

InforMARL: Scalable Multi-Agent Reinforcement Learning through Intelligent Information Aggregation

Siddharth Nayak¹, Kenneth Choi¹, Wenqi Ding¹, Sydney Dolan¹,
Karthik Gopalakrishnan², Hamsa Balakrishnan¹

¹Massachusetts Institute of Technology

²Stanford University

`{sidnayak, kenchoi, wenqi2, sydneyd, hamsa}@mit.edu`

`kgopalakrishnan@stanford.edu`

Background and Motivation

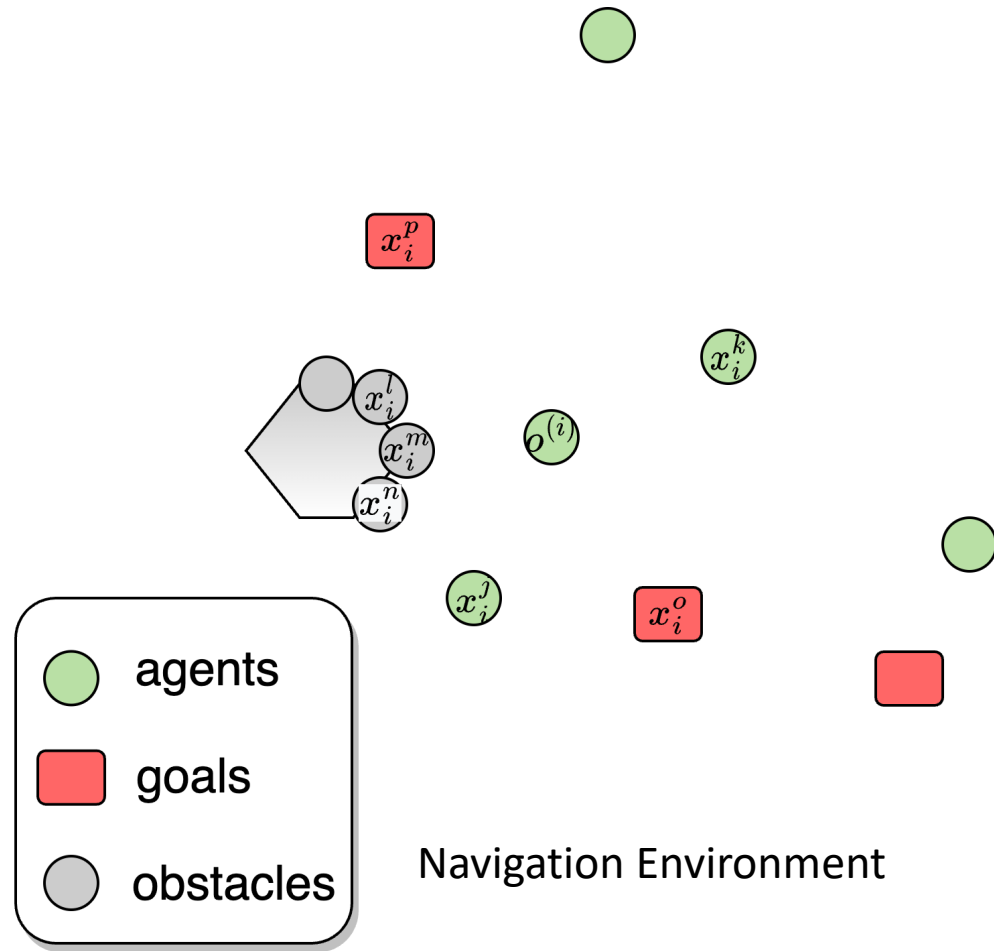


Credit: U.S. Naval Institute

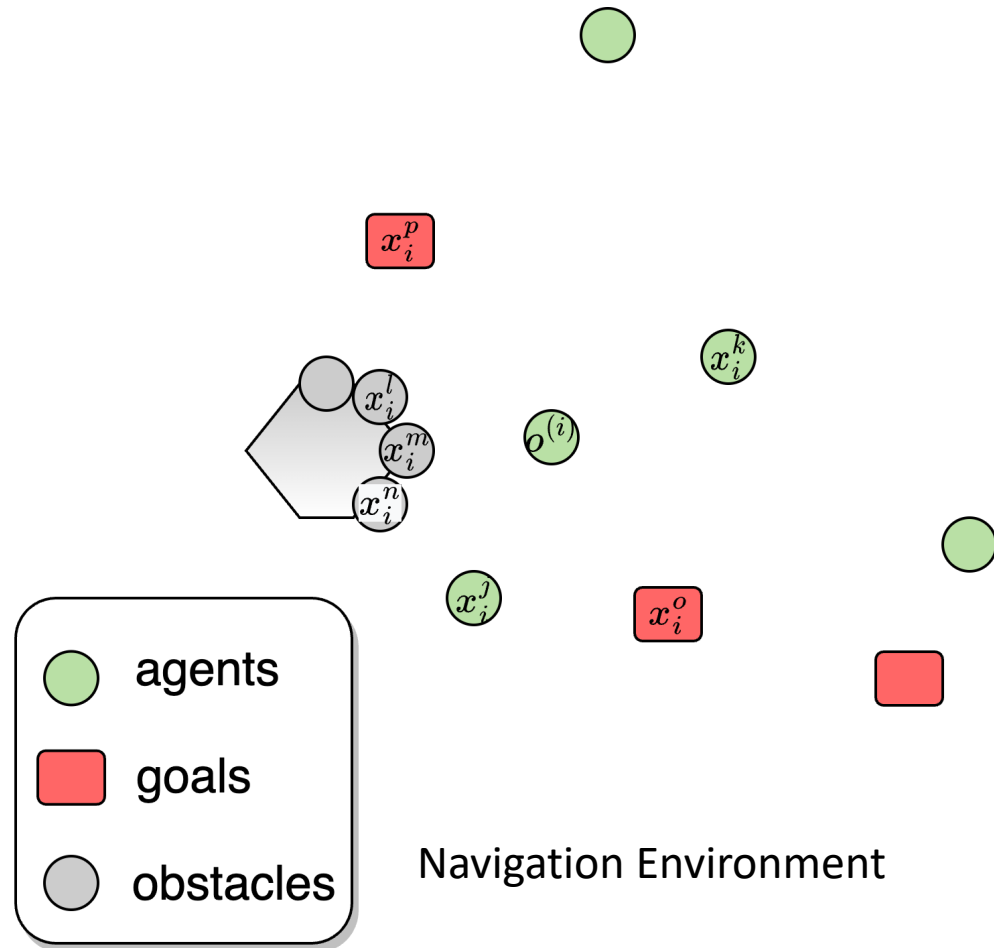


Credit: The Robot Report

Background and Motivation



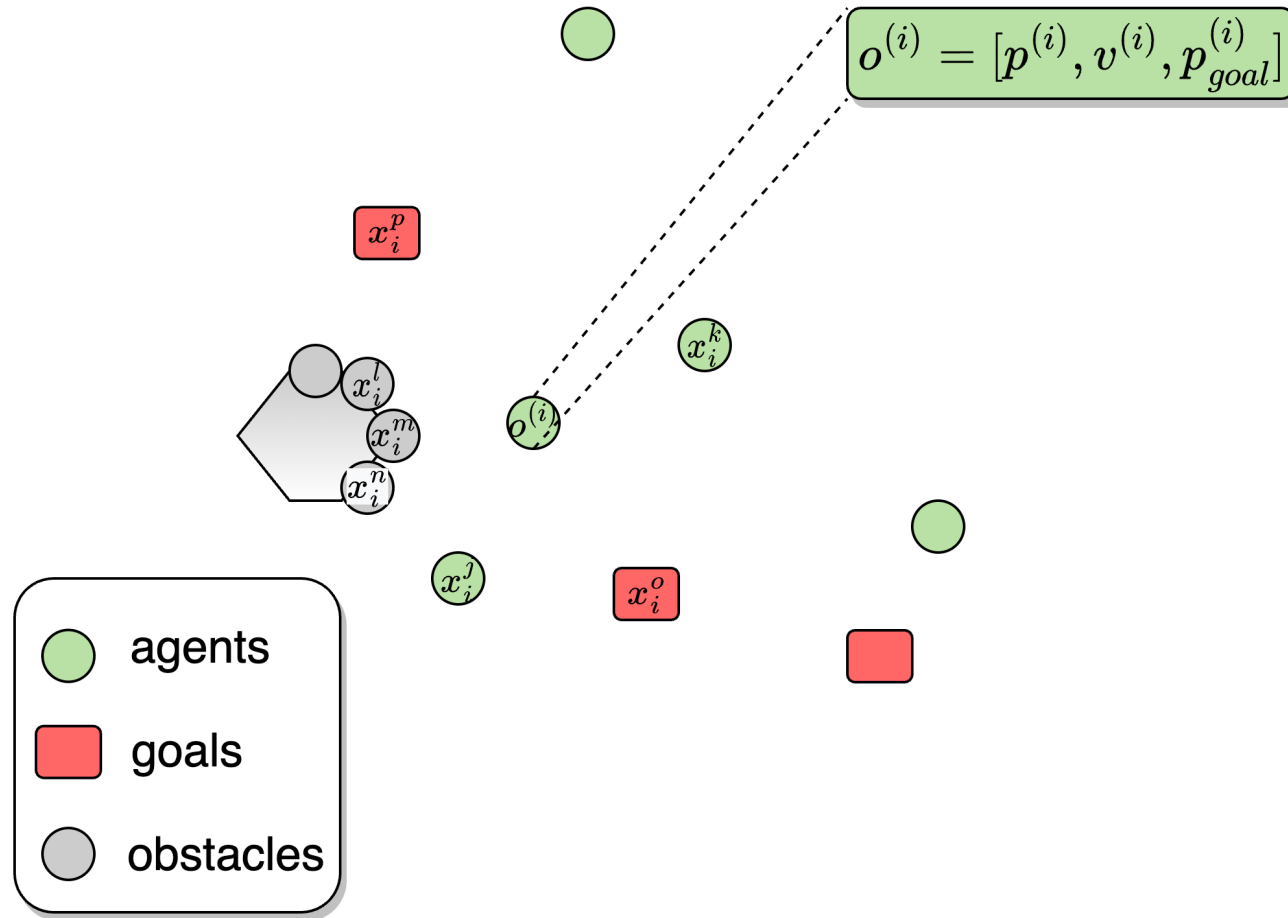
Background and Motivation



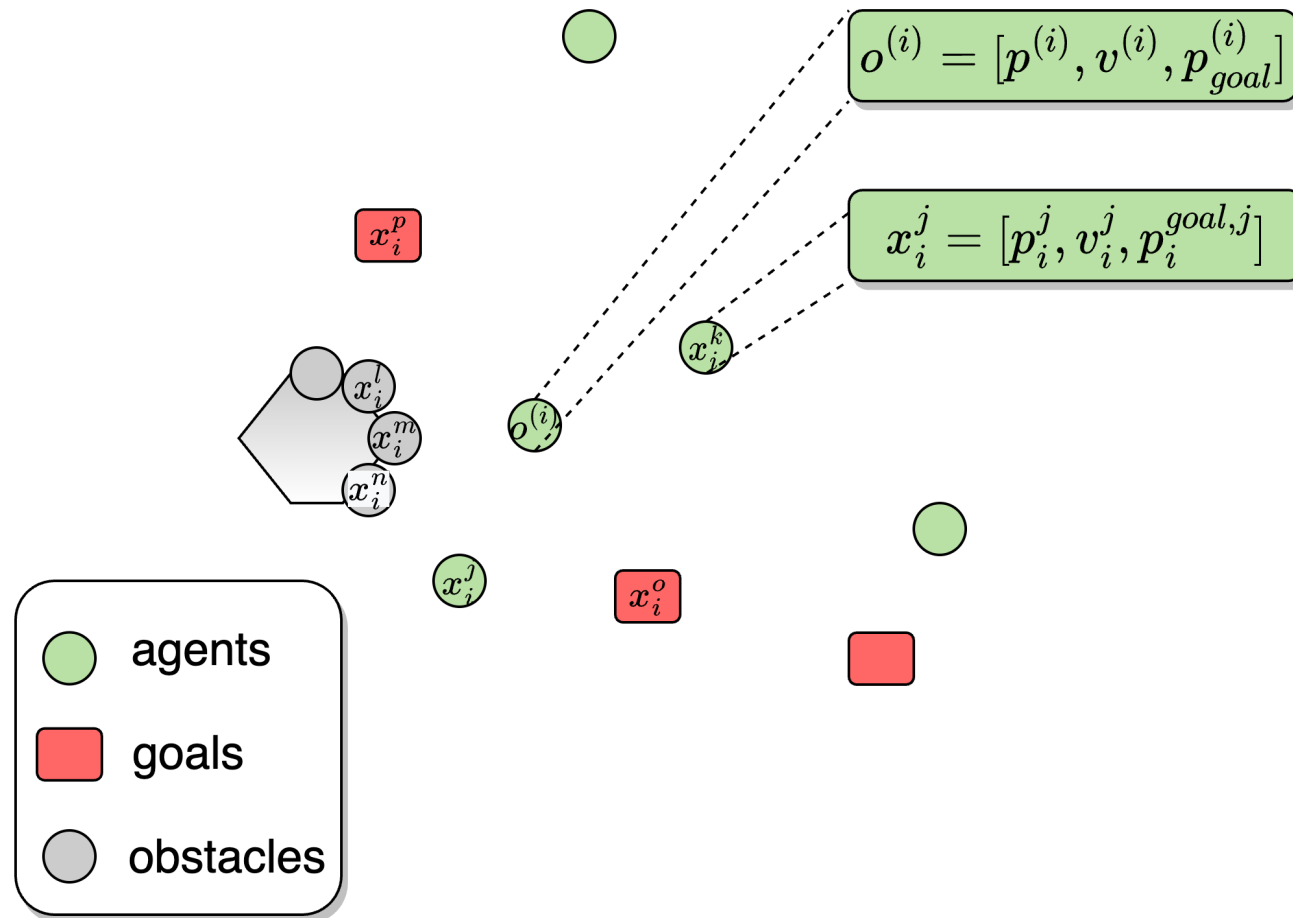
Key features expected from MARL algorithms:

- Decentralized execution
- Scalability
- Efficiency in sample complexity to train

