

Pick up the PACE: A Parameter-Free Optimizer for Lifelong Reinforcement Learning

Aneesh Muppidi

aneeshmuppidi@college.harvard.edu
Harvard College

Zhiyu Zhang

zhiyuz@seas.harvard.edu
Harvard University

Heng Yang

hankyang@seas.harvard.edu
Harvard University

Abstract

A key challenge in lifelong reinforcement learning (RL) is the loss of plasticity, where previous learning progress hinders an agent’s adaptation to new tasks. While regularization and resetting can help, they require precise hyperparameter selection at the outset and environment-dependent adjustments. Building on the principled theory of online convex optimization, we present a parameter-free optimizer for lifelong RL, called PACE, which requires no tuning or prior knowledge about the distribution shifts. Extensive experiments on Procgen and Gym Control environments show that PACE works surprisingly well—mitigating loss of plasticity and rapidly adapting to challenging distribution shifts—despite the underlying optimization problem being nonconvex and nonstationary. Code is available [here](#).

1 Introduction

Spot, the agile robot dog, has been learning to walk confidently across soft, lush grass. But when Spot moves from the grassy field to a gravel surface, the small stones shift beneath her feet, causing her to stumble. When Spot tries to walk across a sandy beach or on ice, the challenges multiply, and her once-steady walk becomes erratic. Spot wants to adjust quickly to these new terrains, but the patterns she learned on grass are not suited to gravel, sand, or ice. Furthermore, she never knows when the terrain will change again and how different it will be, therefore must continually plan for the unknown.

Spot’s struggle exemplifies a critical challenge in real-world decision making, known as *lifelong reinforcement learning* (lifelong RL). The idea is that the learning agent must continually acquire new knowledge to handle the nonstationarity of the environment. At first glance, there appears to be an obvious solution: given a policy gradient oracle, the agent could just keep running gradient descent nonstop. However, recent experiments have demonstrated an intriguing behavior called *loss of plasticity* (Dohare et al., 2021; Lyle et al., 2022; Abbas et al., 2023; Sokar et al., 2023): despite persistent gradient steps, an agent can gradually lose its responsiveness to incoming observations (see Figure 1). There are even extreme cases of loss of plasticity (known as *negative transfer* or *primacy bias*), where prior learning can significantly hamper the performance in new tasks (Nikishin et al., 2022; Ahn et al., 2024).

From the optimization perspective, the above issues might be attributed to the *lack of stability* under gradient descent. That is, the weights of the agent’s parameterized policy can drift far away from the origin (or a good initialization), leading to a variety of undesirable behaviors.¹ Fitting this narrative,

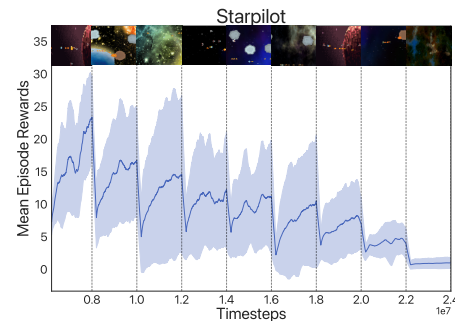


Figure 1: Loss of plasticity in Procgen (Starpilot) shows a steady decline in reward with each distribution shift.

¹Such as the inactivation of many neurons, due to dead ReLU activations (Abbas et al., 2023; Sokar et al., 2023).

it has been shown that simply adding a L_2 regularizer to the optimization objective (Kumar et al., 2023) or periodically resetting the weights (Dohare et al., 2021; Asadi et al., 2023; Sokar et al., 2023; Ahn et al., 2024) can help mitigate the problem. However, a particularly important limitation is their use of *hyperparameters*, such as the magnitude of the regularizer and the resetting frequency. Good performance hinges on the suitable environment-dependent hyperparameter, but how can one confidently choose that *before* interacting with the environment? The classical cross-validation approach would violate the one-shot nature of lifelong RL (and online learning in general; see Chapter 1 of Orabona, 2023), since it is impossible to experience the same environment multiple times. This leads to the contributions of the present work.

Contribution The present work addresses the key challenges in lifelong RL using the principled theory of *Online Convex Optimization* (OCO). Specifically, our contributions are two fold.

- **Algorithm: PACE** Building on a series of results in OCO (Cutkosky & Orabona, 2018; Cutkosky, 2019; Cutkosky et al., 2023; Zhang et al., 2024), we propose a (hyper)-*parameter-free* optimizer for lifelong RL, called **PACE** (Parameter-free Adaption for Continual Environments). Intuitively, the idea is a refinement of regularization: instead of manually selecting the magnitude of regularization beforehand, PACE chooses that in an online, data-dependent manner. From the perspective of OCO theory, PACE is insensitive to its own hyperparameter, which means that no hyperparameter tuning is necessary in practice.
- **Experiment** Using an instantiation of PACE with *Proximal Policy Optimization* (PPO) called PACE PPO, we evaluate on Lifelong settings of Procgen and Gym Control. In settings where existing approaches (Abbas et al., 2023; Kumar et al., 2023) struggle, we find that PACE PPO mitigates loss of plasticity and rapidly adapts when new distribution shifts are introduced. Such findings might be surprising: the theoretical advantage of PACE is motivated by the convexity in OCO, but lifelong RL is *both nonconvex and nonstationary* in terms of optimization.

2 Lifelong RL

As a sequential decision making framework, *reinforcement learning* is commonly framed as a *Markov Decision Process* (MDP) defined by the state space \mathcal{S} , the action space \mathcal{A} , the transition dynamics $P(s_{t+1}|s_t, a_t)$, and the reward function $R(s_t, a_t, s_{t+1})$. In the t -th round, starting from a state $s_t \in \mathcal{S}$, the learning agent needs to choose an action $a_t \in \mathcal{A}$ without knowing P and R . Then, the environment samples a new state $s_{t+1} \sim P(\cdot|s_t, a_t)$, and the agent receives a *reward* $r_t = R(s_t, a_t, s_{t+1})$. From a practical perspective, we measure the agent’s performance by its cumulative reward $\sum_{t=1}^T r_t$.

The standard setting above concerns a *stationary* MDP. The present work studies a nonstationary variant called *lifelong* RL, where the transition dynamics P_t and the reward function R_t can vary over time. Certainly, one should not expect any meaningful “learning” against *arbitrary* unstructured nonstationarity. Therefore, we implicitly assume P_t and R_t to be *piecewise constant* over time, and each piece is called a *task*.

Lifelong RL as online optimization Modern RL approaches, including PPO (Schulman et al., 2017) and others, crucially utilize the idea of *policy parameterization*. We use $\theta_t \in \mathbb{R}^d$ to denote the parameterizing *weight vector*. After sampling a_t and receiving new observations, the agent could define a *loss function* $J_t(\theta)$ that characterizes the “hypothetical performance” of each weight $\theta \in \mathbb{R}^d$. Then, by computing the *policy gradient* $g_t = \nabla J_t(\theta_t)$, one could apply a *first order optimization algorithm*² OPT to obtain the updated weight, $\theta_{t+1} = \text{OPT}(\theta_t, g_t)$.

