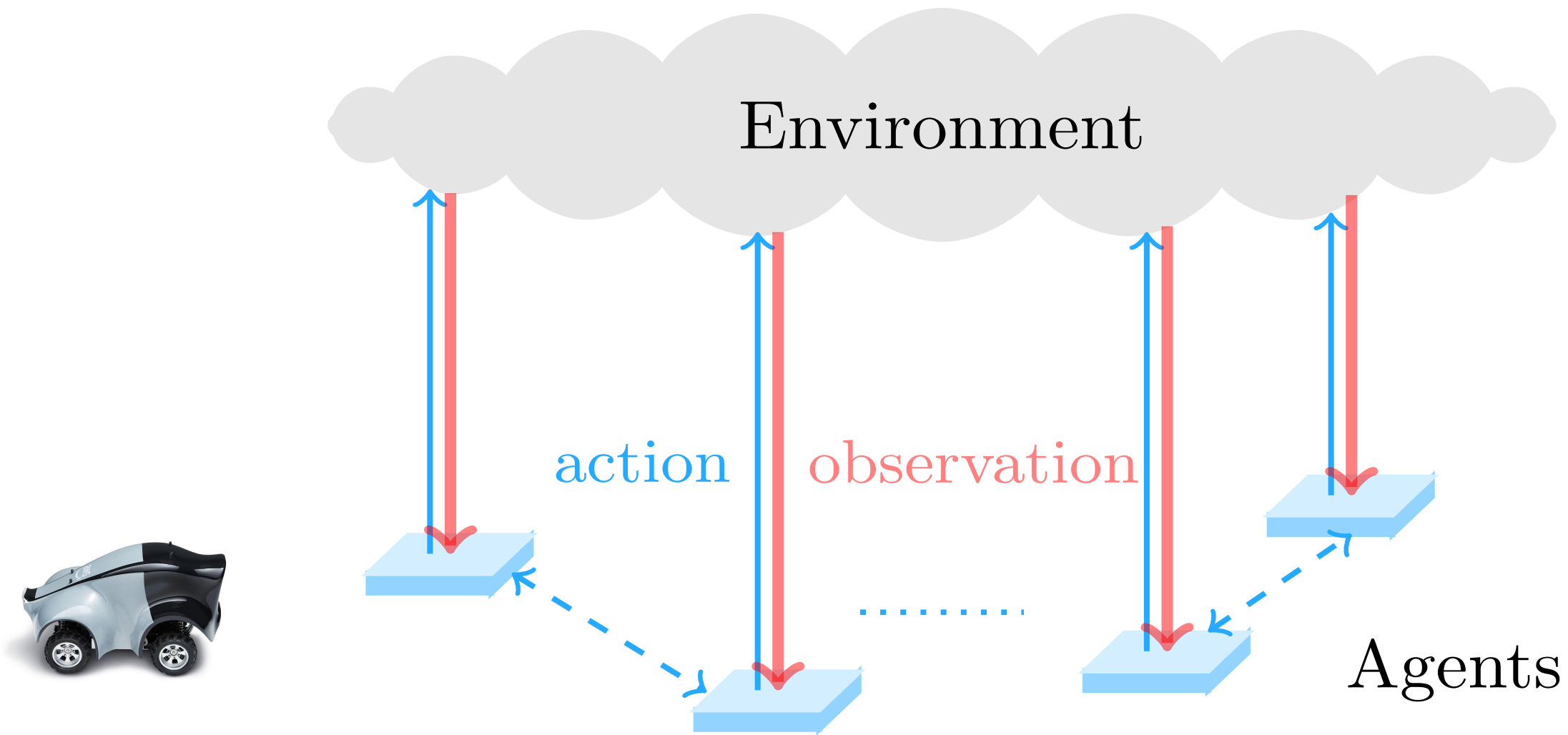


# Provably Efficient Generalized Lagrangian Policy Optimization for Safe MARL

Dongsheng Ding (Penn), Xiaohan Wei (Meta), Zhuoran Yang (Yale),  
Zhaoran Wang (Northwestern), Mihailo R. Jovanović (USC)

