

An aerial photograph of a city, likely Zurich, showing a river flowing through the center. The city is densely packed with buildings of various architectural styles, including modern glass-fronted structures and older stone buildings. A prominent building with a green dome is visible in the lower-left quadrant. The river is surrounded by greenery and a bridge spans across it. The overall scene is bright and clear, suggesting a sunny day.

Robust Knowledge Transfer in Tiered Reinforcement Learning

Jiawei Huang, Niao He



Department of Computer Science, ETH Zurich

Tiered Reinforcement Learning: A Parallel Knowledge Transfer Framework

- **Tiered RL Setting [1]**

- Target/High-Tier task M_{Hi} + Source/Low-Tier task M_{Lo} learning in parallel
- Knowledge transfer from M_{Lo} to M_{Hi}

- **Scenarios in Practice**

- User Interaction Applications [1]
 - Users with higher risk tolerance: 
 - Users with lower risk tolerance: 

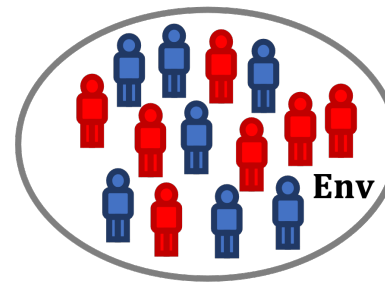


Figure from [1]

- Robotics
 - Multiple robots learning in parallel
 - Some are more vulnerable than others



Figure from [2]

Tiered Reinforcement Learning: A Parallel Knowledge Transfer Framework

- **Main Objective in Tiered RL Setting**

- $\text{Regret}(M_{Lo})$: always near-optimal regret



Source tasks are also important in many cases

- $\text{Regret}(M_{Hi})$:

- **If tasks are similar:** better than optimal regret;



No negative transfer

- **Otherwise:** keep near-optimal

Tiered Reinforcement Learning: A Parallel Knowledge Transfer Framework

- **Main Objective in Tiered RL Setting**

- $\text{Regret}(M_{Lo})$: always near-optimal regret



Source tasks are also important in many cases

- $\text{Regret}(M_{Hi})$:

- **If tasks are similar**: better than optimal regret;



No negative transfer

- **Otherwise**: keep near-optimal

- **Limitation of Existing Knowledge Transfer Frameworks**

	Transfer RL	Multi-Task RL	Parallel Transfer RL (ours; [1])
Guarantees on low-tier/source task?	✗	✓	✓
Tasks learning in parallel?	✗	✓	✓
Distinguish high-tier/target and low-tier/source tasks?	✓	✗	✓

Tiered Reinforcement Learning: A Parallel Knowledge Transfer Framework

- **Main Objective in Tiered RL Setting**

- $\text{Regret}(M_{Lo})$: always near-optimal regret



Source tasks are also important in many cases

- $\text{Regret}(M_{Hi})$:

- **If tasks are similar**: better than optimal regret;



No negative transfer

- **Otherwise**: keep near-optimal

- **Limitation of Existing Knowledge Transfer Frameworks**

	Transfer RL	Multi-Task RL	Parallel Transfer RL (ours; [1])
Guarantees on low-tier/source task?	✗	✓	✓
Tasks learning in parallel?	✗	✓	✓
Distinguish high-tier/target and low-tier/source tasks?	✓	✗	✓

- **Limitation of Existing Tiered RL Literature [1]**

- Strong prior knowledge: $M_{Hi} = M_{Lo}$

Tiered Reinforcement Learning: A Parallel Knowledge Transfer Framework

- **Setting**

- Tabular MDP with finite horizon H
- M_{Hi} shares state-action space with M_{Lo}
- No prior knowledge about similarity between M_{Hi} and M_{Lo}

Tiered Reinforcement Learning: A Parallel Knowledge Transfer Framework

- **Setting**

- Tabular MDP with finite horizon H
- M_{Hi} shares state-action space with M_{Lo}
- No prior knowledge about similarity between M_{Hi} and M_{Lo}