For the rest of this paper, we will work with such an abstraction. The feedback of the environment is treated as a *policy gradient oracle* \mathcal{G} , which maps the time t and the current weight θ_t into a policy gradient $g_t = \mathcal{G}(t, \theta_t)$. Our goal is to design an optimizer OPT well suited for lifelong RL.

Lifelong vs. Continual Abel et al. (2023) characterized *continual reinforcement learning* (CRL) as a never-ending learning process. However, much of the literature cited under CRL, such as (Abbas

²Formally, a dynamical system that given its state θ_t and input g_t outputs the new state $\text{OPT}(\theta_t, g_t)$.

et al., 2023; Ahn et al., 2024), primarily focuses on the problem of *backward transfer* (avoiding catastrophic forgetting). Conversely, the present work addresses the problem of *forward transfer*, which refers to the rapid adaptation to new tasks. Because of this we use “lifelong” rather than “continual” in our exposition, similar to (Thrun, 1996; Abel et al., 2018b; Julian et al., 2020).

3 Method

Inspired by (Cutkosky et al., 2023), we study lifelong RL by exploiting its connection to *Online Convex Optimization* (OCO; Zinkevich, 2003). OCO is a key problem in online learning, with significant efforts to design *parameter-free* algorithms requiring minimal tuning (Streeter & McMahan, 2012; McMahan & Orabona, 2014; Orabona & Pál, 2016; Foster et al., 2017; Cutkosky & Orabona, 2018; Mhammedi & Koolen, 2020; Chen et al., 2021; Jacobsen & Cutkosky, 2022). The surprising observation of Cutkosky et al. (2023) is that several algorithmic ideas closely tied to the convexity of OCO can actually improve nonconvex deep learning training, suggesting certain notions of “near convexity” on its loss landscape. We find that lifelong RL (which is *both nonconvex and nonstationary* in terms of optimization) exhibits a similar behavior, therefore a particularly strong algorithm (named PACE) can be obtained from parameter-free OCO.

Basics of (parameter-free) OCO As a standalone theoretical topic, OCO concerns a sequential optimization problem where the convex loss function l_t can vary arbitrarily over time. In the t -th iteration, the optimization algorithm picks an iterate x_t and then observes a gradient $g_t = \nabla l_t(x_t)$. Motivated by the pursuit of “convergence” in optimization, the standard objective is to guarantee low (i.e., sublinear in T) *static regret*, defined as

$$\text{Regret}_T(l_{1:T}, u) := \sum_{t=1}^T l_t(x_t) - \sum_{t=1}^T l_t(u),$$

where T is the total number of rounds, and u is a *comparator* that the algorithm does not know beforehand. In other words, the goal is to make $\text{Regret}_T(l_{1:T}, u)$ small for *all* possible loss sequence $l_{1:T}$ and comparator u .

For minimizing static regret, classical *minimax* algorithms like gradient descent (Zinkevich, 2003) would assume a small *uncertainty set* \mathcal{U} at the beginning. Then, by setting the hyperparameter (such as the learning rate) according to \mathcal{U} , it is possible to guarantee sublinear *worst case regret*,

$$\max_{(l_{1:T}, u) \in \mathcal{U}} \text{Regret}_T(l_{1:T}, u) = o(T). \quad (1)$$

In contrast, parameter-free algorithms use very different strategies³ to bound $\text{Regret}_T(l_{1:T}, u)$ directly (without taking the maximum) by a function of both $l_{1:T}$ and u . The resulting bound is more refined than Eq.(1) (Orabona, 2023, Chapter 9), and crucially, since there is no need to pick an uncertainty set \mathcal{U} , much less hyperparameter tuning is needed. This is where its name comes from.

PACE for lifelong RL Now back to lifelong RL. As we discussed, a fundamental challenge here is the excessive drifting of the weights θ_t , and this can be fixed by enforcing the proximity to a good reference point θ_{ref} . Different from existing approaches like L_2 regularization and resetting, parameter-free OCO provides a principled solution to this problem without hyperparameter-tuning. Naming this algorithm as PACE, we present its generic template as Algorithm 1, which calls Algorithm 2 (a one-dimensional scale tuner) as the key subroutine.

From the technical perspective, PACE assembles three techniques in parameter-free OCO: the *direction-magnitude decomposition* from (Cutkosky & Orabona, 2018), the *additive aggregation* from (Cutkosky, 2019), and the *erfi potential function* from (Zhang et al., 2024). In addition, we follow (Cutkosky et al., 2023) for the recommended default parameter setup.

Without going deep into the theory, here is an overview of the important ideas.

³The key difference is the use of intricate (non- L_2) regularizers (Fang et al., 2022; Jacobsen & Cutkosky, 2022).

Algorithm 1 PACE: Parameter-free Adaption for Continual Environments.

-
- 1: **Input:** A policy gradient oracle \mathcal{G} ; a first order optimization algorithm BASE; a reference point $\theta_{\text{ref}} \in \mathbb{R}^d$; n discount factors $\beta_1, \dots, \beta_n \in (0, 1]$ (default: $0.9, 0.99, \dots, 0.999999$).
 - 2: **Initialize:** Create n copies of Algorithm 2, denoted as $\mathcal{A}_1, \dots, \mathcal{A}_n$. For each $j \in [1 : n]$, \mathcal{A}_j uses the discount factor β_j . Initialize the algorithm BASE at θ_{ref} . Let $\theta_1 = \theta_{\text{ref}}$.
 - 3: **for** $t = 1, 2, \dots$ **do**
 - 4: Obtain the t -th policy gradient $g_t = \mathcal{G}(t, \theta_t) \in \mathbb{R}^d$.
 - 5: Send g_t to BASE as its t -th input, and get its output $\theta_{t+1}^{\text{Base}} \in \mathbb{R}^d$.
 - 6: For all $j \in [1 : n]$, send $\langle g_t, \theta_t - \theta_{\text{ref}} \rangle$ to \mathcal{A}_j as its t -th input, and get its output $s_{t+1,j} \in \mathbb{R}$.
 - 7: Define the scaling parameter $S_{t+1} = \sum_{j=1}^n s_{t+1,j}$.
 - 8: Update the weight of the policy,

$$\theta_{t+1} = \theta_{\text{ref}} + (\theta_{t+1}^{\text{Base}} - \theta_{\text{ref}}) S_{t+1}.$$

9: **end for**

Algorithm 2 1D Discounted Tuner of PACE.

-
- 1: **Input:** Discount factor $\beta \in (0, 1]$; small value $\varepsilon > 0$ (default: 10^{-8}).
 - 2: **Initialize:** The running variance $v_0 = 0$; the running (negative) sum $\sigma_0 = 0$.
 - 3: **for** $t = 1, 2, \dots$ **do**
 - 4: Obtain the t -th input h_t .
 - 5: Let $v_t = \beta^2 v_{t-1} + h_t^2$, and $\sigma_t = \beta \sigma_{t-1} - h_t$.
 - 6: Select the t -th output

$$s_{t+1} = \frac{\varepsilon}{\text{erfi}(1/\sqrt{2})} \text{erfi}\left(\frac{\sigma_t}{\sqrt{2v_t + \varepsilon}}\right),$$

where erfi is the *imaginary error function* queried from standard software packages.

