# SigmaRL: A Sample-Efficient and Generalizable Multi-Agent Reinforcement Learning Framework for Motion Planning

Jianye Xu,          Pan Hu,          Bassam Alrifaee

Informatik 11 Embedded Software | RWTH AACHEN UNIVERSITY
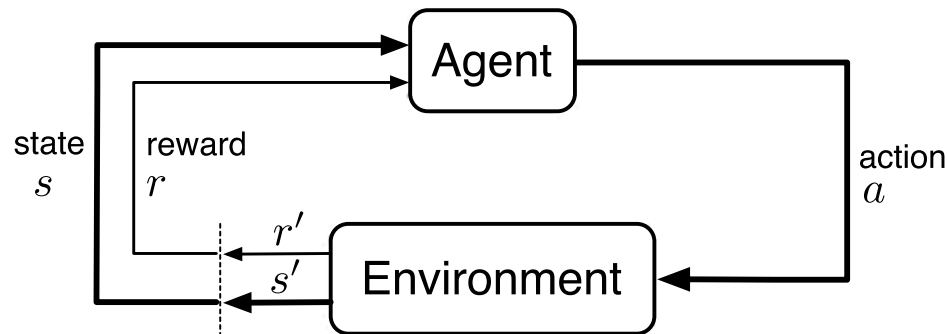
der Bundeswehr Universität München

Preprint | Code | Video

# Introduction

► Multi-Agent Reinforcement Learning (MARL) for motion planning of Connected and Automated Vehicles (CAVs)