- **Main Assumption**

- Optimal Value Dominance:
 - $\forall h, s_h, V_{Lo}^*(s_h) \geq V_{Hi}^*(s_h)$
 - Similar assumptions in [3,4]
 - Theorem 3.1 [Lower bound]: negative transfer is unavoidable if violated

Main Results

- **Bandit & RL Setting with Single Source Task**

- $\text{Regret}(M_{\text{Hi}}, K) = O\left(SH \sum_h \sum_{s_h, a_h \notin \text{Transferable}(M_{\text{Hi}}, M_{\text{Lo}})} \frac{1}{\Delta(s_h, a_h) \vee \frac{\Delta_{\min}}{H}} \log K\right)$

Main Results

- **Bandit & RL Setting with Single Source Task**

- $\text{Regret}(M_{\text{Hi}}, K) = O(SH \sum_h \sum_{s_h, a_h \notin \text{Transferable}(M_{\text{Hi}}, M_{\text{Lo}})} \frac{1}{\Delta(s_h, a_h) \vee \frac{\Delta_{\min}}{H}} \log K)$

Transferable states in single source task setting:

$$\begin{aligned} d_{\text{Lo}}^*(s_h) &> 0; \\ V_{\text{Lo}}^*(s_h) &\leq V_{\text{Hi}}^*(s_h) + O\left(\frac{\tilde{\Delta}}{H}\right); \\ \pi_{\text{Lo}}^*(s_h) &= \pi_{\text{Hi}}^*(s_h) \end{aligned}$$

Main Results

- **Bandit & RL Setting with Single Source Task**

- $\text{Regret}(M_{\text{Hi}}, K) = O\left(SH \sum_h \sum_{s_h, a_h \notin \text{Transferable}(M_{\text{Hi}}, M_{\text{Lo}})} \frac{1}{\Delta(s_h, a_h) \vee \frac{\Delta_{\min}}{H}} \log K\right)$

Transferable states in single source task setting:

$$\begin{aligned} d_{\text{Lo}}^*(s_h) &> 0; \\ V_{\text{Lo}}^*(s_h) &\leq V_{\text{Hi}}^*(s_h) + O\left(\frac{\tilde{\Delta}}{H}\right); \\ \pi_{\text{Lo}}^*(s_h) &= \pi_{\text{Hi}}^*(s_h) \end{aligned}$$

- **Bandit & RL Setting with Multiple Source Tasks**

- W -Source Tasks: $M_{\text{Lo}}^1, \dots, M_{\text{Lo}}^W$

- $\text{Regret}(M_{\text{Hi}}, K) = O\left(SH \sum_h \sum_{s_h, a_h \notin \text{Transferable}(M_{\text{Hi}}, M_{\text{Lo}}^1, \dots, M_{\text{Lo}}^W)} \frac{1}{\Delta(s_h, a_h) \vee \frac{\Delta_{\min}}{H}} \log WK\right)$

Transferable states in single source task setting:

$$\begin{aligned} d_{\text{Lo}}^*(s_h) &> 0; \\ \exists w \in [W] V_{\text{Lo},w}^*(s_h) &\leq V_{\text{Hi}}^*(s_h) + O\left(\frac{\tilde{\Delta}}{H}\right), \\ \pi_{\text{Lo},w}^*(s_h) &= \pi_{\text{Hi}}^*(s_h) \end{aligned}$$

Overview of Robust Knowledge Transfer Mechanism

- **Single Source Task Setting:**

- Key idea: separation between transferable & non-transferable states

- If transferable: $V_{Lo}^*(s_h) \leq V_{Hi}^*(s_h) + O\left(\frac{\tilde{\Delta}}{H}\right) = Q_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right)$

- Otherwise: $V_{Lo}^*(s_h) \geq V_{Hi}^*(s_h) \geq Q_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right) + O\left(\frac{H-1}{H} \Delta_{Hi}(s_h, \pi_{Lo}^*)\right)$