7: **end for**

- First, PACE is a meta-algorithm that operates on top of a “default” optimizer BASE. It can simply be gradient descent with a constant learning rate, or ADAM (Kingma & Ba, 2014) as in our experiments. Applying BASE alone would be equivalent to enforcing the scaling parameter $S_{t+1} \equiv 1$ in PACE, but this would suffer from the drifting of $\theta_{t+1}^{\text{Base}}$ (and thus, the weight θ_{t+1}).
- To fix this issue, PACE uses the tuner (Algorithm 2) to select the scaling parameter S_{t+1} , making it *data-dependent*. Typically S_{t+1} is within $[0, 1]$ (see Figure 10 to 12), therefore essentially, we define the updated weight θ_{t+1} as a *convex combination* of the BASE’s weight θ_t^{Base} and the reference point θ_{ref} ,

$$\theta_{t+1} = S_{t+1} \cdot \theta_{t+1}^{\text{Base}} + (1 - S_{t+1}) \theta_{\text{ref}}.$$

This brings the weight closer to θ_{ref} , which is known to be “safe” (i.e., not overfitting any particular lifelong RL task), although possibly conservative.

- To inject appropriate conservatism without hyperparameter tuning, the tuner (Algorithm 2) uses the erfi function decision rule. This is theoretically optimal in an idealized variant of OCO (Zhang et al., 2022; 2024), but removing the idealized assumptions requires a tiny bit of extra conservatism, which is challenging (and not necessarily practical). For the lifelong RL problem, which considerably deviates from OCO, we apply the erfi decision rule as is.
- Finally, the tuner requires a discount factor β . This crucially controls the strength of regularization, but also introduces a hyperparameter tuning problem. Following (Cutkosky, 2019), we aggregate tuners with different β (on a log-scaled grid) by simply summing up their outputs. This is justified by the *adaptivity* of the tuner itself: in OCO, if we add a parameter-free algorithm \mathcal{A}_1 to any other algorithm \mathcal{A}_2 that already works well, then \mathcal{A}_1 can automatically identify this and “tune down” its aggressiveness, such that $\mathcal{A}_1 + \mathcal{A}_2$ still performs as well as \mathcal{A}_2 . Although PACE requires

a β -grid, constant ε , and the BASE algorithm, it is largely insensitive to these choices. The default values from Cutkosky et al. (2023) suffice.

4 Experiment

Does PACE experience the common pitfalls of loss of plasticity? Does it rapidly adapt to distribution shifts? To answer these questions, we instantiate PPO with two different optimizers: ADAM with constant learning rate for baseline comparison, and PACE for our proposed method (with exactly the same ADAM as the input BASE). We also test ADAM PPO with *concatenated ReLU activations* (CReLU; Shang et al., 2016), previously shown to mitigate loss of plasticity in deep RL (Abbas et al., 2023). More details and experiments can be found in Appendix A, G (including numerical results in Table 1).

OpenAI Procgen We evaluated PACE on OpenAI Procgen, a suite of 16 procedurally generated games (Cobbe et al., 2019). Distribution shifts were introduced by sampling a new level every 2 million steps, treating each level as a distinct task.

In StarPilot, Dodgeball, and Chaser, ADAM PPO and CReLU showed degrading performance with each new level (Figure 2). In contrast, PACE PPO avoided this loss, achieving rapid reward increases. Overall, PACE PPO demonstrated average improvements of 3,212.42% over ADAM PPO and 120.88% over CReLU (Table 1).

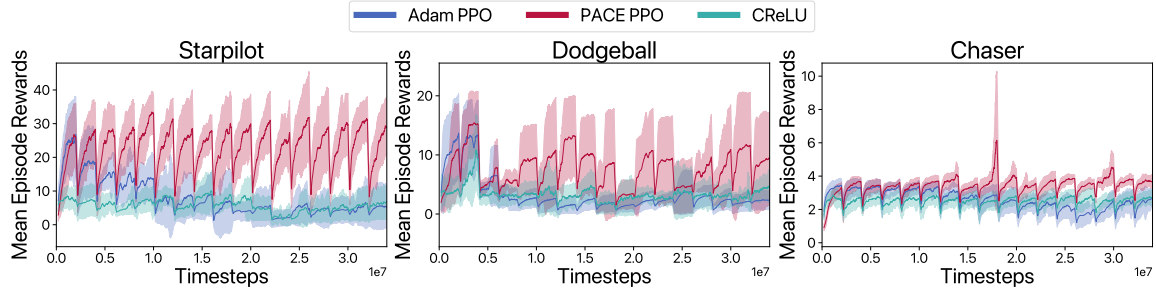


Figure 2: The reward plots show steady loss of plasticity in ADAM PPO and CReLU, while PACE maintains high performance across levels in StarPilot, Dodgeball, and Chaser.

Gym Control We use the CartPole-v1 and Acrobot-v1 environments from the Gym Classic Control suite, along with LunarLander-v2 from Box2d Control. To introduce distribution shifts every 200 steps we perturb each observation dimension with random noise within a range of ± 2 , treating each perturbation phase as a distinct task.

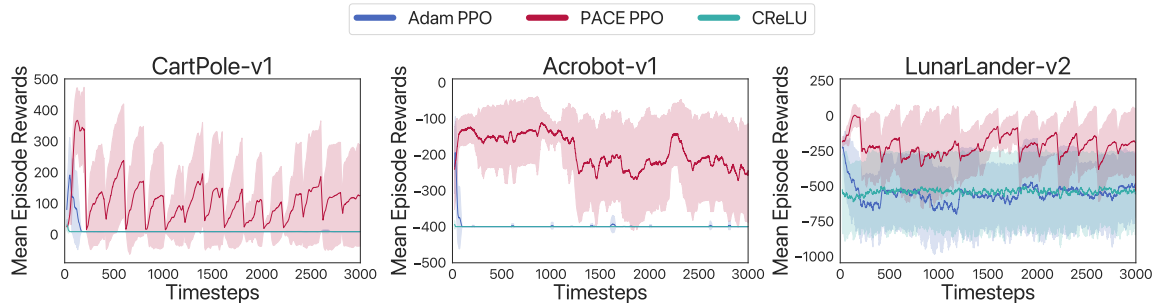


Figure 3: In CartPole, Acrobot, and LunarLander, both ADAM PPO and CReLU fail to recover after the initial distribution shift. PACE PPO rapidly adapts to each extreme distribution shift.