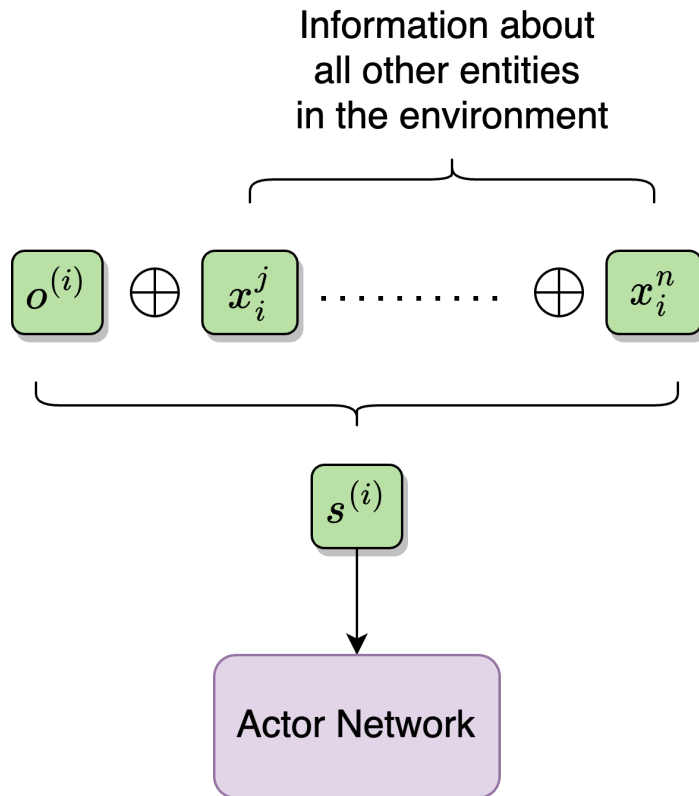
Background and Motivation



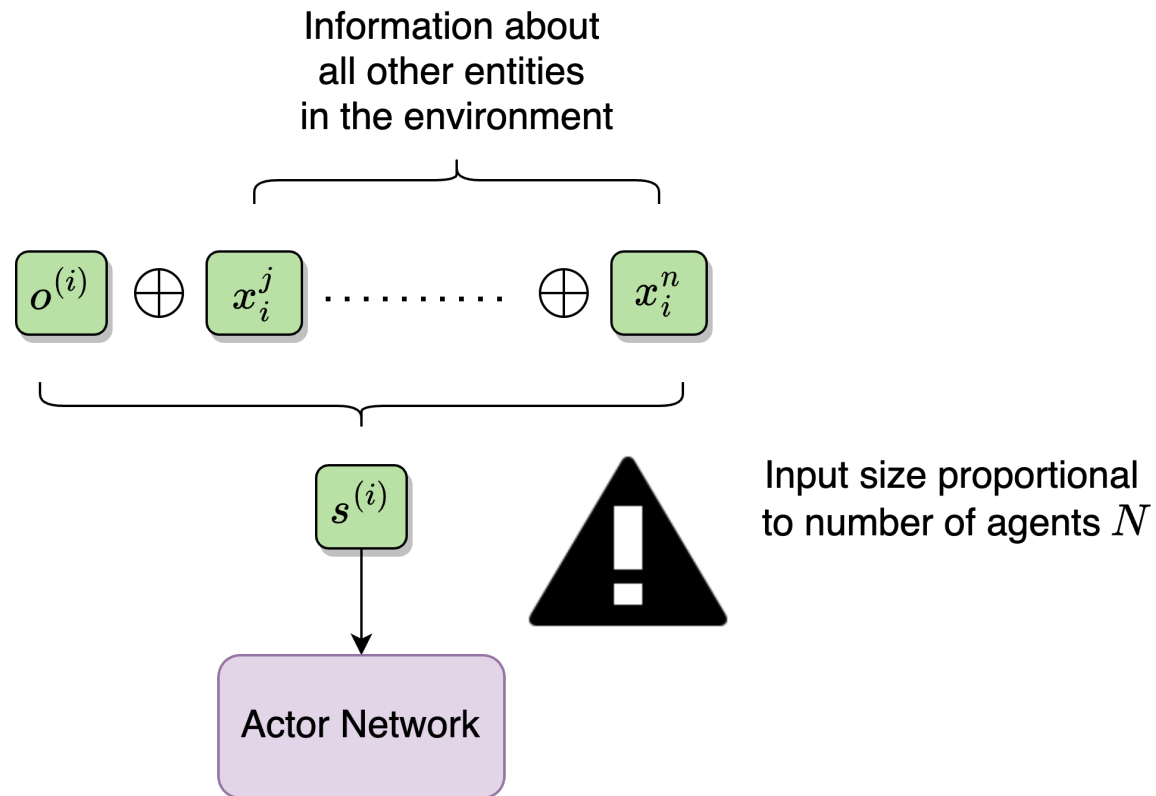
Background and Motivation



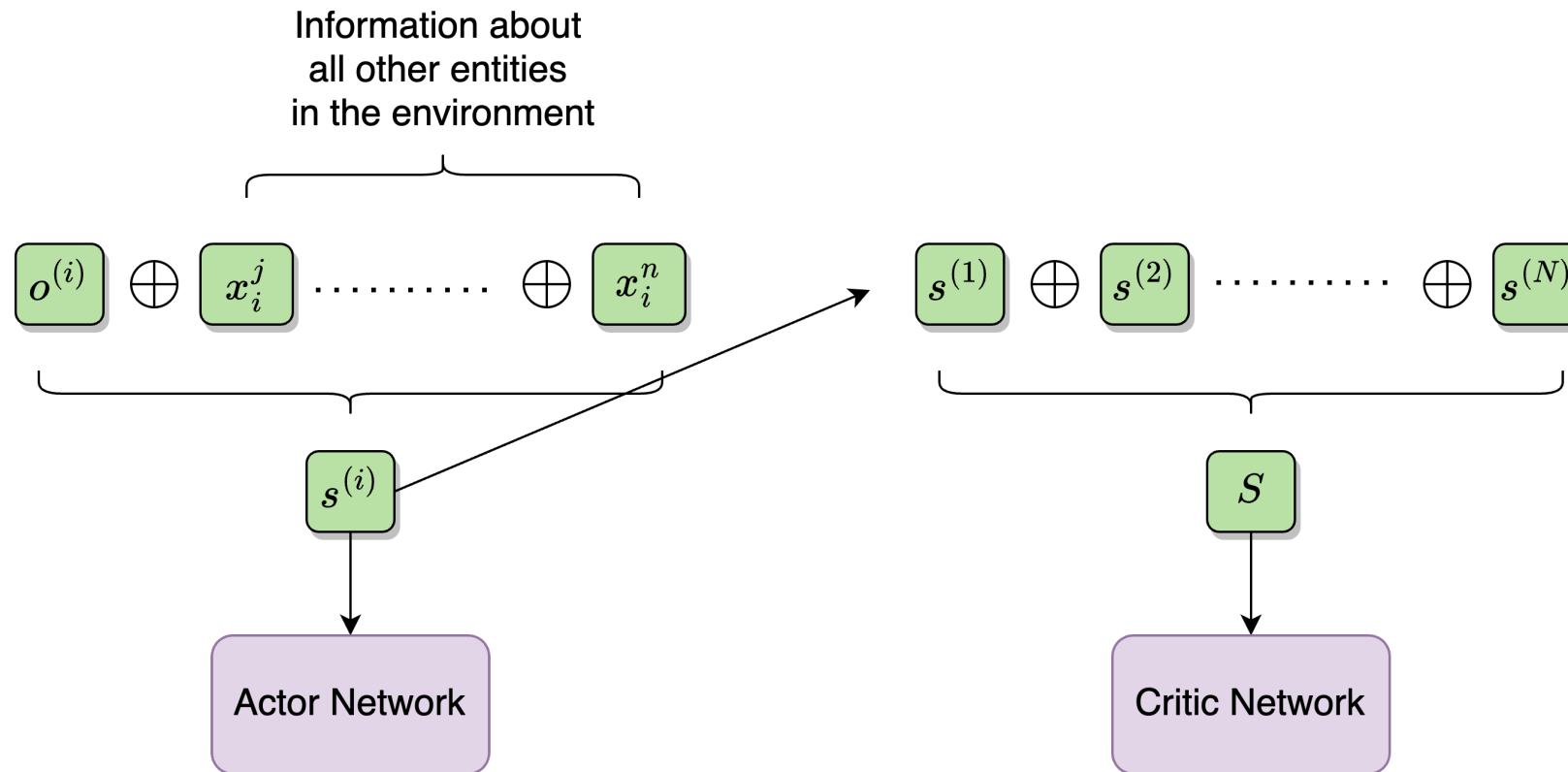
Background and Motivation



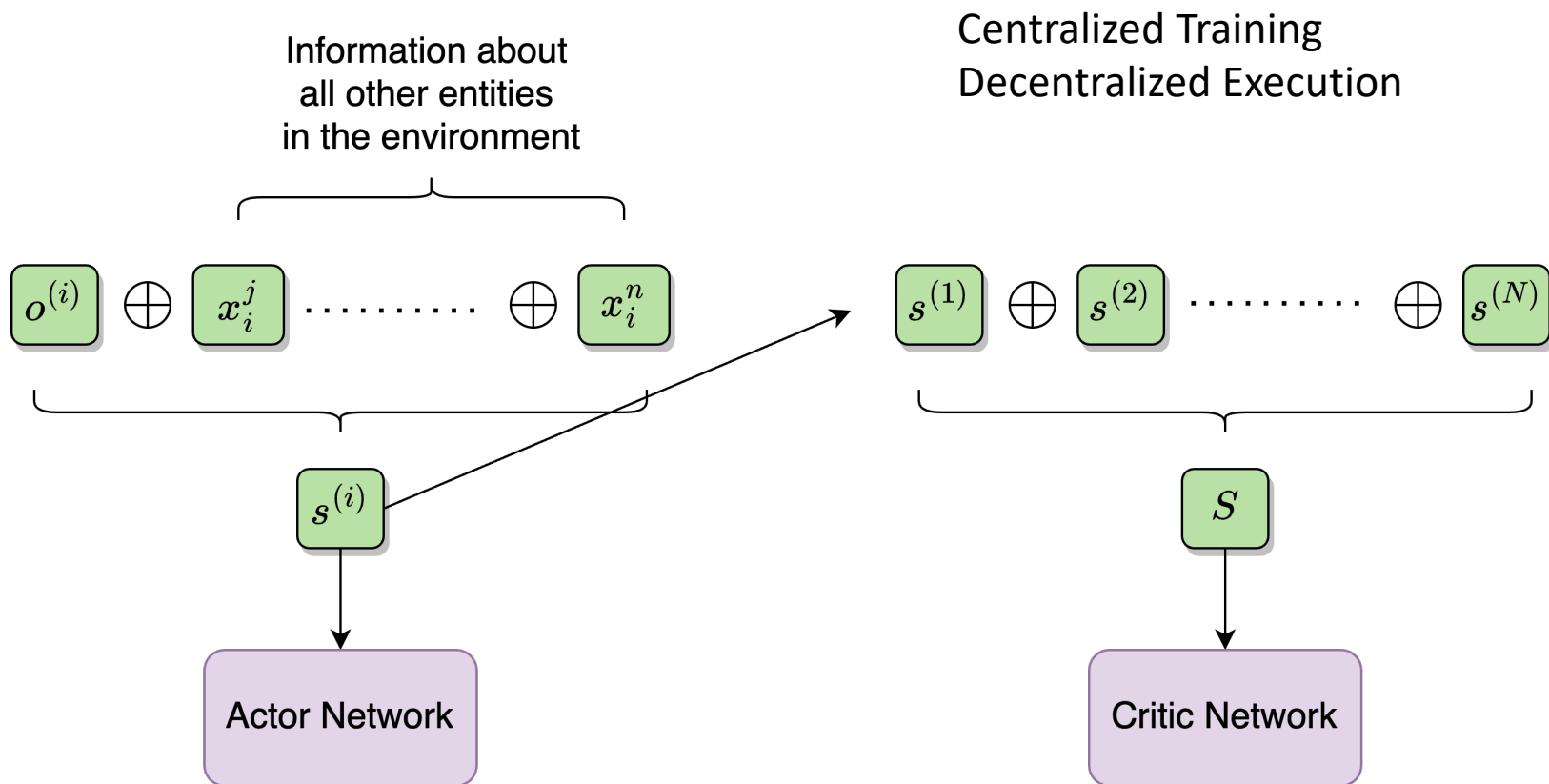
Background and Motivation



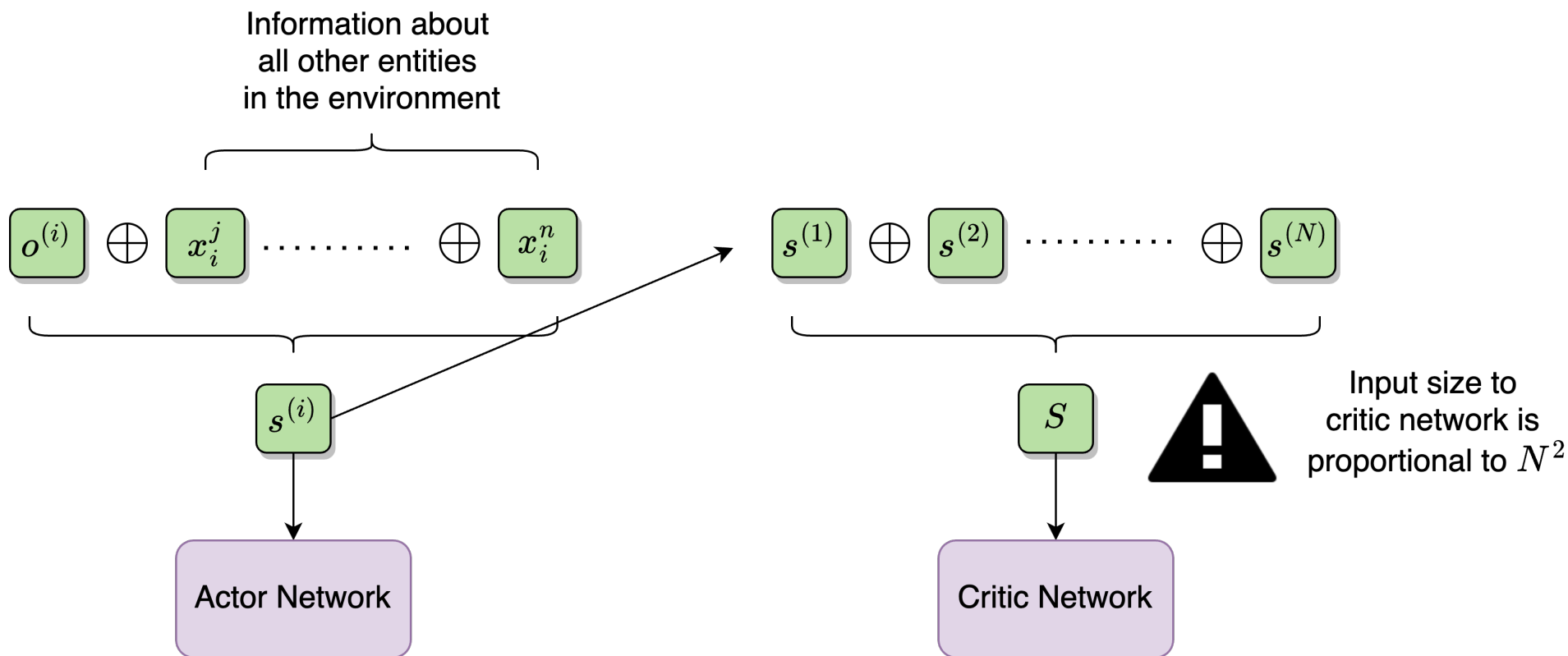
Background and Motivation



Background and Motivation



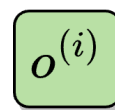
Background and Motivation



Motivating Experiment

Vary the amount of information included in observations for actor

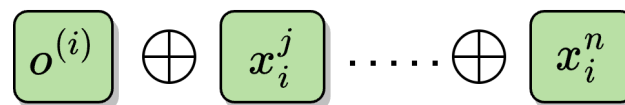
- Local Information Mode:



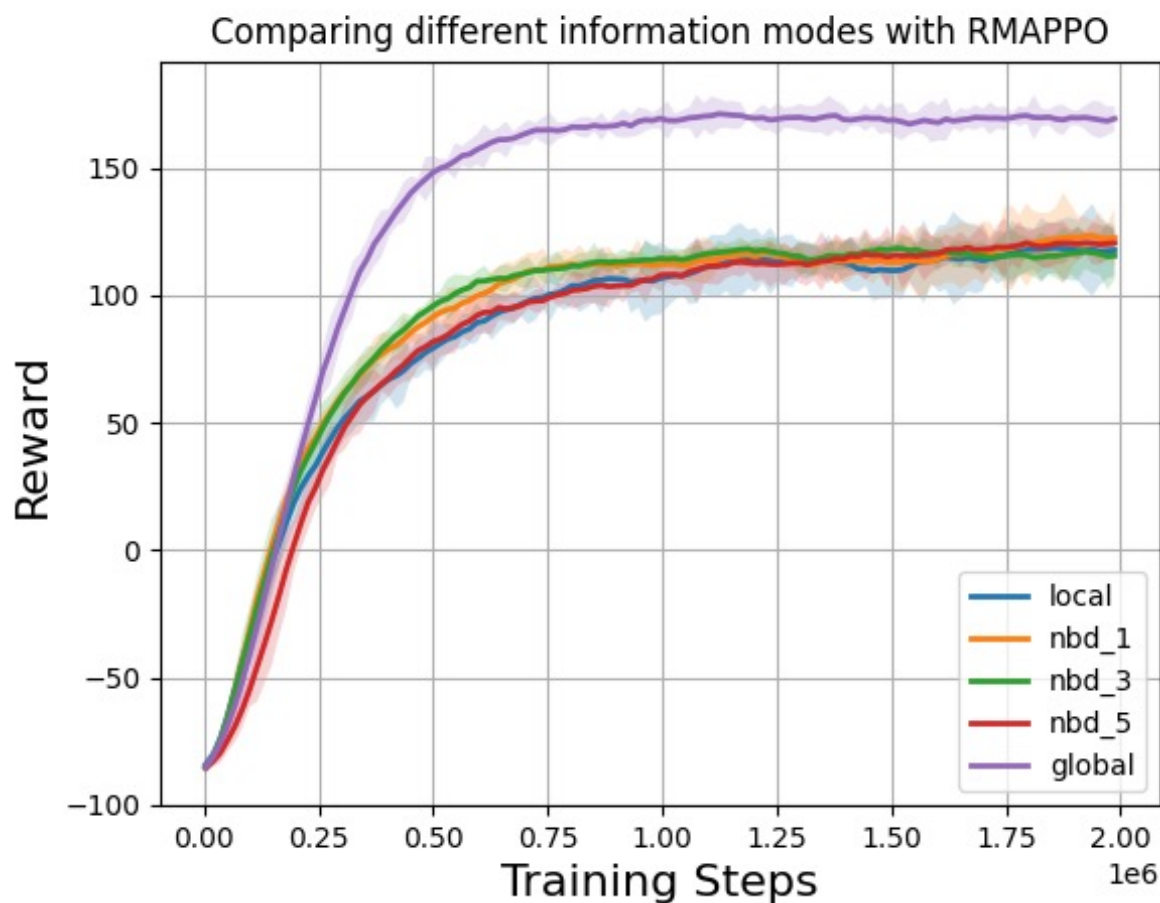
- Global Information Mode:



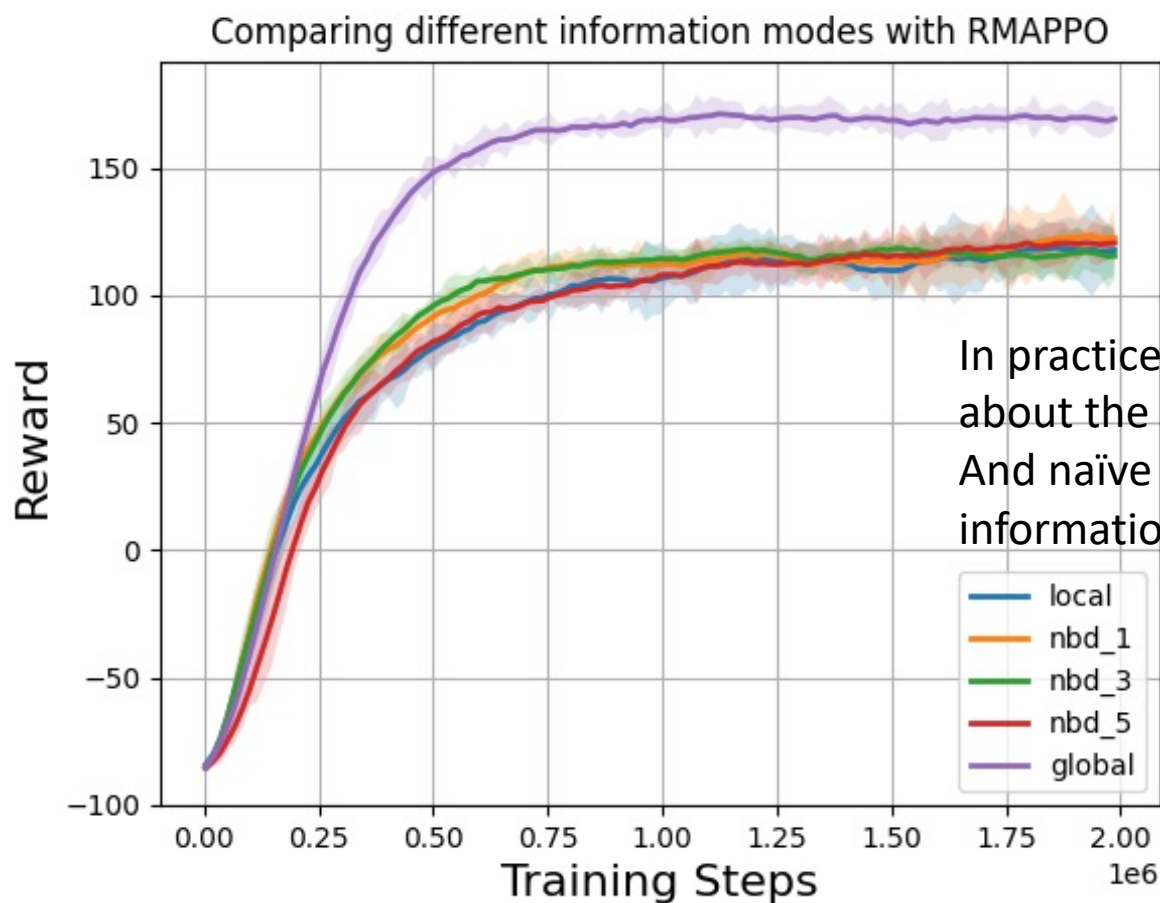
- Neighborhood Information Mode:



