

Takeaways

We proposed a *layered architecture* incorporating *safety-informed reward design for MARL* and CBVF-based safety filters during training via curriculum learning, enabling proactive conflict resolution by avoiding potential conflict zones and improved coordination efficiency.

Conflicting Constraints

Main challenge in composing safe sets in multi-agent interaction: Leaky Corners arising from Conflicting Constraints



Layered Safety

1. Potential Conflict Zone as Soft Constraint for MARL

Based on local observations, GNN-based MARL Policy learns navigation while incentivized to avoid potential conflict zones.

2. Control Barrier-Value Function (CBVF)-based Safety Filtering

Active safety filtering against agent with worst CBVF value, based on the pairwise relative dynamics.



Resolving Conflicting Constraints in Multi-Agent Reinforcement Learning with Layered Safety

Jason J. Choi*, Jasmine Jerry Aloor*, Jingqi Li*, Maria G. Mendoza, Hamsa Balakrishnan†, Claire J. Tomlin†





Potential Conflict Zone:

If **N>2** are within the potential conflict zone, safety of all agents are not guaranteed anymore.



Safety Informed Training

Curriculum learning

- First train without any safety filter or reward penalty Activate safety filter
- Gradually increase safety radius and conflict radius

 $\mathcal{C}_{\text{conflict}} :=$

Combined Reward:

Highlighted Results

1. Air Taxi Operation without Centralized Traffic Control

8 vehicles departing from various locations in Bay Area, merging to land in San Francisco





position of robot 1

position of robot 3

position of robot 2

initial states

(waypoints

Avg. Travel Time: 676 sec. Near Collision: 0.055%

2. Hardware experiments with quadrotors

Three drones going through the same air corridor to get to their landing spots.

Safety-blind MARL



x (m)

Safety-informed MARL





3. An extensive empirical study of different safety-informing methods in MARL

Our method generalizes better than prior works in terms of number of agents. (N = number of agents, M = number of waypoints)

Safety-informed training vs no safety filter or safety-informed reward (Airtaxi scenarios)

	•	•		
Methods	Merging Scenario (N=8, M=5)			
Methous	Travel $t(s)(\downarrow)$	Near collision(%)(\downarrow)	$Conflict(\%)(\downarrow)$	
Safety-blind	675.6	0.055	2.4	
No penalty	617.9	0.042	5.5	
Proposed	450.5	0.021	3.2	
	Intersection Scenario (N=16, M=6)			
Methods	Travel $t(s)(\downarrow)$	Near collision(%)(\downarrow)	Conflict(%)(\downarrow)	
Safety blind	0074	0.050	0.1	
Safety Dilliu	987.4	0.058	2.1	
No penalty	987.4	0.058	3.8	
No penalty Proposed	987.4 780.5 660.8	0.058 0.129 0.056	2.1 3.8 1.6	

N =

N=

References

• Leaky Corners: I. Mitchell. "Scalable calculation of reach sets and tubes for nonlinear systems with terminal integrators." HSCC 2011 • InforMARL: S. Nayak, et al. "Scalable multi-agent reinforcement learning through intelligent information aggregation." ICML 2023. • CBVF: J. Choi, et al. "Robust control barrier-value functions for safety-critical control.", CDC 2021.

• **DG-PPO:** S. Zhang, et al. "Discrete GCBF Proximal Policy Optimization for Multi-agent Safe Optimal Control." ICLR 2025. • Exponential CBF: A. Agrawal and K.Sreenath. "Discrete CBFs for safety-critical control of discrete systems with application to bipedal robot navigation." RSS 2017. Acknowledgement: This work is supported in part by the NASA ULI in Safe Aviation Autonomy, NASA Clean Sheet Airspace Operating Design project, DARPA Assured Autonomy, DARPA ANSR, and ONR Basic Research Challenge in Multibody Control Systems program.



NASA ULI

Safe Aviation Autonom





Penalty for entering potential conflict zone:

		_	$\begin{bmatrix} (ij) \end{bmatrix}$
\sum	$\max\{0, r_{\text{conflict}} - \text{dist}(s^{(ij)})\}$	$ imes \max\left\{0, -\left[x^{(ij)} ight. ight. ight. ight.$	$\left[\begin{array}{c} y^{(ij)} \end{array} ight] \left[egin{matrix} v_x^{(ij)} \ v_x^{(ij)} \end{array} ight] ight\},$
$j ext{dist}(s^{(ij)}) {<} r_{ ext{conflict}} \}$			
		relativ	ve distance change

 $\mathcal{R}_{\text{total}}(o_k^{(i)}, a_k^{(i)}) = \mathcal{R}_{\text{tracking}}(o_k^{(i)}, a_k^{(i)}) + \rho_{\text{goal}}\mathcal{R}_{\text{goal}}(o_k^{(i)}, a_k^{(i)}) - \rho_{\text{conflict}}\mathcal{C}_{\text{conflict}}$

Comparison to other model-based and model-free baselines: (Crazyflie scenarios)

4:	Methods	Goal reach(%)	Near collision(%)
	DG-PPO	96 ± 11.8	0.04 ± 0.16
	Exponential CBF	100 ± 0	0.0 ± 0.0
	Our Method	100 ± 0	0.0 ± 0.0
8:	Methods	Goal reach(%)	Near collision(%)
	DG-PPO	100 ± 0	9.1 ± 2.7
	Exponential CBF	93 ± 8.9	8.8 ± 10.7
	Our Method	100 ± 0	0.0 ± 0.0