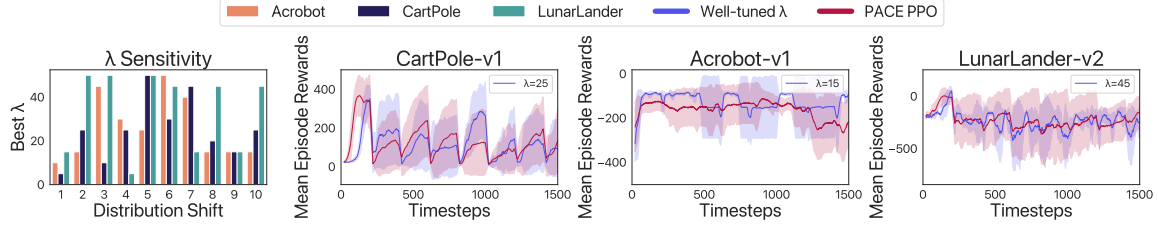


Figure 4: For each Gym Control environment and the initial ten tasks, we identified the optimal λ that maximizes reward for each task’s specific distribution shift. We also determined the best overall λ for each environment: CartPole $\lambda = 25$, Acrobot $\lambda = 15$, and LunarLander $\lambda = 45$.

Here (Figure 3), we notice a peculiar behavior after introducing the first distribution shift in both ADAM PPO and CReLU: policy collapse. We describe this as an *extreme* form of loss of plasticity. Surprisingly, PACE PPO remains resistant to these extreme distribution shifts. Across the three Gym Control environments, PACE PPO shows an average normalized improvement of 204.18% over ADAM PPO and 1044.24% over CReLU (Table 1).

5 Discussion

Related work Combating loss of plasticity has been studied extensively in lifelong RL. A typical challenge for existing solutions is the tuning of their hyperparameters, which requires prior knowledge on the nature of the distribution shift, e.g., (Asadi et al., 2023; Ben-Iwhiwhu et al., 2023; Nikishin et al., 2023; Sokar et al., 2023; Mesbahi et al., 2024). An architectural modification called CReLU is studied in (Abbas et al., 2023), but our experiments suggest that its benefit might be specific to the Atari setup. Besides, Abel et al. (2018a;b) presented a theoretical analysis of skill transfer in lifelong RL, based on value iteration.

Tuning L_2 regularization The success of PACE suggests that adaptive regularization, anchored to θ_{ref} , effectively counters both mild and extreme loss of plasticity. This highlights the limitation of the L_2 regularization approach from (Kumar et al., 2023). It requires selecting a regularization strength parameter λ through cross-validation, which is incompatible with the one-shot nature of lifelong learning settings. However, even when we try the λ -grid suggested by (Kumar et al., 2023), there is no effective λ value within the grid for the lifelong RL environments we consider. All the values are too small. We conduct a hyperparameter search for λ , over various larger values [0.2, 0.8, 1, 5, 10, 15, 20, 25, 30, 35, 40, 45, 50]. We discover that each environment and task responds uniquely to these regularization strengths (Figure 4). This highlights the challenges of tuning λ in a lifelong learning context. In contrast, PACE dynamically adapts to data online, consistently competing with well-tuned λ values in CartPole, Acrobot, and LunarLander (Figure 4).

Near convexity of lifelong RL Our results demonstrate the rapid adaptation of PACE, in lifelong RL problems with complicated function approximation. From the perspective of optimization, the latter requires tackling both nonconvexity and nonstationarity, which is typically regarded intractable in theory. Perhaps surprisingly, when approaching this complex problem using the theoretical insights from OCO, we observe compelling results. This suggests a certain “hidden convexity” in this problem.

6 Conclusion

We introduced PACE, a parameter-free optimizer for lifelong RL using OCO principles. PACE dynamically refines regularization based on data, eliminating hyperparameter tuning. PACE’s results lead to a compelling takeaway: empirical lifelong RL may exhibit more convex properties than previously appreciated, and might inherently benefit from parameter-free OCO approaches.

References

- Zaheer Abbas, Rosie Zhao, Joseph Modayil, Adam White, and Marlos C Machado. Loss of plasticity in continual deep reinforcement learning. *arXiv preprint arXiv:2303.07507*, 2023.
- David Abel, Dilip Arumugam, Lucas Lehnert, and Michael Littman. State abstractions for lifelong reinforcement learning. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 10–19. PMLR, 10–15 Jul 2018a. URL <https://proceedings.mlr.press/v80/abel18a.html>.
- David Abel, Yuu Jinnai, Sophie Yue Guo, George Konidaris, and Michael Littman. Policy and value transfer in lifelong reinforcement learning. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 20–29. PMLR, 10–15 Jul 2018b. URL <https://proceedings.mlr.press/v80/abel18b.html>.
- David Abel, André Barreto, Benjamin Van Roy, Doina Precup, Hado van Hasselt, and Satinder Singh. A definition of continual reinforcement learning, 2023.
- Hongjoon Ahn, Jinu Hyeon, Youngmin Oh, Bosun Hwang, and Taesup Moon. Catastrophic negative transfer: An overlooked problem in continual reinforcement learning, 2024. URL <https://openreview.net/forum?id=o7BwUyXz1f>.
- Kavosh Asadi, Rasool Fakoor, and Shoham Sabach. Resetting the optimizer in deep rl: An empirical study, 2023.
- Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *CoRR*, abs/1207.4708, 2012. URL <http://arxiv.org/abs/1207.4708>.
- Eseoghene Ben-Iwhiwhu, Saptarshi Nath, Praveen K. Pilly, Soheil Kolouri, and Andrea Soltoggio. Lifelong reinforcement learning with modulating masks, 2023.
- Liyu Chen, Haipeng Luo, and Chen-Yu Wei. Impossible tuning made possible: A new expert algorithm and its applications. In *Conference on Learning Theory*, pp. 1216–1259. PMLR, 2021.
- Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. *CoRR*, abs/1912.01588, 2019. URL <http://arxiv.org/abs/1912.01588>.
- Ashok Cutkosky. Combining online learning guarantees. In *Conference on Learning Theory*, pp. 895–913. PMLR, 2019.
- Ashok Cutkosky and Francesco Orabona. Black-box reductions for parameter-free online learning in banach spaces. In *Conference On Learning Theory*, pp. 1493–1529. PMLR, 2018.
- Ashok Cutkosky, Aaron Defazio, and Harsh Mehta. Mechanic: A learning rate tuner. *Advances in Neural Information Processing Systems*, 36, 2023.
- Shibhansh Dohare, Richard S Sutton, and A Rupam Mahmood. Continual backprop: Stochastic gradient descent with persistent randomness. *arXiv preprint arXiv:2108.06325*, 2021.
- Huang Fang, Nicholas JA Harvey, Victor S Portella, and Michael P Friedlander. Online mirror descent and dual averaging: keeping pace in the dynamic case. *Journal of Machine Learning Research*, 23(1):5271–5308, 2022.
- Dylan J Foster, Satyen Kale, Mehryar Mohri, and Karthik Sridharan. Parameter-free online learning via model selection. *Advances in Neural Information Processing Systems*, 30, 2017.
- Andrew Jacobsen and Ashok Cutkosky. Parameter-free mirror descent. In *Conference on Learning Theory*, pp. 4160–4211. PMLR, 2022.
- Ryan Julian, Benjamin Swanson, Gaurav S. Sukhatme, Sergey Levine, Chelsea Finn, and Karol Hausman. Efficient adaptation for end-to-end vision-based robotic manipulation. *CoRR*, abs/2004.10190, 2020. URL <https://arxiv.org/abs/2004.10190>.

- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Saurabh Kumar, Henrik Marklund, and Benjamin Van Roy. Maintaining plasticity via regenerative regularization. *arXiv preprint arXiv:2308.11958*, 2023.
- Clare Lyle, Mark Rowland, and Will Dabney. Understanding and preventing capacity loss in reinforcement learning. In *International Conference on Learning Representations*, 2022.
- Clare Lyle, Zeyu Zheng, Evgenii Nikishin, Bernardo Avila Pires, Razvan Pascanu, and Will Dabney. Understanding plasticity in neural networks, 2023.
- H Brendan McMahan and Francesco Orabona. Unconstrained online linear learning in hilbert spaces: Minimax algorithms and normal approximations. In *Conference on Learning Theory*, pp. 1020–1039. PMLR, 2014.
- Jorge A. Mendez, Boyu Wang, and Eric Eaton. Lifelong policy gradient learning of factored policies for faster training without forgetting, 2020.
- Golnaz Mesbahi, Olya Mastikhina, Parham Mohammad Panahi, Martha White, and Adam White. Tuning for the unknown: Revisiting evaluation strategies for lifelong rl, 2024.
- Zakaria Mhammedi and Wouter M Koolen. Lipschitz and comparator-norm adaptivity in online learning. In *Conference on Learning Theory*, pp. 2858–2887. PMLR, 2020.
- Evgenii Nikishin, Max Schwarzer, Pierluca D’Oro, Pierre-Luc Bacon, and Aaron Courville. The primacy bias in deep reinforcement learning. In *International Conference on Machine Learning*, pp. 16828–16847. PMLR, 2022.
- Evgenii Nikishin, Junhyuk Oh, Georg Ostrovski, Clare Lyle, Razvan Pascanu, Will Dabney, and André Barreto. Deep reinforcement learning with plasticity injection, 2023.
- Francesco Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*, 2023.
- Francesco Orabona and Dávid Pál. Coin betting and parameter-free online learning. *Advances in Neural Information Processing Systems*, 29, 2016.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Wenling Shang, Kihyuk Sohn, Diogo Almeida, and Honglak Lee. Understanding and improving convolutional neural networks via concatenated rectified linear units. In *international conference on machine learning*, pp. 2217–2225. PMLR, 2016.
- Ghada Sokar, Rishabh Agarwal, Pablo Samuel Castro, and Utku Evci. The dormant neuron phenomenon in deep reinforcement learning. *arXiv preprint arXiv:2302.12902*, 2023.
- Matthew Streeter and Brendan McMahan. No-regret algorithms for unconstrained online convex optimization. *Advances in Neural Information Processing Systems*, 25, 2012.
- S. Thrun. *Explanation-Based Neural Network Learning: A Lifelong Learning Approach*. Kluwer Academic Publishers, Boston, MA, 1996.
- Zhiyu Zhang, Ashok Cutkosky, and Ioannis Paschalidis. Pde-based optimal strategy for unconstrained online learning. In *International Conference on Machine Learning*, pp. 26085–26115. PMLR, 2022.
- Zhiyu Zhang, Heng Yang, Ashok Cutkosky, and Ioannis C Paschalidis. Improving adaptive online learning using refined discretization. In *International Conference on Algorithmic Learning Theory*, pp. 1208–1233. PMLR, 2024.
- Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *International Conference on Machine Learning*, pp. 928–936, 2003.

Appendix

A Further Experiments and Numerical Results

Arcade Learning Environment (ALE) Atari The ALE Atari 2600 benchmark tests RL agents across diverse arcade games (Bellemare et al., 2012). We introduce distribution shifts by switching to a new game every 4 million timesteps, posing a greater challenge than OpenAI Procgen by requiring adaptation to changes in both state and reward.

We evaluated two settings with action spaces of 6 and 9. As shown in Figure 5, both ADAM PPO and CReLU often failed in some games, while PACE PPO consistently increased rewards across different games. During the first 12 million steps (3 games) in Atari 6, PACE PPO achieved higher mean rewards and rapid increases. Overall, PACE PPO showed an average improvement of 329.73% over ADAM PPO and 68.71% over CReLU (Table 1). In some instances, like the last 2 million steps of Atari 6, CReLU performed comparably to PACE PPO, aligning with findings that CReLU can prevent plasticity loss in continual Atari setups (Abbas et al., 2023).

Numerical Results Our numerical results for the main experiments (Procgen, Atari, Gym Control) are summarized below in Table 1.

Table 1: Cumulative sum of mean episode reward for PACE PPO, ADAM PPO, and CReLU on Procgen, Atari, and Gym Control environments. Rewards are scaled by 10^5 ; higher is better.

Environment	Adam PPO	CReLU	PACE PPO (Ours)
StarPilot	3.4	3.6	12.5
Dodgeball	1.9	2.3	5.2
Chaser	1.4	1.7	2.2
CartPole	5.1	1.2	39.6
Acrobot	-14.3	-13.9	-12.9
LunarLander	-21.7	-19.4	-8.6
Atari 6	3.1	4.8	10.5
Atari 9	3.9	17.0	20.2

B On the choice of θ_{ref}

In general, the reference point θ_{ref} should be good or “safe” for PACE to perform effectively. One might presume that achieving this requires “warmstarting”, or pre-training using the underlying BASE optimizer. While our experiments validate that such warmstarting is indeed beneficial (Appendix C), our main experiments show that even a random initialization of the policy’s weight serves as a good enough θ_{ref} , even when tasks are similar (Figure 2).

This observation aligns with discussions by Lyle et al. (2023), Sokar et al. (2023), and Abbas et al. (2023), who suggested that persistent gradient steps away from a random initialization can deactivate ReLU activations, leading to activation collapse and loss of plasticity in neural networks. Our results also support Kumar et al. (2023)’s argument that maintaining some weights close to their initial values not only prevents dead ReLU units but also allows quick adaptation to new distribution shifts.

C Warmstarting

In our theoretical framework, we hypothesize that a robust parameter initialization, denoted as θ_{ref} , could enhance the performance of our models, suggesting that empirical implementations might benefit from initializing parameters using a base optimizer such as ADAM prior to deploying PACE. Contrary to this assumption, our experimental results detailed in Section 4 reveal that warmstarting is not essential for PACE’s success. Below, we examine the performance of ADAM PPO and PACE PPO when warmstarted.