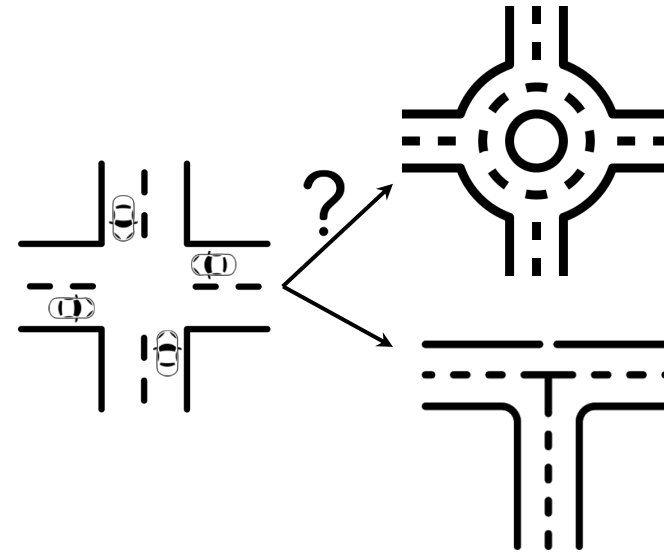**1.  Sample efficiency**

**2.  Generalization**

A $sample \coloneqq (s, a, r', s')$
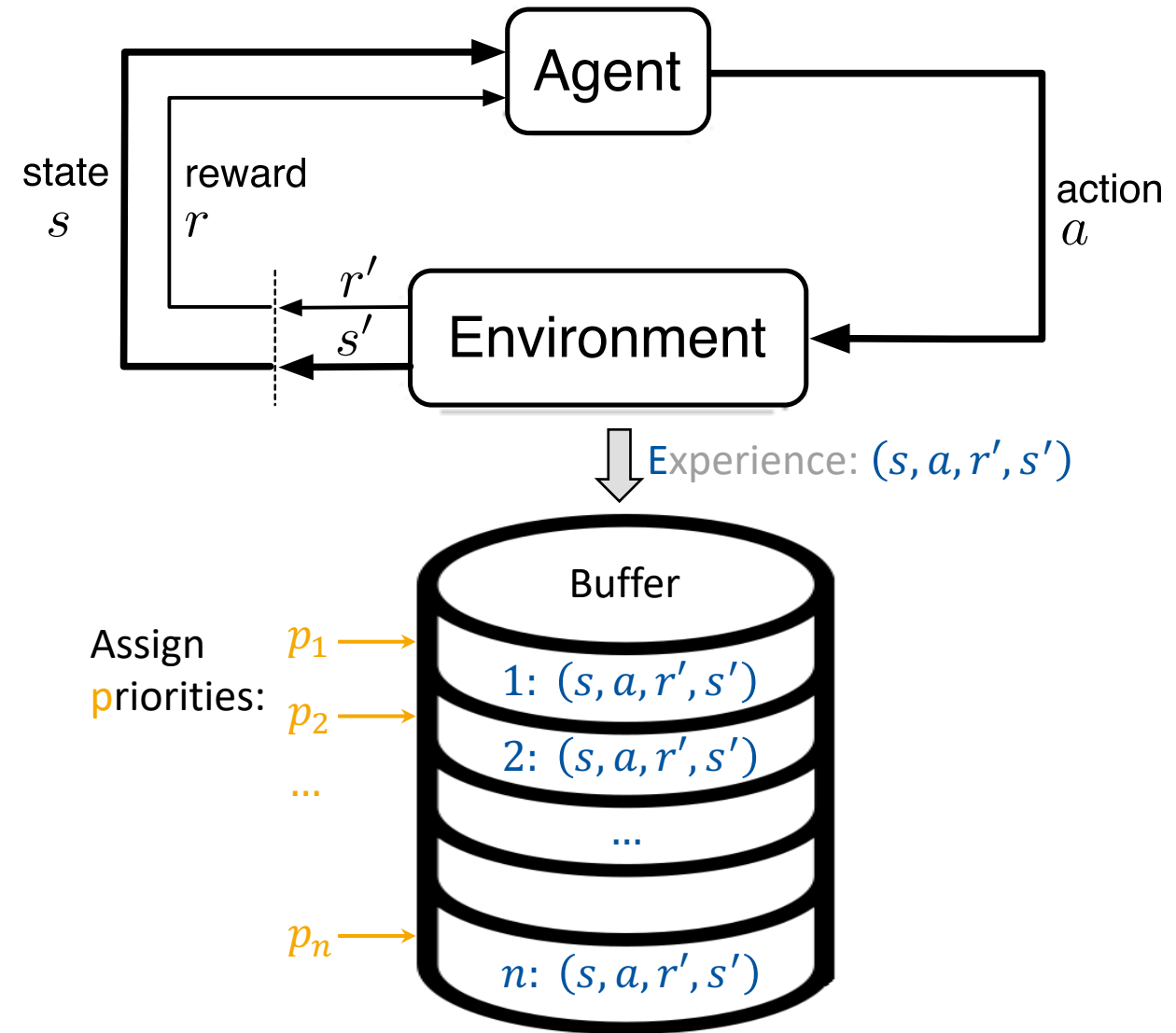


R. S. Sutton et al., Fig. 3.1 modified

# Related Work | Enhance Sample Efficiency

▶ **E**xperience replay buffer [1]

  ▪ Independent and identically distributed (i.i.d.) assumption

▶ **P**rioritized experience replay buffer [2]



[1] L.-J. Lin, "Self-Improving Reactive Agents Based on Reinforcement Learning, Planning and Teaching," *Mach Learn*, 1992.
[2] T. Schaul et al., "Prioritized Experience Replay," *in International Conference on Learning Representations (ICLR)*, 2016.

SigmaRL: A Sample-Efficient and Generalizable Multi-Agent Reinforcement Learning Framework for Motion Planning | ITSC 2024 | September 23
Jianye Xu, Pan Hu, Bassam Alrifaee

# Related Work | Enhance Generalization

▶ Regularization techniques such as dropout and early stopping [1]

▶ Domain randomization [2]

▶ Data augmentation [3]

▶ Better optimization without overfitting [4]

