

DO SYMBOLIC OR BLACK-BOX REPRESENTATIONS GENERALISE BETTER IN LEARNED OPTIMISATION?

Anonymous authors

Paper under double-blind review

ABSTRACT

Until recently, behind *every* algorithmic advance in machine learning was a human researcher. Now, however, algorithms can be *meta-learned automatically*, with little human input. However, to be truly useful, such algorithms must generalise beyond their training distribution. This is especially challenging in reinforcement learning (RL), where transferring algorithms between environments with vastly different dynamics is difficult and training on diverse environments often requires prohibitively expensive large-scale data collection. Learned optimisation is a branch of algorithmic discovery that meta-learns optimiser update rules. Learned optimisers can be classified into two groups: black-box algorithms, where the optimiser is a neural network; or symbolic algorithms, where the optimiser is represented using mathematical functions or code. While some claim that symbolic algorithms generalise better than black-box ones (Chen et al., 2023), testing such assertions is complicated by the fact that symbolic algorithms typically include additional hyperparameters, and thus their evaluation is done *many-shot*. This is an unfair comparison with the *zero-shot* evaluation of black-box optimisers. In this work, we build a pipeline to discover symbolic optimisers which are *hyperparameter-free*, enabling a fair comparison of the generalisation of symbolic optimisers with that of an open-source state-of-the-art black-box optimiser trained for RL¹. Based on our analysis, we propose suggestions to improve the symbolic optimiser discovery pipeline for RL, with an overall objective of reducing the need for hyperparameter tuning to train an agent.

1 INTRODUCTION

Improvements to optimisation algorithms have driven machine learning to new heights over the past few decades. The introduction of components like gradient momentum, second order momentum (Nesterov, 1983; Kingma & Ba, 2017) and adaptive learning rates (Kingma & Ba, 2017; Zhuang et al., 2020) has enabled swifter and more stable convergence, while learning rate annealing has improved the fidelity of converged solutions. Recent evidence (Andrychowicz et al., 2016; Chen et al., 2021; Metz et al., 2022c; Chen et al., 2023; Goldie et al., 2024) suggests that the improvement of optimisers could be automated via *learned optimisation*. In learned optimisation, developing new optimisation algorithms is itself a *meta-learning* process based on data.

Approaches to learned optimisation fall into two camps. Most work (e.g., (Metz et al., 2022a; Kirsch & Schmidhuber, 2022; Andrychowicz et al., 2016; Wichrowska et al., 2017; Goldie et al., 2024)) replaces the optimiser, such as Adam (Kingma & Ba, 2017), with a black-box function using a neural network. In this scenario, the weights of the network are updated in an *outer loop* to maximise the performance of a trained model at the end of an *inner loop*. By contrast, some recent work (Chen et al., 2023; Song et al., 2024a) focuses on discovering *symbolic* optimisation algorithms. In this case, the optimiser is represented by a set of mathematical equations or programming instructions. In general, interest in symbolic algorithm discovery has grown in the past couple of years (Romera-Paredes et al., 2024; Lu et al., 2024a) due to the advent of large language models (OpenAI et al., 2024; Dubey et al., 2024, LLMs). There are arguments in favour of both approaches: black-box algorithms *may* be easier to work with (Goldie et al., 2024), while symbolic optimisers *may* generalise better (Chen et al., 2023). However, there exists little study into the veracity of these claims.

¹Code to be released upon acceptance.

054 Furthermore, direct comparison between the approaches is complicated by the fact that they target
 055 subtly different problems; black-box optimisers are typically evaluated zero-shot without any tune-
 056 able hyperparameters, whereas symbolic optimisers such as Lion (Chen et al., 2023) tune hyperpa-
 057 rameters *per-task*, making evaluation *many-shot*. Therefore, it is hard to compare these different
 058 paradigms like-for-like based on current literature.

059 The need for general optimisation algorithms is exacerbated in reinforcement learning (Sutton &
 060 Barto, 2018, RL) due to its many idiosyncratic issues which make optimisation challenging. In par-
 061 ticular, RL is very sensitive to hyperparameters (Eimer et al., 2023) which can cause catastrophic
 062 instability if they are not correctly tuned. This instability may stem from the fact that RL often
 063 **uses** algorithms imported from supervised learning, motivating the development of RL-specific ap-
 064 proaches (Henderson et al., 2018; Sarigül & Avci, 2017). For instance, many conventional optimis-
 065 ers, like Adam Kingma & Ba (2017), are designed for stationary learning tasks and are thus ill-suited
 066 for the non-stationarity of RL (Igl et al., 2021; Bengio et al., 2021). Learned optimisers tailored for
 067 RL show promise in addressing these issues (Lan et al., 2024; Goldie et al., 2024).

