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

**Siddharth Nayak**
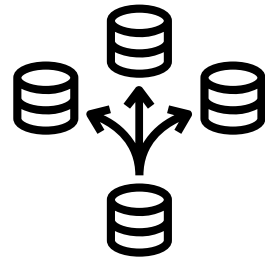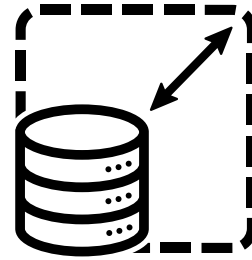
sidnayak@mit.edu

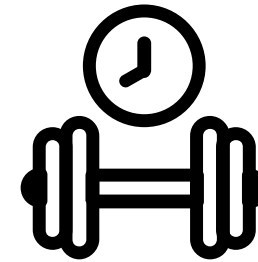# Background and Motivation

# Background and Motivation

Key Features Expected from MARL Algorithms
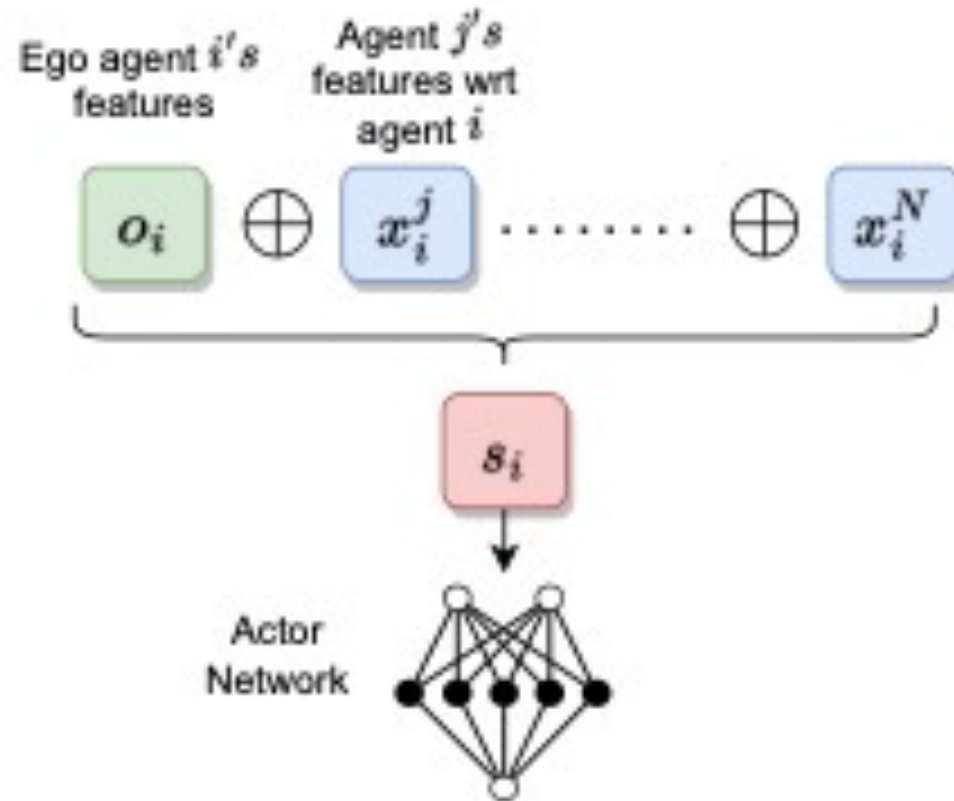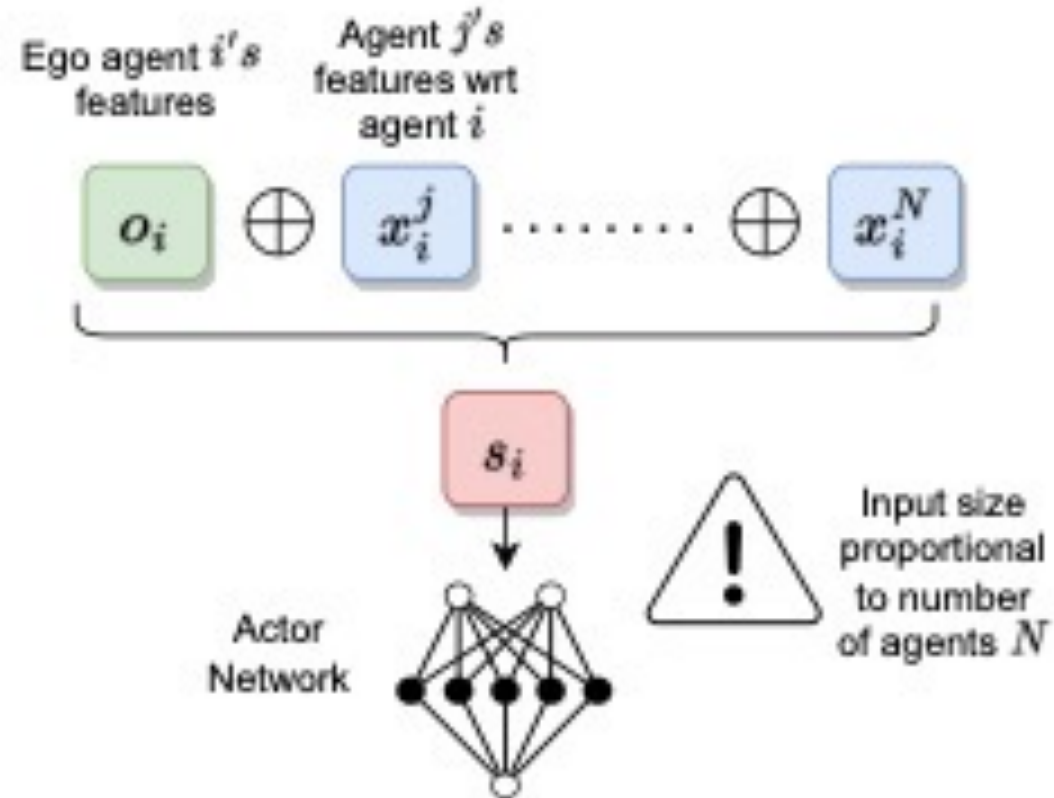