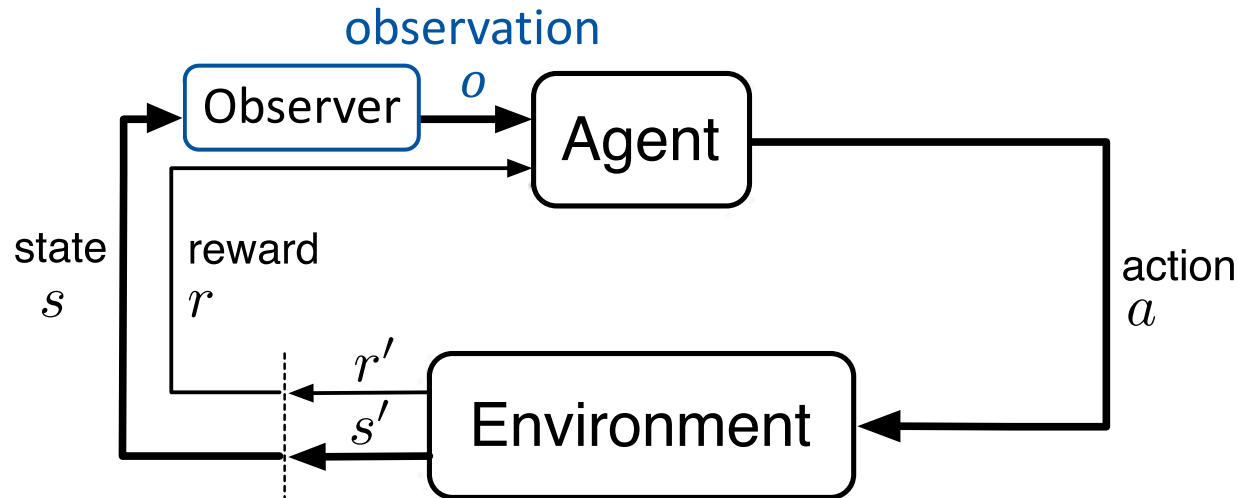[1] J. Farebrother et al., "Generalization and Regularization in DQN," *arXiv*, 2020.
[2] Tobin, Josh, et al., "Domain randomization for transferring deep neural networks from simulation to the real world," *IEEE/RSJ IROS*, 2017.
[3] Connor Shorten, and Taghi M. Khoshgoftaar, "A survey on image data augmentation for deep learning," *Journal of big data*, 2019.
[4] Cobbe, Karl W., et al. "Phasic policy gradient." *International Conference on Machine Learning*, 2021.