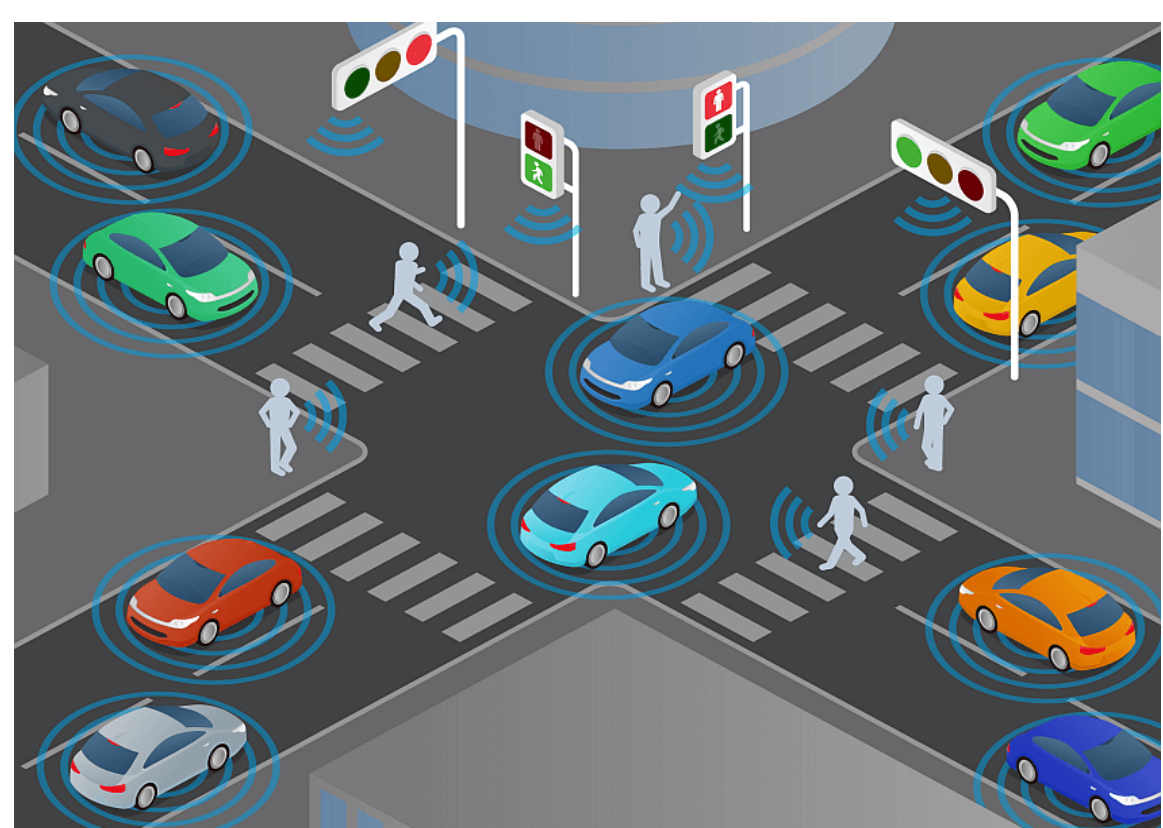
## MOTIVATION

### Multi-agent sequential decision making



**Trade-off** {reward, profit, ...} vs. {safety, budget, fairness, ...}

### Constraint-rich multi-agent systems



automated vehicles



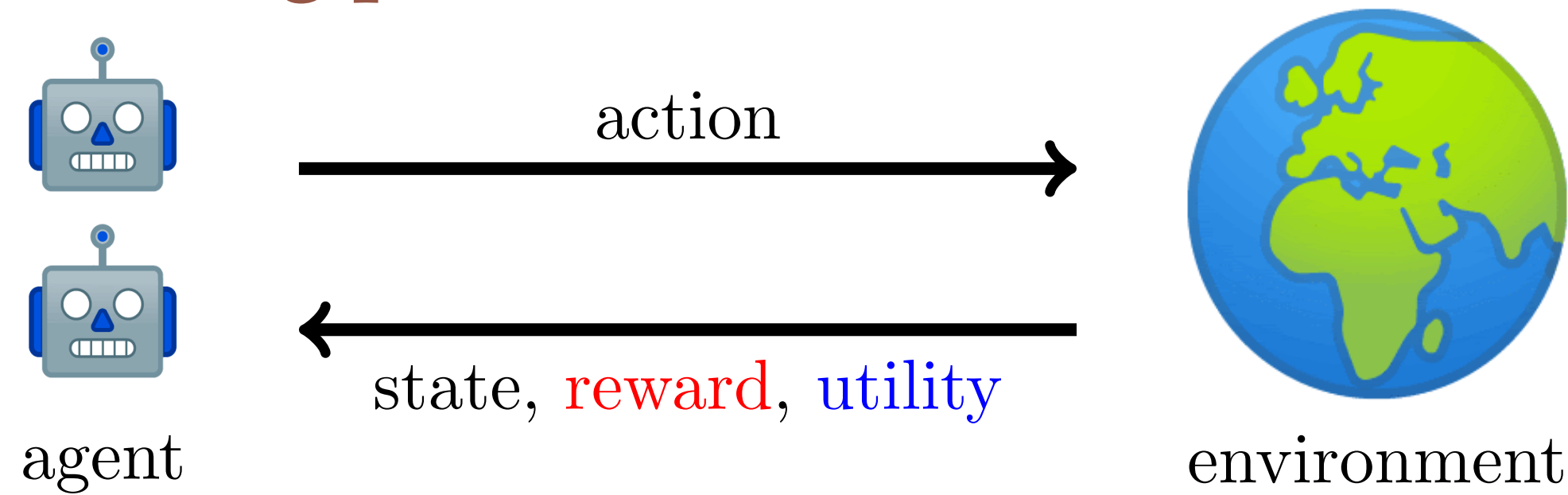
satellite communication

### Challenges

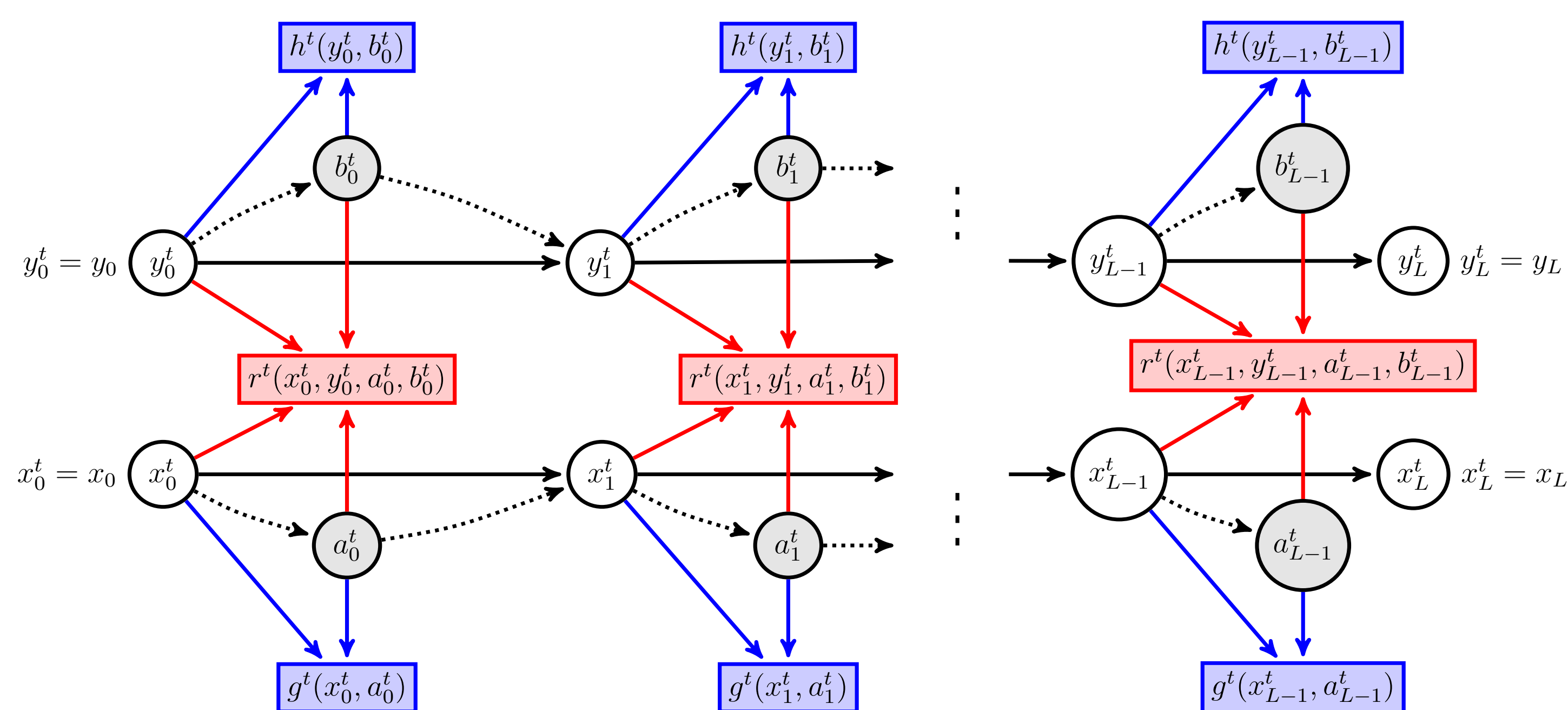
- safe exploration
- efficiency

## PROBLEM FORMULATION

### Episodic learning protocol



- $\ell = 0, \dots, L$  – horizon
- $t = 0, \dots, T-1$  – episode
- $(x_\ell^t, y_\ell^t), (a_\ell^t, b_\ell^t), r_\ell^t, (g_\ell^t, h_\ell^t)$  – state, action, reward, utilities
- $a_\ell^t \sim \pi^t(\cdot | x_\ell^t), b_\ell^t \sim \mu^t(\cdot | y_\ell^t)$  – policies
- $x_{\ell+1}^t \sim P_1(\cdot | x_\ell^t, a_\ell^t), y_{\ell+1}^t \sim P_2(\cdot | y_\ell^t, b_\ell^t)$  – independent dynamics



