# Scalable Decision-Theoretic Planning in Open and Typed Multiagent Systems

Adam Eck[1], Maulik Shah[2], Prashant Doshi[2], and Leen-Kiat Soh[3]

[1]Oberlin College, [2]University of Georgia, [3]University of Nebraska
aeck@oberlin.edu, mns28652@uga.edu, pdoshi@uga.edu, lksoh@cse.unl.edu

OBERLIN COLLEGE & CONSERVATORY

UNIVERSITY OF GEORGIA

UNIVERSITY OF NEBRASKA Lincoln

## Problem: Scalable Planning in Open Environments

**Many-Agent Environments**: environments with dozens to thousands of agents

**Open Environments**: agents join and leave the environment (temporarily or permanently) over time

- *Real-world Examples*: wildfire fighting, autonomous ridesharing, cybersecurity
- <u>Problem</u>: Openness requires agents to not only predict what actions their peers will take in order to choose a best response, but also *first estimate which peers will even be present to take actions*



- Tracking the presence of neighbors results in a *more complicated problem model*. The increase in the problem model size is:
  - *exponential* if agent presence is considered in the environment state, which is intractable, especially for many-agent environments
  - only *linear* if agent presence is considered within the mental models of each agent in an I-POMDP-Lite problem model (or other I-POMDP variants)

- Larger planning problem *affects scalability* as the number of agents increases (to many-agents)
  - MCTS algorithms: each trajectory requires sampling actions for all neighbors, so more agents results in fewer sampled trajectories in a fixed time budget (for responsive reasoning)
  - With openness, time is also spent updating estimates of presence of each neighbor, further reducing the number of trajectories possible within a time budget
  - **Frame-action Anonymity** offers some relief by replacing joint actions with counts of actions called **configurations $C$** when environment dynamics *do not depend on which agents take which actions*.
    - <u>Key observation</u>: estimating *counts* of actions might not require estimating actions for all agents *individually*

## Solution: Many-Agent Planning under Openness



**Main principle**: *only model some neighbors* to counter the increased complexity caused by openness

- In polling and survey theory, social scientists only survey a small randomly sampled portion of the population to estimate the behaviors and opinions of all people

Subject agent estimates action counts for the entire neighborhood by *extrapolating* the estimated actions of modeled neighbors $\hat{N}_\theta(i)$ using the **multinomial distribution**

$$P(C_\theta|s^t, M^t) \sim Multi\left(|N_\theta(i)|, \left\{\hat{p}_{a_1,\hat{N}_\theta(i)}, \ldots, \hat{p}_{a_{|A|},\hat{N}_\theta(i)}\right\}\right)$$

$$\hat{p}_{a,\hat{N}_\theta(i)} = \frac{\hat{n}_{\pi(s^t)=a,\hat{N}_\theta(i)}}{|\hat{N}_\theta(i)|}$$

**I-POMCP$_O$ Many-Agent MCTS Algorithm**: adapts POMCP algorithm to many-agent open environments

- Estimate actions of a few neighbors then extrapolate to all neighbors' behaviors
  - *Better time complexity* than previous I-POMCP MCTS algorithm (Hua et al., 2015).
- Comparable to Dec-POMDP MCTS algorithms (Amato & Oliehoek, 2015; Best et al., 2019) *but does not require* the subject agent to observe other agents' actions nor their observations

**Algorithm 1** I-POMCP$_O$: Open Many-Agent MCTS

**Note:** $T$ is the tree (initially empty), $p$ is a path from the root of the tree (with $p = \emptyset$ signifying the root), $B_p$ is the particle filter signifying the set of state-model pairs encountered at the node at $p$ in the tree, $PF$ is the root particle filter, $N$ is count of the number of visits to each node in the tree initialized to some constant $\nu \geq 0$, $Q$ is the Q function initialized to 0, $c$ a constant from UCB-1.