▶ Observation design: under-explored
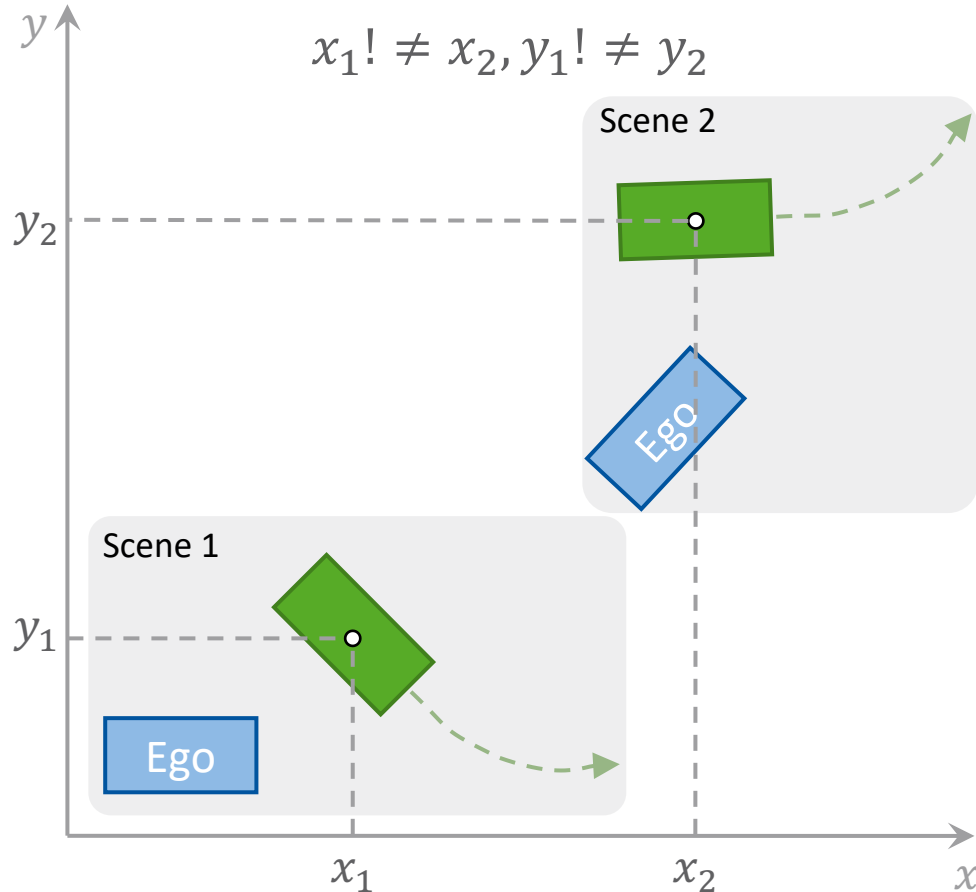
▶ Observer (also called observation function )
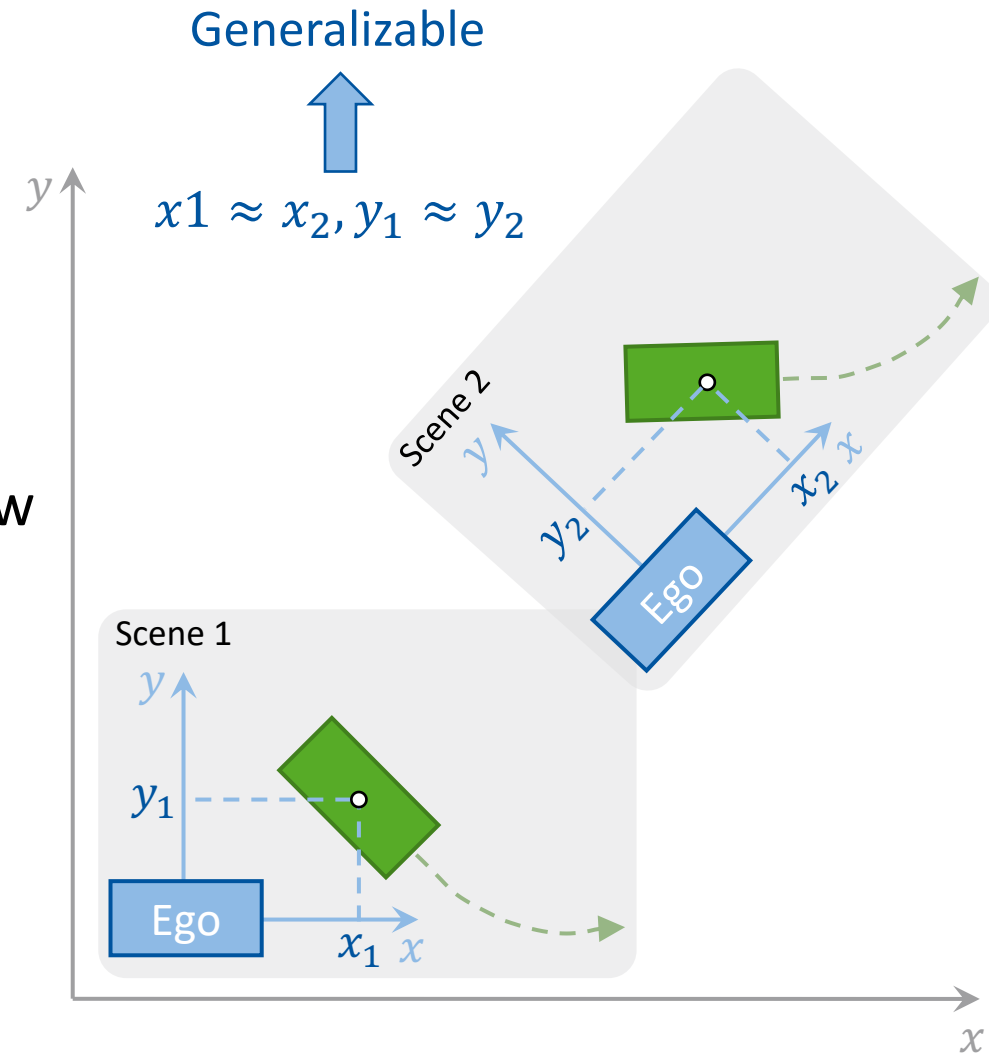


➤ Five observation-design strategies

▶ Use **ego view** instead of bird-eye view



Generalizable

$$x1 \approx x_2, y_1 \approx y_2$$

$$x_1! \neq x_2, y_1! \neq y_2$$

Bird-eye view:

Ego View (**our**):

