

# Learning Cooperation for Partially Observable Multi-Agent Path Finding

M.Sc. Thesis Defence

---

Qiushi Lin

Simon Fraser University

# Motivations



Automated Warehouses<sup>1</sup>



Sorting Centers<sup>2</sup>



Video Games<sup>3</sup>



Drone Swarm Systems<sup>4</sup>

---

<sup>1</sup><https://www.amazon.science/latest-news/how-amazon-robots-navigate-congestion>

<sup>2</sup><https://www.wired.com/story/amazon-warehouse-robots>

<sup>3</sup>[https://upload.wikimedia.org/wikipedia/commons/0/02/3\\_cossacks\\_european\\_wars.jpg](https://upload.wikimedia.org/wikipedia/commons/0/02/3_cossacks_european_wars.jpg)

<sup>4</sup><https://futureoflife.org/wp-content/uploads/2019/04/Why-ban-lethal-AI-1030x595.jpg>

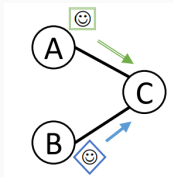
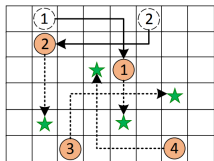
- Planning paths for multiple agents are usually modelled as **NP-hard** combinatorial search problems in discretized worlds with discretized time steps
- Search-based planning algorithms are **computationally expensive** and **cannot generalize well among instances**
- Learning-based methods provide solvers that are computationally efficient but usually have **poor solution quality**
- Most multi-agent reinforcement learning (MARL) methods can only deal with instances with small numbers of agents (usually 5-10 agents) and are usually **not scalable** when it comes to congested environments
- It is not clear in the literature how various types of cooperation and objectives among agents can be learned

- Introduction
  - Problem Definitions
  - Related Works
- Single-Objective Cooperation
  - Methodology (SACHA)
  - Empirical Evaluation
- Bi-Objective Cooperation
  - Methodology (MFC-EQ)
  - Empirical Evaluation
- Conclusion and Future Work

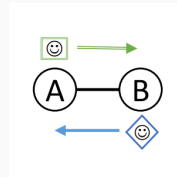
## Introduction

---

# Problem Definitions: Single-Objective Cooperation



Vertex Collision



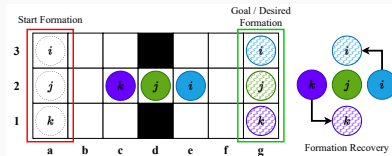
Edge Collision

## Multi-Agent Path Finding (MAPF)<sup>5</sup>

- a connected and undirected graph and a set of  $M$  agents
- Each agent has a unique start vertex and a unique goal vertex.
- For each time step, each agent can either **move** to one of its adjacent vertices or **wait** at its current vertex.
- Collisions: vertex collisions and edge collisions
- The goal is to find a set of **collision-free** path, one for each agent, while minimizing the **flowtime** (i.e., sum of all path length)

<sup>5</sup>Stern et al., "Multi-agent pathfinding: Definitions, variants, and benchmarks". In SoCS, 2019.

# Problem Definitions: Bi-Objective Cooperation



## Moving Agents in Formation (MAiF)<sup>6</sup>

- maintain close adherence to a designated formation (specified by goal locations) while moving towards the goals
- The formation deviation  $\mathcal{F}_t$  quantifies the least effort required to transform from the current formation to the desired formation:

$$\mathcal{F}_t := \min_{\Delta} \sum_{i=1}^M \| \mathbf{u}^i - (\mathbf{v}^i + \Delta) \|_1,$$

where  $\Delta$  is the element-wise median of  $\{ \mathbf{u}^i - \mathbf{v}^i \}_{i \in [M]}$

- The goal is to minimize (i) both the total (average) **formation deviation** (summing up  $\mathcal{F}_t$  over  $t$ ) and the **makespan** (i.e., maximum of all path length) or (ii) a linear combination of them

<sup>6</sup>J. Li et al., "Moving agents in formation in congested environments". In AAMAS, 2020