Decentralized Execution

Scalability

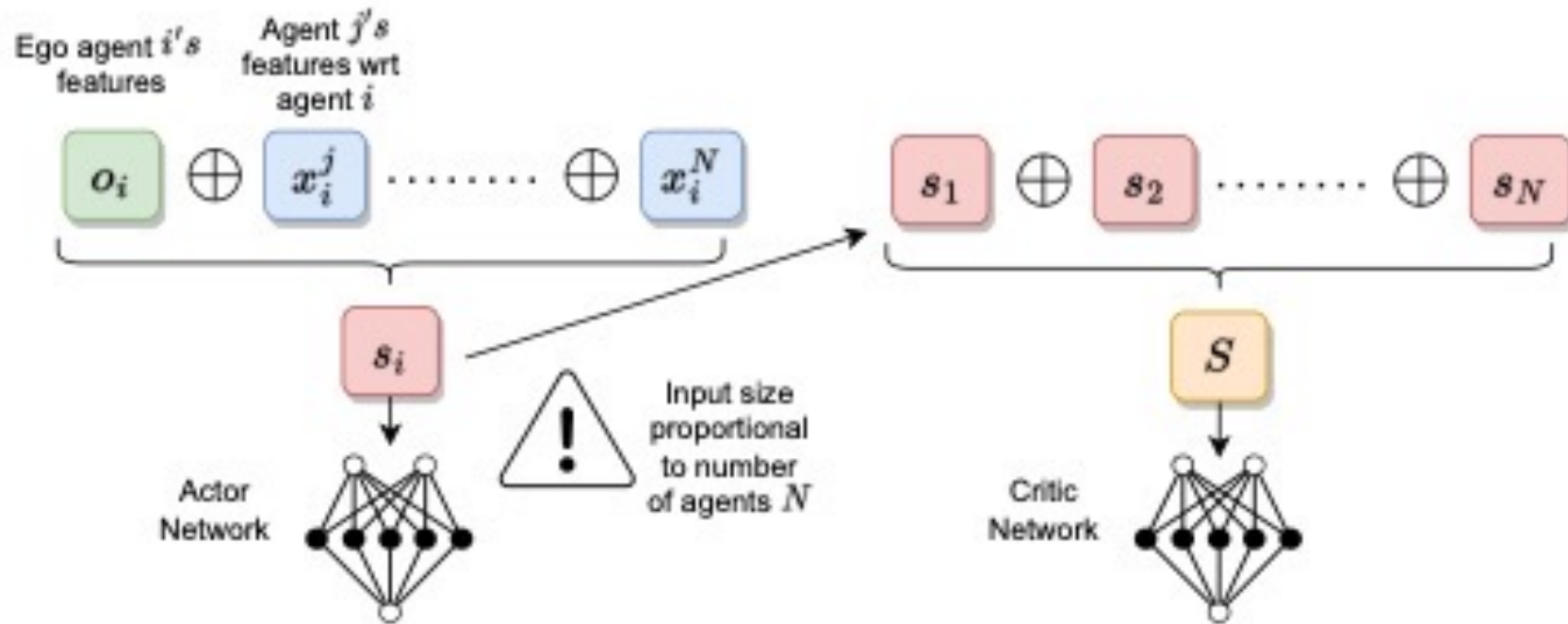Efficiency in training sample complexity
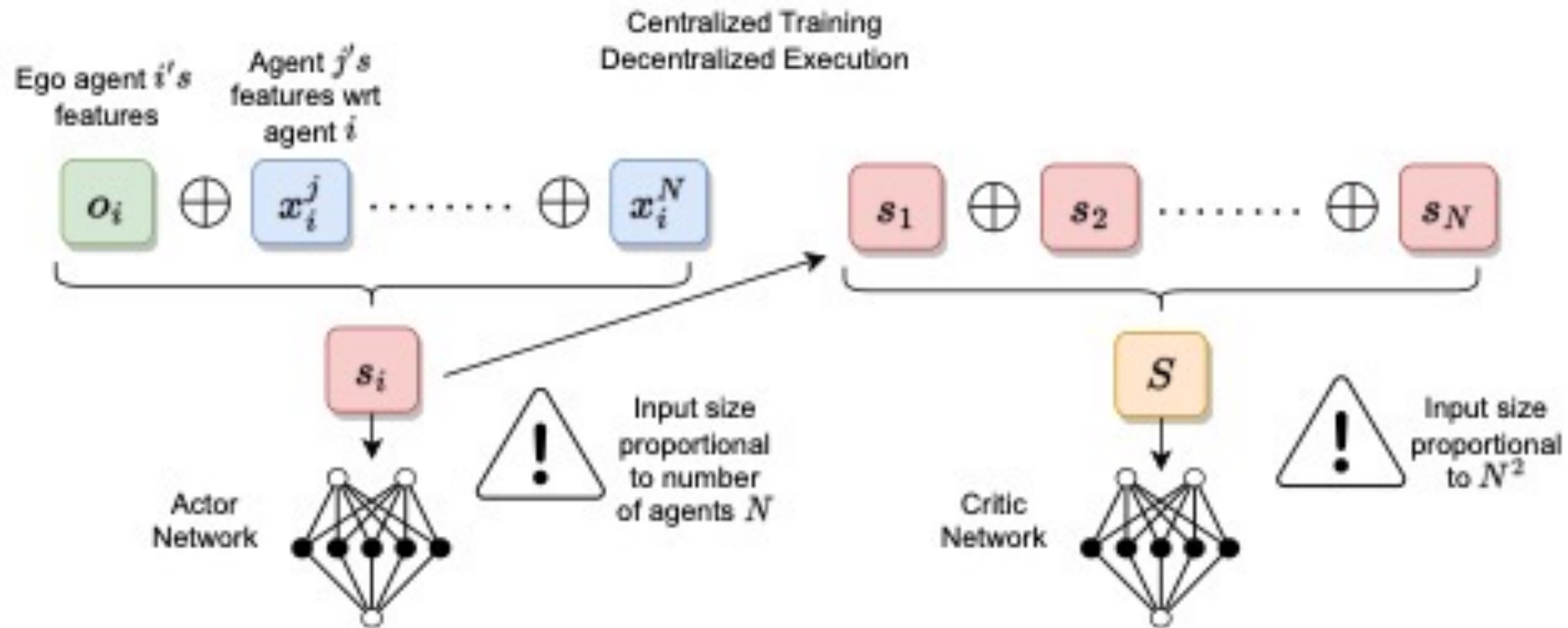
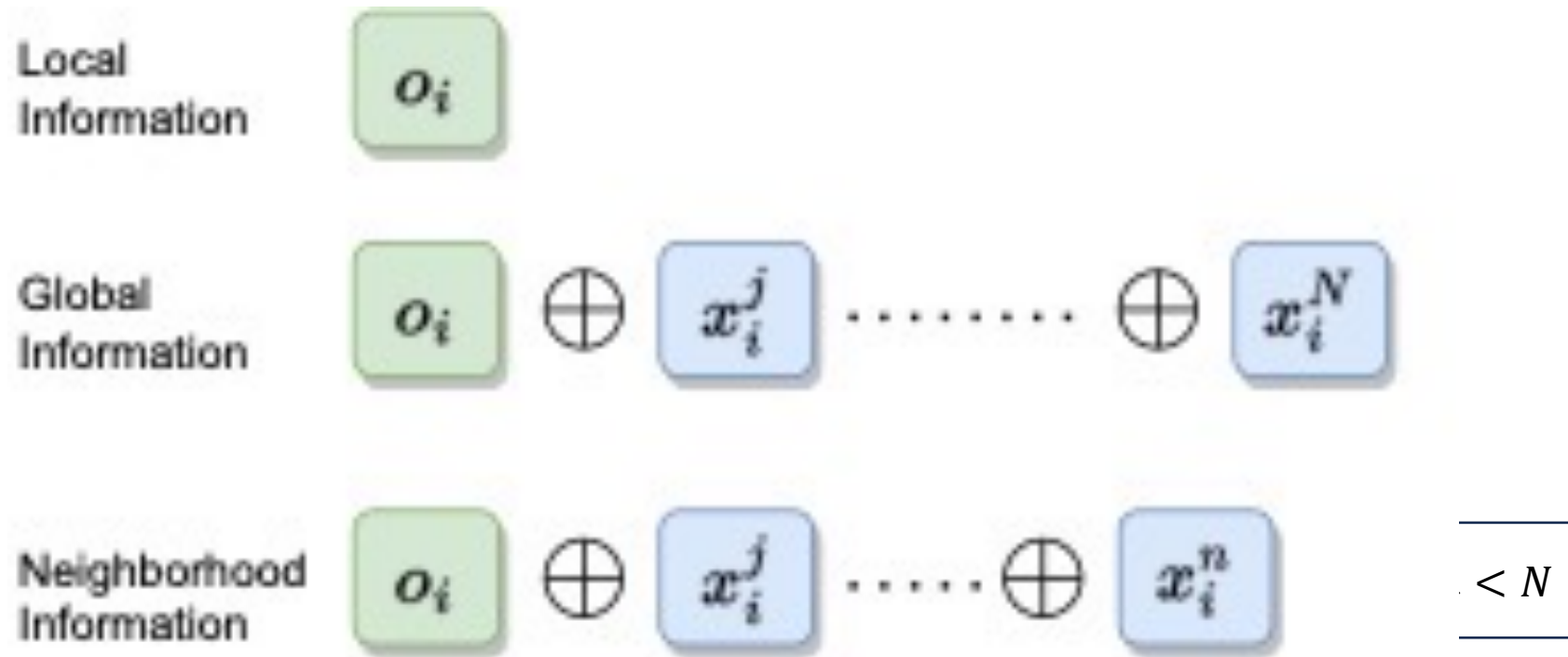# Motivation: Prior Approaches