1: **procedure** I-POMDP-MCTS($PF, \tau$)
2:    $time \leftarrow 0$
3:    **while** $time < \tau$ **do**
4:      $s^0, M^0 \leftarrow SampleParticle(PF)$
5:      UpdateTree($s^0, M^0, 0, \emptyset$)
6:      Increment $time$
7:    **return** $\arg\max_{a \in A_i} Q(\emptyset, a)$

1: **procedure** UPDATETREE($s^t, M^t, t, p$)
2:    **if** $t \geq H$ **then**
3:      **return** 0
4:    $B_p \leftarrow B_p \cup \{(s^t, M^t)\}$
5:    **if** $p \notin T$ **then**
6:      $T \leftarrow T + leafnode(p)$
7:      **return** Rollout($s^t, M^t, t$)
8:    **else**
9:      $C^t \leftarrow$ SampleConfiguration($s^t, M^t$)
10:      $a_i^t \leftarrow \arg\max_{a \in A_i} Q(p, a) + c\sqrt{(\log N_p)/N_{p \to a}}$
11:      $s^{t+1}, M^{t+1}, o_i^t, r_i^t \leftarrow$ Simulate($s^t, M^t, a_i^t, C^t$)
12:      $N_p \leftarrow 1 + N_p$
13:      $N_{p \to a_i^t} \leftarrow 1 + N_{p \to a_i^t}$
14:      $p' \leftarrow p + (a_i^t, o_i^t)$
15:      $R \leftarrow r_i^t + \gamma$·UpdateTree($s^{t+1}, M^{t+1}, t+1, p'$)
16:      $Q(p, a_i^t) \leftarrow Q(p, a_i^t) + R - Q(p, a_i^t)/N_{p \to a_i^t}$
17:      **return** $R$

1: **procedure** ROLLOUT($s^t, M^t, t$)
2:    $R \leftarrow 0, t' \leftarrow t$
3:    **while** $t < H$ **do**
4:      $C^t \leftarrow$ SampleConfiguration($s^t, M^t$)
5:      $a_i^t \leftarrow$ SampleAction($A_i$)
6:      $s^{t+1}, M^{t+1}, o_i^t, r_i^t \leftarrow$ Simulate($s^t, M^t, a_i^t, C^t$)
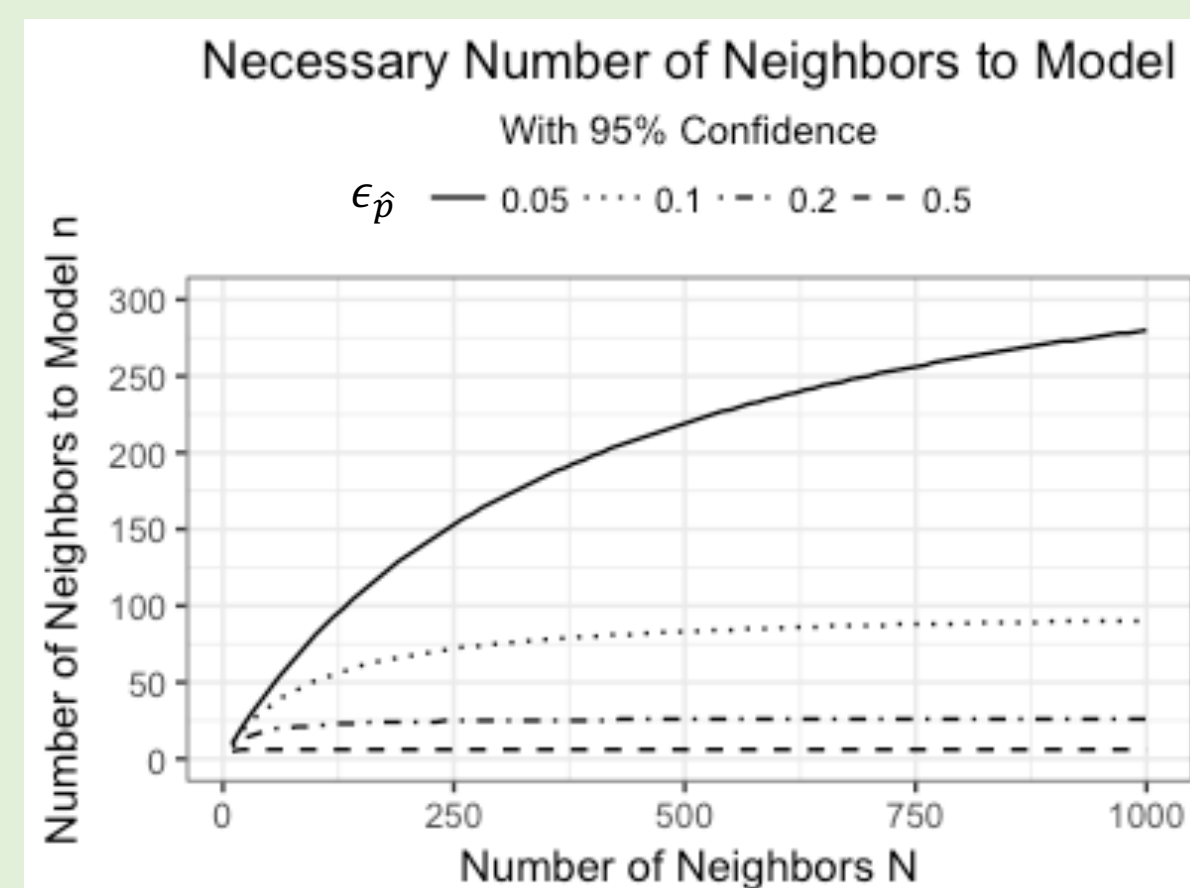7:      $R \leftarrow R + \gamma^{t-t'} \cdot r_i^t, t \leftarrow t+1$
8:    **return** $R$