- $\langle q_1^t \cdot q_2^t, r^t \rangle := \mathbb{E}[\sum_{\ell=0}^{L-1} r^t(x_\ell, a_\ell, a_\ell, b_\ell)]$  – reward value
- $\langle q_1^t, g^t \rangle := \mathbb{E}[\sum_{\ell=0}^{L-1} g^t(x_\ell, a_\ell)]$  – utility value; also for  $\langle q_2^t, h^t \rangle$

### Constrained zero-sum Markov game

$$\begin{aligned} & \text{maximize}_{q_1 \in \Delta(P_1)} & \text{minimize}_{q_2 \in \Delta(P_2)} & \sum_{t=0}^{T-1} \langle q_1 \cdot q_2, r^t \rangle \\ & \text{subject to} & & \langle q_1, g \rangle + \langle q_2, h \rangle \leq b \end{aligned}$$

- $r^t$  – adversarial
- $g, h$  – expectations of stochastic  $g^t, h^t$
- $(q_1^*, q_2^*)$  – constrained Nash equilibrium

## PERFORMANCE MEASURE

$$\text{Regret}(K) := \sum_{t=0}^{T-1} (\langle q_1^t \cdot q_2^*, r^t \rangle - \langle q_1^* \cdot q_2^t, r^t \rangle)$$

$$\text{Violation}(K) := \sum_{k=1}^K (\langle q_1^k, g^k \rangle + \langle q_2^k, h^k \rangle - b)$$

- $q_1^t / q_2^t$  – occupancy measures induced by policies  $\pi^t / \mu^t$

## ALGORITHM DESIGN

### One-episode constrained minimax problem

$$\begin{aligned} & \text{maximize}_{q_1 \in \Delta(P_1)} & \text{minimize}_{q_2 \in \Delta(P_2)} & \langle q_1 \cdot q_2, r^{t-1} \rangle \\ & \text{subject to} & & \langle q_1, g^{t-1} \rangle + \langle q_2, h^{t-1} \rangle \leq b \end{aligned}$$

- $L^t(q_1, q_2; \lambda)$  – generalized Lagrangian
- =  $\langle q_1 \cdot q_2, r^{t-1} \rangle$  minimax obj.
- +  $\lambda (\langle q_1, g^{t-1} \rangle + \langle q_2, h^{t-1} \rangle - b)$  min's vio.
- $\lambda (\langle \hat{q}_1^t, g^{t-1} \rangle + \langle q_2, h^{t-1} \rangle - b)$  max's vio.

