



UMassAmherst

Manning College of Information  
& Computer Sciences

# RL<sup>3</sup>: Boosting Meta Reinforcement Learning via RL inside RL<sup>2</sup>

Abhinav Bhatia, Samer B. Nashed, Shlomo Zilberstein

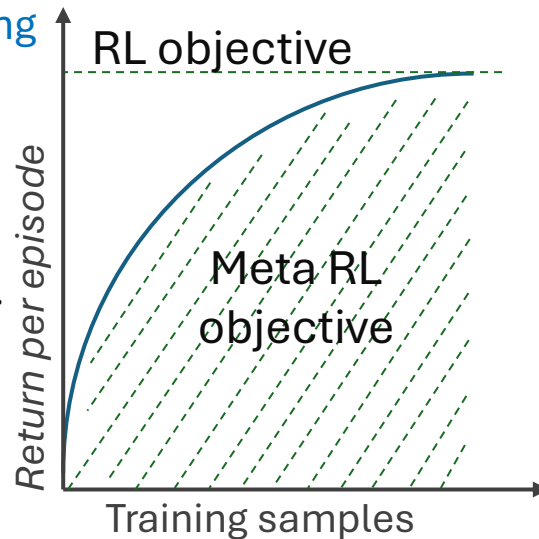
Aug 6, 2025

# Overview

An RL algorithm: a mapping from experience history to actions.

## Standard RL

- Given: an MDP
- Project: learn a *state-to-action mapping* to maximize *cumulative reward per episode*.
- Output: “Policy”
- Involves *value functions* to distill data.
- Standard RL Pros & cons:
  - Data inefficient (not the objective)
  - General-Purpose (achieves objective on any MDP)
  - Asymptotically optimal