Motivating Experiment

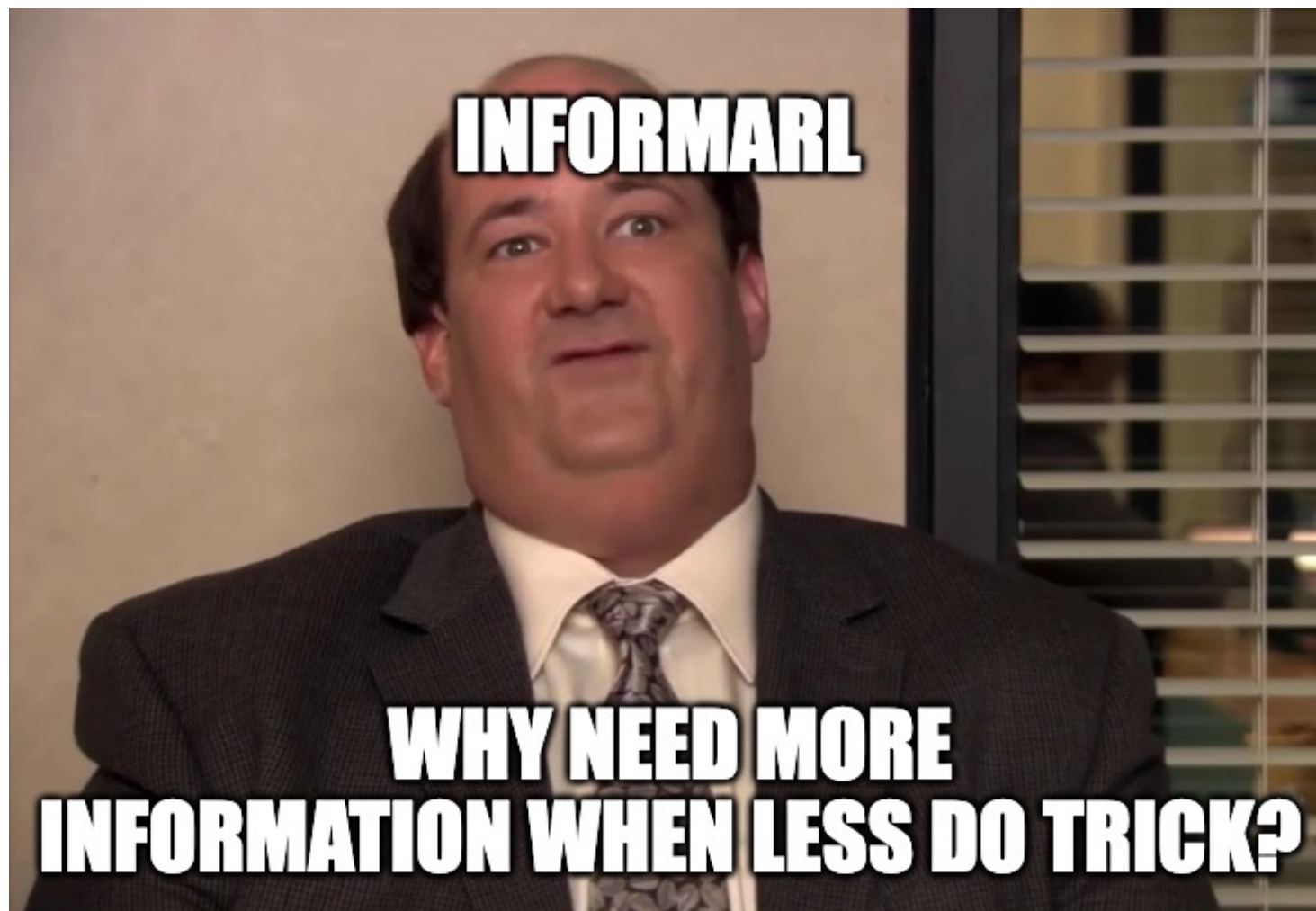


Motivating Experiment

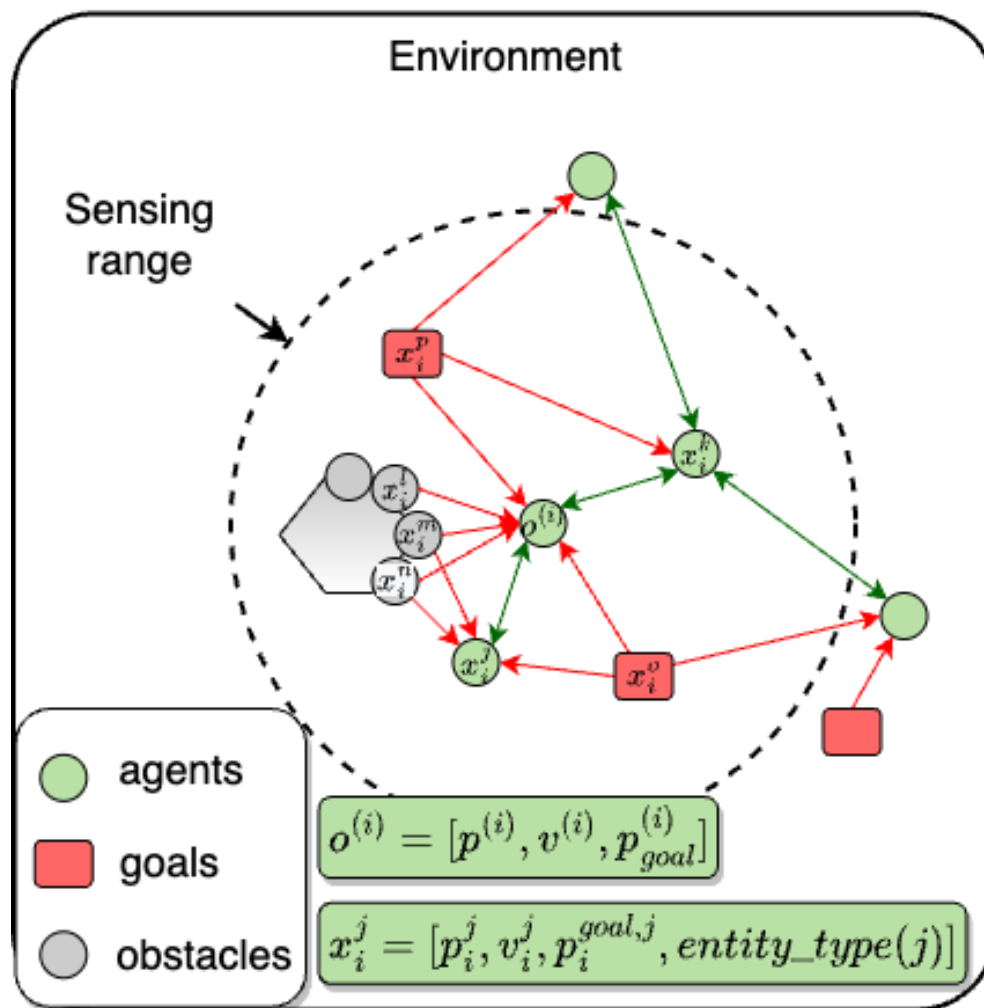


In practice, we just have local information about the neighborhood
And naïve concatenation of neighborhood information doesn't work

