

Convergence Rates for Localized Actor-Critic in Networked Markov Potential Games

Zhaoyi Zhou¹, Zaiwei Chen², Yiheng Lin², Adam Wierman²

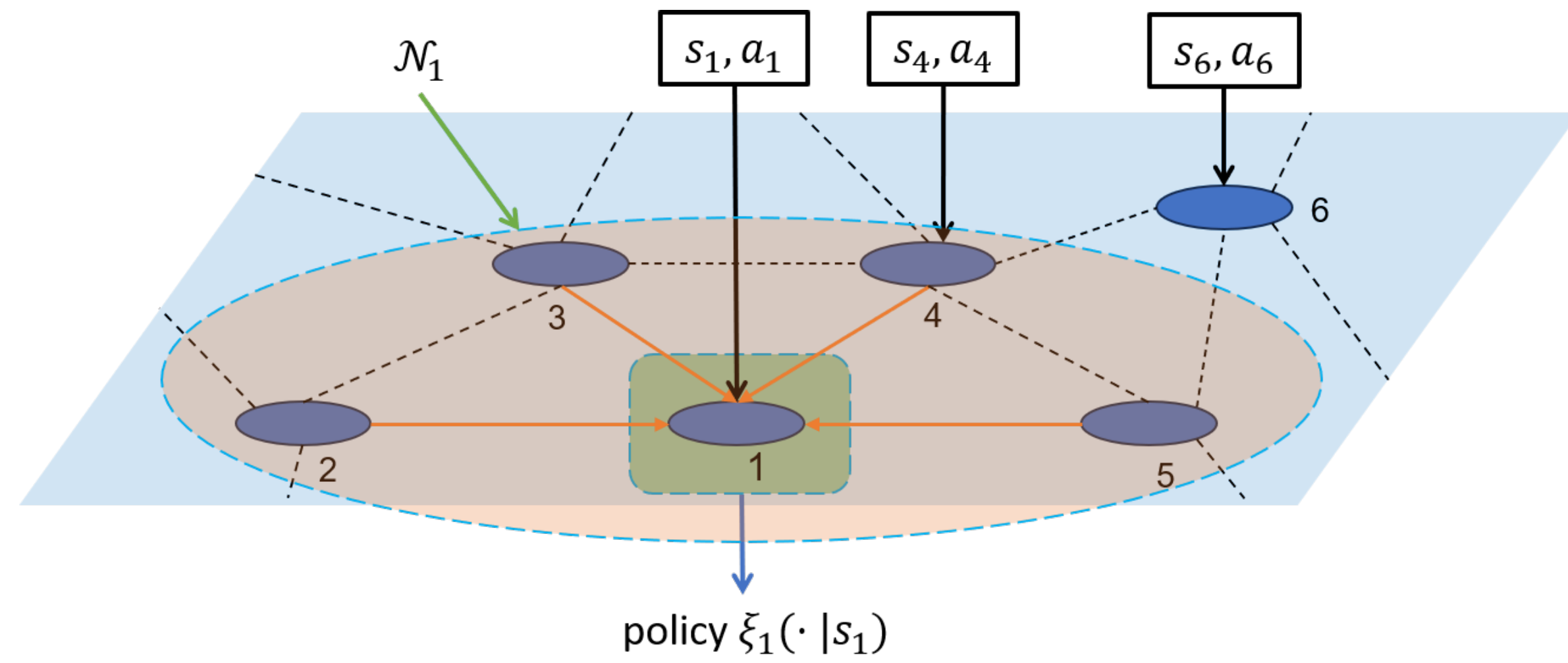
¹ Tsinghua University ² California Institute of Technology

39th Conference on Uncertainty in Artificial Intelligence

Setting

We consider n agents associated with an undirected graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ (“communication network”). Each agent i has its local state $s_i \in \mathcal{S}_i$ and local action $a_i \in \mathcal{A}_i$. Global state/action space can be decomposed as

$$\mathcal{S} = \mathcal{S}_1 \times \mathcal{S}_2 \times \cdots \times \mathcal{S}_n, \text{ and } \mathcal{A} = \mathcal{A}_1 \times \mathcal{A}_2 \times \cdots \times \mathcal{A}_n.$$



The next local state depends on local action and neighboring local states:

$$\mathcal{P}(s(t+1) | s(t), a(t)) = \prod_{i=1}^n \mathcal{P}_i(s_i(t+1) | s_{\mathcal{N}_i}(t), a_i(t)).$$

Local reward depends on states and actions of agents within κ_r -graph distance, i.e., $r_i(s, a) = r_i(s_{\mathcal{N}_i^{\kappa_r}}, a_{\mathcal{N}_i^{\kappa_r}}), \forall i \in \mathcal{N}$.