Overview of Robust Knowledge Transfer Mechanism

- **Single Source Task Setting:**

- Key idea: separation between transferable & non-transferable states

- If transferable: $V_{Lo}^*(s_h) \leq V_{Hi}^*(s_h) + O\left(\frac{\tilde{\Delta}}{H}\right) = Q_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right)$
- Otherwise: $V_{Lo}^*(s_h) \geq V_{Hi}^*(s_h) \geq Q_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right) + O\left(\frac{H-1}{H} \Delta_{Hi}(s_h, \pi_{Lo}^*)\right)$
- Checking condition:
 - $Q_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right) \geq V_{Lo}^*(s_h)$

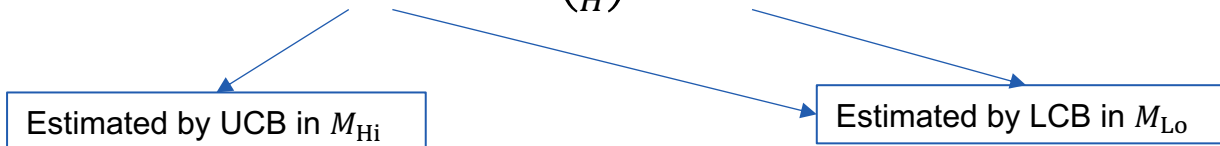
Overview of Robust Knowledge Transfer Mechanism

- **Single Source Task Setting:**

- Key idea: separation between transferable & non-transferable states

- If transferable: $V_{L_0}^*(s_h) \leq V_{H_i}^*(s_h) + O\left(\frac{\tilde{\Delta}}{H}\right) = Q_{H_i}^*(s_h, \pi_{L_0}^*) + O\left(\frac{\tilde{\Delta}}{H}\right)$
- Otherwise: $V_{L_0}^*(s_h) \geq V_{H_i}^*(s_h) \geq Q_{H_i}^*(s_h, \pi_{L_0}^*) + O\left(\frac{\tilde{\Delta}}{H}\right) + O\left(\frac{H-1}{H} \Delta_{H_i}(s_h, \pi_{L_0}^*)\right)$
- Checking condition:

- $\bar{Q}_{H_i}^*(s_h, \pi_{L_0}^*) + O\left(\frac{\tilde{\Delta}}{H}\right) \geq \underline{V}_{L_0}^*(s_h)$



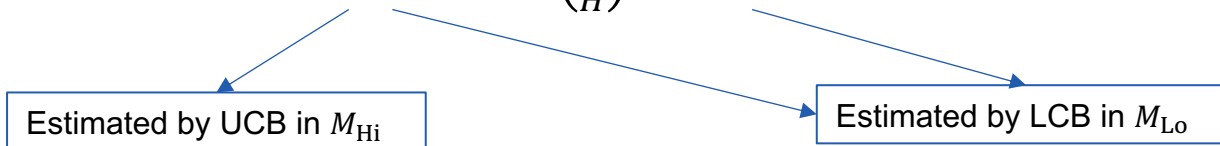
Overview of Robust Knowledge Transfer Mechanism

- **Single Source Task Setting:**

- Key idea: separation between transferable & non-transferable states

- If transferable: $V_{Lo}^*(s_h) \leq V_{Hi}^*(s_h) + O\left(\frac{\tilde{\Delta}}{H}\right) = Q_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right)$
- Otherwise: $V_{Lo}^*(s_h) \geq V_{Hi}^*(s_h) \geq Q_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right) + O\left(\frac{H-1}{H} \Delta_{Hi}(s_h, \pi_{Lo}^*)\right)$
- Checking condition:

- $\bar{Q}_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right) \geq \underline{V}_{Lo}^*(s_h)$



- Avoid negative transfer
 - Every negative transfer will result in tighter estimation of Q_{Hi}^* and V_{Lo}^*

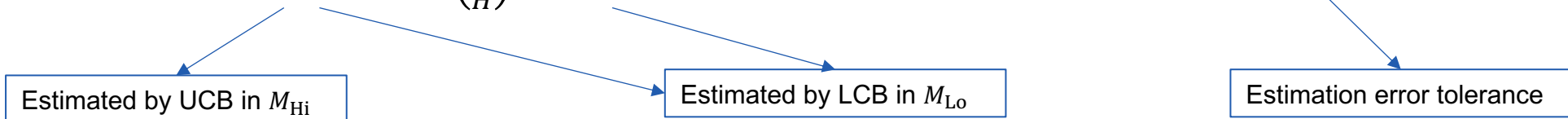
Overview of Robust Knowledge Transfer Mechanism

