

# Learning Fair Policies in Decentralized Cooperative Multi-Agent Reinforcement Learning

Matthieu Zimmer\*, Claire Glanois\*, Umer Siddique and Paul Weng

Univ. of Michigan-Shanghai Jiao Tong Univ. Joint Institute

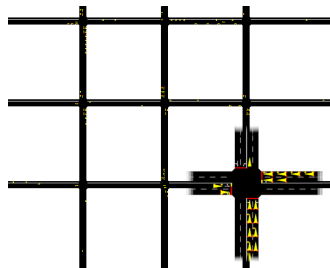
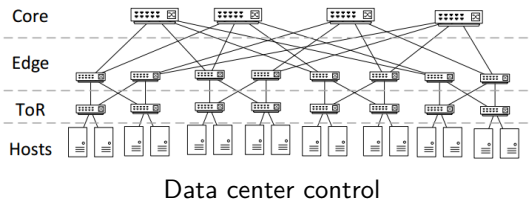
ICML - July 2021



\*Equal contribution

# Motivations

## Examples



## Framework

Dec-POMDP  $(\mathcal{S}, \mathcal{A} = (\mathcal{A}_i)_{i \in [M]}, (\mathcal{O}_i)_{i \in [M]}, P, (\Omega_i)_{i \in [M]}, \mathbf{r}, \gamma)$

- ▶ vectorial reward function  $\mathbf{r} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^D$  where  $D$  is the number of users
- ▶ partial reward observability: an agent  $i$  observes  $\mathbf{r}_{l_i} = (r_k)_{k \in l_i}$  where  $l_i \subseteq [D]$

# Propositions

## Optimizing Social Welfare Functions

$$\max_{\theta} \mathfrak{J}(\theta) = \max_{\theta} \phi(\mathbf{J}_{I_1}(\theta_1), \dots, \mathbf{J}_{I_N}(\theta_N))$$

where  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}$  is a social welfare function, such as:

- ▶ Generalized Gini social welfare function:  $G_{\mathbf{w}}(\mathbf{u}) = \sum_{k \in [D]} w_k u_k^{\uparrow}$
- ▶  $\alpha$ -fairness:  $\phi_{\alpha}(\mathbf{u}) = \sum_{k \in [D]} \frac{u_k^{1-\alpha}}{1-\alpha}$

## Fairness properties

- ▶ Impartiality, Efficiency, Pigou-Dalton principle

## Advantage sharing

- ▶ No need for a centralized critic
- ▶ Less communication needed

# Theoretical Analysis

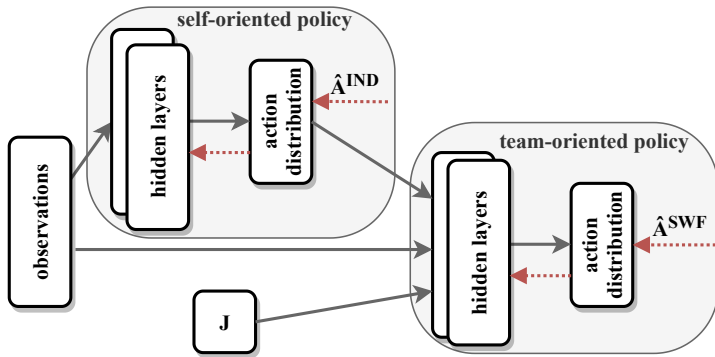
## Theorem

*Under standard assumptions, the SWF objective  $\mathfrak{J}(\boldsymbol{\theta}^k)$  converges almost surely and with a sub-linear convergence rate within a radius of convergence  $\tilde{\epsilon}$  of the optimal value  $\mathfrak{J}^*$  where  $\tilde{\epsilon}$  depend on the approximation errors of (a) estimating  $\mathbf{J}$ , (b) estimating  $\mathbf{A}(\mathbf{o}, \mathbf{a})$ , and (c) ignoring the effects of one agent's action over other agents.*

- ▶ Corollary providing a high-probability bound on the number of iterations before convergence
- ▶ Reducing (b) by learning two critics per agents

# Self-Oriented Team-Oriented (SOTO) Architecture

- ▶ Transfer learning with advice taking: the self-oriented policy advises the team-oriented policy
- ▶ Learning from two losses from two critics
- ▶ Progressively switch from the self-oriented policy to the team-oriented one



# Conclusion

## Fair optimization in multi-agent reinforcement learning

- ▶ Scalable (no centralized critic nor centralized policy)
- ▶ Evaluation on two scenarios and 5 domains
  - ▶ Centralized learning with decentralized execution
  - ▶ Fully decentralized
- ▶ Convergence proof

## Future directions

- ▶ Learning the communications
- ▶ Relaxation of impartiality