

# Randomized Exploration in Cooperative Multi-Agent Reinforcement Learning

Hao-Lun Hsu\*, [Weixin Wang\\*](#), Miroslav Pajic, Pan Xu

Department of Electrical and Computer Engineering



- **Thompson Sampling (TS)**
  - Outperform Upper Confidence Bound (UCB) empirically
  - **Not easily scalable to large environments** (multi-agent scenarios)
- **Randomized Exploration**
  - Effective in bandit and single-agent RL
  - **Remain underexplored** in Cooperative MARL

## Parallel MDPs

- $M$  agents interact independently with their respective MDPs
- Share the **same but independent state and action spaces**
- Each agent **might have its unique reward functions and transition kernels**
- Agents and server can communicate to share data

# Unified Algorithm Framework

## Algorithm Unified Algorithm Framework for Randomized Exploration in Parallel MDPs

```
1: for episode  $k = 1, \dots, K$  do
2:   for agent  $m \in \mathcal{M}$  do
3:     Receive initial state  $s_{m,1}^k$  and  $V_{m,H+1}^k(\cdot) \leftarrow 0$ .
4:      $\{Q_{m,h}^k(\cdot, \cdot), V_{m,h}^k(\cdot, \cdot)\}_{h=1}^H \leftarrow$  Randomized Exploration ◁ PHE or LMC
5:     for step  $h = 1, \dots, H$  do
6:        $a_{m,h}^k \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} Q_{m,h}^k(s_{m,h}^k, a)$ .
7:       Receive  $s_{m,h+1}^k$  and  $r_{m,h}$ , then update local data.
8:       if Condition then
9:         SYNCHRONIZE  $\leftarrow$  True.
10:      end if
11:    end for
12:  end for
13:  if SYNCHRONIZE then
14:    All the agents upload their newly collected local data to the server.  $U_{m,h}^{\text{loc}}(k)$ 
15:    The server gathers all information and sends it back to each agent.  $U_h^{\text{ser}}(k)$ 
16:  end if
17: end for
```

# Synchronization Conditions

## Synchronization Conditions

- 1) synchronize at a constant frequency
- 2) synchronize at an exponential frequency
- 3) synchronize when the following **feature mapping  $\phi(\cdot, \cdot)$ -based** condition is satisfied

$$\log \frac{\det(\text{ser} \mathbf{\Lambda}_h^k + \text{loc} \mathbf{\Lambda}_{m,h}^k + \lambda \mathbf{I})}{\det(\text{ser} \mathbf{\Lambda}_h^k + \lambda \mathbf{I})} \geq \frac{\gamma}{(k - k_s)},$$

where  $\phi(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ ,  $\text{ser} \mathbf{\Lambda}_h^k = \sum_{U_h^{\text{ser}}(k)} \phi(s^l, a^l) \phi(s^l, a^l)^\top$ ,  
 $\text{loc} \mathbf{\Lambda}_{m,h}^k = \sum_{U_{m,h}^{\text{loc}}(k)} \phi(s^l, a^l) \phi(s^l, a^l)^\top$ .

# Perturbed-History Exploration

- Regression loss with added random Gaussian noises  $\epsilon_h^{k,l,n}$  and  $\xi_h^{k,n}$  to **perturb reward and regularizer**

$$\tilde{L}_{m,h}^{k,n}(\mathbf{w}) = \sum_{l=1}^{\mathcal{K}(k)} L((r_h^l + \epsilon_h^{k,l,n}) + V_{m,h+1}^k(s^l), f(\mathbf{w}; \phi^l)) + \lambda \|\mathbf{w} + \xi_h^{k,n}\|^2.$$

- Unified Algorithm Framework + PHE  $\Rightarrow$  **CoopTS-PHE**

---

## Algorithm Perturbed-History Exploration

---

- 1: **for** step  $h = H, \dots, 1$  **do**
  - 2:   **for**  $n = 1, \dots, N$  **do**
  - 3:     Sample  $\{\epsilon_h^{k,l,n}\}_{l \in [\mathcal{K}(k)]} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$  and  $\xi_h^{k,n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$  independently.
  - 4:     Obtain the perturbed estimated parameter  $\tilde{\mathbf{w}}_{m,h}^{k,n} = \operatorname{argmin} \tilde{L}_{m,h}^{k,n}(\mathbf{w})$ .
  - 5:   **end for**
  - 6:    $Q_{m,h}^k \leftarrow \min \{ \max_{n \in [N]} f(\tilde{\mathbf{w}}_{m,h}^{k,n}; \phi), H - h + 1 \}^+$ .
  - 7:    $V_{m,h}^k(\cdot) \leftarrow \max_{a \in \mathcal{A}} Q_{m,h}^k(\cdot, a)$ .
  - 8: **end for**
  - 9: Output:  $\{Q_{m,h}^k(\cdot, \cdot), V_{m,h}^k(\cdot, \cdot)\}_{h=1}^H$ .
-