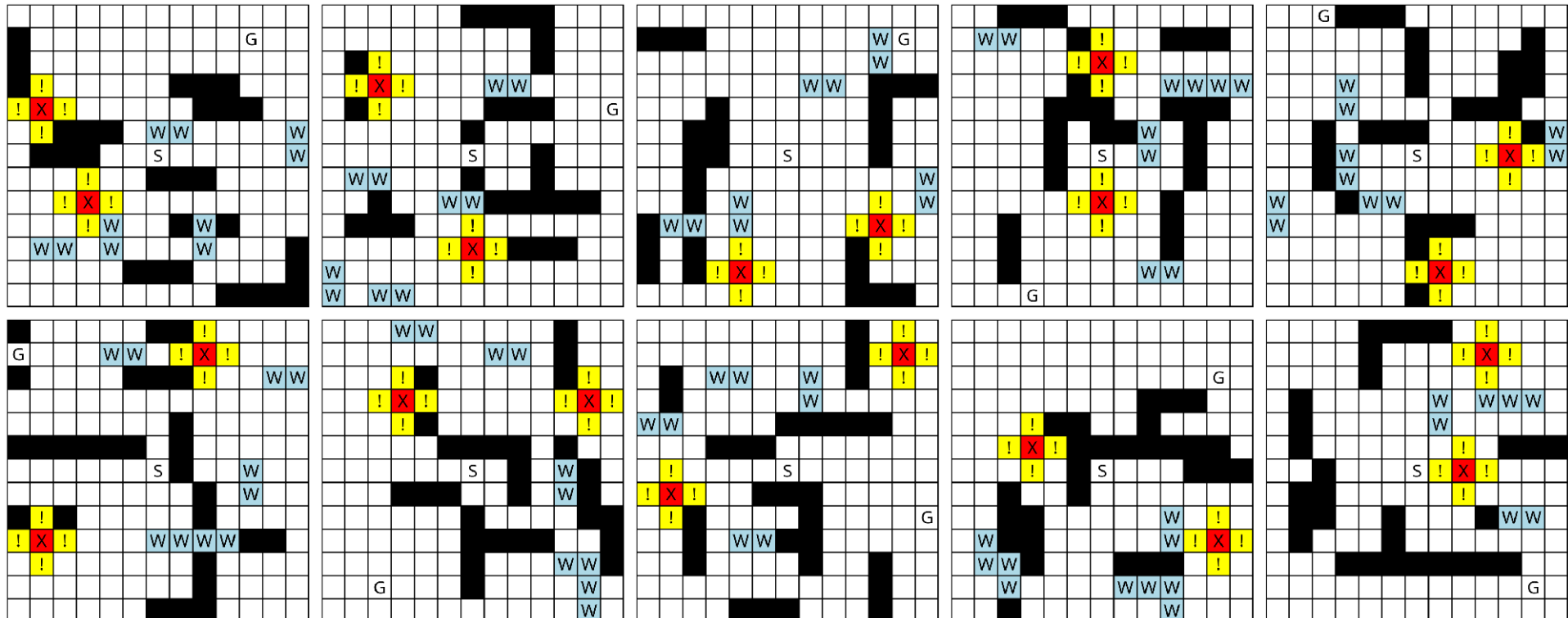
## Meta RL

- Given: a distribution of MDPs / “tasks”
- Project: learn a *data-to-action mapping* maximize *cumul. reward over entire interaction period (fixed)*. Min regret.
- “Meta-RL policy” or “Fast RL”
- Involves a *sequence model (RNN/Transformer)* to ingest data (and may be an incremental task-inference module). E.g.,  $RL^2$ , VariBAD.
- Meta RL Pros & cons:
  - Data-efficient (objective = min regret)
  - OOD issue (not trained for that)
  - Long-context problem (compute bottleneck, context truncation, gradient problems, long range credit assignment difficulty, compounding errors)

$RL^3$ : Injects RL into Meta-RL: Augment experience inputs with  $Q^*$  estimates

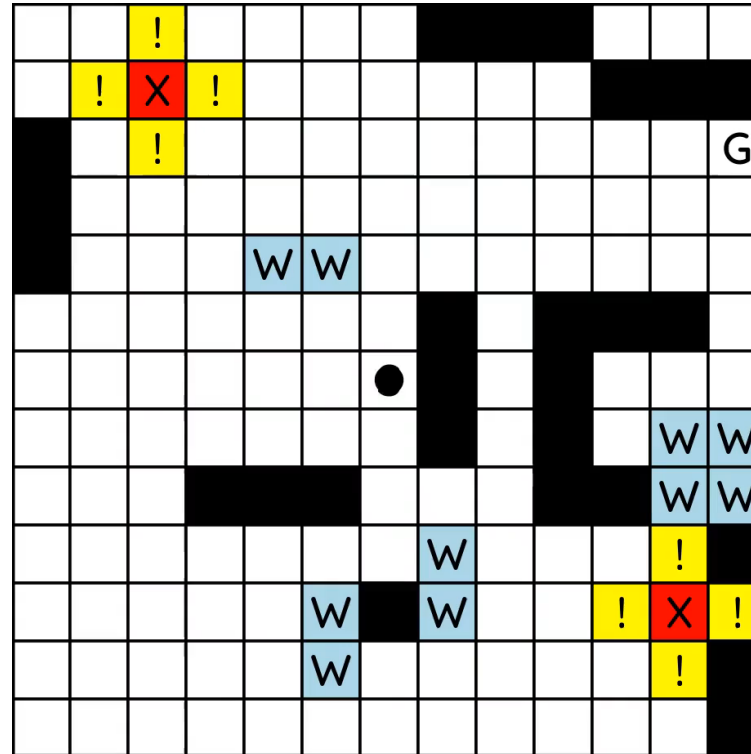
# Distribution of MDPs

- Same state, action space. Different reward, transition functions, with commonalities.





