# Tiered Reinforcement Learning: Pessimistic in the Face of Uncertainty and Constant Regret

Jiawei Huang, Li Zhao, Tao Qin, Wei Chen, Nan Jiang, Tie-Yan Liu
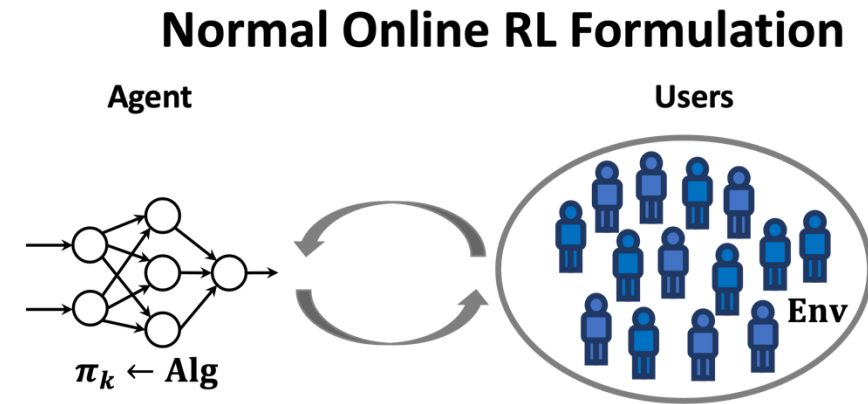
ETH *zürich*

UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN

Microsoft

# Introduction

- RL has been applied in many applications with user interaction:
  - Medical Treatment

  - Recommendation System

  - Other Online Application Services



**Normal Online RL Formulation**

Agent

$\pi_k \leftarrow$ Alg

Users

Env

- The normal learning protocol (Fig. RHS)
  - Repeatedly:
    - **Policy Improvement:** Learn a policy from collected data
    - **Collect New Data from Env:** User comes; generate trajectories during interaction.

# Introduction

> However, users may have **' Tiered Structure'**:
> *The users can be divided into 2 (or more) groups by their different preference on exploration risk.*

- RL has been applied in many applications with user interaction:
  - Medical Treatment
    Paid Volunteers v.s. Conservative Patients
  - Recommendation System
    Paid Tester v.s. General Customers
  - Other Online Application Services
    Customers Using Free Services v.s. Paid VIP Customers



**Normal Online RL Formulation**

Agent

$\pi_k \leftarrow$ **Alg**

Users

Env

User $\in G^O$   User $\in G^E$

$G^O$ (short for Group$^{Online}$): risk-tolerant user group
$G^E$ (short for Group$^{Exploit}$): risk-averse user group
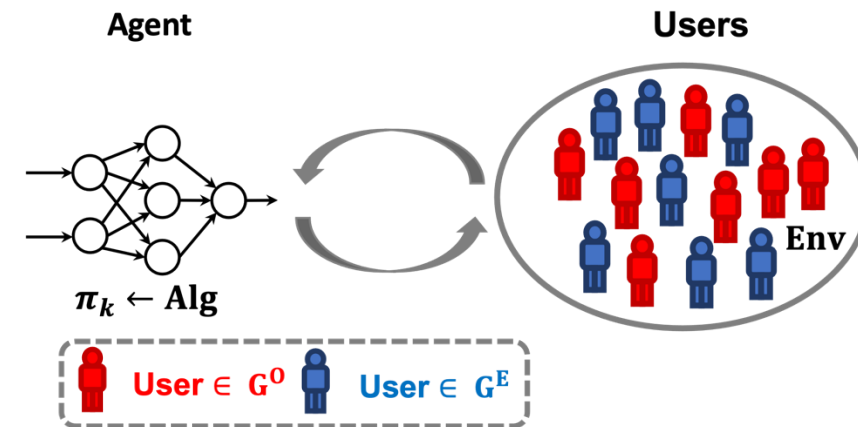
- The normal learning protocol (Fig. RHS)
  - Repeatedly:
    - **Policy Improvement:** Learn a policy from collected data
    - **Collect New Data from Env:** User comes; generate trajectories during interaction.

# Introduction

- RL has been applied in many applications with user interaction:
  - Medical Treatment
    Paid Volunteers v.s. Conservative Patients
  - Recommendation System
    Paid Tester v.s. General Customers
  - Other Online Application Services
    Customers Using Free Services v.s. Paid VIP Customers

**Normal Online RL Formulation**



$G^O$ (short for $Group^{Online}$): risk-tolerant user group
$G^E$ (short for $Group^{Exploit}$): risk-averse user group