SigmaRL: A Sample-Efficient and Generalizable Multi-Agent Reinforcement Learning Framework for Motion Planning | ITSC 2024 | September 23
Jianye Xu, Pan Hu, Bassam Alrifaee

▶ Observe **vertices** of surrounding agents instead of poses and dimensions

Poses and dimensions:



Vertices (**our**):

▶ Observe **distances** to surrounding agents

▶ Observe **distances** to lane boundaries instead of sampled points

Sample points:

$(x_1, y_1)$ ... $(x_n, y_n)$

Distances (**our**):

$d_{LB} := \min\{ d_1, d_2, \quad d_3 - \frac{w}{2}, \quad d_4, d_5 \}$

$w/2$

Ego

Ego

$d_{RB} := \min\{ d'_1, d'_2, \quad d_3 - \frac{w}{2}, \quad d'_4, d'_5 \}$

$(x'_1, y'_1)$ ... $(x'_n, y'_n)$

▶ Observe **distances** to lane center lines

▶ Own speed $v$

▶ Reference path

▶ Velocity of surrounding agents



Sampled points from reference path

► Integrated our observation-design strategies into multi-agent PPO [1-2]



[1] J. Schulman *et al.*, "Proximal Policy Optimization Algorithms," *arXiv*, 2017.
[2] R. Lowe *et al.*, "Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments,"
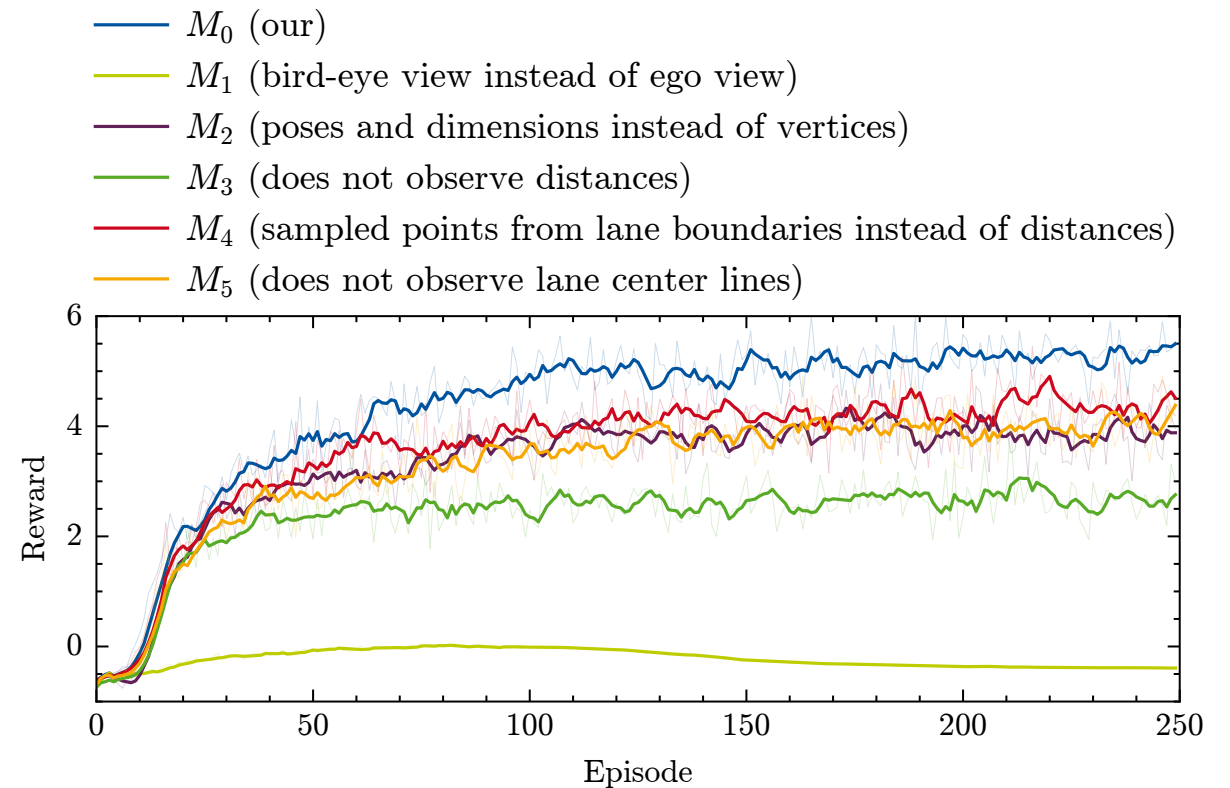in *Advances in Neural Information Processing Systems*, Curran Associates, Inc., 2017