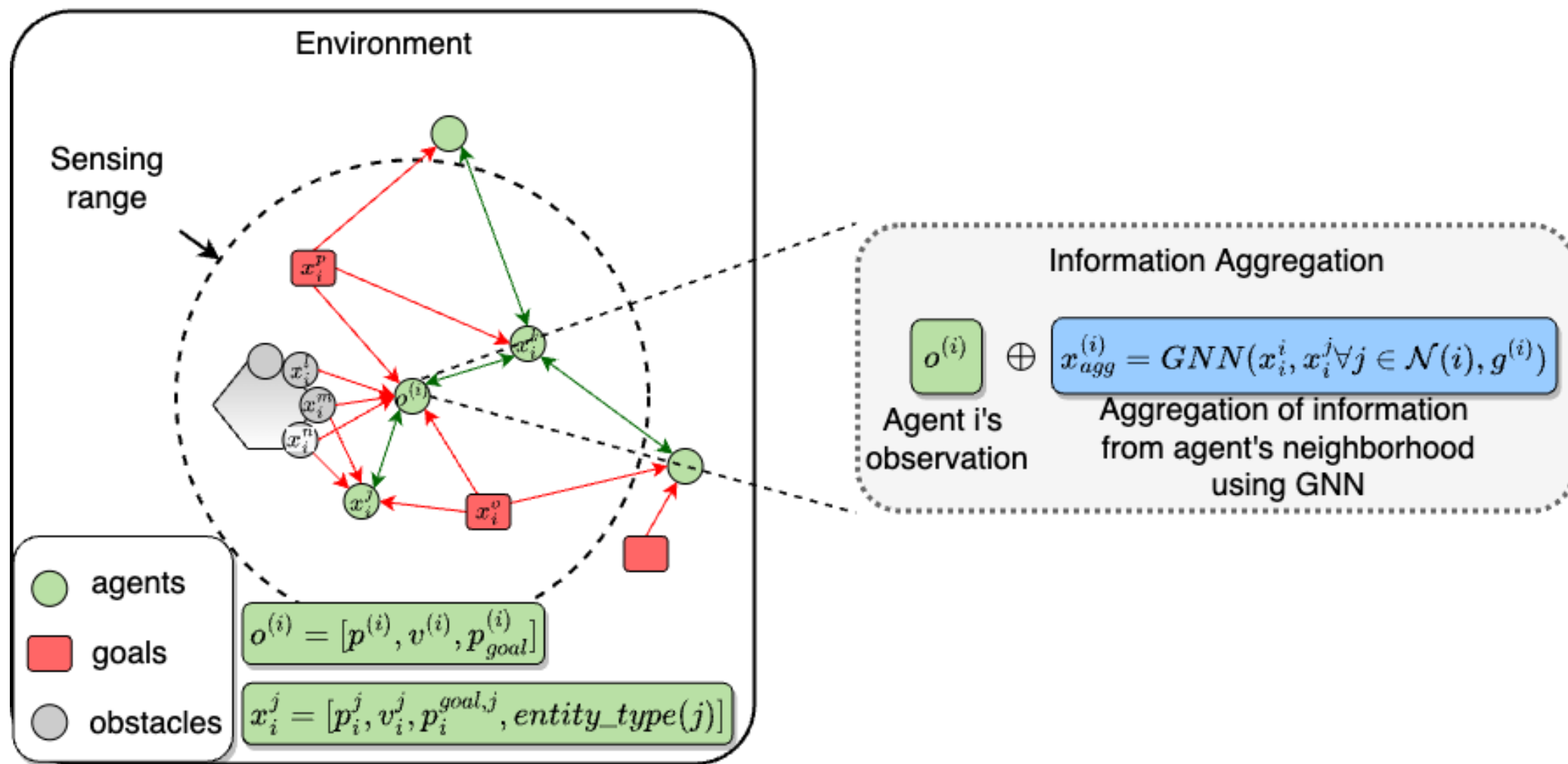
Motivating Experiment



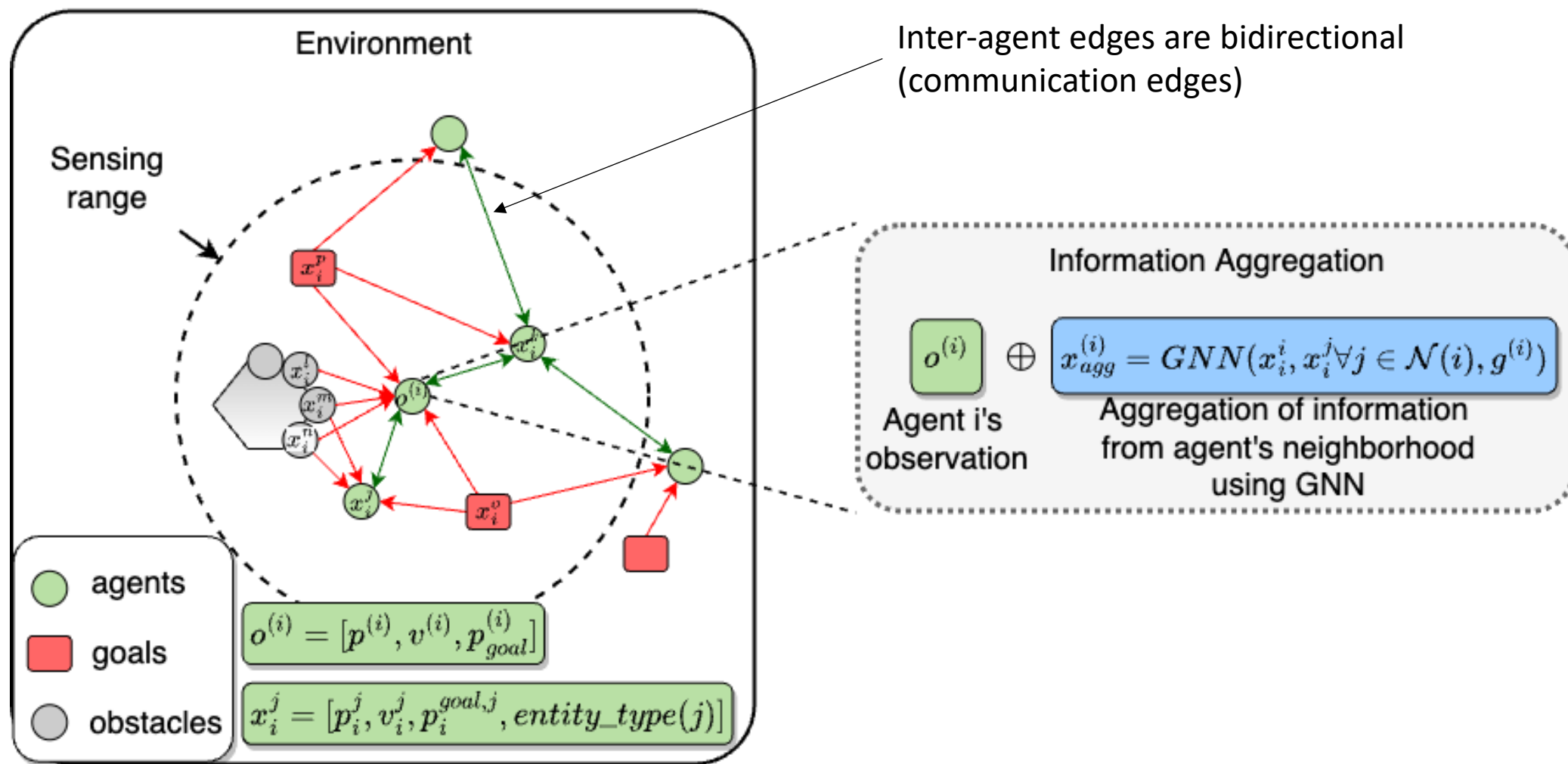
Method



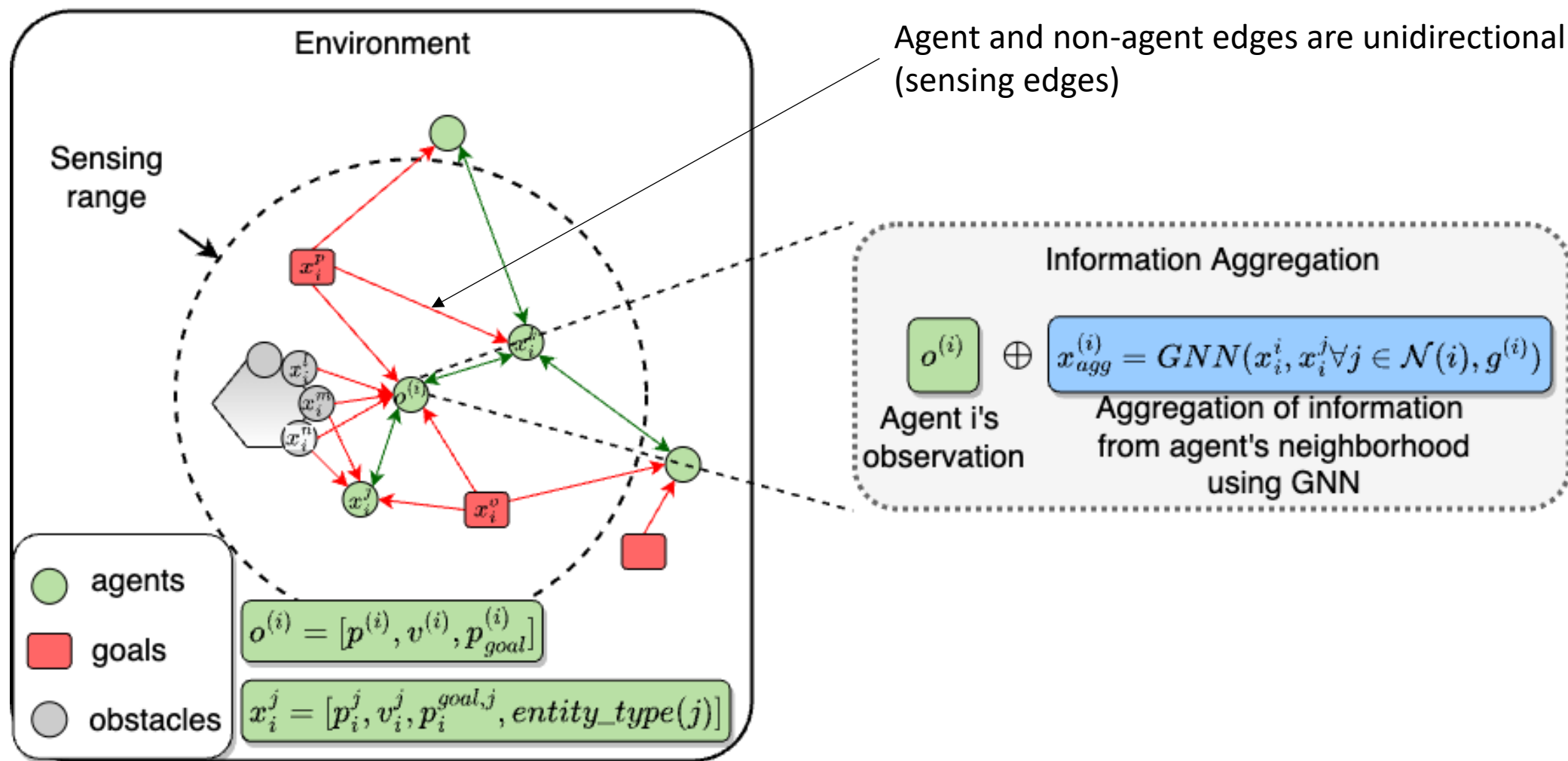
Method



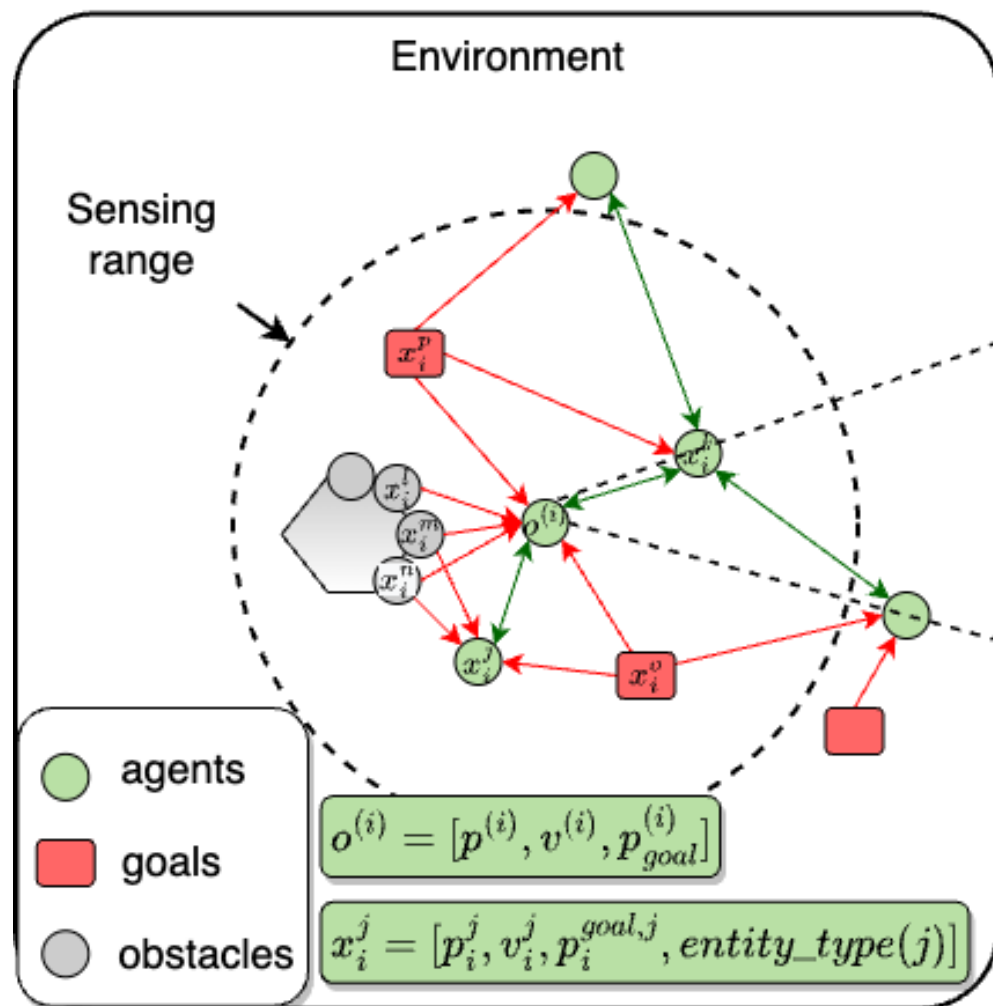