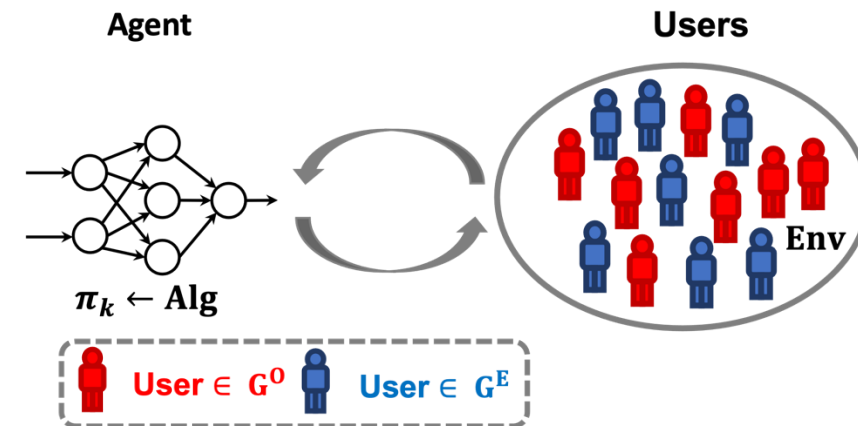- The normal learning protocol (Fig. RHS)
  - Repeatedly:
    - **Policy Improvement:** Learn a policy from collected data
    - **Collect New Data from Env:** User comes; generate trajectories during interaction.
  - Users in different groups will be treated equivalently and suffer similar loss…

# Tiered RL Framework

Initialize $D_1 \leftarrow \{\}$.
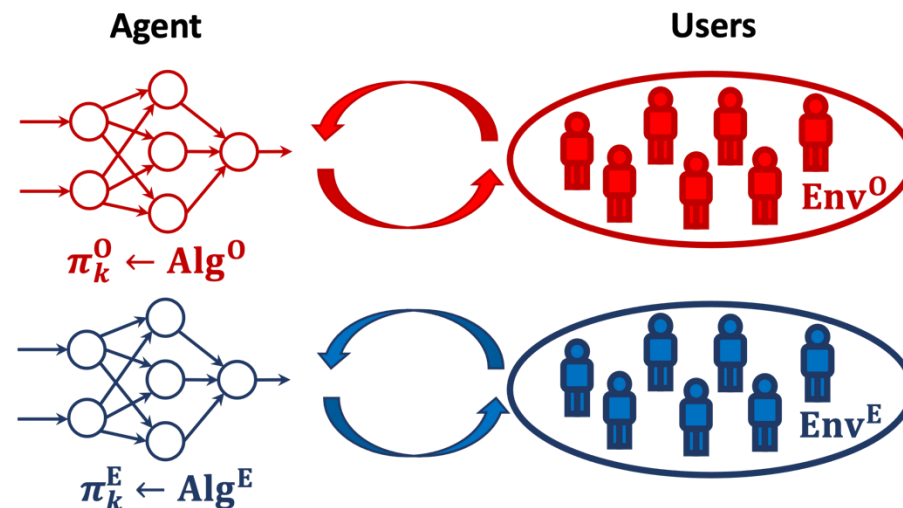**for** $k = 1, 2, ..., K$ **do**

    $\pi_{O,k} \leftarrow \text{Alg}^O(D_k); \pi_{E,k} \leftarrow \text{Alg}^E(D_k)$.

    $\pi_{O,k}/\pi_{E,k}$ interacts with users in exploration/exploitation tier, and collect data $\tau_{O,k}/\tau_{E,k}$.

    $D_{k+1} = D_k \cup \{\tau_{O,k}\} \cap \{\tau_{E,k}\}$.

**end**

## Tiered Online RL Formulation (Ours)



Assume $\text{Env}^O = \text{Env}^E$;
Relaxtion of this assumption leave to future work.

# Tiered RL Framework

Initialize $D_1 \leftarrow \{\}$.
**for** $k = 1, 2, ..., K$ **do**
$\quad \pi_{O,k} \leftarrow \text{Alg}^O(D_k);\ \pi_{E,k} \leftarrow \text{Alg}^E(D_k)$.
$\quad \pi_{O,k}/\pi_{E,k}$ interacts with users in exploration/exploitation tier, and collect data $\tau_{O,k}/\tau_{E,k}$.
$\quad D_{k+1} = D_k \cup \{\tau_{O,k}\} \cap \{\tau_{E,k}\}$.
**end**

- Objective
  - Consider the pseudo-regret $\textbf{Regret}_K(\cdot)$:
    - $\textbf{Regret}_K(\textbf{Alg}^E) := \mathbb{E}\left[\sum_{k=1}^K V^*(s_1) - V^{\pi_k^E}(s_1)\right]$;　　$\textbf{Regret}_K(\textbf{Alg}^O) := \mathbb{E}\left[\sum_{k=1}^K V^*(s_1) - V^{\pi_k^O}(s_1)\right]$.
  - Is it possible for $\textbf{Regret}_K(\textbf{Alg}^E)$ **to be strictly lower than any online learning algorithms** in certain scenarios, while **keeping $\textbf{Regret}_K(\textbf{Alg}^O)$ near-optimal**?

Benefits for $G^E$ under our framework.

Not too much additional cost for $G^O$.