### Online mirror descent primal-dual step

$$\hat{q}^t = \arg \min_{q_1 \in \hat{\Delta}_1} \arg \max_{q_2 \in \hat{\Delta}_2} L^t(q_1, q_2, \lambda^{t-1}) + \frac{1}{\eta} \underbrace{D_{\text{KL}}(q, \hat{q}^{t-1})}_{\text{KL regularization}}$$

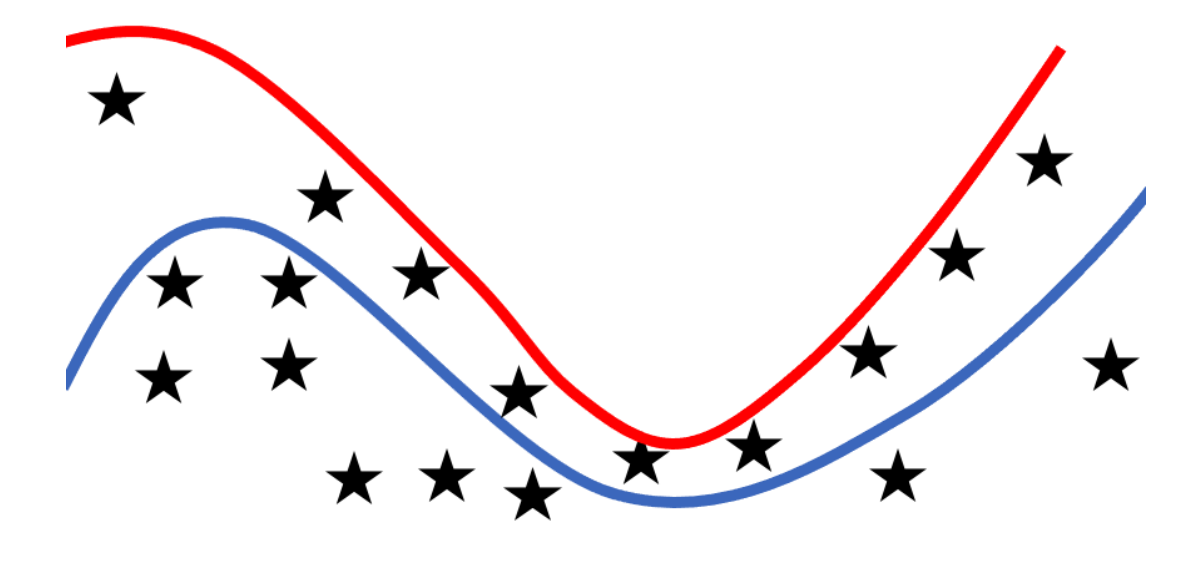
$$\lambda^t = \max \left( \lambda^{t-1} + \underbrace{(\langle \hat{q}_1^t, g^{t-1} \rangle + \langle \hat{q}_2^t, h^{t-1} \rangle - b)}_{\text{violations}}, 0 \right)$$

- compete for rewards
- cooperate for constraints

### Optimistic estimation of $\hat{\Delta}_1, \hat{\Delta}_2$

$$\begin{aligned} \hat{P}_1 & \leftarrow \bar{P}_1 + \text{UCB}_1 \\ \hat{P}_2 & \leftarrow \bar{P}_2 + \text{UCB}_2 \end{aligned}$$

exploit      explore

$$\hat{\Delta}_i = \text{Linear constraint}(\bar{P}_i, \text{UCB}_i)$$


## THEORETICAL GUARANTEE

### Constrained Markov games with independent dynamics

$$\text{Regret}(K), \text{Violation}(K) = \tilde{O}((|X| + |Y|)L \sqrt{T(|A| + |B|)})$$

- $T$  – # episodes;
- $L$  – horizon length
- $|X| + |Y|, |A| + |B|$  – state/action space sizes
- no sampling assumptions & adversarial reward function
- applicable to side constraint case and single-controller case

## REFERENCE

- [1] D. Ding, X. Wei, Z. Yang, Z. Wang, M. Jovanovic, "Provably Efficient Generalized Lagrangian Policy Optimization for Safe Multi-Agent Reinforcement Learning", arXiv:2306.00212 (a long version with appendices).