### MAPF (Single-Objective Cooperation)

- Conflict-Based Search<sup>7</sup>(CBS): A two-level search algorithm. The high level constructs a **constraint search tree** by adding constraints to different nodes, while the low-level plans paths w.r.t these constraints via **A\* search**
- Priority-Based Search<sup>8</sup>(PBS): A two-level search algorithm. The high level constructs a **priority search tree** by adding **partial orderings** to different nodes, while the low-level plans paths w.r.t these partial orderings via **prioritized planning**

### MAiF (Bi-Objective Cooperation)

- SWARM-MAPF<sup>9</sup>(SWARM): A two-phase method combining swarm-based formation control with MAPF algorithms. Phase 1 selects a **leader**, plans its path, and then partitions it into segments; Phase 2 runs **conflict-based search** to plan paths for other agents to follow the leader in those congested segments
- Scalarized Prioritized Planning (SPP)<sup>10</sup>: Plan agents' paths one by one using a specific (or random) ordering and optimize over the scalarized objective

---

<sup>7</sup>Sharon et al., "Conflict-based search for optimal multi-agent pathfinding". In Artificial Intelligence, 2015.

<sup>8</sup>H. Ma et al., "Searching with consistent prioritization for multi-agent path finding". In AAAI, 2019.

<sup>9</sup>J. Li et al., "Moving agents in formation in congested environments". In AAMAS, 2020.

<sup>10</sup>Silver, "Cooperative pathfinding". In AAAI, 2015.

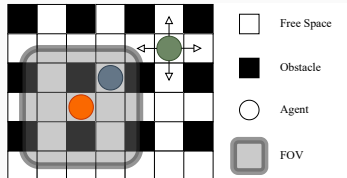


## Related Works: Learning Environments

There are some universal settings that learning-based methods follow which we adopt as our learning environments

Reward Function Design

Action	Reward
Move (up / down / left / right)	-0.075
Wait (on goal, away goal)	0, -0.075
Collision (obstacles or agents)	-0.5
Reaching Goal	3



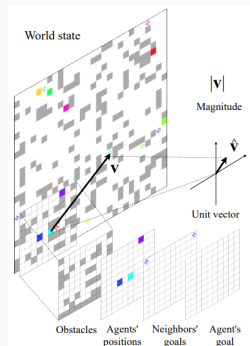
- 2-Dimensional 4-neighbor Grid Worlds: 2D grids environments where agents are only allowed to move along 4 cardinal directions
- Reward Design: we penalize each move with a small negative reward to incentivize agents to reach their goals as fast as possible
- Partially Observable Environments: each agent can only observe its surrounding  $\mathcal{L} \times \mathcal{L}$  area, namely field-of-view (FOV)
- Homogeneous Multi-Agent Systems: each agent shares the same policy but makes different decisions based on their observations

### PRIMAL<sup>11</sup>

- Mixture of A3C (RL) and Behavior Cloning (IL)
- Only observe **goal direction** as path planning guidance without considering the obstacles
- Use neighbouring agents' **goal directions** as cooperative guidance (not informative for complex cooperation)

### DHC<sup>12</sup> and DCC<sup>13</sup>

- Independent Q-Learning (IQL)
- Embed **communication** models
- Utilize **single-agent heuristic maps** as path planning guidance (equivalent to obstacle avoidance)
- No explicit guidance for cooperation



System design from PRIMAL

<sup>11</sup>Sartoretti et al., "Primal: Pathfinding via reinforcement and imitation multi-agent learning". In IEEE RA-L, 2019.

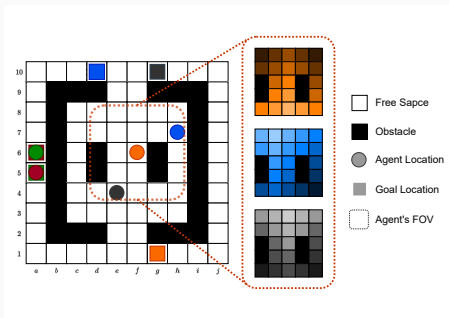
<sup>12</sup>Z. Ma, Luo, and H. Ma, "Distributed heuristic multi-agent path finding with communication". In ICRA, 2021.

<sup>13</sup>Z. Ma, Luo, and Pan, "Learning selective communication for multi-agent path finding". In IEEE RA-L, 2021.