Method



Method

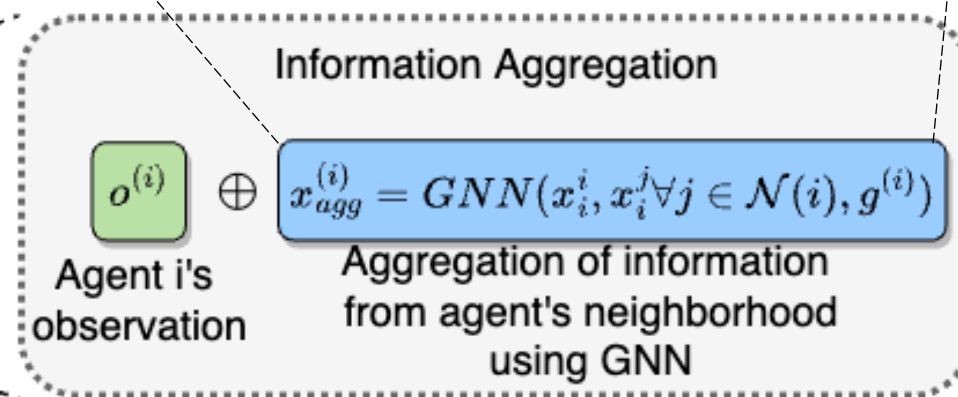


Method

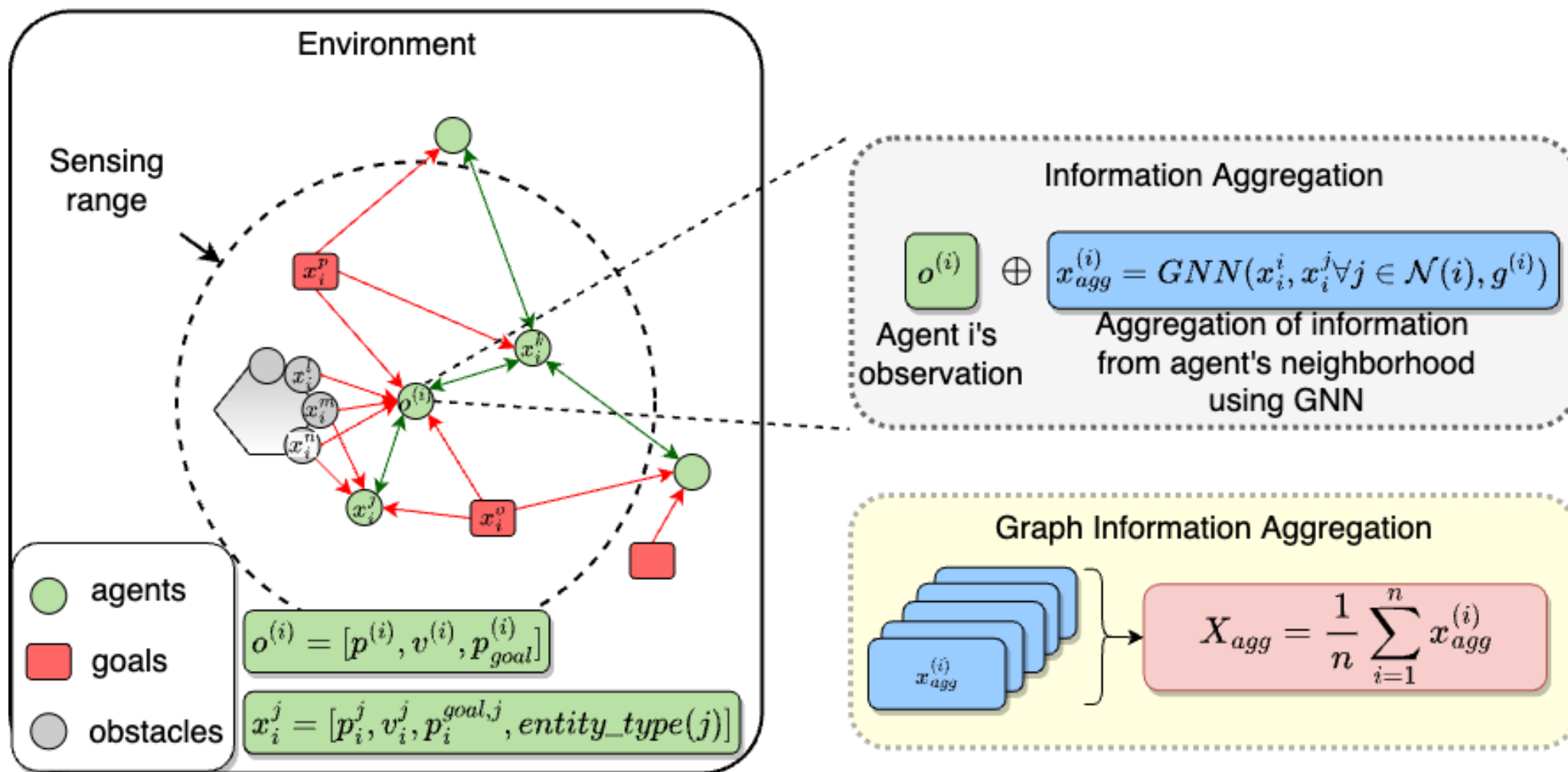


$$x'_i = W_1 \cdot x_i + \sum_{j \in \mathcal{N}(i)} \alpha_{i,j} W_2 \cdot x_j$$