1: **procedure** SAMPLECONFIGURATION($s^t, M^t$)
2:    $C(a, \theta) \leftarrow 0, \hat{n}_{\pi(s^t)=a,\hat{N}_\theta(i)} \leftarrow 0 \quad \forall a, \theta$
3:    **for** $\mathcal{M}_{j,l-1} \in M^t$ **do**
4:      $a \sim \pi_{j,l-1}(s^t)$
5:      $\hat{n}_{\pi(s^t)=a,\hat{N}_{\theta_j}(i)} \leftarrow \hat{n}_{\pi(s^t)=a,\hat{N}_{\theta_j}(i)} + 1$
6:    **for** $\theta \in \Theta$ **do**
7:      **for** $a \in A$ **do**
8:        $\hat{p}_{a,N(i)} \leftarrow \hat{n}_{\pi(s^t)=a,\hat{N}_\theta(i)} / |\hat{N}_\theta(i)|$
9:    **for** $j \in N_\theta(i)$ **do**
10:      $a \sim Cat(\hat{p}_{a_1,\hat{N}_\theta}, \hat{p}_{a_2,\hat{N}_\theta}, \ldots, \hat{p}_{a_{|A|},\hat{N}_\theta})$
11:      $C(a, \theta) \leftarrow C(a, \theta) + 1$
12:    **return** $C$

## Theoretical Results

**Theorem 1** With confidence $1 - \alpha$, the *error* incurred by the subject agent in its estimate *of the proportion $\hat{p}_{a,\hat{N}_\theta(i)}$* of its neighbors of frame $\theta$ that will perform action $a$ is *bounded by the given $\epsilon_{\hat{p}}$* so long as it models the following number of neighbors:

$$n_\theta = |\hat{N}_\theta(i)| \geq \frac{N\left(\frac{t_{n-1,\frac{\alpha}{2}}}{\epsilon_{\hat{p}}}\right)^2}{N - 1 + \left(\frac{t_{n-1,\frac{\alpha}{2}}}{\epsilon_{\hat{p}}}\right)^2}$$