## Single-Objective Cooperation

---

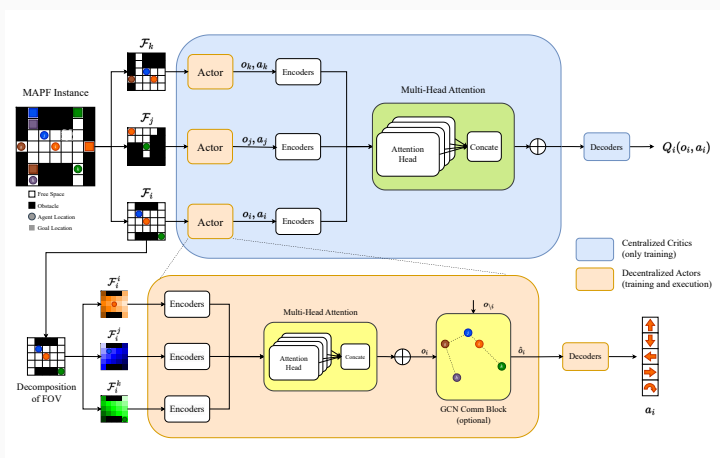
# SACHA: Multi-Agent Heuristic Maps



Example of multi-agent heuristic maps. A darker shade means a larger heuristic.

- Make use of heuristic maps from not only the center agent but also its neighboring agents
- Each cell holds a heuristic value proportional to its shortest path distance to the goal, and these heuristic maps inform each agent about its paths and neighboring agents' potential plans
- These heuristics can be pre-computed before execution with polynomial-time algorithms (e.g., Dijkstra) and will stay constant during execution

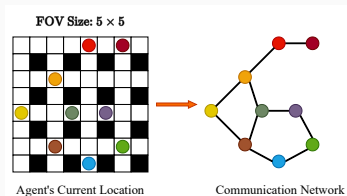
# SACHA: Soft Actor-Critic and Heuristic-Based Attention



## Centralized Training and Decentralized Execution (CTDE) and Multi-Agent Actor-Critic

- Policy network with heuristic-based attention for greater cooperative potentials
- Partially centralized attention critic network for better credit assignment

# SACHA: Graph-Based Communication



We also proposed another communication-based variant, named SACHA(C)

- We establish a dynamic communication network  $G_t$  that depends on agents' current positions: each vertex represents one agent, and each edge means two connecting agents lie within each other's FOV
- We run two-layer GCN<sup>14</sup> to encode, re-normalize, and pass messages along the communication network
- It has been shown by other works (e.g., Li et al.<sup>15</sup>) that this technique can encourage communication among agents

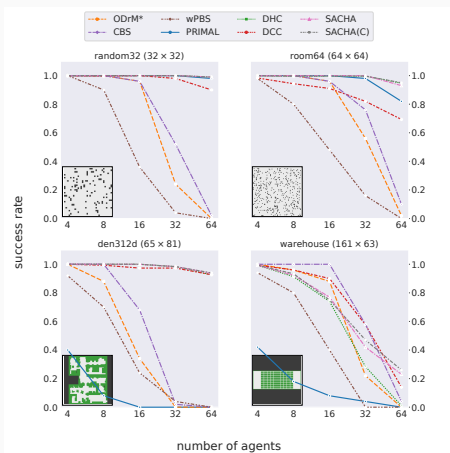
<sup>14</sup>Kipf and Welling, "Semi-supervised classification with graph convolutional networks". In arXiv, 2016.

<sup>15</sup>Li et al., "Graph neural networks for decentralized multi-robot path planning". In IROS, 2020.

# SACHA: Empirical Evaluation

- Besides learning-based methods (PRIMAL, DHC, and DCC), we also compare our methods with some planning algorithms with runtime limits: CBS (120s), PBS(120s), and OrDM\*(20s)