# Langevin Monte Carlo Exploration

- **Langevin Monte Carlo update:** for iterate  $j = 1, \dots, J_k$ , the update is given by

$$\mathbf{w}_{m,h}^{k,j,n} = \mathbf{w}_{m,h}^{k,j-1,n} - \eta_{m,k} \nabla L_{m,h}^k(\mathbf{w}_{m,h}^{k,j-1,n}) + \sqrt{2\eta_{m,k}\beta_{m,k}^{-1}} \epsilon_{m,h}^{k,j,n}.$$

- Unified Algorithm Framework + LMC  $\Rightarrow$  **CoopTS-LMC**

---

## Algorithm Langevin Monte Carlo Exploration

---

- 1: **for** step  $h = H, \dots, 1$  **do**
- 2:   **for**  $n = 1, \dots, N$  **do**
- 3:      $\mathbf{w}_{m,h}^{k,0,n} = \mathbf{w}_{m,h}^{k-1, J_{k-1}, n}$ .
- 4:     **for**  $j = 1, \dots, J_k$  **do**
- 5:       Sample  $\epsilon_{m,h}^{k,j,n} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$  and obtain  $\mathbf{w}_{m,h}^{k,j,n}$  through **LMC update**.
- 6:     **end for**
- 7:   **end for**
- 8:    $Q_{m,h}^k \leftarrow \min \{ \max_{n \in [N]} f(\mathbf{w}_{m,h}^{k, J_k, n}; \phi), H - h + 1 \}^+$
- 9:    $V_{m,h}^k(\cdot) \leftarrow \max_{a \in A} Q_{m,h}^k(\cdot, a)$ .
- 10: **end for**

**Linear MDP** (linear reward and transition functions) An MDP  $(\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r)$  is a **linear MDP** with feature map  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ , if for any  $h \in [H]$ , there exist  $d$  unknown measures  $\mu_h = (\mu_h^1, \dots, \mu_h^d)$  over  $\mathcal{S}$  and an unknown vector  $\theta_h \in \mathbb{R}^d$  such that for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\mathbb{P}_h(\cdot | s, a) = \langle \phi(s, a), \mu_h(\cdot) \rangle, \quad r_h(s, a) = \langle \phi(s, a), \theta_h \rangle.$$

**Cumulative Group Regret** The learning goal is to minimize the **cumulative group regret** among  $M$  agents after  $K$  episodes, which is defined as

$$\text{Regret}(K) = \sum_{m \in \mathcal{M}} \sum_{k=1}^K [V_{m,1}^*(s_{m,1}^k) - V_{m,1}^{\pi_m^k}(s_{m,1}^k)].$$

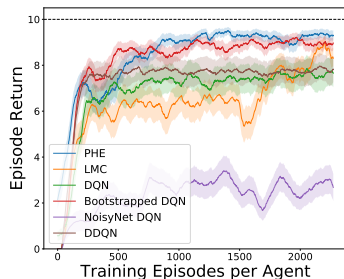


# Theoretical Results

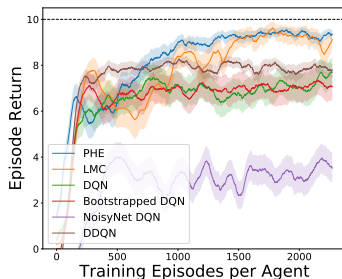
**Table:** Comparison on episodic, non-stationary, linear MDPs. **We define the average regret as the cumulative regret divided by the total number of samples used by the algorithm.** Here  $d$  is the feature dimension,  $H$  is the episode length,  $K$  is the number of episodes, and  $M$  is the number of agents in a multi-agent setting.