Necessary Number of Neighbors to Model
With 95% Confidence
$\epsilon_{\hat{p}}$ — 0.05 ·· 0.1 -· 0.2 - 0.5



**Corollary 1** The *maximum error in the multinomial distribution* of the configuration of other agents' actions $P(C|s^t, M^t)$ is bounded by:

$$\epsilon_{P(C)} = |P^*(C|s^t, M^t) - P(C|s^t, M^t)|$$

$$< \frac{\prod_\theta |N_\theta(i)|!}{\prod_{\theta,a} C(a,\theta)!}\left[\prod_{\theta,a}(\hat{p}_{a,\theta} - \epsilon_{\hat{p}})^{C(a,\theta)} - \prod_{\theta,a}\hat{p}_{a,\theta}^{C(a,\theta)}\right]$$

when the subject agent models only $n_\theta$ neighbors (from Theorem 1) to achieve at most $\epsilon_{\hat{p}}$ error in its estimates of the action probabilities parameterizing the multinomial distribution.

**Theorem 2** The *regret* of the subject agent from modeling only a subset of its neighbors *is bounded by*:
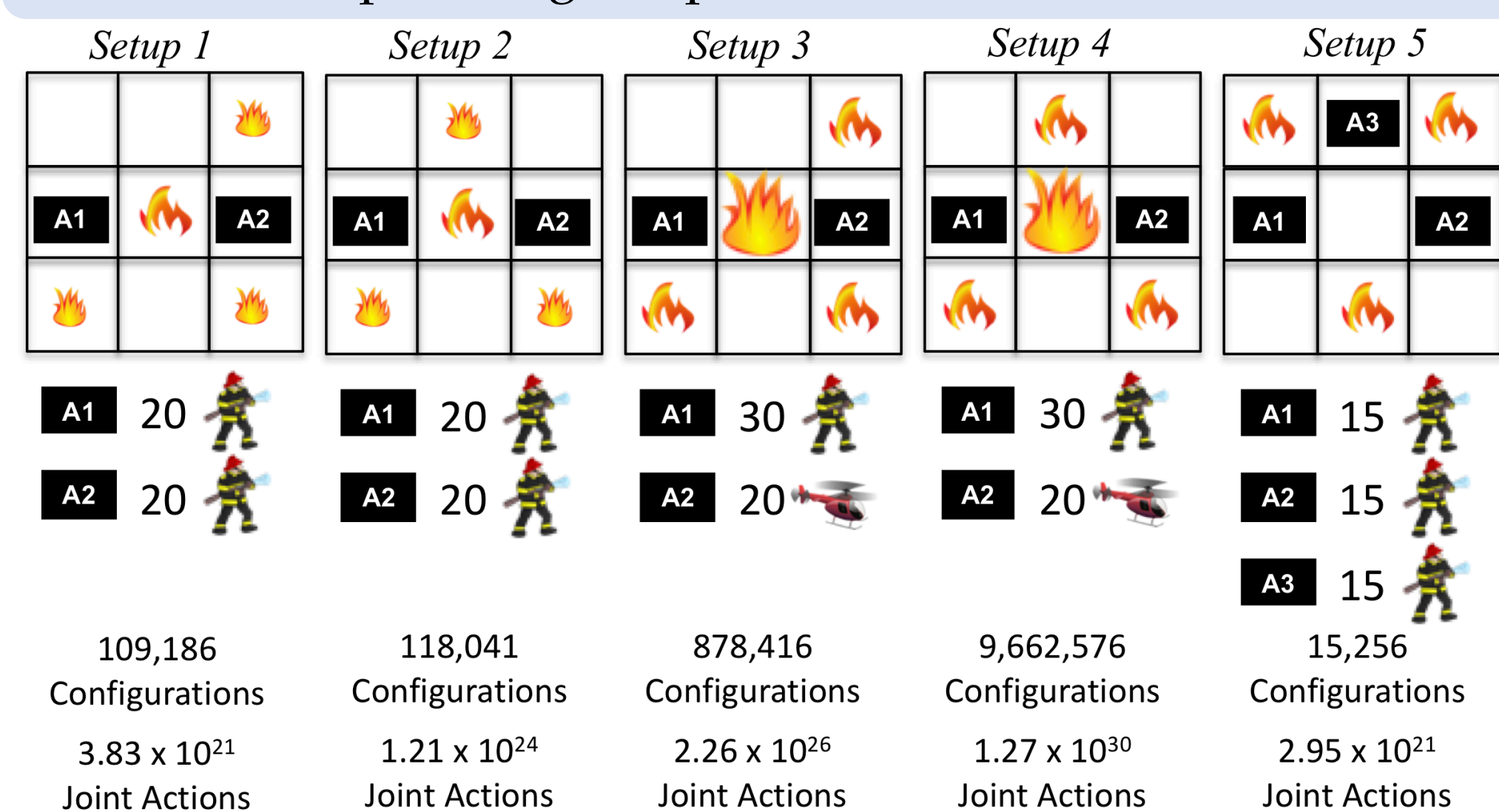
$$\|V_{i,k}^* - J_{i,k}\|_\infty$$