068 However, simply relying on a large meta-task diversity to enable generalisation across RL is im-
 069 practical. For anything beyond simple environments, sampling in RL is expensive. Therefore, find-
 070 ing learned optimisation strategies which demonstrate generalisation, whilst maintaining a limited
 071 meta-training cost, would significantly improve the practicality of RL. In this work, we compare
 072 the generalisation capabilities of a pretrained, black-box optimiser for RL (Goldie et al., 2024) with
 073 a roughly equivalent symbolic optimiser discovered using an evolutionary process based around
 074 LLMs. We focus on a regime in which optimisers can only be learned from a small number of en-
 075 vironments; we believe this represents a scenario of greater interest than training in a distribution of
 076 gridworlds, which has been a previous focus for generalisation (Goldie et al., 2024; Lan et al., 2024)
 077 but does not transfer well to the modern LLM-driven discovery pipeline. In doing so, we explore
 078 the question of whether black-box or symbolic optimisers are *actually* best for generalisation across
 079 a number of axes, including to different environments and to longer training lengths. We use these
 080 findings to recommend promising directions for future work in this field, thus providing a pathway
 081 to unlock truly general learned optimisation algorithms.

082 2 BACKGROUND

083 **Optimisation** Optimisation is ubiquitous throughout machine learning. Given a general training
 084 objective $f_{\theta}(\cdot)$, there is an extensive set of optimisation algorithms whose goal is to guide θ , a
 085 model’s parameters, to the optimal θ^* . Most fundamental of optimisers is gradient descent, where θ
 086 is updated iteratively towards negative gradient as $\theta_{t+1} \leftarrow \theta_t - \eta \nabla_{\theta} f(\cdot)$, using a step-size η .
 087

088 A number of augmentations are frequently applied to gradient descent to enable quicker conver-
 089 gence, less noisy updates or improved asymptotic performance. For instance, modern optimisers
 090 like Adam (Kingma & Ba, 2017) and RMSProp (Tieleman et al., 2012) use *momentum*, a time-
 091 based moving average of gradients or updates which provides more consistent updates over training.
 092 Similarly, learning rate *annealing* or *warmup* change the step size over time to provide closer con-
 093 vergence to the optimum by the end of training, or improved stability at the beginning of training,
 094 respectively (Robbins, 1951; Gotmare et al., 2018).

095 **Reinforcement Learning** Reinforcement learning focuses on Markov Decision Processes (Sutton
 096 & Barto, 2018, MDPs), defined as $\langle \mathcal{A}, \mathcal{S}, T, R, \rho, \gamma \rangle$. The *agent* learns a policy $\pi(\cdot|s_t) \in \Pi$ and,
 097 at each discrete timestep t , samples an action $a_t \in \mathcal{A}$ based on the current state $s_t \in \mathcal{S}$ (where
 098 $s_0 \sim \rho$). After sampling an action, the agent transitions to the next state $s_{t+1} \in \mathcal{S}$ according
 099 to a transition distribution $T(s_{t+1}|s_t, a_t)$ and receives a reward according to the reward function
 100 $R(s_t, a_t)$. The policy is trained to maximise the *discounted expected return*, J^{π} , based on the
 101 discount factor $\gamma \in [0, 1)$, which is defined over a fixed length episode as

$$102 \quad J^{\pi} := \mathbb{E}_{a_0, \infty \sim \pi, s_0 \sim \rho, s_1: \infty \sim T} \left[\sum_{t=0}^T \gamma^t R_t \right]. \quad (1)$$

105 Sample complexity is a major issue in reinforcement learning. Due to the potential cost of inter-
 106 acting with the environment, it can often be prohibitively expensive to collect large datasets. One
 107 opportunity to reduce sample complexity is to remove the reliance on hyperparameters intrinsic to
 RL. Learned optimisers without hyperparameters could help to unlock this capability.

