we include the following CLF constraint in the QP:

$$(v_x^i - v_{x,\text{des}})a_x^i + (v_y^i - v_{y,\text{des}})a_y^i \leq -\alpha \left( (v_x^i - v_{x,\text{des}})^2 + (v_y^i - v_{y,\text{des}})^2 \right),$$

1175 where  $\alpha > 0$  is a tunable rate parameter.

1176 In the Target Bicycle environment, each agent's velocity is represented as a scalar speed  $v^i \in \mathbb{R}$ , and  
 1177 the longitudinal acceleration  $a^i \in \mathbb{R}$  serves as the control input directly affecting the speed via  $\dot{v}^i = a^i$ .  
 1178 To drive  $v^i$  to a desired speed  $v_{\text{des}} \in \mathbb{R}$ , we define the CLF as:

$$V(v^i) = \frac{1}{2} (v^i - v_{\text{des}})^2.$$

1179 The time derivative of this CLF, using  $\dot{v}^i = a^i$ , is:

$$\dot{V} = (v^i - v_{\text{des}})a^i.$$

1180 To enforce convergence of the velocity to the target, the following CLF constraint is imposed in the  
 1181 QP:

$$(v^i - v_{\text{des}})a^i \leq -\alpha (v^i - v_{\text{des}})^2,$$

1182 where  $\alpha > 0$  is a tunable parameter determining the exponential rate of convergence.

1183 **Skills Description.** We define four interpretable low-level skills for each agent in the LiDAR suite  
 1184 environments, which are: i. speeding up, ii. slowing down, iii. turning left, and iv. turning right.

1185 *Speeding Up.* This skill allows an agent to increase its translational velocity to navigate more  
 1186 efficiently through open space or move rapidly towards its assigned goal. Let the velocity at the start  
 1187 of the skill be  $\mathbf{v}_0 \in \mathbb{R}^2$ ; the skill is considered complete once the velocity reaches  $\mathbf{v}_T = \mathbf{v}_0 + \Delta \mathbf{v}$ ,  
 1188 where  $\Delta \mathbf{v} \in \mathbb{R}^2$ , with  $\|\Delta \mathbf{v}\| > 0$ . The CLF constraint associated with velocity tracking is enforced

1189 by setting the desired velocity  $\mathbf{v}_{\text{des}} = \mathbf{v}_0 + \Delta \mathbf{v}$  in the QP-based low-level control policy. This skill  
 1190 facilitates faster progression in low-density scenarios while maintaining safety guarantees via CBFs.

1191 *Slowing Down.* This skill enables the agent to reduce its speed when approaching dynamic obstacles,  
 1192 other agents, or goal regions where precise navigation is required. Let the initial velocity be  $\mathbf{v}_0 \in \mathbb{R}^2$ ,  
 1193 and define the skill termination condition as reaching  $\mathbf{v}_T = \mathbf{v}_0 - \Delta \mathbf{v}$ , where  $\Delta \mathbf{v} \in \mathbb{R}^2$  and  $\|\Delta \mathbf{v}\| > 0$ .  
 1194 To enforce this behavior, the velocity CLF is specified using  $\mathbf{v}_{\text{des}} = \mathbf{v}_0 - \Delta \mathbf{v}$ . This skill is critical for  
 1195 smooth and safe deceleration, reducing the risk of collisions and improving maneuverability near  
 1196 goals.

1197 *Turning Left.* This skill enables the agent to redirect its heading counterclockwise to avoid obstacles  
 1198 or navigate toward an alternate subgoal. Let the initial heading be  $\theta_0$ , and the desired heading be  
 1199  $\theta_{\text{des}} = \theta_0 + \Delta\theta$ , where  $\Delta\theta > 0$ . Instead of using a heading CLF, we convert the desired heading into a  
 1200 velocity vector using a fixed target magnitude  $v_{\text{mag}}$ , and set  $\mathbf{v}_{\text{des}} = v_{\text{mag}}[\cos(\theta_{\text{des}}), \sin(\theta_{\text{des}})]^\top$ . This  
 1201 desired velocity is then tracked via the standard velocity CLF in the QP. The turning-left skill is  
 1202 useful for curvature adjustment in cluttered or multi-agent scenes.