# Meta-Episode Example




Action:  Reward: 0.0

- Many episode on one MDP  $\leftrightarrow$  1 meta-episode

# Meta Reinforcement Learning

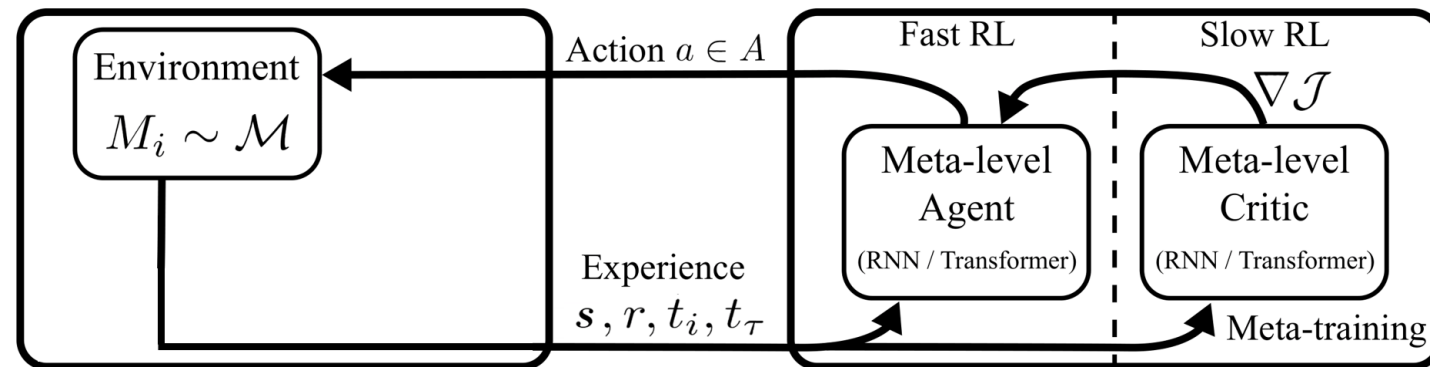
- Objective: Learn a data-to-action mapping to maximizes cumulative reward

$$\mathcal{J}(\pi) = \mathbb{E}_{M_i \sim \mathcal{M}} \left[ \sum_{t=0}^H \gamma^t \mathbb{E}_{(s_t, a_t) \sim \rho_{i,t}^{\pi}} [R_i(s_t, a_t)] \right]$$

- As a meta-level Markov decision process:
  - Each **meta-episode: sample a new MDP**, or “task”, play for  $H$  interactions.
  - Optimal meta policy maximizes cumulative reward. Least regret given  $H$  interactions .
  - Dynamics different across meta-episodes, because every meta-episode is a new MDP.
  - It's **POMDP** at meta-level, where **hidden variable is task identity**, aka **BAMDP**.
  - Beliefs over tasks capture history sufficiently.
  - Conditioning meta-policy on  $(s_t, b_t)$  is sufficient.

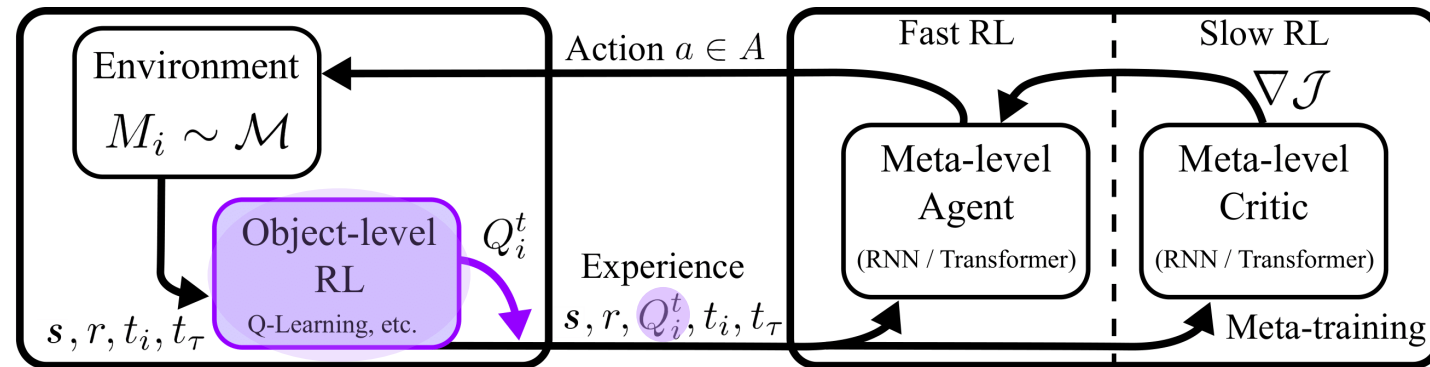
# $RL^2$ : Fast RL using Slow RL (Duan et al. 2016)

- $RL^2$  maps raw-data to actions directly (A simple black-box method).  
(Note: Some approaches map data-to-beliefs first e.g., VeriBAD (Zintgraf et al., 2019))
- Trained with standard “slow” deep RL.
- No general-purpose components.