# Motivation: Prior Approaches

# Motivation: Prior Approaches

# Motivation: Prior Approaches

# Motivating Experiment



Local Information: $o_i$

Global Information: $o_i \oplus x_i^j \cdots\cdots \oplus x_i^N$

Neighborhood Information: $o_i \oplus x_i^j \cdots \oplus x_i^n$    $< N$

# Motivating Experiment



Comparing different information modes with RMAPPO

- local
- nbd_1
- nbd_3
- nbd_5
- global

• In practice, we just have local information about the neighborhood

• And naïve concatenation of neighborhood information doesn't work

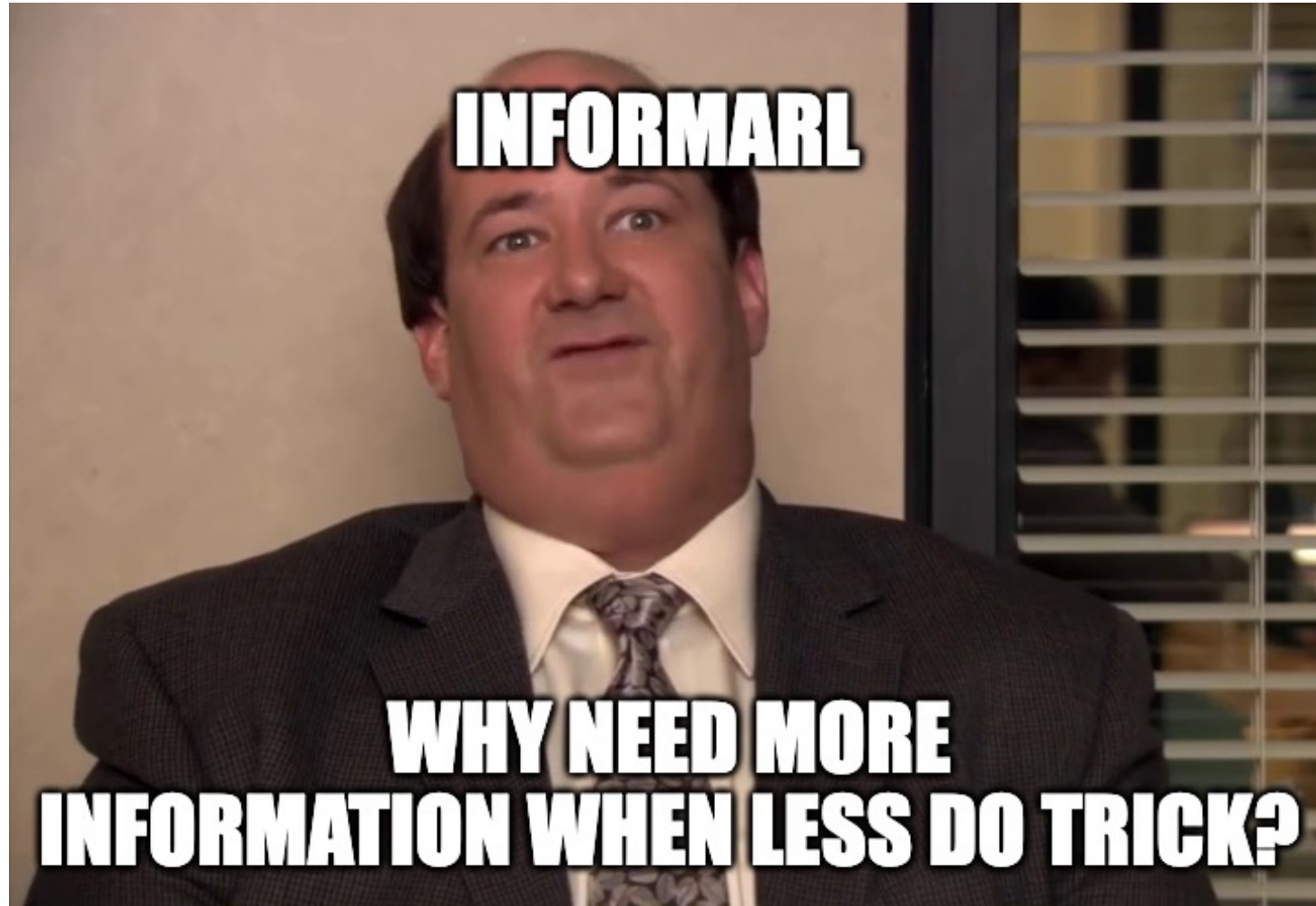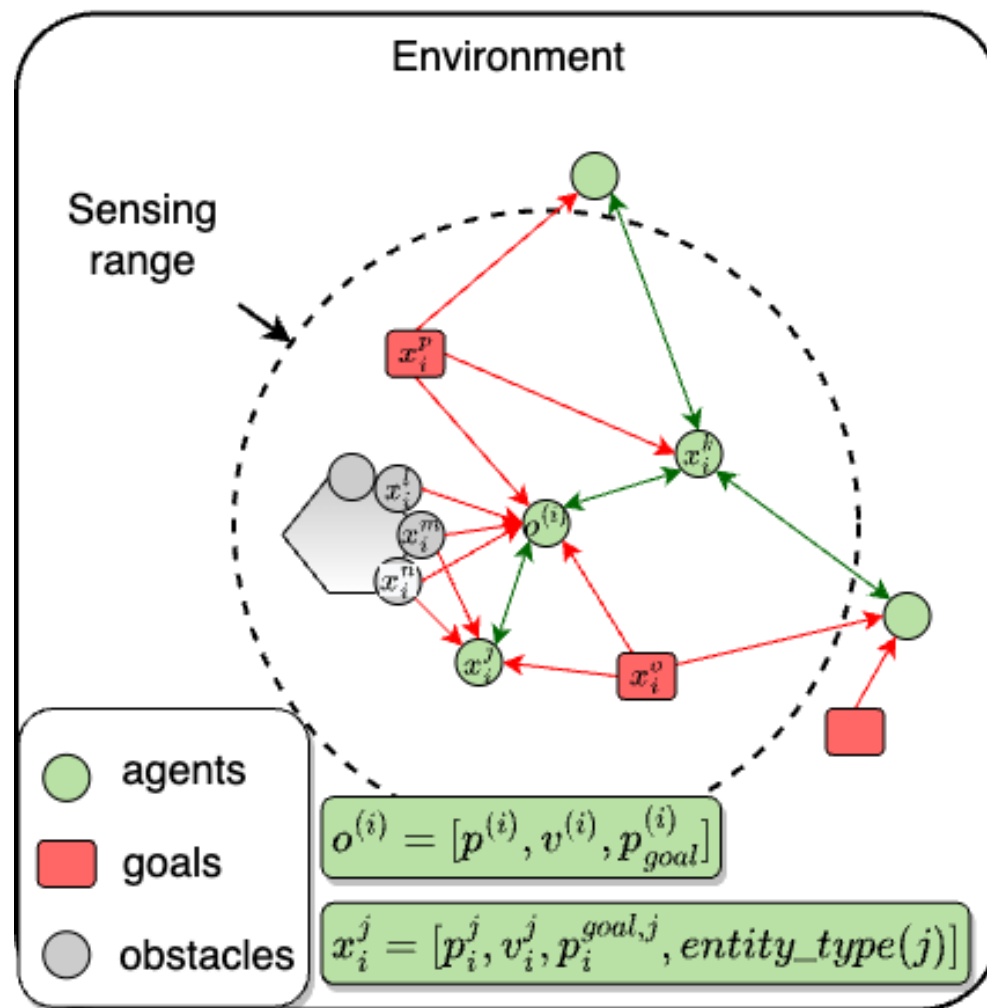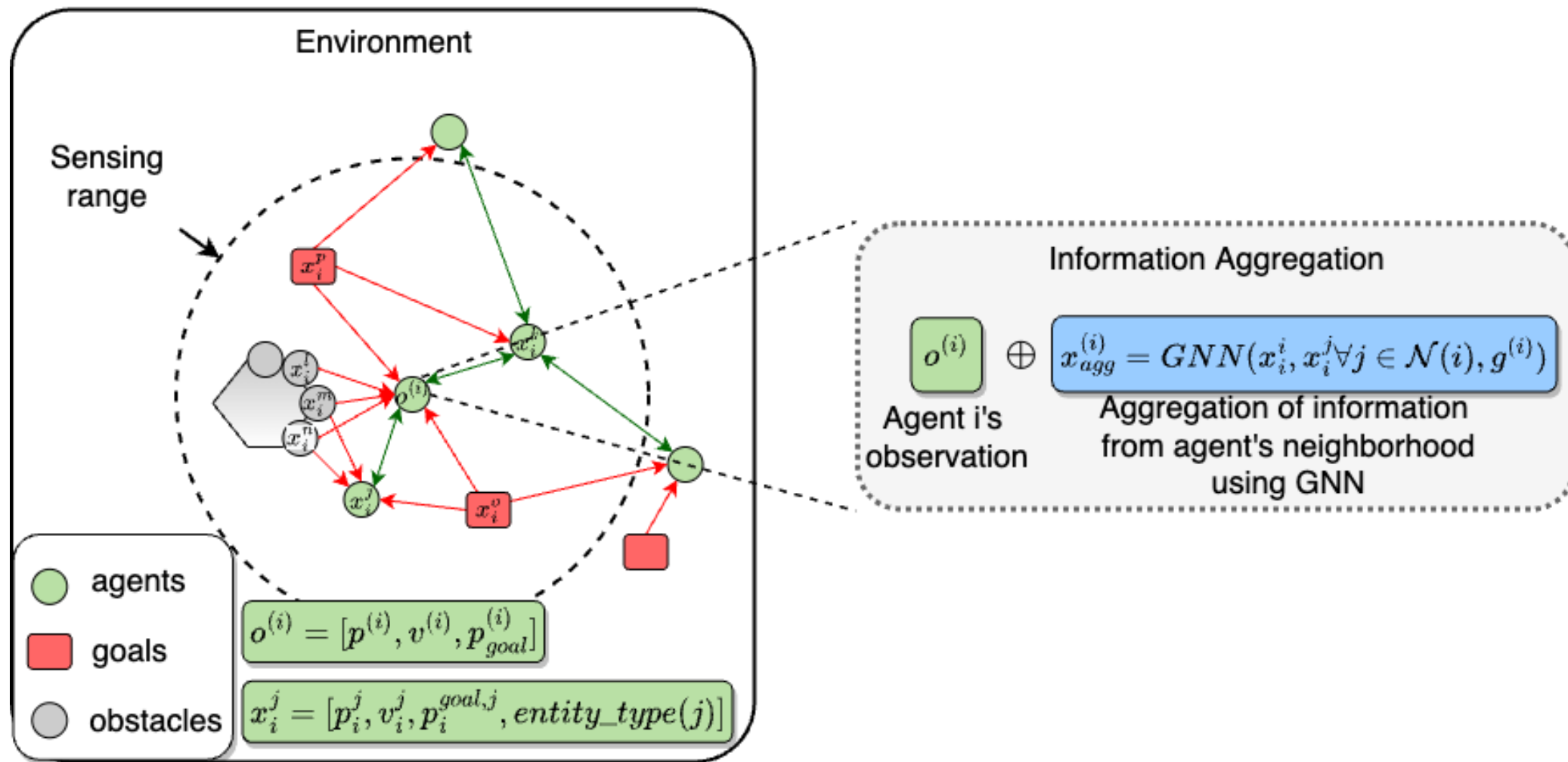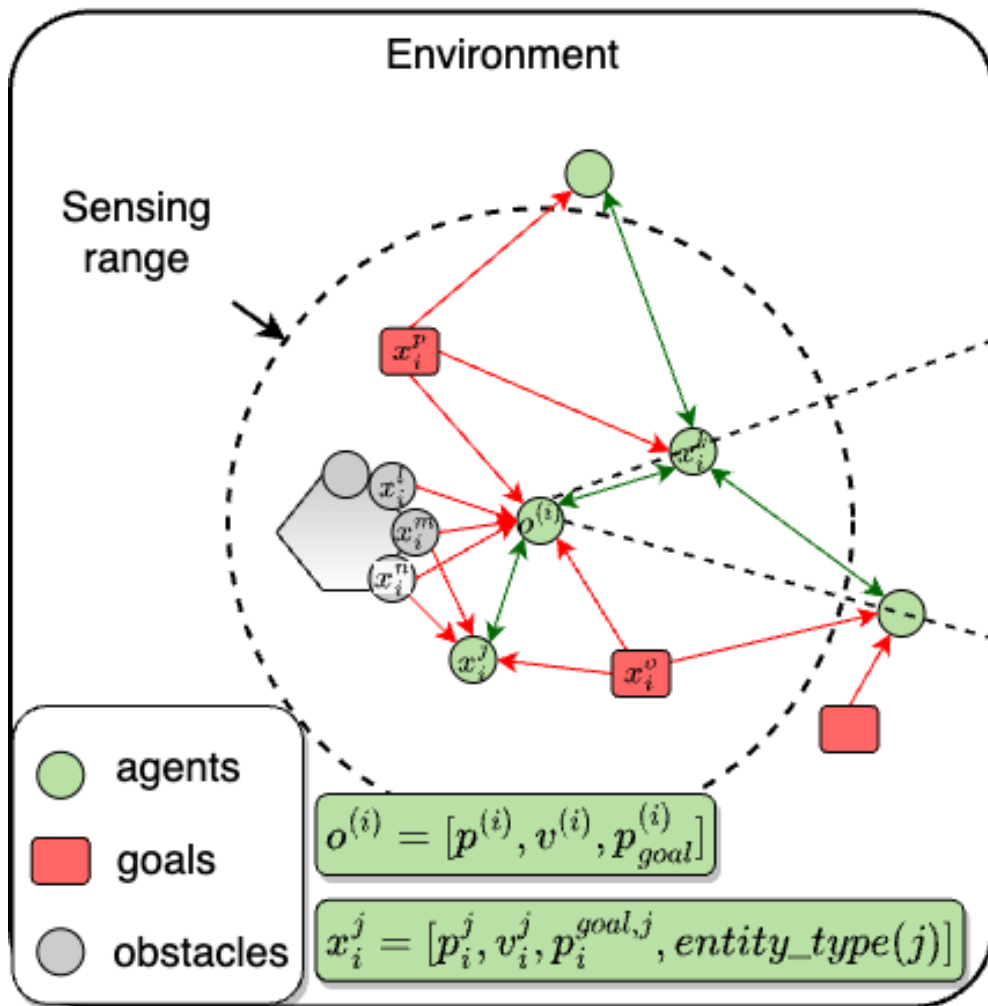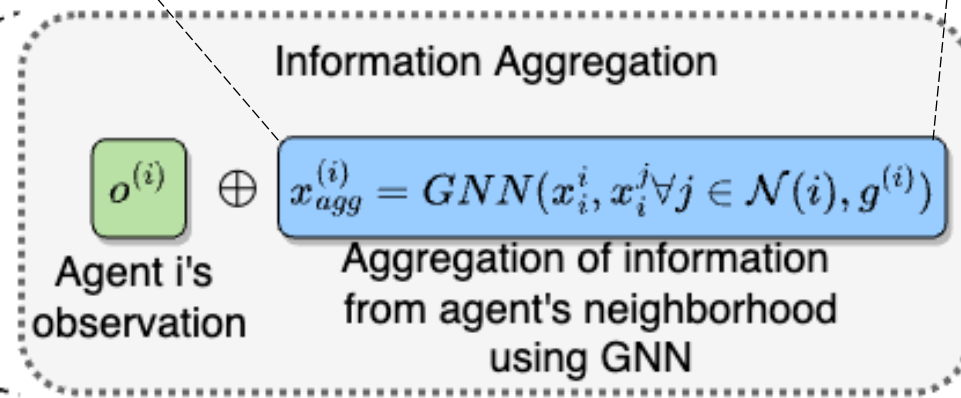# Motivating Experiment