$$\leq 2\epsilon_{P(C)} \cdot |C| \cdot R_{max}\left[\gamma^{k-1} + \frac{1}{1-\gamma}\left(1 + 3\gamma\frac{|\Omega_i|}{1-\gamma}\right)\right]$$

which is *linear* in the error $\epsilon_{P(C)}$ in the agent's estimation of configuration likelihoods caused by modeling some neighbors only and *proportional to only $\sqrt{1/n_\theta}$ in the worst case* (due to the fact that $\epsilon_{P(C)}$ is at worst linear in $\epsilon_{\hat{p}}$ and $\epsilon_{\hat{p}}$ is proportional to $\sqrt{1/n_\theta}$)
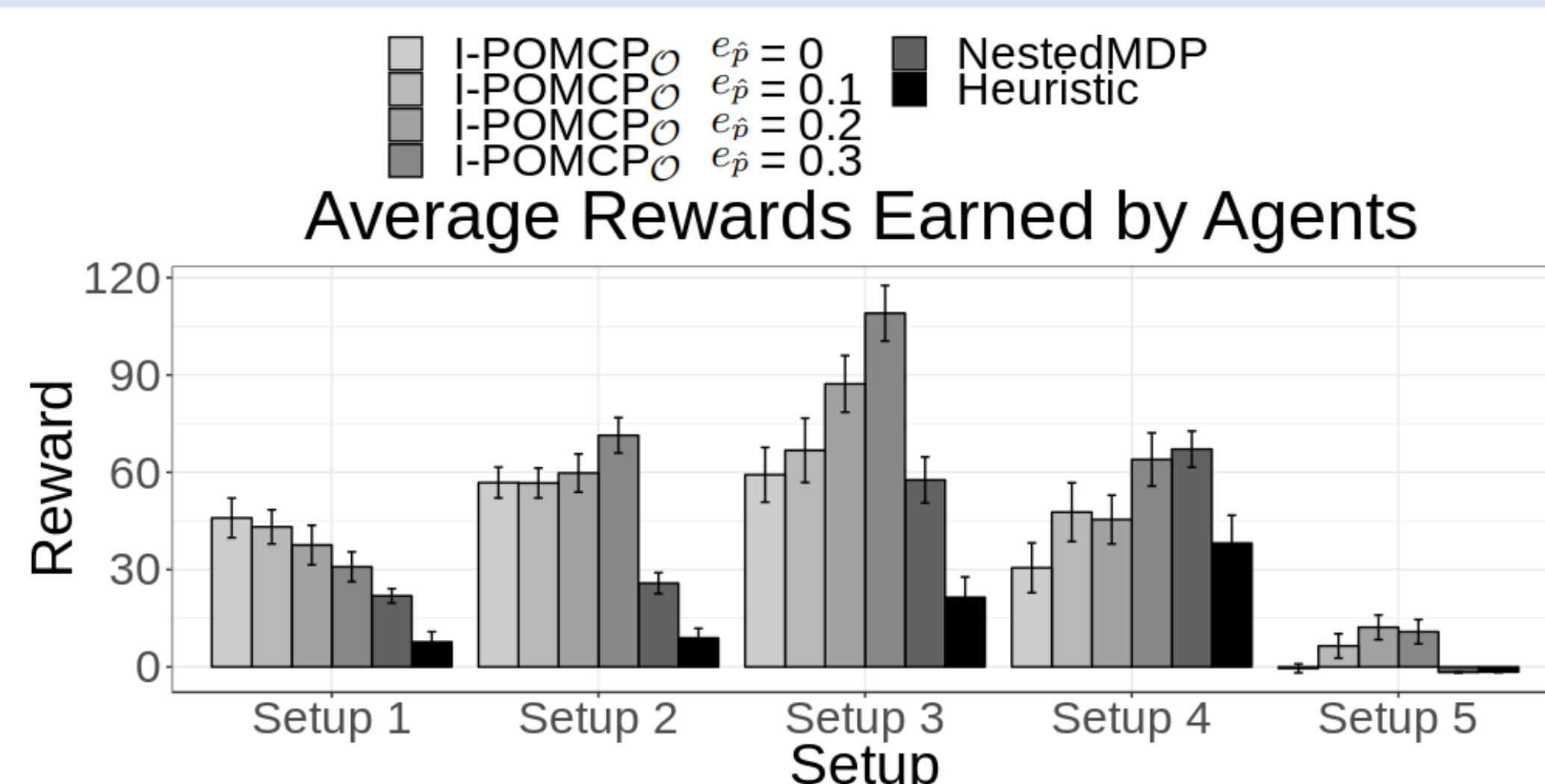
## Experimental Results