# RL<sup>3</sup>: Inject RL into RL<sup>2</sup>

- Insert a RL subroutine: **estimate Q\*-values** e.g., use Q-learning.
- **Provide to meta-RL**. (Provide action-counts too).
- Meta-RL decides how to use.
- Over time, Q-values improve.



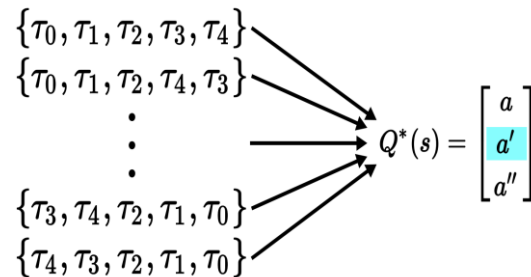
# RL<sup>3</sup>: Inject RL into RL<sup>2</sup>: But why?

Claim: Q-injection  improves OOD generalization and long-context reasoning.

ep-greedy uct  
exploration count-based  
curiosity-driven ucb  
sac  
boltzman dqn  
ddpg

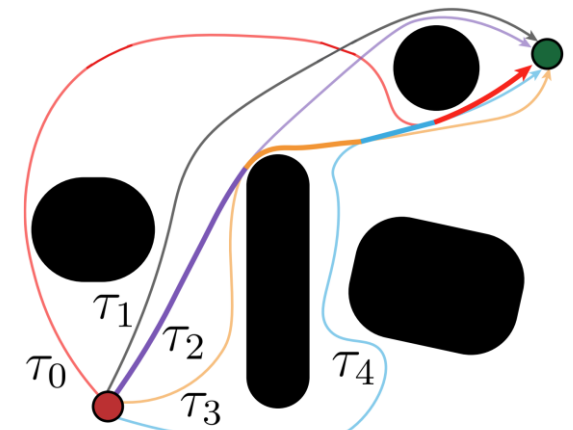
## Inherent generality:

Key component in general-purpose RL



**Summarization:** Many-to-one mapping.  
Order is irrelevant.

Lossy, but “remembers” key details



## Actionability:

Handles long-range credit assignment,  
“stitching trajectories”, optimal policy given data.

Eventually, Can ignore history, just exploit

Bottom line: Over time, data becomes overwhelming, Q-estimates more useful.



# RL<sup>3</sup>: Inject RL into RL<sup>2</sup>: But Why?

## Theoretical Reasons

### **Excellent task discriminators / identifiers:**

Rare for MDPs to have same Q-value function, or same evolution of Q-estimates.

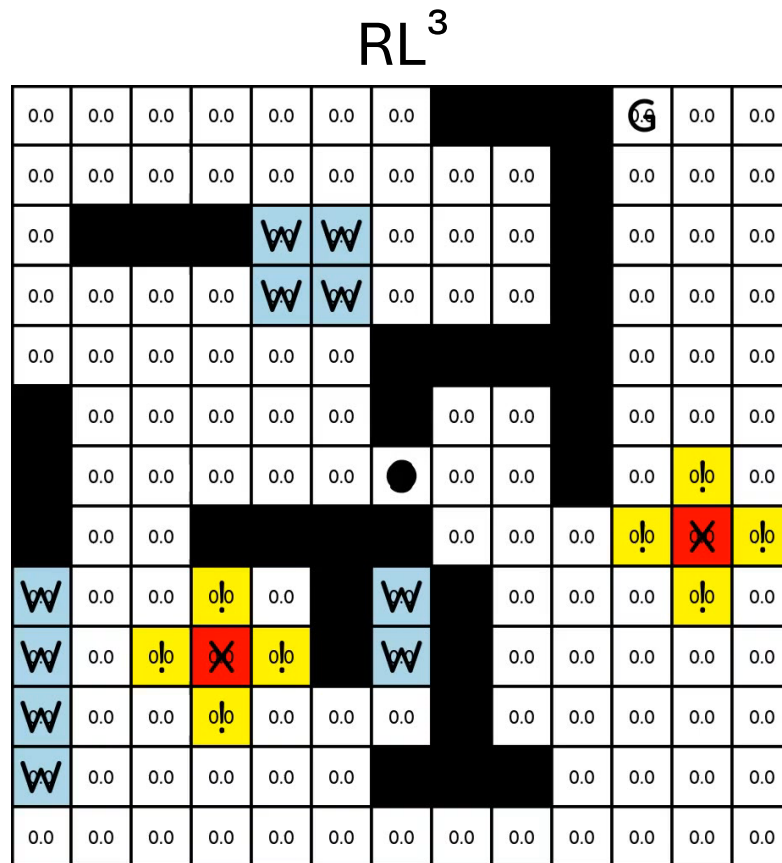
### **Related to meta-value function:**

Equivalent in the limit!