Comparison of success rate in different maps



Comparison of solution quality

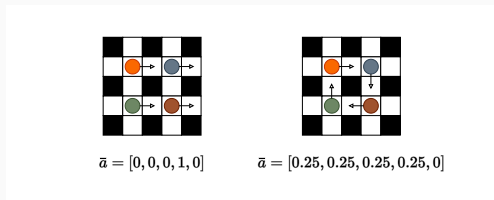
Map	Agents	Average Step Per Agent				
		PRIMAL	DHC	DCC	SACHA	SACHA(C)
random32	4	32.96	35.70	32.83	<b>29.93</b>	31.03
	16	45.12	48.67	43.56	41.71	<b>41.30</b>
	64	69.40	<b>66.05</b>	88.79	76.47	74.48
random64	4	67.82	71.04	70.80	<b>65.47</b>	67.10
	16	89.22	94.22	102.27	83.74	<b>82.17</b>
	64	105.12	120.68	154.72	99.02	<b>96.42</b>
den312d	4	196.54	86.56	82.99	<b>78.33</b>	81.43
	16	<del>256.00</del>	109.24	108.29	97.86	<b>96.74</b>
	64	<del>256.00</del>	153.17	145.21	<b>140.79</b>	142.97
warehouse	4	355.80	146.12	135.89	<b>131.43</b>	134.59
	16	492.04	281.37	208.72	<b>192.30</b>	198.72
	64	<del>512.00</del>	<del>512.00</del>	473.92	449.83	<b>437.29</b>

- PRIMAL has the worst performance since behavior cloning from experts (planning algorithms) hinders the trained model from generalizing toward unseen environments
- Our methods can outperform DHC and DCC in most cases



## Bi-Objective Cooperation

---



## Mean Field Reinforcement Learning<sup>16</sup>

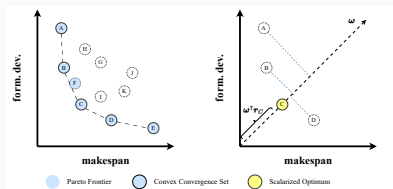
- Approximate interactions within agents by those between every single agent and the average effect from the overall population

$$Q^j(s, \{a^k\}_{k \in [M]}) = \frac{1}{M} \sum_{k=1}^M Q^j(s, a^j, a^k) \approx Q^j(s, a^j, \bar{a}),$$

where  $\bar{a} = \sum_{k=1}^M a^k$  is the mean action

- Avoid the exponential growth of agents' interactions (the curse of dimensionality) and thus enhance scalability
- Most importantly, **mean action can reflect on the formation change**

<sup>16</sup>Y. Yang et al., "Mean field multi-agent reinforcement learning". In ICML, 2018

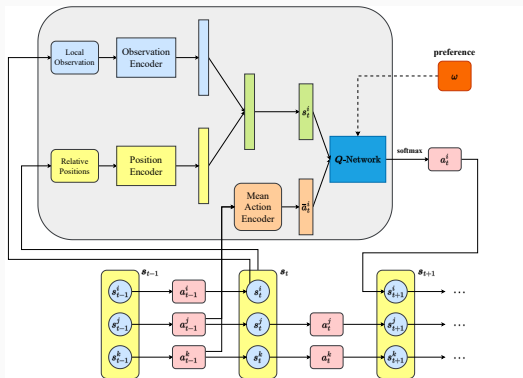


- Reward:  $r_t^j = (c_t^j, \mathcal{F}_t^j)^\top$ , where  $c_t^j$  is the moving cost and  $\mathcal{F}_t^j$  is agent  $j$ 's contribution to the formation deviation
- Goal: learn a universal model to minimize  $\sum_{t=0}^T \gamma^t \omega^\top (\sum_{j \in [M]} r_t^j)$  for **any given linear preference**  $\omega = (\lambda, 1 - \lambda)^\top \in \Omega$ ,
- We adopt the Envelope Q-Learning<sup>17</sup> to tackle the bi-objective optimization in the multi-agent settings (optimistic approach)

$$(\mathcal{T}Q)(s, a, \omega) := r(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} \arg_Q \left\{ \max_{\omega' \in \Omega} \max_{a'} \omega'^\top Q(s', a', \omega') \right\}$$

where  $\arg_Q$  takes the  $Q$ -value that corresponds to the maximal  $\omega^\top Q$

<sup>17</sup>R. Yang, Sun, and Narasimhan, "A generalized algorithm for multi-objective reinforcement learning and policy adaptation". In NeurIPS, 2019.