Setting	Algorithm	Cumulative Group Regret	Average Regret	Randomized Exploration	Generalizable to Deep RL	Communication Complexity
single-agent	OPT-RLSVI [Zanette et al., 2020]	$\tilde{O}(d^2 H^{5/2} \sqrt{K})$	$\tilde{O}(d^{3/2} H^{3/2} \sqrt{1/K})$	Yes	No	–
	LSVI-UCB [Jin et al., 2020]	$\tilde{O}(d^{3/2} H^2 \sqrt{K})$	$\tilde{O}(d^{3/2} H \sqrt{1/K})$	No	No	–
	LSVI-PHE [Ishfaq et al., 2021]	$\tilde{O}(d^{3/2} H^2 \sqrt{K})$	$\tilde{O}(d^{3/2} H \sqrt{1/K})$	Yes	Yes	–
	LMC-LSVI [Ishfaq et al., 2024]	$\tilde{O}(d^{3/2} H^2 \sqrt{K})$	$\tilde{O}(d^{3/2} H \sqrt{1/K})$	Yes	Yes	–
multi-agent	Coop-LSVI [Dubey & Pentland, 2021]	$\tilde{O}(d^{3/2} H^2 \sqrt{MK})$	$\tilde{O}(d^{3/2} H \sqrt{1/MK})$	No	No	$\tilde{O}(dHM^3)$
	Asyn-LSVI [Min et al., 2023]	$\tilde{O}(d^{3/2} H^2 \sqrt{K})$	$\tilde{O}(d^{3/2} H \sqrt{1/K})$	No	No	$\tilde{O}(dHM^2)$
	<b>CoopTS-PHE (Ours)</b>	$\tilde{O}(d^{3/2} H^2 \sqrt{MK})$	$\tilde{O}(d^{3/2} H \sqrt{1/MK})$	Yes	Yes	$\tilde{O}(dHM^2)$
	<b>CoopTS-LMC (Ours)</b>	$\tilde{O}(d^{3/2} H^2 \sqrt{MK})$	$\tilde{O}(d^{3/2} H \sqrt{1/MK})$	Yes	Yes	$\tilde{O}(dHM^2)$

# N-Chain Experiments



(a)  $m=2$



(b)  $m=3$

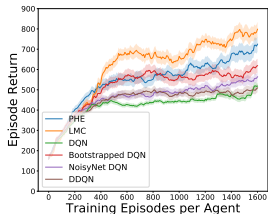
Figure:  $N$ -chain task with  $N = 25$  states and  $m = 2, 3$  agents.

- When  $m = 2$ , **PHE** achieve **higher average returns** and **LMC** eventually catches up. When  $m = 3$ , **PHE** and **LMC** outperform baselines. **PHE** shows less fluctuation, which supports theoretical results in misspecified settings.

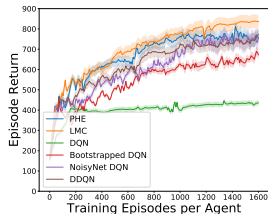
# Super Mario Bro Experiments



(a) Illustration



(b) Parallel MDP

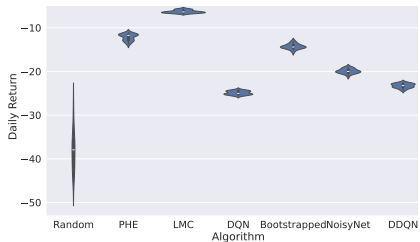


(c) Federated Learning

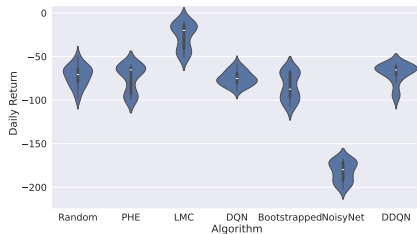
**Figure:** Super Mario Bros task with  $m = 4$  agents: (a) illustration; (b) parallel MDP (**communicate whole data**); (c) federated learning (**only communicate value function**).

- **PHE** and **LMC** outperform baselines in parallel MDP and federated learning settings. **LMC** consistently outperforms **PHE**, as the added noise in **PHE** may not always accurately reflect true posterior in practice.

# Real-world Experiments in Thermal Control



(a) Tampa (hot humid)



(b) Great Falls (cold dry)

**Figure:** Evaluation performance at Tampa and Great Falls in building energy systems.

- Experiments, trained with parallel data sharing across cities and varying weather, **aim to meet temperature specifications while minimizing electricity use**. **LMC** shows **higher mean returns** in the violin plots.

# Summary of Main Contributions

- A unified framework for parallel MDPs + TS-related exploration strategies **PHE** and **LMC** → **CoopTS-PHE** and **CoopTS-LMC**.
  - **CoopTS-PHE**: perturb reward and regularizer (equivalent to TS)
  - **CoopTS-LMC**: perform noisy gradient descent (converge to TS)

• **Regret Upper Bound**:  $\tilde{O}(d^{3/2}H^2\sqrt{M}(\sqrt{dM\gamma} + \sqrt{K}))$

• **Communication Complexity**:  $\tilde{O}((d + K/\gamma)MH)$

- Extend to **misspecified settings** (approximately linear reward and transition functions)
- Extensive experiments on various benchmarks
  - **N-chain** (require deep exploration)
  - **Super Mario Bros** (misspecified setting; federated learning)
  - **Thermal control problem in building energy systems**

Outperform existing DQN-based baselines

# END

**Thank you very much for your valuable time!**