1203 *Turning Right.* This skill enables clockwise heading correction using the same velocity CLF formula-  
 1204 tion. Given initial heading  $\theta_0$  and desired adjustment  $\Delta\theta > 0$ , the target heading is  $\theta_{\text{des}} = \theta_0 - \Delta\theta$ ,  
 1205 which is converted to a velocity vector as  $\mathbf{v}_{\text{des}} = v_{\text{mag}}[\cos(\theta_{\text{des}}), \sin(\theta_{\text{des}})]^\top$ . The agent’s translational  
 1206 motion is then regulated by the velocity CLF, enforcing convergence to the desired heading direction  
 1207 through velocity tracking. This skill is particularly effective in obstacle-rich or competitive zones  
 1208 where sharp directional changes are required.

## 1209 A.6 Additional Results

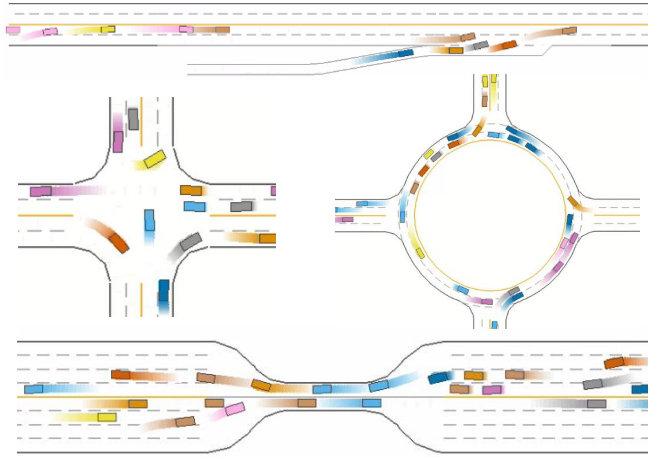


Figure 10: Visualization of agents’ trajectories under the trained model: for each vehicle, its ten most recent positions are shown in a unique hue, with brightness fading from light (earlier points) to dark (later points) to indicate temporal progression.

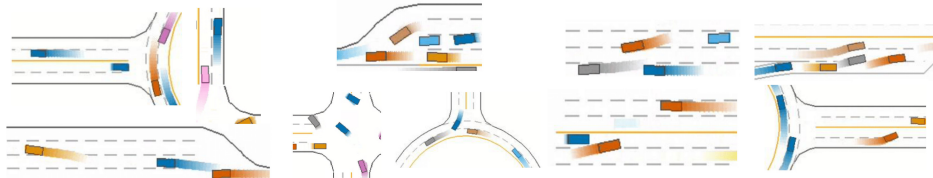


Figure 11: Visualization of skills under the trained model. From left to right: (a) acceleration (speed-up), (b) yielding (slow-down), and (c) lane changing

1210 Here, we present results to further demonstrate the merit of our proposed algorithm. Specifically, we  
 1211 examine the trajectories in detail to analyze the distinct skills exhibited by the agents, offering insight

1212 into how our method sustains an extremely high success rate while also outperforming the baselines  
1213 in terms of energy efficiency and episode duration. To this end, we employ a visualization scheme in  
1214 which each agent’s ten most recent positions are depicted in a distinct color, with brightness gradually  
1215 diminishing from light (earlier points) to dark (later points) to convey temporal progression. An  
1216 example of this visualization is depicted in Figure 10. By looking more closely, we can identify  
1217 some of the skills that the agents exhibit during the episode. For instance, as agents approach their  
1218 destinations (i.e., exit points of the environment), they execute speed up skill only when safety  
1219 constraints are satisfied, thereby enhancing road throughput and reducing average travel time. Several  
1220 such examples are shown in Figure 11a. Yielding (i.e., controlled slow down/deceleration) is another  
1221 common skill agents employ to maintain a safe distance from other vehicles. By anticipating the need  
1222 to yield, agents avoid abrupt braking, which both enhances safety and reduces overall episode energy  
1223 consumption. Several such examples are shown in Figure 11b. Another skill exhibited by agents  
1224 is lane changing: when executed safely at appropriate locations, it enhances road throughput and  
1225 consequently reduces the average episode duration. An example of such lane changing is depicted in  
1226 11c where the agent with brown color is changing its lane— creating a safe gap into which the grey  
1227 agent then merges onto the main road.