# Method

# Method

# Method



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

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

Environment

Sensing range

agents

goals

obstacles

$$o^{(i)} = [p^{(i)}, v^{(i)}, p_{goal}^{(i)}]$$

$$x_i^j = [p_i^j, v_i^j, p_i^{goal,j}, entity\_type(j)]$$

Information Aggregation

$o^{(i)}$ $\oplus$ $x_{agg}^{(i)} = GNN(x_i^i, x_i^j \forall j \in \mathcal{N}(i), g^{(i)})$

Agent i's observation

Aggregation of information from agent's neighborhood using GNN

DINaMo

# Method

# Method

# Method

# Method

# Experiments: Environments



Target

Coverage

Formation

Line Formation

# Experiments: Sample complexity

# Experiments: Scalability

$\uparrow$ - higher better
$\downarrow$ - lower better

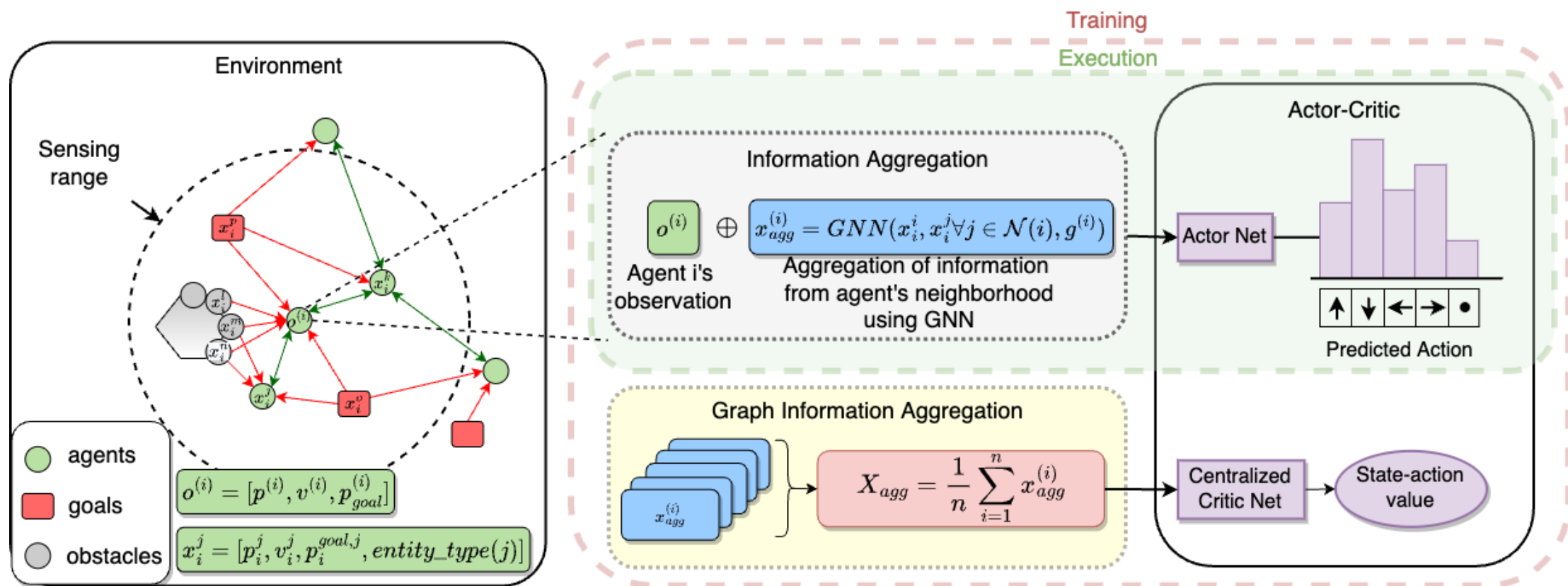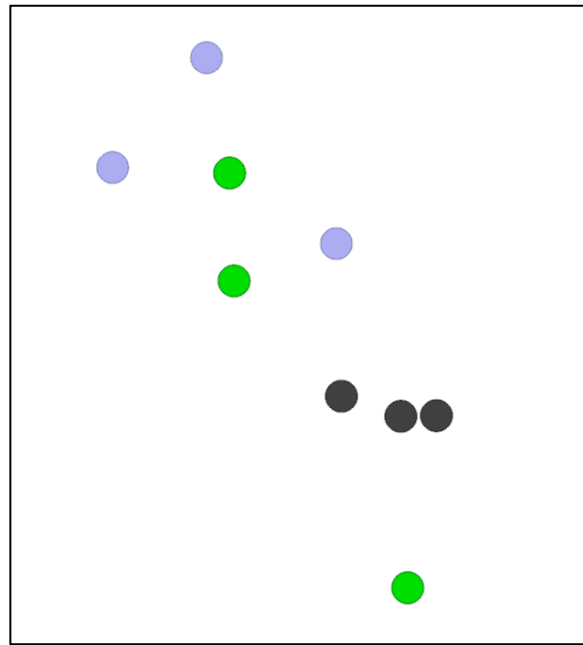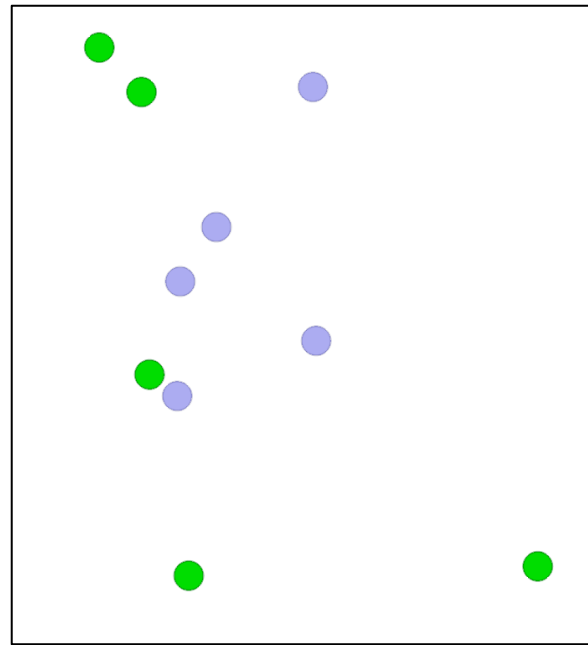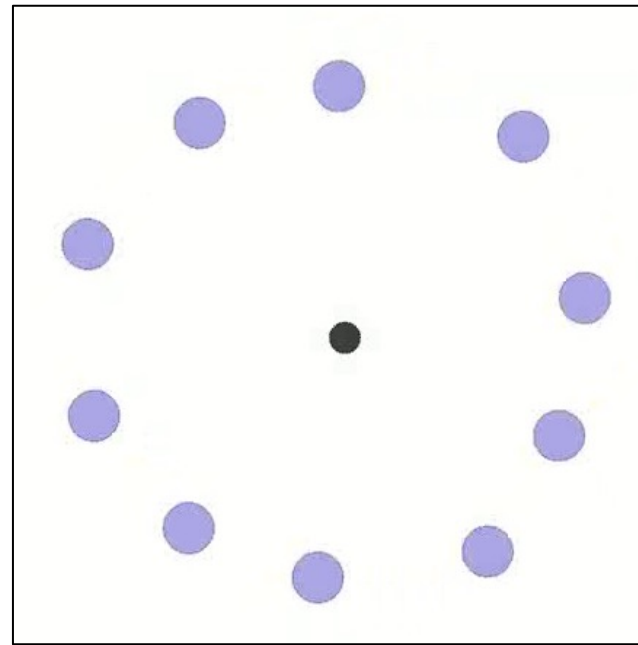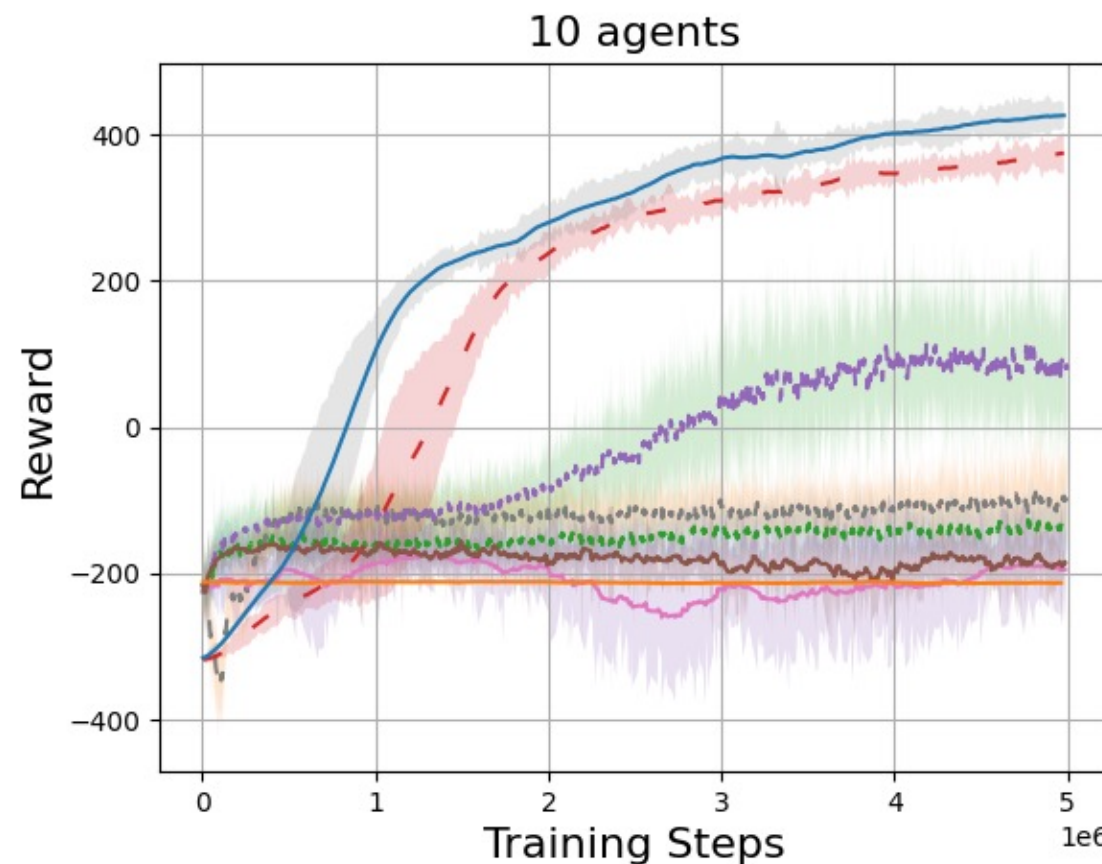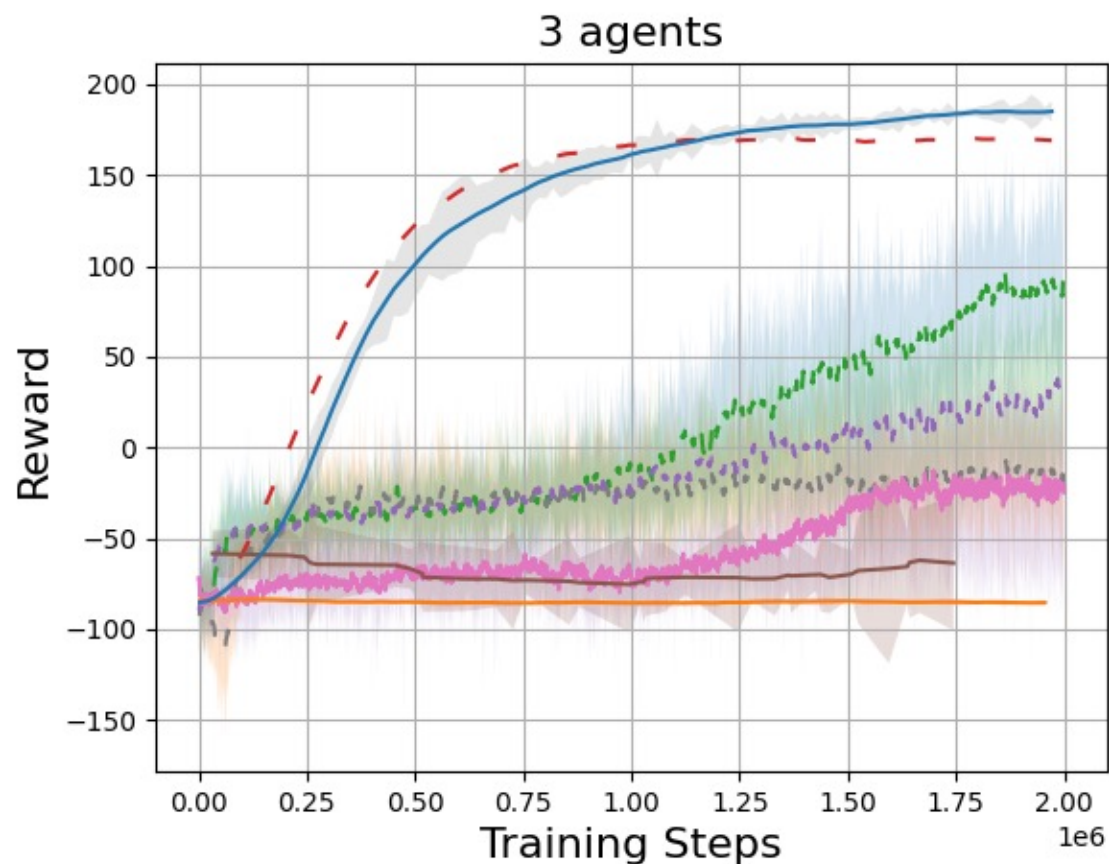| Testing / Training | | $n$=3 | $n$=7 | $n$=10 |