for any  $\epsilon > 0$ , there exists  $\kappa \in \mathbb{N}$  such that for  $t \geq \kappa$ ,

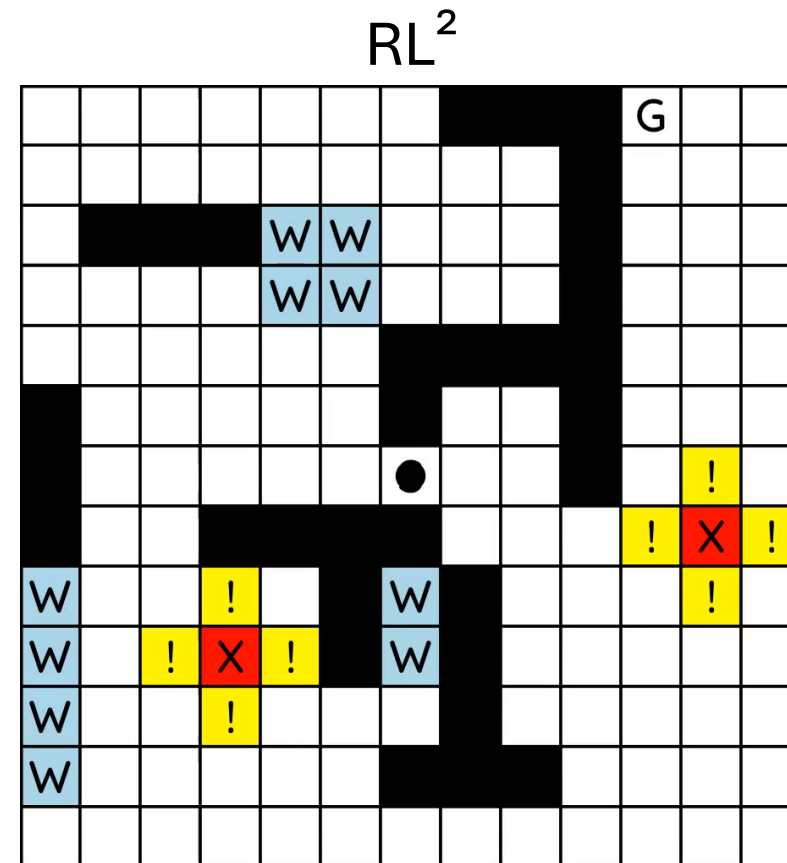
$$\left| \max_{a \in A} \left[ Q_i^t(s, a) \right] - \bar{V}^*(\bar{b}) \right| \leq \epsilon \quad \forall s \in S.$$