https://github.com/cas-lab-munich/sigmarl

SigmaRL: A Sample-Efficient and Generalizable Multi-Agent Reinforcement Learning Framework for Motion Planning | ITSC 2024 | September 23
Jianye Xu, Pan Hu, Bassam Alrifaee

# Evaluations | Training

► An intersection with four agents

► ~1 million samples

► < 1 h of training time with a single CPU

► Ablation studies with six models

- $M_0$: All five observation-design strategies
- $M_{i \in \{1,...,5\}}$: Omits the $i_{\text{th}}$ strategy



Legend:
- $M_0$ (our)
- $M_1$ (bird-eye view instead of ego view)
- $M_2$ (poses and dimensions instead of vertices)
- $M_3$ (does not observe distances)
- $M_4$ (sampled points from lane boundaries instead of distances)
- $M_5$ (does not observe lane center lines)

► Four unseen scenarios

► 32 one-minute simulations per scenario per model

► Collision rate of each simulation: proportion of time steps in which collisions occur



(a) With 15 agents

(b) With 6 agents

(c) With 8 agents

(d) With 8 agents

SigmaRL: A Sample-Efficient and Generalizable Multi-Agent Reinforcement Learning Framework for Motion Planning | ITSC 2024 | September 23
Jianye Xu, Pan Hu, Bassam Alrifaee

$M_0$ (our)

$M_1$ (do not use an ego view)

$M_2$ (do not observe vertices of surrounding agents)

$M_3$ (do not observe distances to surrounding agents)

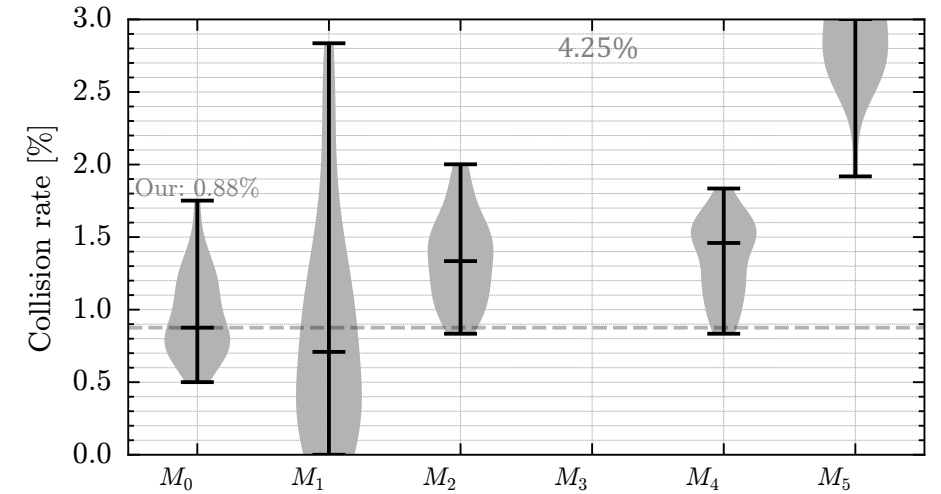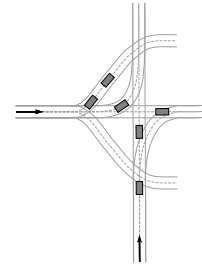$M_4$ (do not observe distances to lane boundaries)

$M_5$ (do not observe distances to lane center lines)

SigmaRL: A Sample-Efficient and Generalizable Multi-Agent Reinforcement Learning Framework for Motion Planning | ITSC 2024 | September 23
Jianye Xu, Pan Hu, Bassam Alrifaee