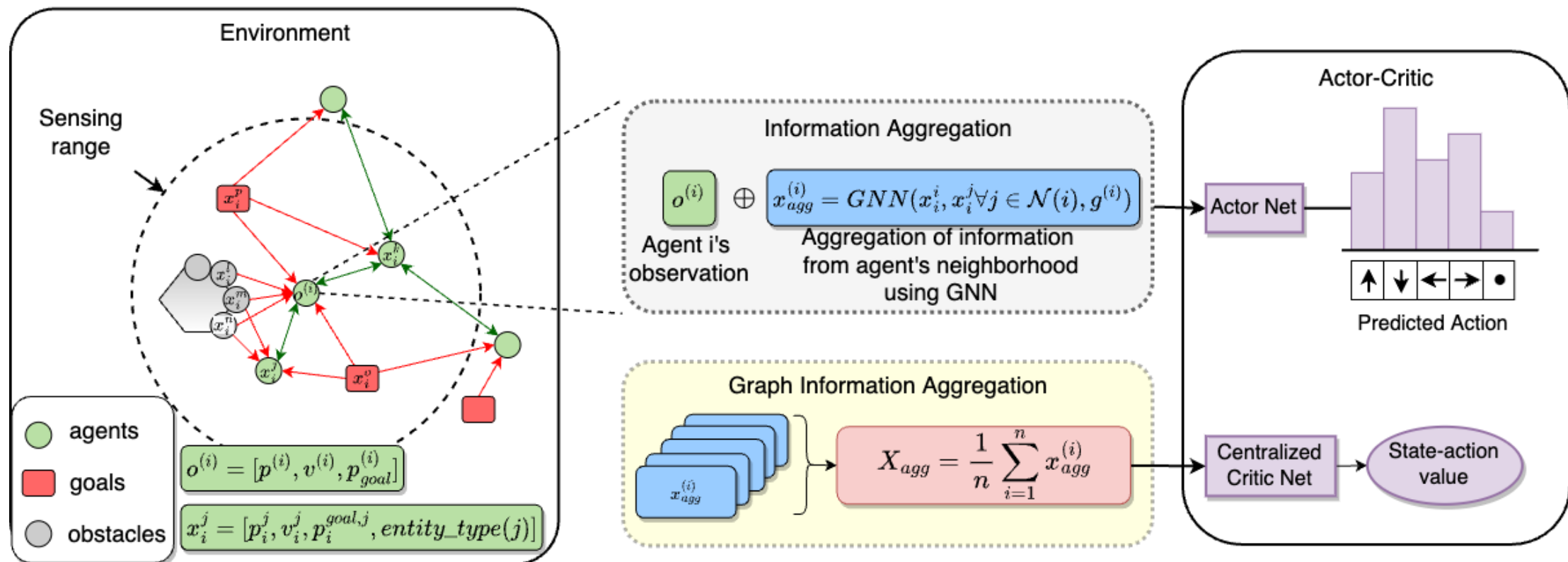
$$\alpha_{i,j} = \text{softmax} \left(\frac{(W_3 \cdot x_i)^T (W_4 \cdot x_j + W_5 \cdot e_{i,j})}{\sqrt{c}} \right)$$