# RL<sup>3</sup> vs RL<sup>2</sup> - Gridworlds Results Demo



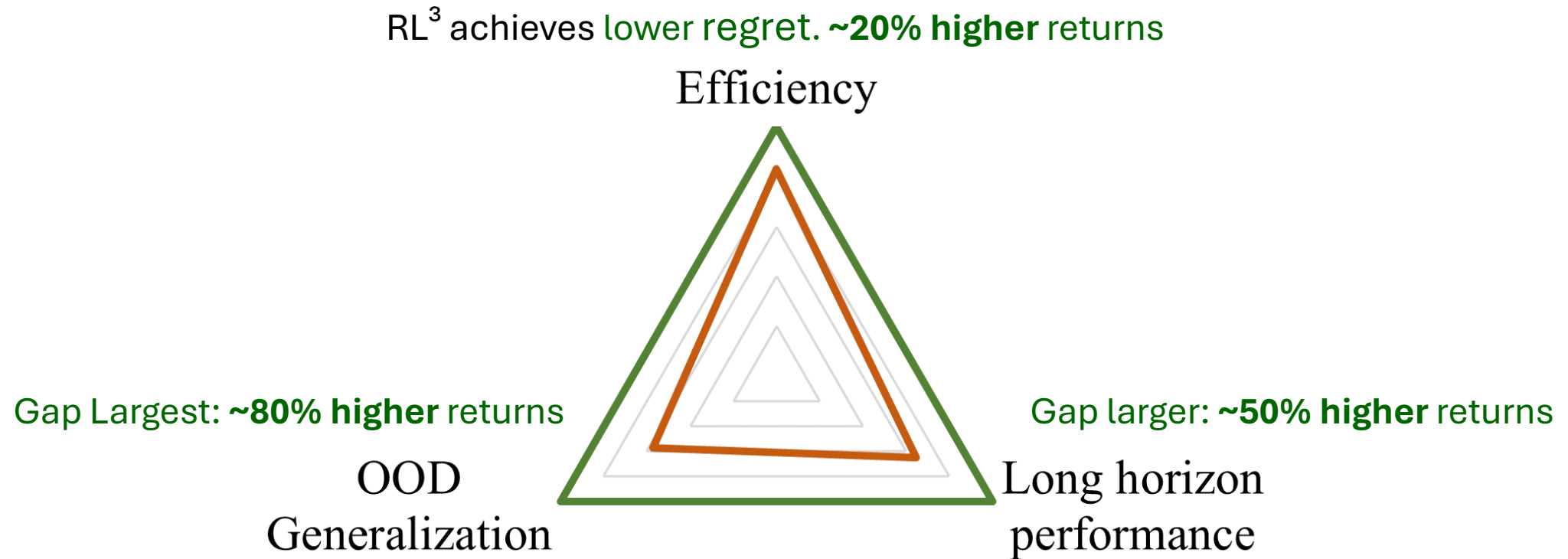
Action:  Reward: 0.0

[https://youtu.be/eLA\\_S1BQUYM](https://youtu.be/eLA_S1BQUYM)



Action:  Reward: 0.0

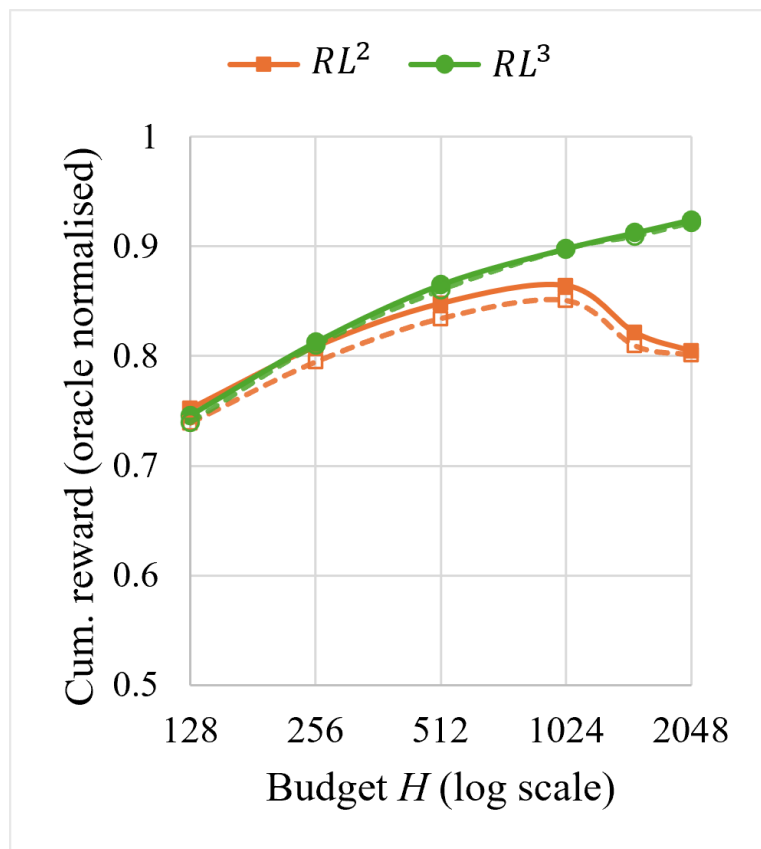
# RL<sup>3</sup> vs RL<sup>2</sup> - Gridworlds Results