# Highlight of Main Results

Normal Tabular RL Setting

Tabular RL with Strictly Positive Gap
$$\forall h, s, a: \Delta_h(s,a) = 0 \text{ or } \Delta_h(s,a) \geq \Delta_{min} > 0$$
$$\text{where } \Delta_h(s,a) := V_h^*(s) - Q_h^*(s,a)$$

- No benefits by comparing with standard online RL (from minimax optimality perspective)

  - $\min_{\text{Alg}^O, \text{Alg}^E} \max_{\text{MDP}} \text{Regret}(\text{Alg}^E) \geq O(\sqrt{H^3 SAK})$

Minimax lower bound of normal online RL setting

- By choosing:

  - **P**essimistic **V**alue **I**teration (PVI) as $\text{Alg}^E$,

  - Arbitrary online algorithm with near-optimal regret as $\text{Alg}^O$

- Guarantee

  - $\text{Regret}_K(\text{Alg}^O)$ keeps near-optimal.

  - $\text{Regret}_K(\text{Alg}^E)$ is constant/independent of K.

In contrast with $O(\log K)$ lower bound in online setting

# Why Pessimistic Value Iteration?

- Why Pessimism?
  - Key property [Jin et. al., 2021]:
    - The more optimal trajectories occurs in $D_k$, the smaller $V^* - V^{\pi_k^{\text{PVI}}}$ would be

# Why Pessimistic Value Iteration?

- Why Pessimism?
    - Key property [Jin et. al., 2021]:
        - The more optimal trajectories occurs in $D_k$, the smaller $V^* - V^{\pi_k^{\mathrm{PVI}}}$ would be




    - Coincide with the optimality constraint of **Alg$^{\mathbf{O}}$**
        - Low regret of **Alg$^{\mathbf{O}}$** $\Leftrightarrow$ Faster accumulation of optimal trajectories in $D_k$
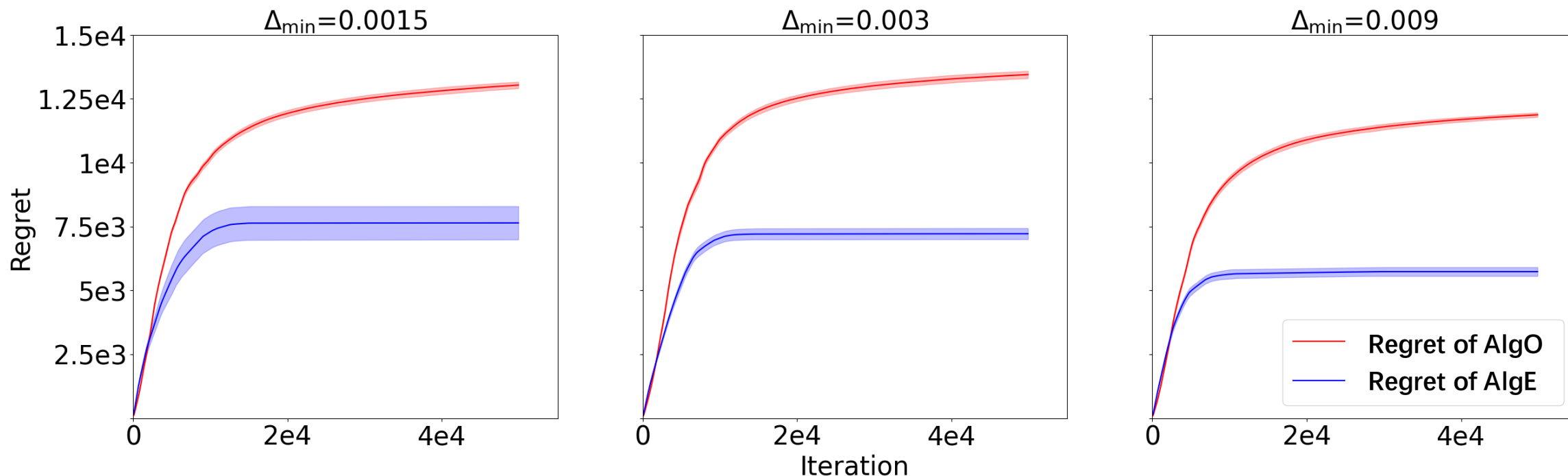
# Why Pessimistic Value Iteration?

- Why Pessimism?
  - Key property [Jin et. al., 2021]:
    - The more optimal trajectories occurs in $D_k$, the smaller $V^* - V^{\pi_k^{\text{PVI}}}$ would be

  - Coincide with the optimality constraint of **Alg$^O$**
    - Low regret of **Alg$^O$** $\Leftrightarrow$ Faster accumulation of optimal trajectories in $D_k$

Thanks to strictly positive gap, $V^* - V^{\pi_k^{\text{PVI}}}$ will be zero after certain steps, which implies constant regret.

# Verification Experiments in Tabular MDP

- S=A=H=5. Random generated transition/reward functions.

- $\mathrm{Alg^O}$: StrongEulder [2]; $\mathrm{Alg^E}$: PVI with Adaptive Bonus Term in [2]



[2] NeurIPS 2019, Max Simchowitz and Kevin Jamieson, Non-Asymptotic Gap-Dependent Regret Bounds for Tabular MDPs.

# Thanks