

TL;DR: Supplementing Meta-RL inputs with  $Q$ -values learned online improves long-horizon performance and OOD generalization.

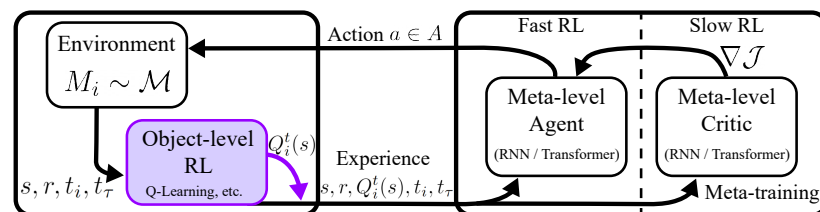
## Meta-Reinforcement Learning

|                 | RL   | Meta-RL  |
|-----------------|--|--|
| Objective       | Maximize return <i>within an episode</i> . | Maximize return over a <i>meta-episode</i> (whole interaction window) via online adaptation.   |
| Scope           | Single MDP (asymptotically optimal).       | Family of MDPs (or <i>tasks</i> ) from a known distribution.   |
| Uses experience | To learn value functions.                  | To map history $\rightarrow$ action with sequence models like RNN/Transformer (e.g., RL <sup>2</sup> ), and/or incremental task-inference as an intermediate step. |

- Meta-RL Challenges:** i) OOD generalization and ii) long horizons (truncated gradients / long-range credit assignment / compounding inference errors).

## RL<sup>3</sup>: RL inside RL<sup>2</sup>

- RL<sup>2</sup> input:** experience sequence  $(s, a, r)$ , into a Transformer (our impl.).
- RL<sup>3</sup> input:** experience sequence + online  $Q_t(s)$  (value estimates for the current MDP).



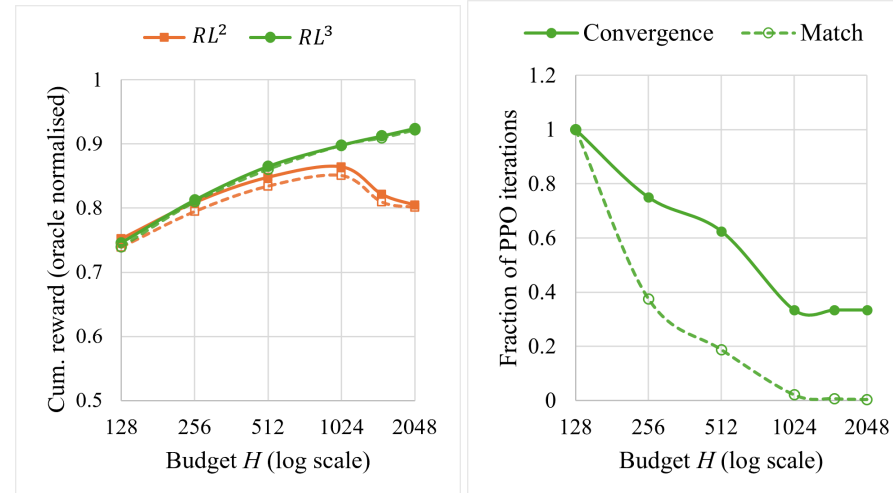
Overview: purple = additions to RL<sup>2</sup>.

- Insight:** Online  $Q_t$  adds **task-agnostic inductive bias** and an **actionable summary of experience** that *improves* as data accumulates, easing the sequence model's long-range credit assignment burden.
- Outcome:** Improved OOD generalization and higher long-horizon returns.

|                       | RL                | RL <sup>2</sup> | RL <sup>3</sup> |
|-----------------------|-------------------|-----------------|-----------------|
| Short-Term Efficiency | ✗                 | ✓               | ✓               |
| Long-Term Performance | ✓                 | ✗               | ✓               |
| OOD Generalization    | ✓                 | ✗               | ✓               |
|                       | (General Purpose) |                 | (Improved)      |