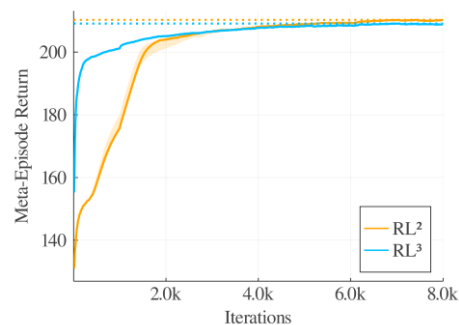
RL<sup>3</sup> with state-abstractions: *RL<sup>3</sup>-coarse*: 2x fast, 90% of RL<sup>3</sup>.

Bonus: RL<sup>3</sup> meta-training: **~30% more sample efficient**

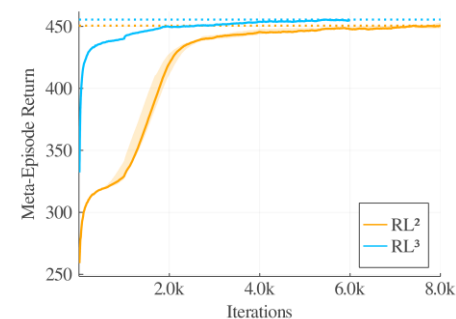
# RL<sup>3</sup> vs RL<sup>2</sup> - Random MDPs Results



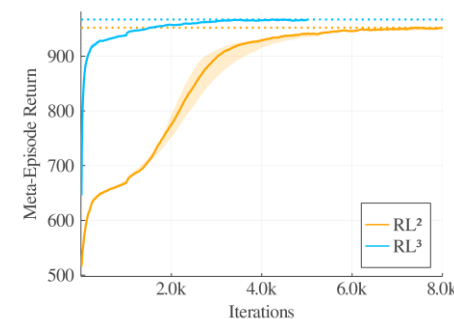
Meta-training graphs



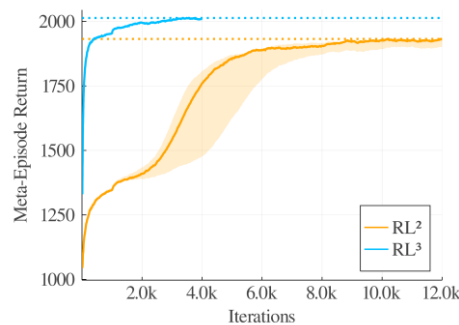
128



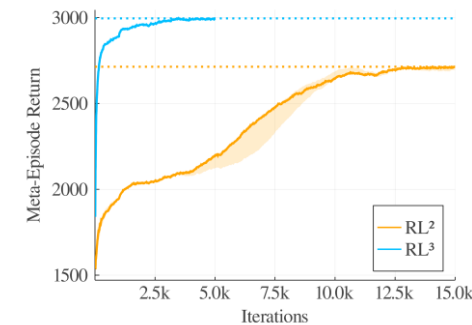
256



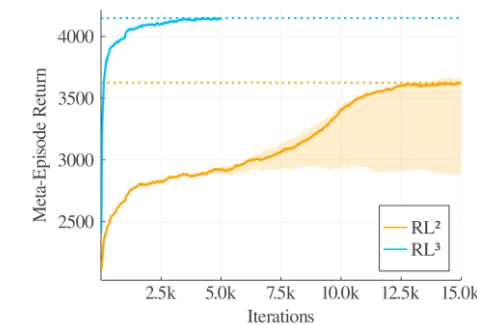
512



1024



1500



2048

# Conclusion

- We introduced  $RL^3$ , aiming to combine best of RL and  $RL^2$  – to achieve good efficiency (minimize regret), better long-term reasoning, better OOD generalization.
- Intuitions: Universality, summarization, actionability, long-term credit assignment and with helps task identification. With time, data gets overwhelming, Q-estimates useful, eventually sufficient.
- Key takeaways:
  - $RL^3$  retains short-term efficiency of  $RL^2$  on all domains
  - $RL^3$  benefits increase with horizon.
  - $RL^3$  benefits increase with distribution shift.
  - Bonus result: meta-training efficiency.
- Plug & play