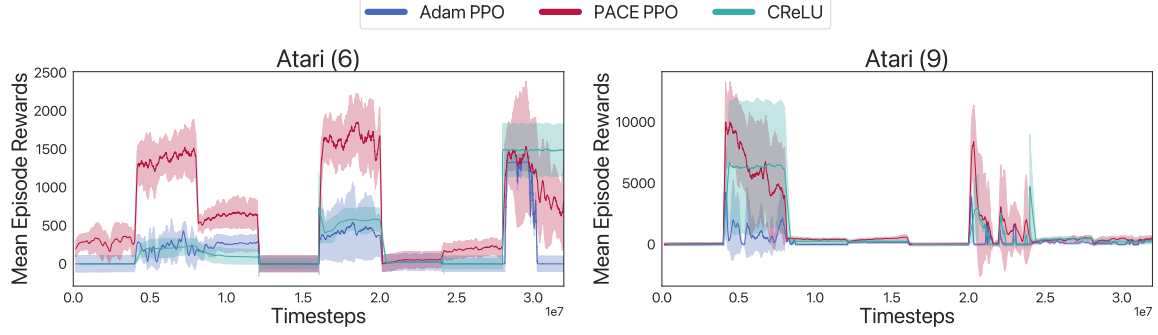


Figure 5: Reward in the lifelong Atari environments, across games with action spaces of 6 and 9. These plots demonstrate that PACE PPO rapidly adapts to new tasks, in contrast to the ADAM PPO and CReLU which struggle to achieve high reward, indicating mild loss of plasticity.

Both PACE PPO and ADAM PPO were warmstarted using ADAM for the initial 150,000 steps in all games for the Atari and Procgen environments, and for the first 30 steps in the Gym Control experiments. As seen in Figure 6, in games like Starpilot and Dodgeball, PACE PPO surpasses ADAM PPO in the first “level” of the online setup, with its performance closely matching that of ADAM PPO in Chaser. Importantly, PACE PPO continues to circumvent the loss of plasticity encountered by ADAM PPO, even when both are warmstarted. This makes sense since all of the distributions share some foundational game dynamics; the initial learning phases likely explore these dynamics, so leveraging a good parameter initialization to regularize in this early region can be beneficial for PACE—we observe that forward transfer occurs somewhat in later level distribution shifts as the reward does not drop back to zero where it initially started from.

Our findings indicate that warmstarting does not confer a significant advantage in the Atari games. This makes sense because a parameter initialization that is good in one game setting is likely a random parameterization for another setting, which is equivalent to the setup without warmstarting where PACE regularizes towards a random parameter initialization. In the Gym Control experiments although warmstarted PACE PPO manages to avoid the extreme plasticity loss and policy collapse seen in warmstarted ADAM PPO, it does not perform as well as non-warmstarted PACE PPO. This variability underscores that the efficacy of warmstarting is environment-specific and highlights the challenge in predicting when ADAM PPO may achieve a parameter initialization that is advantageous for PACE PPO to regularize towards. From an overall perspective, warmstarting PACE PPO in every setting still shows substantial improvement over ADAM PPO (Table 2)

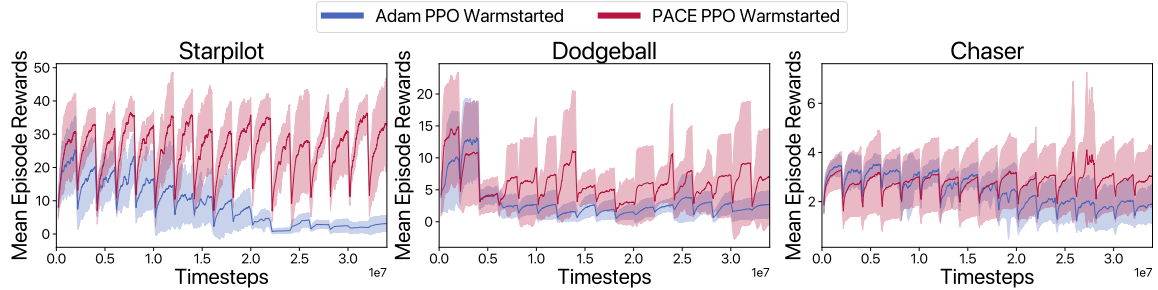


Figure 6: Comparison of reward in the lifelong Procgen environments for StarPilot, Dodgeball, and Chaser with warmstarted PACE PPO and warmstarted ADAM PPO. Initial performance of PACE PPO is improved with warmstarting and continues to avoid loss of plasticity.

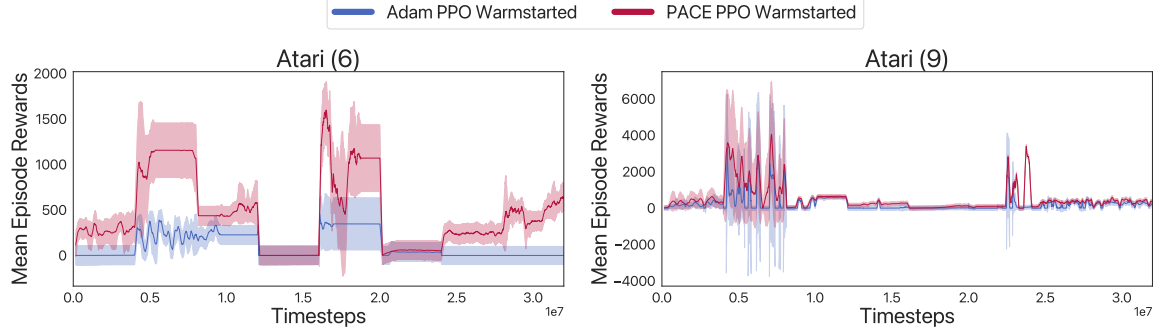


Figure 7: Comparison of reward in the lifelong Atari environments with warmstarted PACE PPO and warmstarted ADAM PPO. No significant benefit is found by warmstarting PACE compared to not warmstarting it.

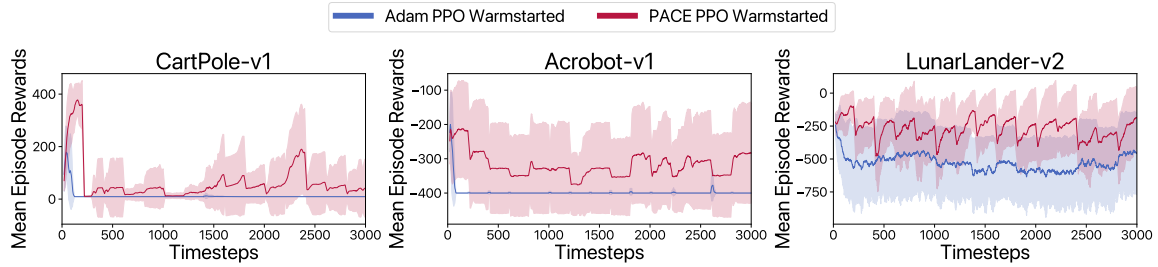


Figure 8: Comparison of reward in the lifelong Gym control environments for CartPole-v1, Acrobot-v1, and LunarLander-v2 with warmstarted PACE PPO and warmstarted ADAM PPO.

