

# Novelty-Guided Proximal Curriculum Learning

Jan Malte Töpperwien

16.09.2024

- need lots of exploration
- sparse rewards -> infrequent learning signal
- curriculum learning for appropriate task difficulty

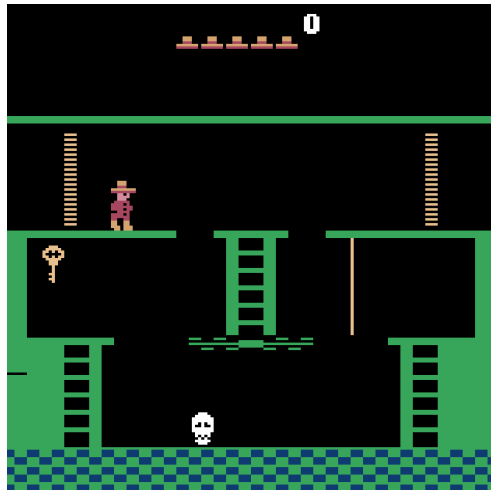


Figure: <https://images.openai.com/blob/2c736f64-38dc-4c65-a1dd-aabe2ceb8ddf/learning-montezumas-revenge-from-a-single-demonstration.png>

- need lots of exploration
- sparse rewards -> infrequent learning signal
- curriculum learning for appropriate task difficulty
- *How to set proper tasks?*
- current approaches often rigid and demonstration-based

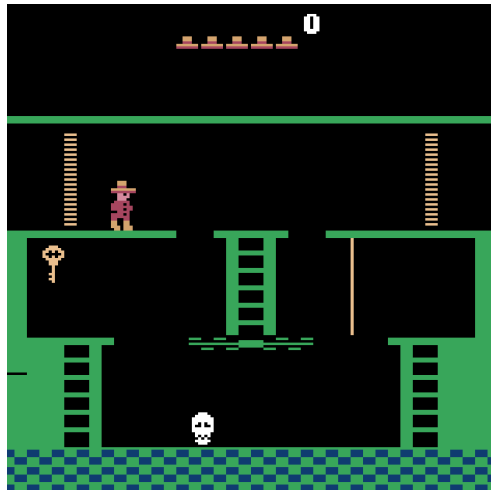


Figure: <https://images.openai.com/blob/2c736f64-38dc-4c65-a1dd-aabe2ceb8ddf/learning-montezumas-revenge-from-a-single-demonstration.png>

# Proximal Curriculum Learning (PCL) [Tzannetos et al. 2023]

- starting state based on *probability of success*  $PoS$
- $PoS \approx 0.5$
- distribution over  $S_{init}$
- Approximate  $PoS$  via agents value function  $V$  (scaled to  $[0, 1]$ )

Introduction

Method

Results

Discussion

References

# Proximal Curriculum Learning (PCL) [Tzannetos et al. 2023]

- starting state based on *probability of success*  $PoS$
- $PoS \approx 0.5$
- distribution over  $S_{init}$
- Approximate  $PoS$  via agents value function  $V$  (scaled to  $[0, 1]$ )

## Problems:

- $V$  initialized randomly
- $V$  inaccurate for seldomly seen states
- states may get  $PoS$  of 0 or 1 and never get chosen

# State Novelty

Exploration may help us fix the problem.

- state novelty  $\sim V$  inaccuracy
- explore seldomly seen states
- set starting state based on state novelty

# Novelty-Guided Proximal Curriculum Learning (NGPCL)

Introduction

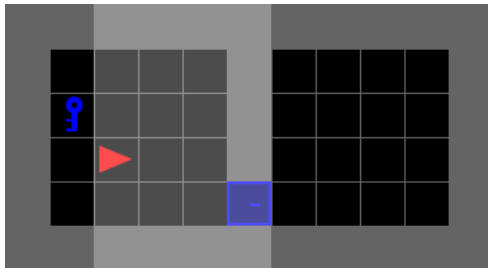
Method

Results

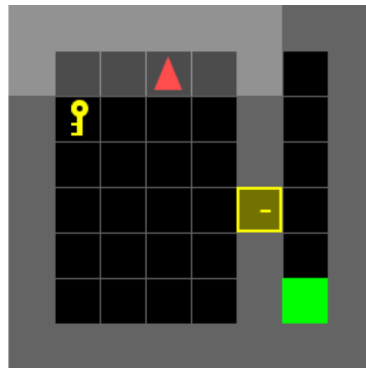
Discussion

References

- incorporate novelty by creating distribution over  $S_{init}$
- faster  $V$  convergence
- skip environment steps needed for intrinsic reward
- overlay both distributions by using weighted sum