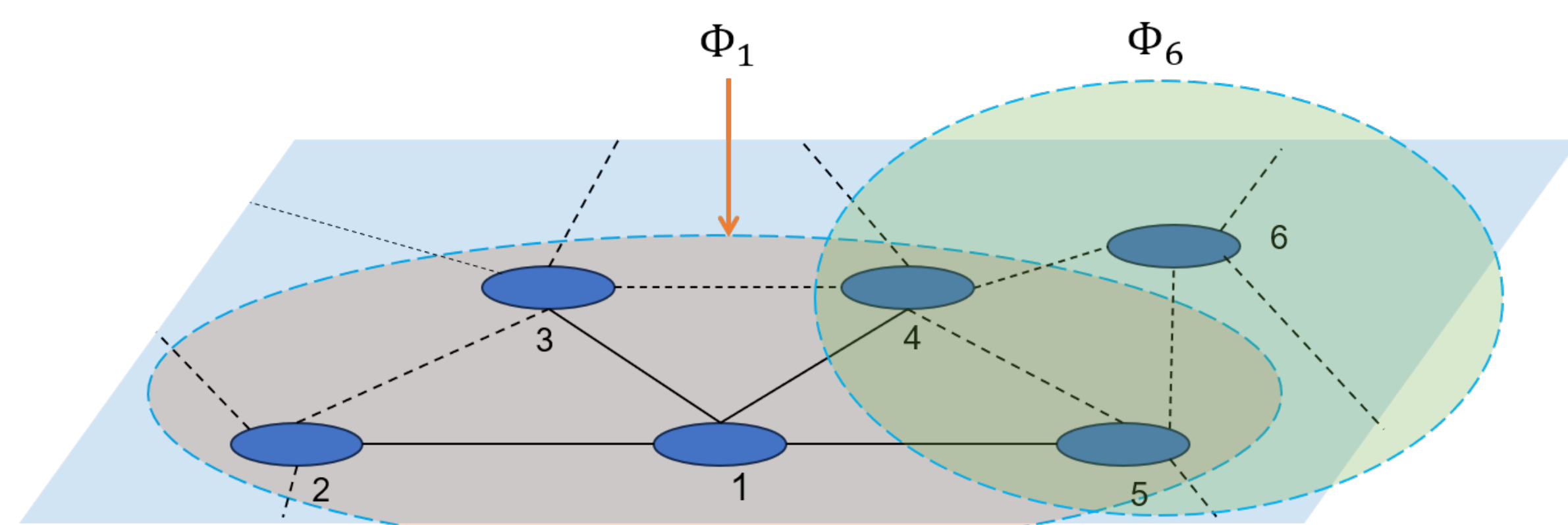
Each agent i adopts a localized policy $\xi_i(\cdot | s_i)$, which maps local state space \mathcal{S}_i to $\Delta(\mathcal{A}_i)$ using softmax parameterization. Every agent i tries to maximize its own expected γ -discounted return, denoted by $J_i(\xi)$.

Networked Markov Potential Game (NMPG)

Relax MPG by restricting the impact of potential to “nearby” agents.

Def 1 (NMPG). A multi-agent Markov game is called a κ_G -NMPG if there exists a set of **local** potential functions $\{\Phi_i\}_{i \in \mathcal{N}}$, such that Φ_i tracks policy deviation of any agent j **within κ_G -graph-distance of i** :

$$J_j(\xi'_j, \xi_{-j}) - J_j(\xi_j, \xi_{-j}) = \Phi_i(\xi'_j, \xi_{-j}) - \Phi_i(\xi_j, \xi_{-j}), \forall i, \forall j \in \mathcal{N}_i^{\kappa_G}, \forall \xi_j, \xi_{-j}, \xi'_j.$$



Our Contributions:

(1) **Networked** Markov potential games;

(2) **Localized** Actor-Critic;

(3) **Finite-Sample** Analysis.

Localized Actor-Critic

Hyperparameter κ_c controls trade-off between communication distance and truncation accuracy.

Algorithm 1 Localized Actor-Critic

```

1: Let  $\theta_i(0) = 0, \Delta_i^0(m) = 0$  for all  $i$  and  $m$ .
2: for  $m = 0, 1, 2, \dots, M-1$  do
3:   Construct  $\epsilon$ -exploration policy  $\hat{\xi}$  from  $\xi^{\theta(m)}$ .
4:   Use  $\hat{\xi}$  to collect samples for  $K$  iterations:
   Each agent  $i$  records  $\{(s_{N_i^{\kappa_c}}(k), a_i(k), r_i(k))\}_{0 \leq k \leq K-1}$ .
5:   Each agent  $i$  estimates  $\kappa_c$ -truncated averaged  $Q$ -function
   by linear function approximation, with weight vector  $w_i^m$ .
6:   for  $t = 0, 1, \dots, T-1$  do
7:     Use  $\xi^{\theta(m)}$  to collect samples for  $H$  iterations:
     Each agent  $i$  records  $\{(s_{N_i^{\kappa_c}}^t(k), a_i^t(k))\}_{0 \leq k \leq H-1}$ .
8:     Every agent  $i$  uses its own record and  $w_i^m$  to obtain
      $\Delta_i^T(m) \approx \nabla_i J_i(\theta(m))$ .
9:     Independent Policy Gradient: Every agent  $i$  performs
      $\theta_i(m+1) = \theta_i(m) + \beta \Delta_i^T(m)$ .
10:   end for
11: end for