**Wildfire Problem:** ground firefighter and helicopter agents with different capabilities cooperate to put out nearby fires.
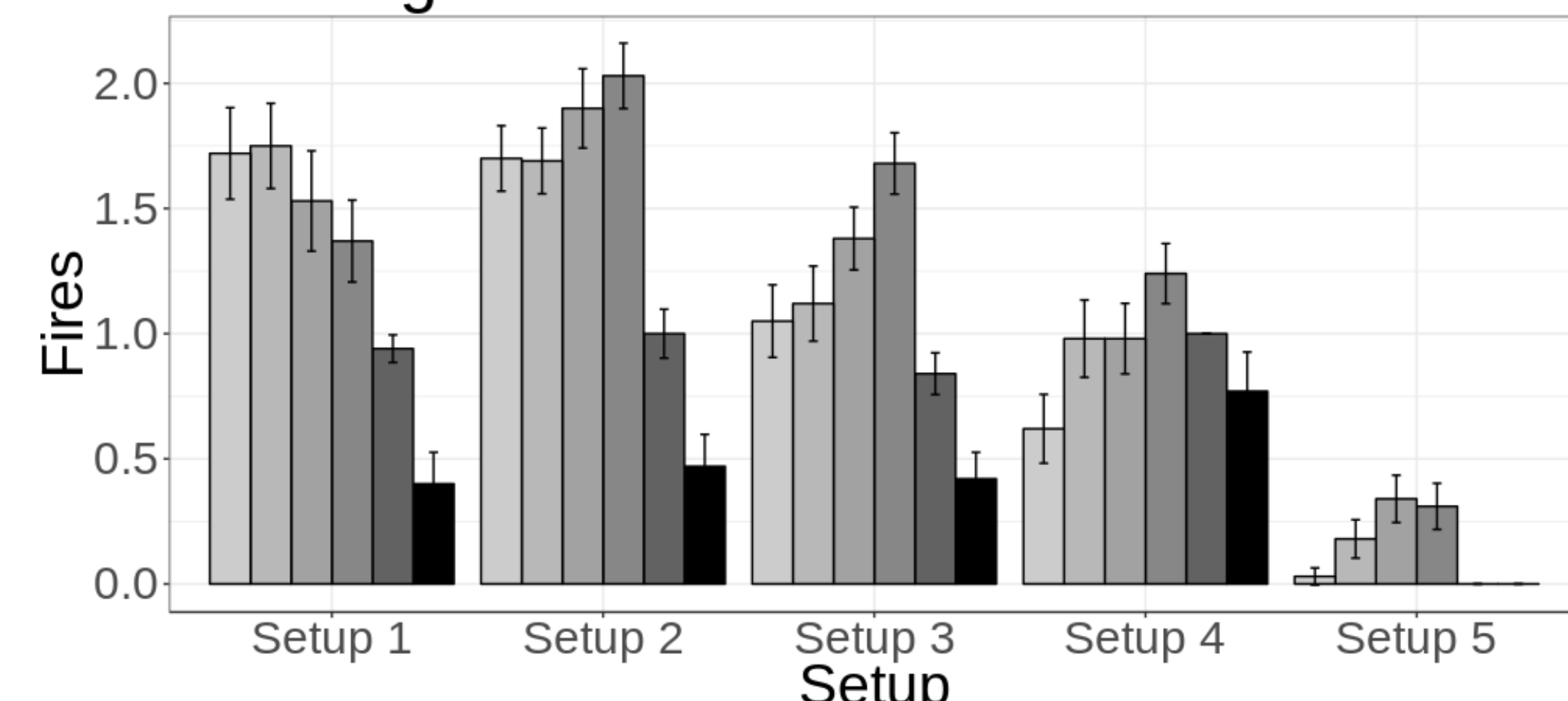
- **Goal:** maximize rewards earned for putting out fires (and minimize costs for fires burning out)
- Agents can *run out of suppressant*, requiring them to temporarily leave the environment to recharge.
- Consider an order of magnitude more agents than prior studies on planning in open environments