(a) Unlock

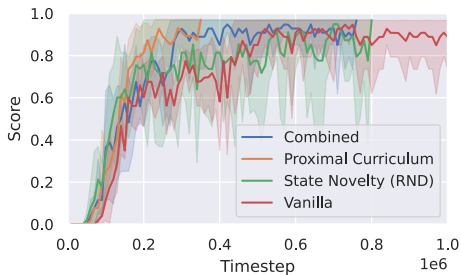


(b) DoorKey-8x8

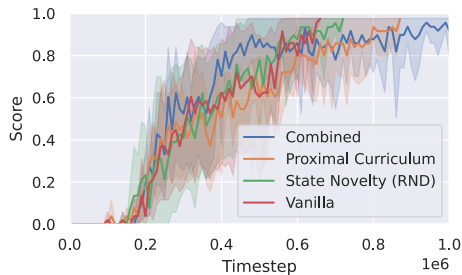
**Figure:** Environments used for experiments [Chevalier-Boisvert et al. 2023]. Key was added to observation.



- $S_{init}$  over all states, not evolving
- Random Network Distillation (RND) [Burda et al. 2018] as state novelty implementation
- Hyperparameters and architectures optimized by bayesian optimization using SMAC [Lindauer et al. 2022]

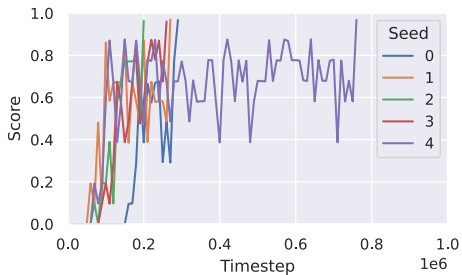


(a) Unlock

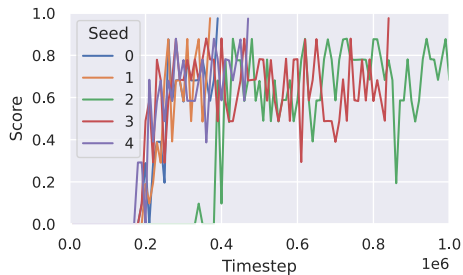


(b) Doorkey-8x8

**Figure:** Rewards over training steps starting from usual starting state (NGPCL). 5 seeds with 95% CI.



(a) Unlock



(b) Doorkey-8x8

Figure: Rewards over the training steps given seed (NGPCL).

## Results:

- some seeds showed fast learning
- ... others failed to solve the environment
- performance differs per environment

## Possible causes:

- distribution overlaying may be destructive
- Hyperparameters could be set poorly (especially for RND)

# Future Work

## Things to try:

- schedule for overlay parameter
- interleave PCL with RND
- try other state novelty approaches
- (dynamic)  $S_{init}$  determination

## Environments to evaluate:

- dense rewards
- big/continuous state- and/or action-spaces
- reward space not being convex

## References



Burda, Yuri et al. (2018). “Exploration by Random Network Distillation”. In: *CoRR* abs/1810.12894. arXiv: 1810.12894. URL: <http://arxiv.org/abs/1810.12894>.



Chevalier-Boisvert, Maxime et al. (2023). “Minigrid & Miniworld: Modular & Customizable Reinforcement Learning Environments for Goal-Oriented Tasks”. In: *CoRR* abs/2306.13831.



Lindauer, Marius et al. (2022). “SMAC3: A Versatile Bayesian Optimization Package for Hyperparameter Optimization”. In: *Journal of Machine Learning Research* 23.54, pp. 1–9. URL: <http://jmlr.org/papers/v23/21-0888.html>.



Tzannetos, Georgios et al. (2023). “Proximal Curriculum for Reinforcement Learning Agents”. In: *Trans. Mach. Learn. Res.* 2023. URL: <https://openreview.net/forum?id=8WUyeeMxMH>.