# Evaluations | Collision Rate

SigmaRL: A Sample-Efficient and Generalizable Multi-Agent Reinforcement Learning Framework for Motion Planning | ITSC 2024 | September 23
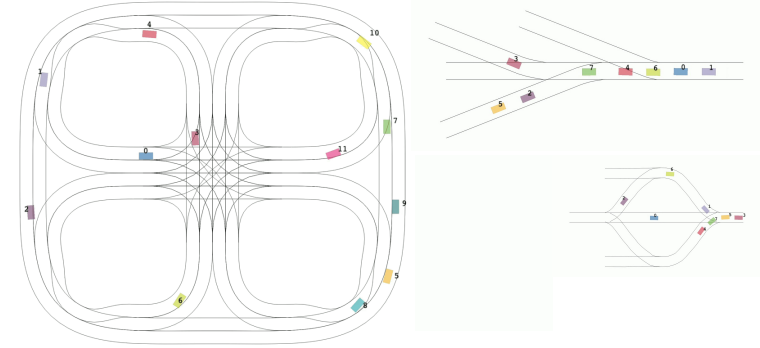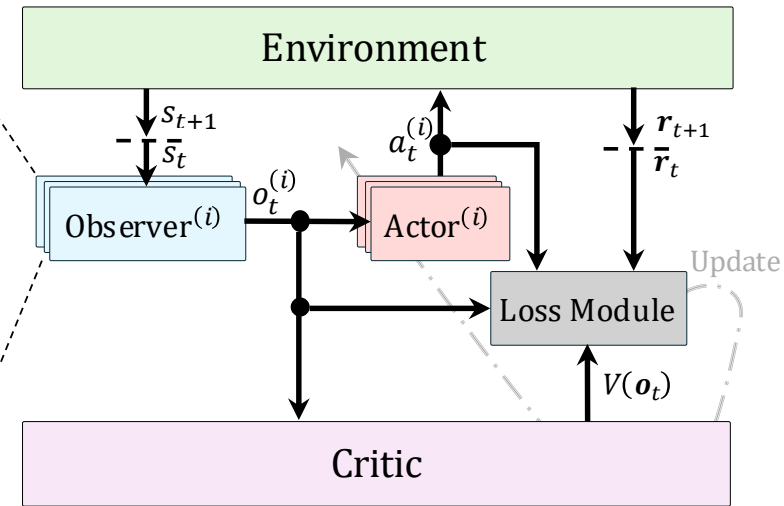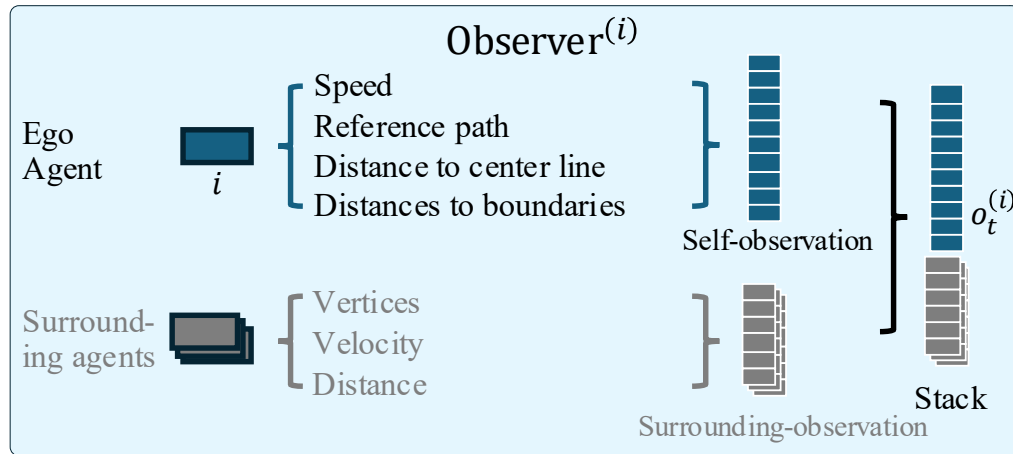Jianye Xu, Pan Hu, Bassam Alrifaee