| Setup 1 | Setup 2 | Setup 3 | Setup 4 | Setup 5 |
|---|---|---|---|---|
| A1 20 | A1 20 | A1 30 | A1 30 | A1 15 |
| A2 20 | A2 20 | A2 20 | A2 20 | A2 15 |
| | | | | A3 15 |
| 109,186 Configurations | 118,041 Configurations | 878,416 Configurations | 9,662,576 Configurations | 15,256 Configurations |
| $3.83 \times 10^{21}$ Joint Actions | $1.21 \times 10^{24}$ Joint Actions | $2.26 \times 10^{26}$ Joint Actions | $1.27 \times 10^{30}$ Joint Actions | $2.95 \times 10^{21}$ Joint Actions |

**Agent Reasoning Models:** Heuristic randomly chooses active fires, NestedMDP is principled reasoning with value iteration at level 1, and I-POMCP$_O$ is principled reasoning at level 2 with MCTS modeling $n_\theta$ neighbors from Theorem 1 based on $\epsilon_{\hat{p}}$ maximum error

Legend:
- I-POMCP$_O$ $\epsilon_{\hat{p}} = 0$
- I-POMCP$_O$ $\epsilon_{\hat{p}} = 0.1$
- I-POMCP$_O$ $\epsilon_{\hat{p}} = 0.2$
- I-POMCP$_O$ $\epsilon_{\hat{p}} = 0.3$
- NestedMDP
- Heuristic



Average Rewards Earned by Agents



Average Number of Fire Locations Put Out

- I-POMCP$_O$ *outperformed all baselines* (statistically significantly higher rewards in 4 of 5 setups) due to *better coordination* between agents by working together to put out fires in more locations

- In the more complicated Setups 2-4, *modeling fewer neighbors $n_\theta$* due to higher allowable error $\epsilon_{\hat{p}}$ *led to improved performance* due to more sampled trajectories in the fixed time budget. This implies that the approximation error caused by modeling only some neighbors *was less than* the approximation error from sampling fewer trajectories in MCTS.