Optimisation Difficulties in RL Goldie et al. (2024) discuss three optimisation difficulties present in RL: plasticity loss (Lyle et al., 2023; 2022), a phenomenon in which neural networks *lose* the ability to learn when given new data; exploration, where the optimiser must escape local optima from the agent being trapped in a localised state-action space; and non-stationarity (Igl et al., 2021), which arises as the input and output distributions in RL are continuously changing. OPEN incorporates a number of features to tackle each individual problem. To be specific:

- For plasticity, OPEN conditions on neuron dormancy (Sokar et al., 2023), a metric which measures what proportion of a layer’s activation comes from a specific neuron. Near-zero dormancy neurons are dormant and need to be reactivated. OPEN also learns separate update rules for each layer by conditioning on *layer proportion*.
- For nonstationarity, OPEN conditions on two timescales: *batch proportion*, or progress through epochs with the current batch of data; and *training proportion* (Jackson et al., 2023a), meaning how far through the training horizon optimisation is.
- To boost exploration, OPEN introduces stochasticity of a learned variance to the update. This enables similar exploration behaviour to parameter space noise (Plappert et al., 2018) or noisy nets (Fortunato et al., 2019) while also incidentally helping with dormancy.

3 RELATED WORK

Meta-Learning Algorithms Meta-learning intends to replace handcrafted algorithms with ones learned from data. Though some approaches use *meta-gradients* which are backpropagated through training episodes (e.g., (Lan et al., 2024; Oh et al., 2020)), this is impractical in our setting. Firstly, meta-learning in RL requires long horizon rollouts, where untruncated backpropagation experiences exploding or vanishing gradients but truncating biases towards greedy algorithms (Wu et al., 2018; Metz et al., 2022b; Lu et al., 2022b). Secondly, with a *symbolic* optimiser, it is not obvious how to project gradients on to the non-numerical symbols of our algorithm, **requiring more complex techniques** (e.g. (Kuang et al., 2024; Chen et al., 2024)).

Evolutionary methods (Rechenberg, 1973; De Jong, 2006) provide an alternative. These are derivative-free optimisation methods which mutate and evaluate a populations of candidates. Common evolutionary methods include genetic algorithms (Such et al., 2018), covariance matrix adaptation (Hansen & Ostermeier, 2001), evolution strategies (Salimans et al., 2017) or, in the symbolic case, genetic programming (Koza, 1992). Evolution involves sequentially sampling population members, randomly changing their parameters and evaluating the final performance of the candidate. By optimising based on the final evaluation, rather than backpropagating *through* the rollout, evolutionary methods avoid many of the issues with meta-gradients.

Since the advent of LLMs, a new form of symbolic evolution has emerged (Romera-Paredes et al., 2024). Rather than applying *random* mutations, recent methods have replaced the evolutionary system with LLMs that suggest edits and reason about performance to guide search (Lu et al., 2024a; Meyerson et al., 2024; Lehman et al., 2022; Shojaee et al., 2024). This uses an LLM’s prior knowledge to make ‘intelligent’ changes, in effect limiting the search to reasonable if not limited edits. Despite its recent invention, this technique has led to impressive results in function discovery (Romera-Paredes et al., 2024) or solving symbolic regression tasks (Shojaee et al., 2024).

Learned Optimisation Learning to optimise (Metz et al., 2020; 2022c;a; Chen et al., 2023; Goldie et al., 2024, L2O) automates the discovery of better optimisers by *meta-learning* the algorithms. Generally, L2O replaces the optimiser with a neural network which conditions on the gradient, and potentially extra features, and outputs an update for *each parameter* in the training model. This method has proven effective in supervised and unsupervised learning (Metz et al., 2022c), but naïvely fails to transfer to RL. Due to the opportunity of learning *specialised* optimisation algorithms, OPEN (Goldie et al., 2024) and Optim4RL (Lan et al., 2024) L2O directly for RL. This is justified by many works suggesting RL-specific algorithms are warranted (Henderson et al., 2018; Bengio et al., 2021; Sarigül & Avci, 2017). Whereas Optim4RL attempts to L2O in RL by constraining the structure of the update, OPEN targets a number of difficulties present only in RL. Unfortunately, while these works have demonstrated signs of life for generalisation, there is little work exploring whether black-box optimisation is the best route to discover truly generalist optimisers.