Table 2: Cumulative sum of mean episode reward over all distributions for ADAM PPO warmstarted and PACE PPO warmstarted on Procgen, Atari, and Gym Control environments. Rewards are scaled by 10^5 ; higher is better.

Environment	Adam PPO	Pace PPO (Ours)
Starpilot	3.0	10.2
Dodgeball	1.2	2.5
Chaser	1.3	1.6
CartPole	4.6	22.8
Acrobot	-142.9	-114.5
LunarLander	-190.7	-97.3
Atari6	16.7	72.2
Atari9	34.6	80.6

D Gravity Based Distribution Shifts

One method to introduce distribution changes in reinforcement learning environments is by altering the dynamics [Mendez et al. \(2020\)](#), such as adjusting the gravity in the CartPole environment. In this set of experiments, we manipulate the gravity by a magnitude of ten, randomly adding noise for one distribution shift, and then inversely, dividing by ten and adding random noise for the next shift. This process continues throughout the experiment.

Our observations suggest that ADAM PPO is robust to such dynamics-based distribution shifts, as shown in Figure 9. This indicates that while ADAM PPO implicitly models the dynamics of the

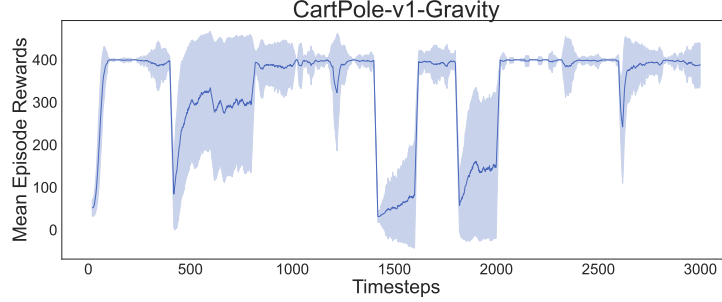


Figure 9: Mean Episode Reward for ADAM PPO on CartPole-v1 with varying gravity. ADAM PPO demonstrates robust policy recovery across most gravity-based distribution shifts.

environment well—where changes in dynamics minimally impact performance—it struggles more with adapting to out-of-distribution observations such as seen in the main experiments (Figure 3) and in the warmstarting experiments (Figure 8).

E Scaling-Value Convergence

As discussed in the algorithm section (see Section 3), PACE operates as a meta-algorithm on top of a standard optimizer, denoted as BASE. The crucial component of PACE involves the dynamic adjustment of the scaling parameter S_{t+1} , managed by the tuner algorithm (Algorithm 2). This parameter is data-dependent and typically ranges between $[0, 1]$. The weight update θ_{t+1} is consequently defined as a convex combination of the current optimizer’s weight θ_t^{BASE} and a predetermined reference point θ_{ref} .

This section presents the convergence behavior of the scaling parameter S_{t+1} across different environments, analyzed through the mean values over multiple seeds.

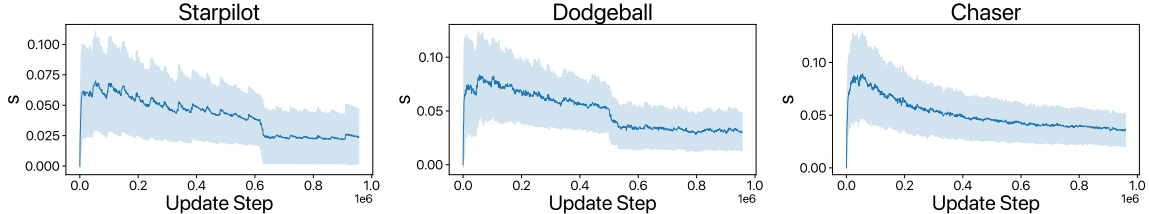


Figure 10: Convergence of the scaling parameter S_{t+1} in the ProcGen environments.

The convergence of the scaling parameter S_{t+1} observed across the Procgen and Gym Control environments, as depicted in Figures 10 and 12, reflects a good scaling value that effectively determines the strength of regularization towards the initialization points, yielding robust empirical outcomes in lifelong RL settings. Interestingly, in Procgen environments, this converged scaling value exhibits consistency across various games, typically hovering between 0.02 and 0.03, as shown in Figure 10. In contrast, in Control environments, the scaling values are notably lower, ranging between 0.005 and 0.01, as illustrated in Figure 12.

F Comparison to Mechanic

Our algorithm PACE builds on a long line of works on parameter-free OCO (see Section 3). In particular, we are inspired by the MECHANIC algorithm. Compared to MECHANIC, PACE improves the scale tuner there (which is based on the *coin-betting* framework; Orabona & Pál, 2016) by the

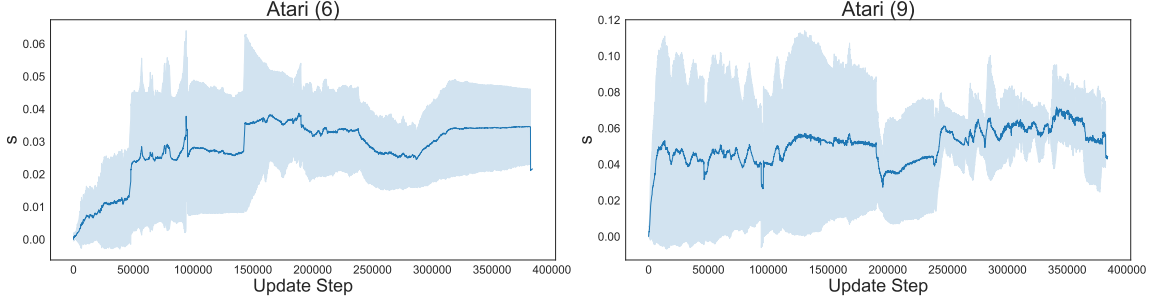


Figure 11: Evolution of the scaling parameter S_{t+1} in the Atari environments. Here we don't see a meaningful convergence of S_{t+1} .

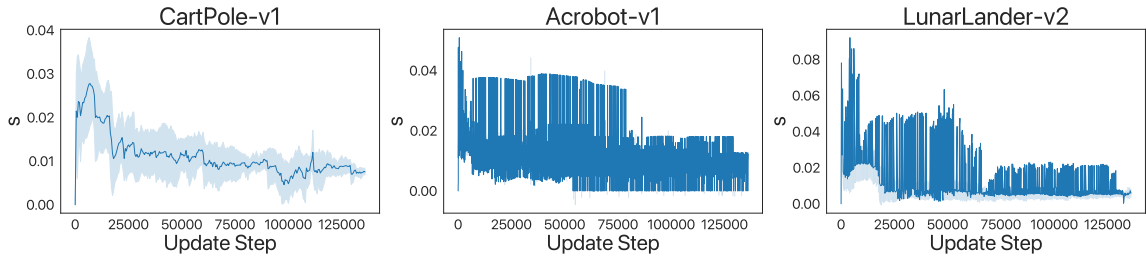


Figure 12: Convergence of the scaling parameter S_{t+1} in the Gym Control environments.

erfi algorithm that enjoys a better OCO performance guarantee. We empirically compare PACE and MECHANIC in the Starpilot game from the Procgen suite (Figure 13). The results indicate that while MECHANIC effectively mitigates plasticity loss and adapts quickly to new distribution shifts, it slightly underperforms in comparison to PACE. This suggests potential for the effectiveness of the general “parameter-free” principle in lifelong RL.