# Summary

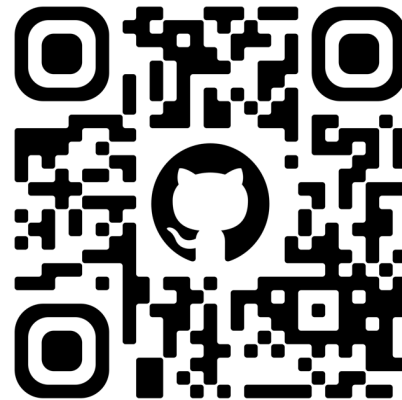▶ **SigmaRL**, a <u>sample-efficient</u> and <u>generalizable</u> MARL Framework for motion planning of CAVs

< 1 h training time on a single CPU    Completely unseen scenarios

Five observation-design strategies

Preprint | Code | Video

# Backup
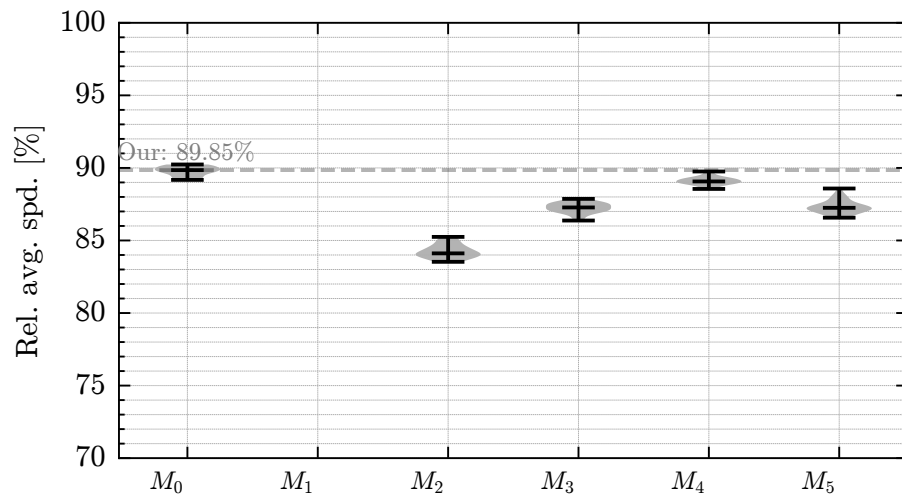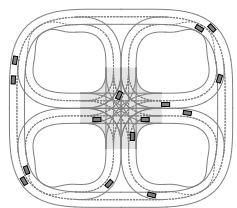
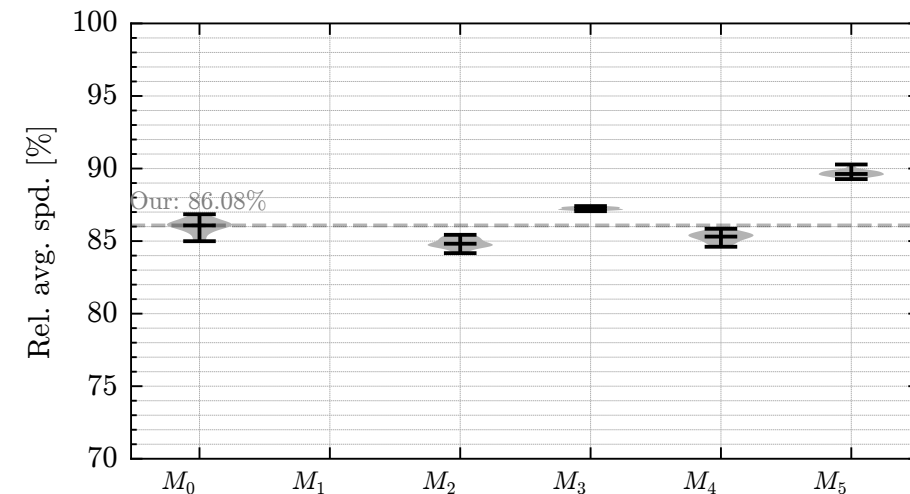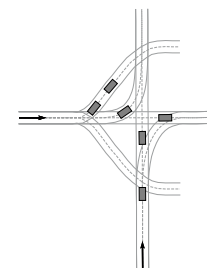# Evaluations | Average Speed

## Scenario (a)



## Scenario (b)



## Scenario (c)



## Scenario (d)

SigmaRL: A Sample-Efficient and Generalizable Multi-Agent Reinforcement Learning Framework for Motion Planning | ITSC 2024 | September 23
Jianye Xu, Pan Hu, Bassam Alrifaee