```

Truncated averaged Q -function:

- Can be approximated using information in κ_c -hop neighborhood.
- Truncation error decays exponentially with κ_c .

$$\begin{array}{ccc} \text{Averaged } Q\text{-function} & \xrightarrow{\text{Use } \kappa_c\text{-hop information}} & \kappa_c\text{-truncated averaged } Q\text{-function} \\ \bar{Q}_i(s, a_i) & \xrightarrow{\text{Use } \kappa_c\text{-hop information}} & \bar{Q}_i^{\kappa_c}(s_{N_i^{\kappa_c}}, a_i) \end{array}$$

$$\xrightarrow{\text{Further reduce parameter dimension}} \hat{Q}_i^{\kappa_c}(s_{N_i^{\kappa_c}}, a_i; w_i)$$

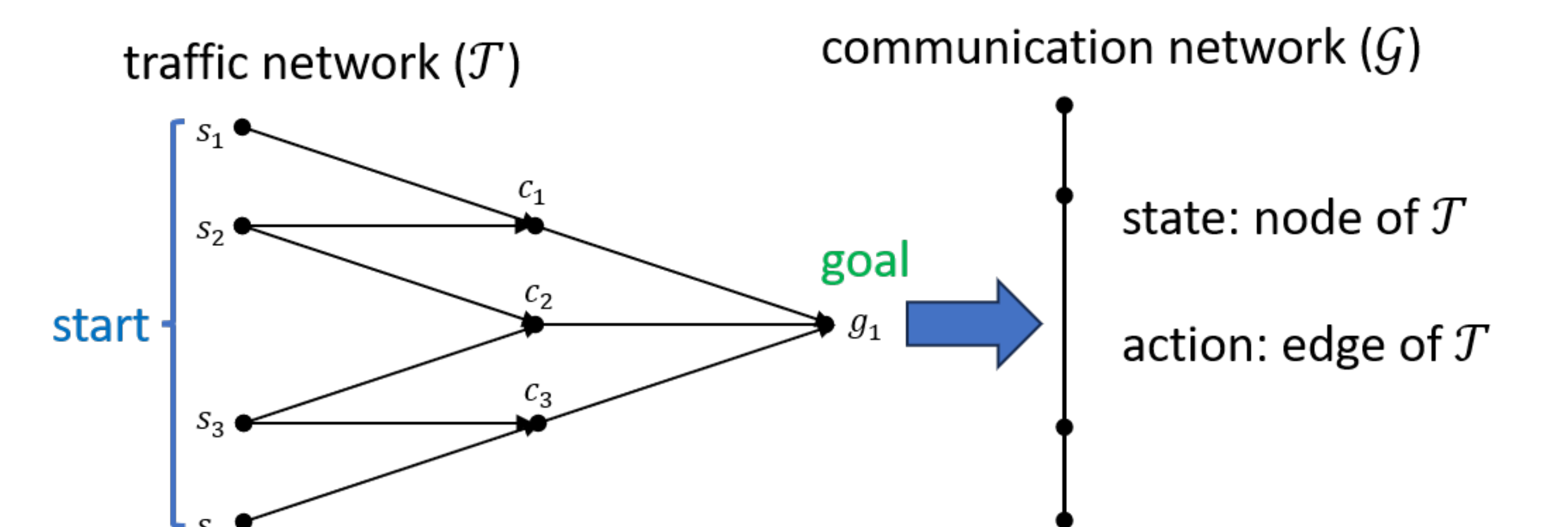
Main Results

We prove that Localized Actor-Critic has $\tilde{O}(\tilde{\epsilon}^{-4})$ sample complexity.

Thm 1 (Informal). With $\tilde{O}(\tilde{\epsilon}^{-4})$ samples, the expectation of averaged Nash regret is upper bounded by $\tilde{\epsilon}$ plus following additional errors:

- 1) Localization error, decaying exponentially with hyperparameter κ_c ;
- 2) Exploration error, depending on exploration coefficient ϵ ;
- 3) Function approximation error.

Example: Markov Congestion Game



Each agent tries to commute from its start node to its destination. Two agents share an edge in the communication network iff they can reach the same node in traffic network.



Take a picture to download the full paper