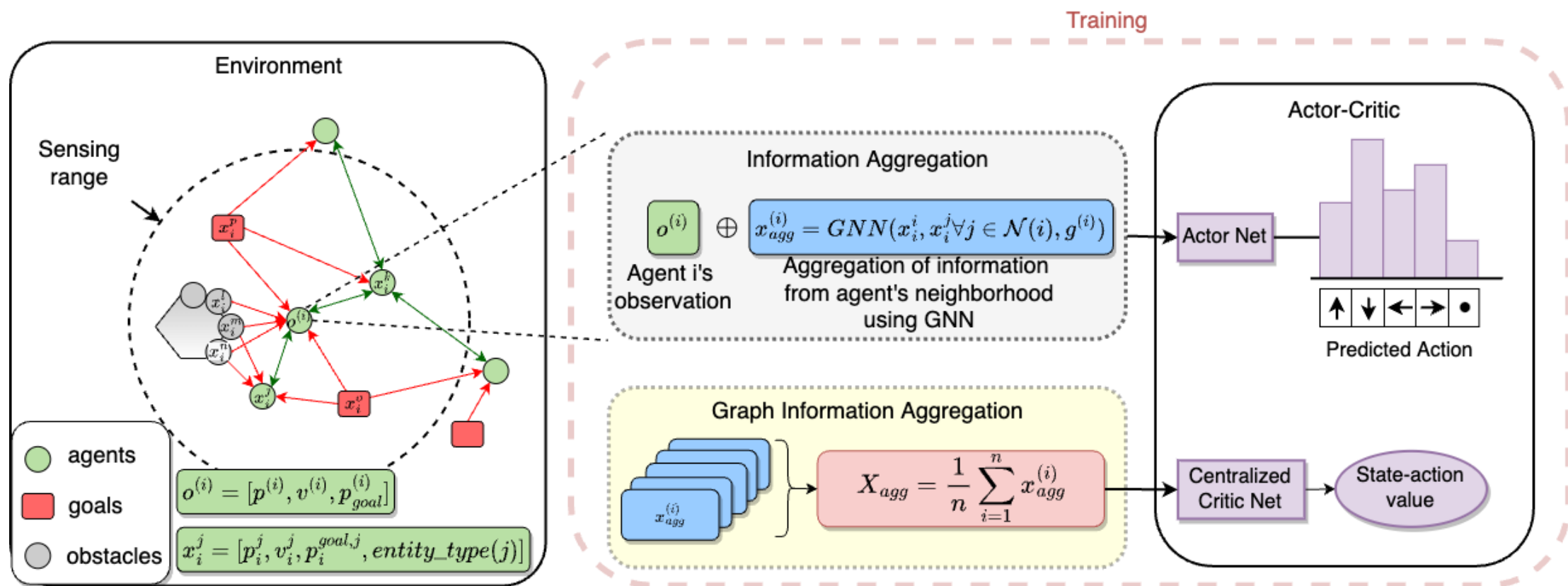
Method



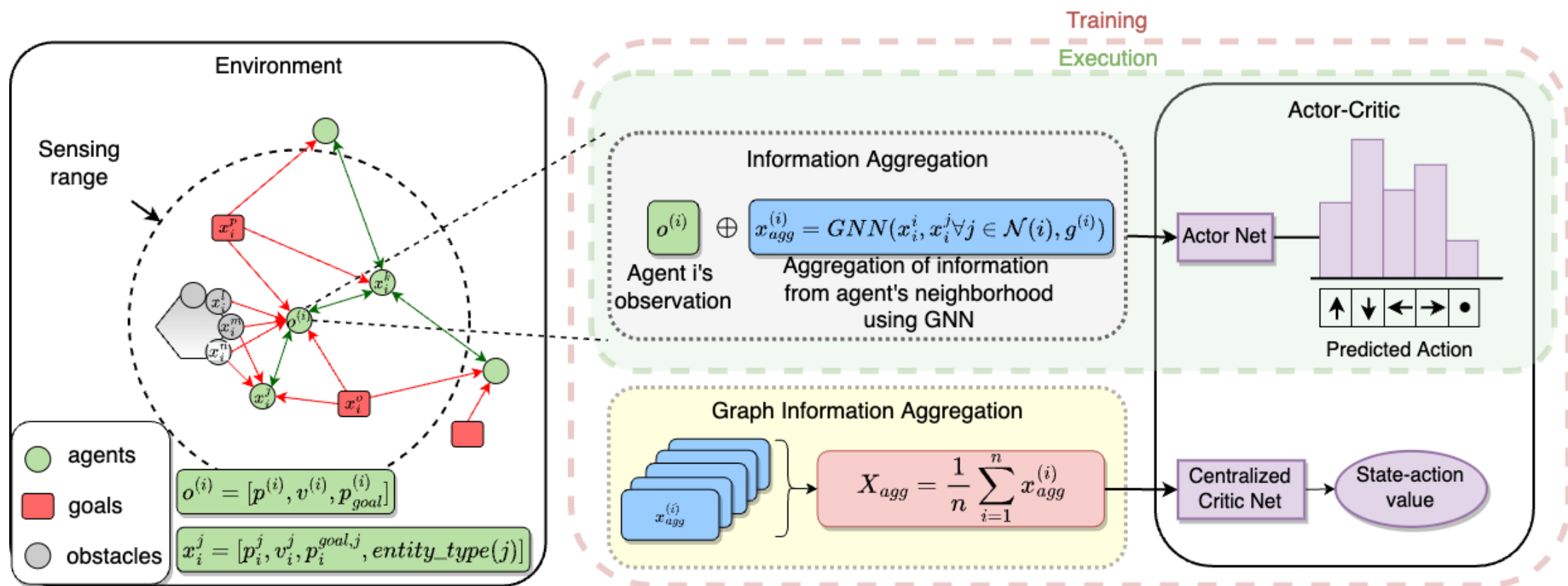
Method



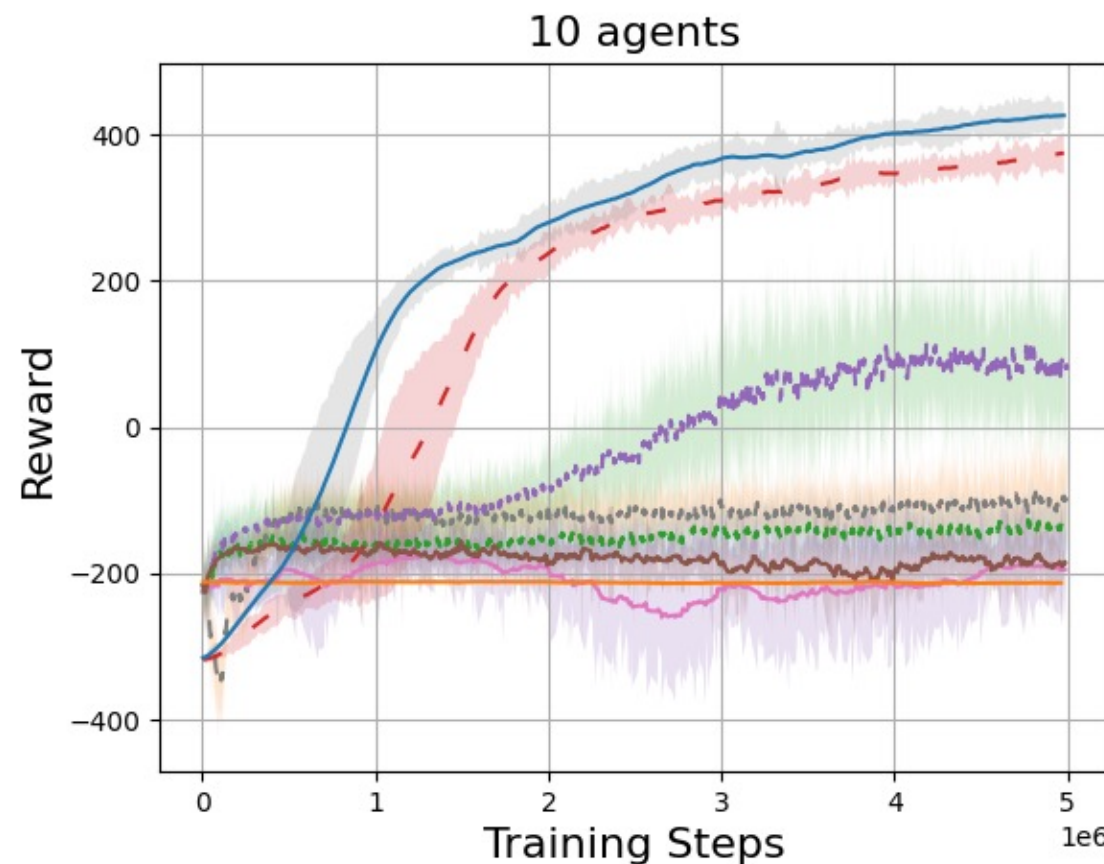
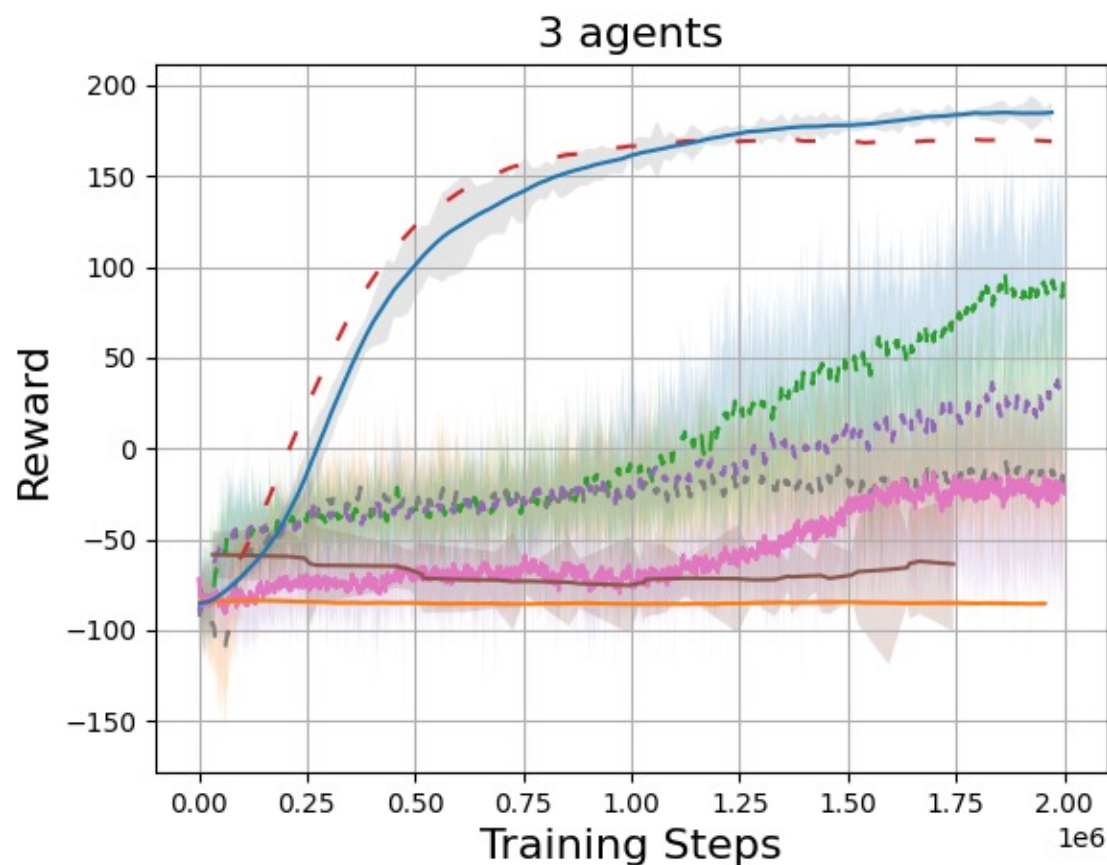
Method



Method



Experiments: Sample complexity



Global Information

Local Information

-- RMATD3

-- GPG (dynamic)

-- RQMIX

-- DGN + ATOC

-- RVDN

-- EMP

-- RMAPPO

-- InforMARL



DINaMo

Experiments: Scalability

Train \ Test	$m=3$			$m=10$			$m=15$		
	R/m	# col/ m	$S\%$	R/m	# col/ m	$S\%$	R/m	# col/ m	$S\%$
$n = 3$	68.31	0.48	100	58.59	1.60	100	53.19	2.24	100
$n = 7$	58.30	0.61	100	53.25	1.43	99	46.39	2.31	99
$n = 10$	57.27	0.64	99	52.10	1.68	100	48.15	2.20	99

Experiments: Scalability

Train \ Test		<i>m</i> =3			<i>m</i> =10			<i>m</i> =15		
		<i>R/m</i>	# col/ <i>m</i>	<i>S</i> %	<i>R/m</i>	# col/ <i>m</i>	<i>S</i> %	<i>R/m</i>	# col/ <i>m</i>	<i>S</i> %
<i>n</i> = 3		68.31	0.48	100	58.59	1.60	100	53.19	2.24	100
<i>n</i> = 7		58.30	0.61	100	53.25	1.43	99	46.39	2.31	99
<i>n</i> = 10		57.27	0.64	99	52.10	1.68	100	48.15	2.20	99

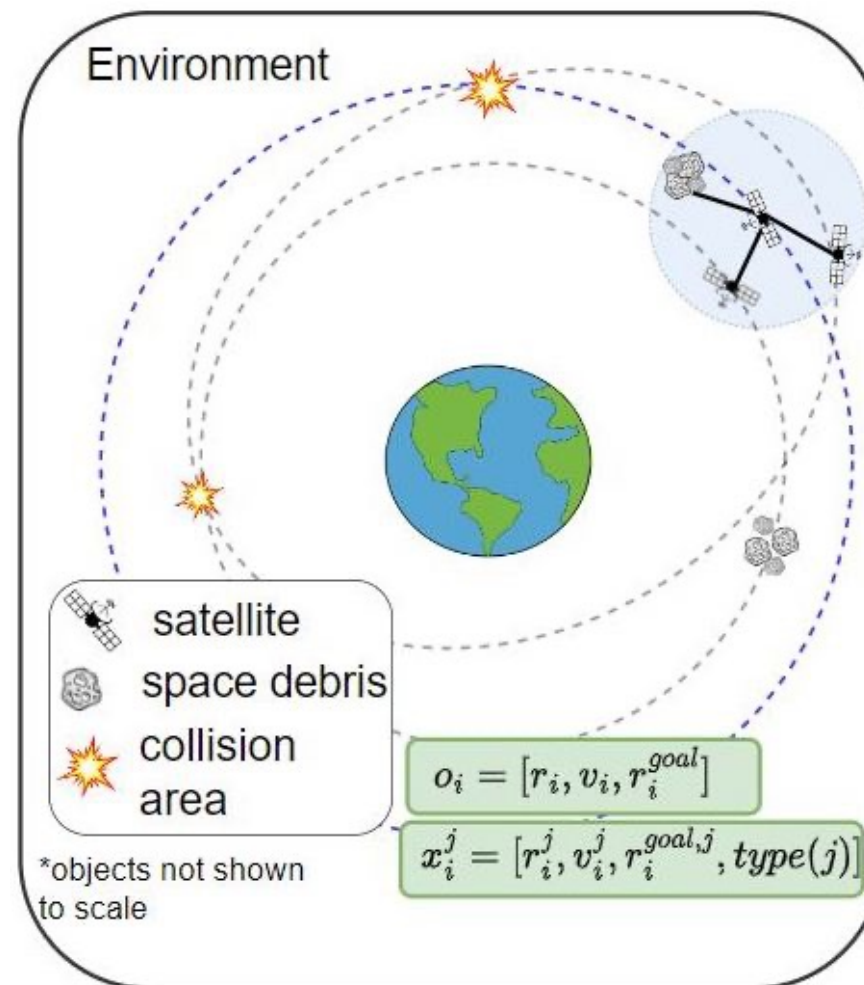
Reward per agent

Number of collisions per agent

Success Rate

Satellite Environment

- Concurrent work in leveraging transfer learning for satellite environment
- More complex non-linear dynamics for all entities in the environment



Conclusions and Future Work

- Introduced a graph-based algorithm for scaling standard MARL algorithms to arbitrary scenarios.
- Uses just neighborhood information instead of global information required by previous methods.



Project Website