|---|---|---|---|---|
| $m$=3 | Reward/agent $\uparrow$ | 63.21 | 63.25 | 62.87 |
| | Avg. completion time $\downarrow$ | 0.39 | 0.40 | 0.40 |
| | Avg. #collisions/agent $\downarrow$ | 0.40 | 0.46 | 0.49 |
| | Completion rate $\uparrow$ | 100% | 100% | 99% |
| $m$=7 | Reward/agent $\uparrow$ | 61.16 | 62.23 | 61.32 |
| | Avg. completion time $\downarrow$ | 0.38 | 0.40 | 0.40 |
| | Avg. #collisions/agent $\downarrow$ | 0.74 | 0.66 | 0.70 |
| | Completion rate $\uparrow$ | 100% | 100% | 100% |
| $m$=10 | Reward/agent $\uparrow$ | 58.59 | 58.23 | 58.67 |
| | Avg. completion time $\downarrow$ | 0.38 | 0.40 | 0.39 |
| | Avg. #collisions/agent $\downarrow$ | 0.95 | 0.88 | 0.87 |
| | Completion rate $\uparrow$ | 100% | 99% | 100% |
| $m$=15 | Reward/agent $\uparrow$ | 53.19 | 53.46 | 54.21 |
| | Avg. completion time $\downarrow$ | 0.39 | 0.40 | 0.40 |
| | Avg. #collisions/agent $\downarrow$ | 1.28 | 1.21 | 1.20 |
| | Completion rate $\uparrow$ | 100% | 99% | 99% |