- Observation: encode agents' observation inside the FOV
- Position: encode agents' relative positions with others for formation control
- Mean Action:  $\bar{a}$  from mean field reinforcement learning
- Preference:  $\omega = (\lambda, 1 - \lambda)^\top$  from envelope Q-learning

# MFC-EQ: Empirical Evaluation

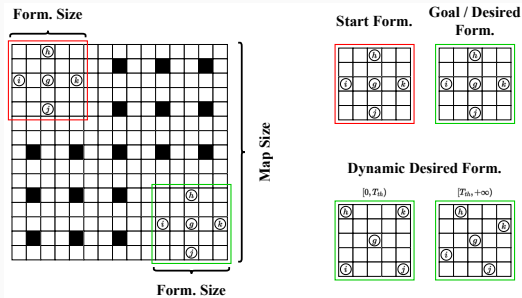


Illustration of experimental settings for MAiF

- The start position is located at the top-left corner, and the goal / desired formation lies at the bottom-right corner
- Agents travel from the top-left corner to the bottom-right corner while minimizing makespan and formation deviation
- We also conduct the dynamic formation experiment in which we alter the desired formation after a certain time threshold  $T_{th}$

## MFC-EQ: Empirical Evaluation

- Define the *MIX* metric:  $MIX(\lambda) = \lambda T + (1 - \lambda) \cdot \frac{\sum_{t=0}^T \mathcal{F}_t}{M}$
- We first test the model by setting  $\omega = (0.5, 0.5)^T$

Comparison of solution quality in different sizes of maps

Map Size	$M$	Success Rate			MIX(0.5)		
		SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ
32	10	1.00	1.00	1.00	39.05	32.70	<b>32.30</b>
×	20	1.00	0.99	0.99	40.73	37.78	<b>35.21</b>
32	30	0.79	0.96	0.90	46.90	39.87	<b>37.46</b>
48	10	1.00	0.99	0.99	67.18	53.14	<b>49.56</b>
×	20	0.95	0.99	0.96	76.26	66.27	<b>62.89</b>
48	30	0.74	0.94	0.88	90.48	<b>68.91</b>	72.30
64	10	1.00	0.99	0.99	105.65	79.77	<b>76.79</b>
×	20	1.00	0.97	0.93	114.04	94.64	<b>84.80</b>
64	30	0.22	0.98	0.90	111.62	<b>99.98</b>	103.47

- Test for adaptability towards various linear preferences

$\omega(\lambda)$	Makespan	Form. Dev.	MIX(0.1)	MIX(0.3)	MIX(0.5)	MIX(0.7)	MIX(0.9)
0.1	106.33	14.67	23.84	42.17	60.50	78.83	97.16
0.3	101.14	15.37	23.95	41.10	58.26	75.41	92.56
0.5	98.64	16.84	25.02	41.38	57.74	74.10	90.46
0.7	96.74	19.16	26.92	42.43	57.95	73.47	88.98
0.9	96.42	21.75	29.22	44.15	59.09	74.02	88.95

- Dynamic Formations: require agents to change formation midway

$M$	Success Rate			MIX(0.5)		
	SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ
10	1.00	0.98	0.96	88.05	115.70	<b>80.47</b>
15	1.00	1.00	1.00	90.50	137.34	<b>85.95</b>
20	0.97	1.00	1.00	95.41	135.94	<b>88.75</b>
25	0.72	1.00	1.00	100.09	133.46	<b>92.25</b>
30	0.90	0.98	0.93	98.41	123.20	<b>96.12</b>
35	0.48	0.94	0.87	107.58	123.76	<b>98.89</b>
40	0.25	0.81	0.74	110.66	109.59	<b>101.31</b>

## **Conclusion and Future Work**

---



## Conclusions

- **SACHA** and **SACHA(C)**<sup>18</sup> address the issues of learning single-objective cooperation in partial observable multi-agent path finding (**MAPF**), with a focus on **generalizability** among different environments
- **MFC-EQ**<sup>19</sup> tackles the challenges of learning bi-objective cooperation in decentralized moving agents in formation (**MAiF**), with a focus on **scalability** towards large-scale instances and **adaptability** towards any linear preferences

## Future Work

- To improve generalizability, a more generalized scheme is to consider meta-learning, in which one can pre-train a model as initialization and fine-tune to different environments
- To solve multi-object path planning tasks, we can apply multi-objective RL algorithms (e.g., Pareto  $Q$ -learning) to directly approach the Pareto frontier
- Design algorithms/frameworks for more sophisticated types of cooperation in multi-agent path planning

---

<sup>18</sup>published paper in IEEE Robotics and Automation Letters 2023

<sup>19</sup>submitted paper under review

Thank you for listening!