G Experimental Setup

Procgen and Atari Vision backbone For both the Atari and Procgen experiments, the Impala architecture was used as the vision backbone. The Impala model had 3 Impala blocks, each containing a convolutional layer followed by 2 residual blocks. The output of this is flattened and connected to a fully connected layer. The impala model parameters are initialized using Xavier uniform initialization.

Policy and Value Networks Across all experiments—including Control, Atari, and Procgen—the policy and value functions are implemented using a multi-layer perceptron (MLP) architecture. This architecture processes the input features into action probabilities and state value estimates. The MLP comprises several fully connected layers activated by ReLU. The output from the final layer uses a softmax activation.

PACE PACE, for all experiments, was implemented using the same experiment-specific baseline architectures and baseline optimizer. For the Procgen and Atari experiments, the base ADAM optimizer was configured as the same as baseline, with a learning rate of 0.001, and for Control experiments, a learning rate of 0.01 was used. Both learning rates were tested for all experiments and found to have negligible differences in performance outcomes. Other than the learning rate, we

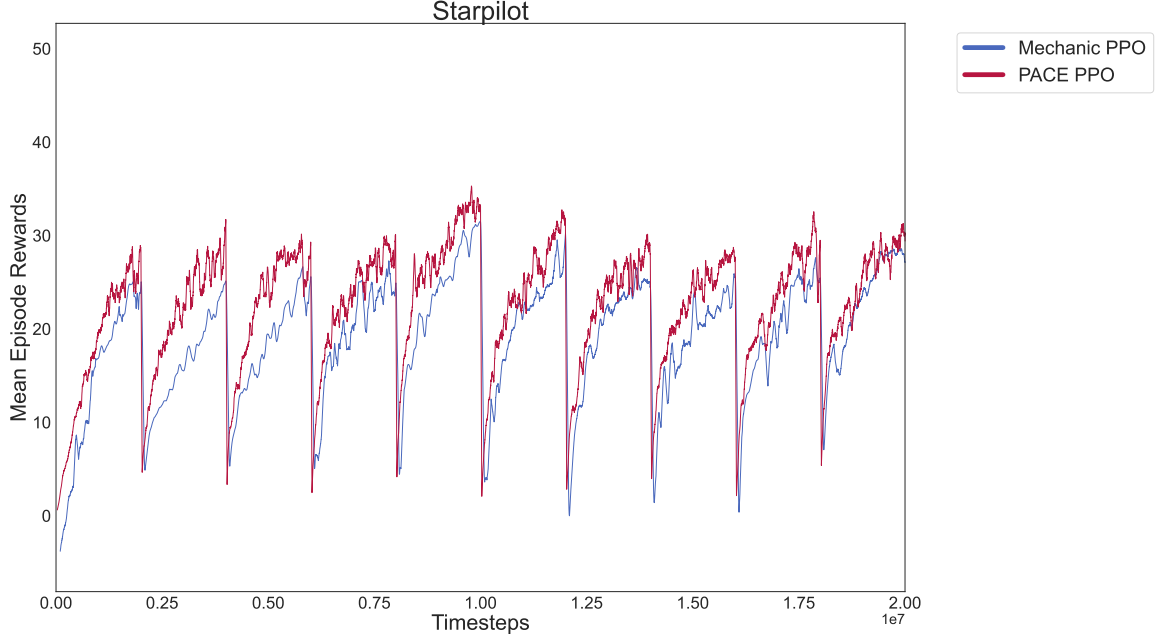


Figure 13: Comparison of reward in the lifelong StarPilot environment with both PACE PPO and MECHANIC PPO. MECHANIC PPO performs similarly to PACE PPO, although slightly underperforming PACE PPO.

Table 3: PPO Parameters for Atari, Procgen, and Control Experiments

Parameter	Atari	Procgen	Control
Steps per update	2,000	1,000	800 (2 episodes with 400 steps)
Batch size	250	125	32
Epochs per update	3	3	5
Epsilon clip for PPO	0.2	0.2	0.2
Value coefficient	0.5	0.5	0.5
Entropy coefficient	0.01	0.01	0.01
Base Optimizer	ADAM (LR: 0.001)	ADAM (LR: 0.001)	ADAM (LR: 0.01)
Architecture	Impala + MLP	Impala + MLP	MLP

use the default ADAM parameters, including weight decay and betas, followed by the specifications outlined in the PyTorch Documentation.⁴

The setup for PACE included β values for adaptive gradient adjustments: 0.9, 0.99, 0.999, 0.9999, 0.99999, and 0.999999. Both S_t and ε were initially set to (1×10^{-8}) . Modifications were made to a PyTorch error function library, which accepts complex inputs to accommodate the necessary computations for the imaginary error function. This library can be found at Torch Erf GitHub.⁵

Distribution Shifts In the Atari experiments, game environments were switched every 4 million steps. The sequence for games with an action space of 6 included “BasicMath”, “Qbert”, “SpaceInvaders”, “UpNDown”, “Galaxian”, “Bowling”, “Demonattack”, “NameThisGame”, while games with an action space of 9 included “LostLuggage”, “VideoPinball”, “BeamRider”, “Asterix”, “Enduro”, “CrazyClimber”, “MsPacman”, “Koolaid”.

⁴<https://pytorch.org/docs/stable/generated/torch.optim.Adam.html>

⁵https://github.com/redsnic/torch_erf

For Procgen experiments, individual game levels were sampled using a seed value as the *start_level* parameter, which was incremented sequentially to generate new levels. Each new environment was introduced every 2 million steps.

In the Gym Control experiments, each observation dimension was randomly perturbed by a value ranging from 0 to 2. This perturbation was constant for 200 timesteps, after which a new perturbation was applied, effectively switching the environmental conditions every 200 steps.

Statistical Significance The Procgen and Atari experiments were conducted with 8 seeds/runs, while the Gym Control experiments utilized 25 seeds/runs. The exception was in the L_2 initialization experiments, which used 15 seeds/runs per regularization strength. In Figures 2, 5, 3, 4, 6, 7, 8, 9, the plotted lines represent the mean of all of the Mean Episode Rewards from the different seeds/runs, and the shaded error bands indicate the standard deviation of all of the Mean Episode Rewards from the different seeds/runs.

Compute Resources For the Procgen and Atari experiments, each was allocated a single A100 GPU, typically running for 3-4 days to complete. The Gym Control experiments were conducted using dual-core CPUs, generally concluding within a few hours. In both scenarios, an allocation of 8GB of RAM was sufficient to meet the computational demands.