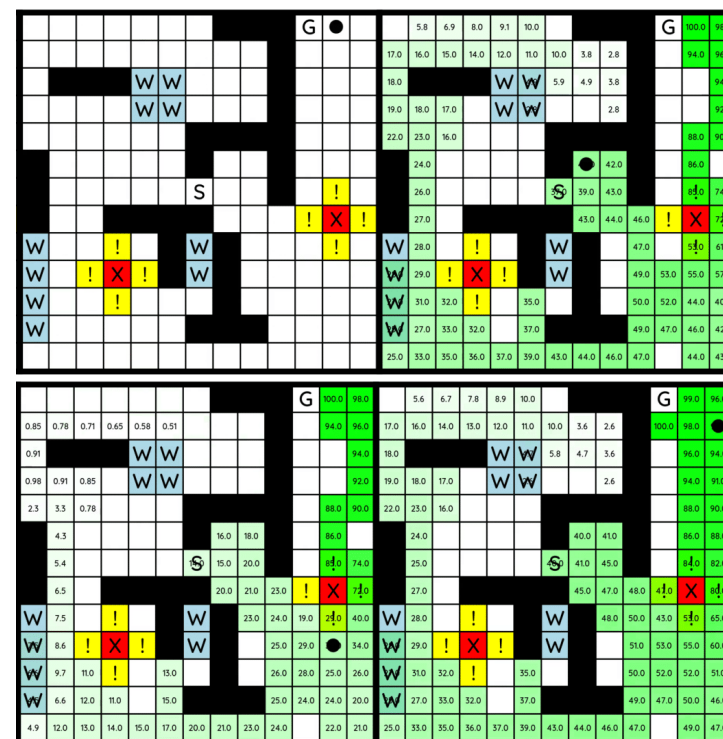
## Results

**Random MDPs:** Stochastic MDPs drawn from a fixed distribution; parameters are varied for OOD tests. RL<sup>3</sup> preserves RL's asymptotic scaling, maintains strong OOD performance, and meta-training is drastically more efficient.



Left: Return vs. interaction budget  $H$ . The gap widens as  $H$  increases 128  $\rightarrow$  2048 (solid: in-dist., dashed: OOD). Right: RL<sup>3</sup> meta-training converges or matches RL<sup>2</sup> with a fraction of the samples.

**Gridworlds:** Procedurally generated grids with obstacles (black), goals (G), hazards ('X', '!'), and slippery tiles ('W'). RL<sup>3</sup> **averages +50% return vs. RL<sup>2</sup>**; and **+80% under OOD shifts**, where we varied obstacle density, stochasticity, number of hazard tiles, etc. Interestingly, even with value estimates on a **2 $\times$  coarser state abstractions**, return drops only  $\approx$ 10% while compute is  $\approx$ 2 $\times$  faster. Finally, **Meta-training** is 30% more sample-efficient.



Example meta-episode. After exploration in the first episode (top left), the RL<sup>3</sup> agent uses the estimated value function (text and green color gradients inside tiles) in future episodes. Demonstrably, RL<sup>3</sup> avoids relying only on the RNN/transformer to plan the shortest path.

## Why RL<sup>3</sup> Works (Intuition)

RL<sup>3</sup> exploits properties of the action-value function  $Q(s, a)$ :

- Inductive bias:**  $Q$  is a ubiquitous RL signal; adding it provides task-agnostic structure to the meta-learner.
- Compression:** summarizes arbitrarily long, out-of-order experience into a fixed-size representation.
- Actionable:** a policy greedy w.r.t.  $Q$  approaches optimal as  $Q \rightarrow Q^*$ , shortening long-range credit assignment.
- Refinement:** DP/TD updates contract toward  $Q^*$ , so estimates typically improve with more interaction.
- Task identification:** the  $Q$ -landscape is task-specific, helping disambiguate the current MDP.
- Separation of concerns:** the RL module estimates/updates  $Q$ ; the meta-policy explores and decides when/how to trust  $Q_t$  (at convergence, meta- and task-level value align).

## Summary

RL<sup>3</sup> augments RL<sup>2</sup> with online  $Q$ -values for the current task, adding a **task-agnostic inductive bias** and a **compact, actionable history summary**. This reduces **long-horizon credit assignment** load on the sequence model, so the meta-learner **adapts more effectively over long-horizons** and **generalizes OOD**; the idea extends to function-approximate  $Q$  in richer domains.

## Additional Information

Find the [link to the paper](#) using the QR code. For additional information contact [abhinavbhati@umass.edu](mailto:abhinavbhati@umass.edu)

## Acknowledgments

This work was supported in part by the National Science Foundation grant numbers 1954782, 2205153, and 2321786.