- **Single Source Task Setting:**

- Key idea: separation between transferable & non-transferable states

- If transferable: $V_{Lo}^*(s_h) \leq V_{Hi}^*(s_h) + O\left(\frac{\tilde{\Delta}}{H}\right) = Q_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right)$
- Otherwise: $V_{Lo}^*(s_h) \geq V_{Hi}^*(s_h) \geq Q_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right) + O\left(\frac{H-1}{H} \Delta_{Hi}(s_h, \pi_{Lo}^*)\right)$
- Checking condition:

- $\bar{Q}_{Hi}^*(s_h, \pi_{Lo}^*) + O\left(\frac{\tilde{\Delta}}{H}\right) \geq \underline{V}_{Lo}^*(s_h)$



- Avoid negative transfer
 - Every negative transfer will result in tighter estimation of Q_{Hi}^* and V_{Lo}^*

Overview of Robust Knowledge Transfer Mechanism

- **Multiple Source Tasks Setting:**

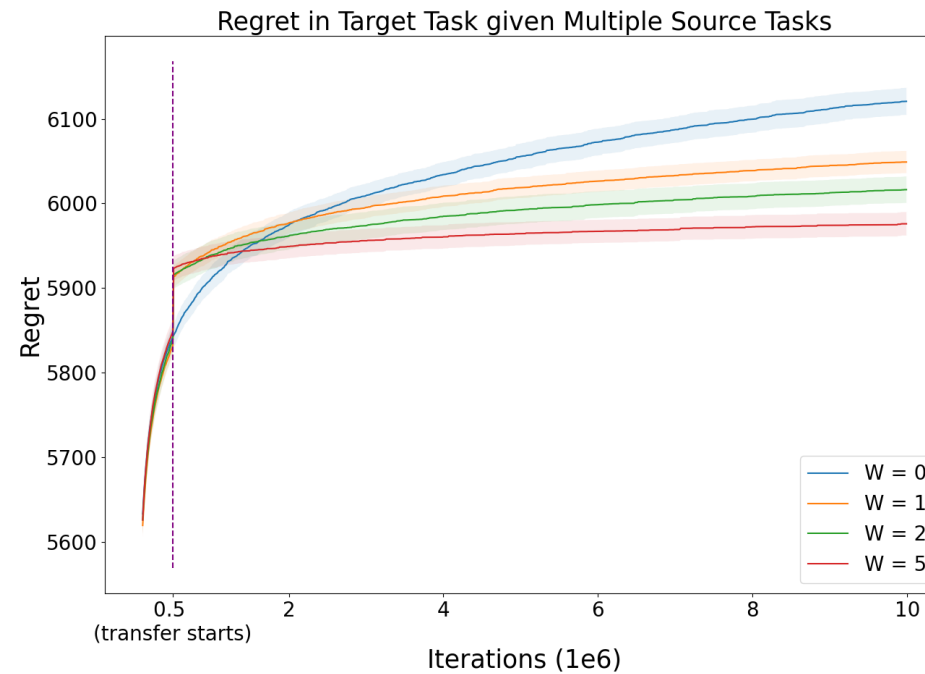
- **New issue:** how to select transferable tasks from task set?
- **Solution:** A novel task selection mechanism: “*Trust till Failure*”
 - For each state:
 - Maintain a feasible task set \mathcal{M}_{sh}
 - Pick $M_{\text{Trust}} \in \mathcal{M}_{sh}$ to trust until it is no longer feasible
 - When selecting the next task to trust:
 - Priorly select the feasible task recommending the same action

Experiments

- **Setting**

- Toy tabular MDP example;
- 5 source tasks at most;
- Different tasks created by permuting transition matrix

- **Results**



Thank you!