DINaMo

# Experiments: Other environments

↑ - higher better
↓ - lower better

| Task environment | $m$ | Metric | RMAPPO (global info) | InforMARL (local info) |
|---|---|---|---|---|
| Coverage | $m$=3 | Avg. completion time ↓ | 0.34 | 0.36 |
| | | Completion rate ↑ | 100% | 100% |
| | $m$=7 | Avg. completion time ↓ | 0.42 | 0.43 |
| | | Completion rate ↑ | 100% | 99% |
| Formation | $m$=3 | Avg. completion time ↓ | 0.31 | 0.30 |
| | | Completion rate ↑ | 100% | 100% |
| | $m$=7 | Avg. completion time ↓ | 0.47 | 0.43 |
| | | Completion rate ↑ | 100% | 100% |
| Line | $m$=3 | Avg. completion time ↓ | 0.24 | 0.21 |
| | | Completion rate ↑ | 100% | 100% |
| | $m$=7 | Avg. completion time ↓ | 0.38 | 0.36 |
| | | Completion rate ↑ | 100% | 100% |

DINaMo

# Conclusions

- InforMARL uses a graph neural network (GNN)-based architecture for **scalable** multi-agent RL in a **decentralized** fashion.

- InforMARL is **transferable** to scenarios with a different number of entities in the environment than what it was trained on.

- InforMARL has **better sample complexity** than most other standard MARL algorithms with global observations