An alternative approach is Lion (Chen et al., 2023), an optimiser discovered by *symbolic evolution*. However, to enable comparison between black-box and symbolic optimisation, we make a number

of key design changes from Lion. Firstly, our method searches in a **code**, rather than mathematical, parameterisation. This enables a richer space of functions by allowing conditional statements, like (if , $>$, $<$). Secondly, by building on modern LLM-based methods, we diverge from Lion’s naïve mutation operation. Since we attempt to *directly* compare against OPEN, whose inputs expands the algorithm design space drastically, the prior knowledge of an LLMs limits search to grounded mutations, thus preventing an excessive computation budget. Finally, we direct our search towards hyperparameter-free optimisers for RL to enable a fair comparison with OPEN.

LLM-Guided Research LLMs have increasingly been used for evolution-like optimisation recently (Song et al., 2024b). FunSearch (Romera-Paredes et al., 2024) demonstrated the validity of this approach by prompting an LLM to write functions for specific tasks. Like FunSearch, many works have synthesised the expressiveness of code with the creativity of LLMs: Hu et al. (2024) use LLMs to *design* agents for complex problems; DiscoPOP (Lu et al., 2024a) finds new objectives for preference optimisation in LLMs; and Lehman et al. (2022) incorporate Quality-Diversity approaches (Mouret & Clune, 2015) to produce different robot morphologies. While a common thread exists between these works and ours – using LLMs as a mutation operator for evolution – our discovery pipeline differs in its end-goal of learning an *optimisation algorithm*. We also consider how an LLM can be used to handle additional inputs, defined by OPEN, with natural language descriptions. Finally, we are approaching this setting from a purely analytical perspective.

4 MOTIVATION

To motivate our study into the generalisation capabilities of symbolic and black-box optimisers, we briefly compare the two in terms of potential advantages, grounded in both literature and intuition.

Black-Box Optimisers Since black-box optimisers are principally neural networks, they have a number of inherent advantages. Firstly, since they typically use *small* networks, they can easily be trained with evolution (Salimans et al., 2017) to avoid issues of short-term bias from truncated meta-gradients (Wu et al., 2018; Lu et al., 2022b). This does, however, have the issue of high memory usage and training sample complexity since each meta-update needs a number of full training loops equal to the population size. Though GPU-vectorisation (Bradbury et al., 2018) helps speed up this training dramatically (Lu et al., 2022b), it can require both high-end hardware and easy-to-sample environments which may not be practical.

Also, the simplicity of introducing additional inputs to black-box optimisers was demonstrated by OPEN, as well as an ease to learn interactions between input variables. This ability to easily scale with inputs may make black-box optimisers the best option in some settings.

Finally, due to their iterative meta-learning process, black-box optimisers can converge Goldie et al. (2024). This is in contrast to symbolic optimisers, which may not converge due to the mechanisms of symbolic evolution. This convergence can have advantages – training is predictable and usually stable – but can also lead to the optimiser being trapped in subpar optima.

Symbolic Optimisers Though symbolic discovery of optimisers is relatively unexplored, it has a number of *potential* advantages. It is worth noting, however, that we focus on a novel evaluation regime which aligns symbolic and black-box optimisation. Whereas Lion (Chen et al., 2023) needed *tuning* for its hyperparameters, black-box optimisers are applied zero-shot to new environments. Therefore, we concern ourselves with symbolic algorithms which *do not use hyperparameters*.

In this paper, we assess how black-box and symbolic optimisation algorithms generalise. Chen et al. (2023) suggest, without justification, that symbolic algorithms *should* generalise better, which seems intuitive. Symbolic optimisers are usually simpler; whereas Lion is 8 lines of code, OPEN uses up to ~ 4000 parameters, increasing the opportunity for overfitting. Also, symbolic optimisers must start from *something*, meaning they can be initialised from pre-existing optimisers.

A key advantage of symbolic algorithm discovery is that LLMs can interface into the discovery pipeline to improve the search efficiency, leaning on their vast knowledge-base to find new algorithms Lu et al. (2024a); Romera-Paredes et al. (2024). This also gives a large amount of control to the human-in-the-loop. As a researcher can describe design specifications in natural language, the search can be biased towards algorithms based on design requirements. We find this can help with including additional inputs to the algorithms, such as those from OPEN, even if the inputs are not included in the LLM’s training data.

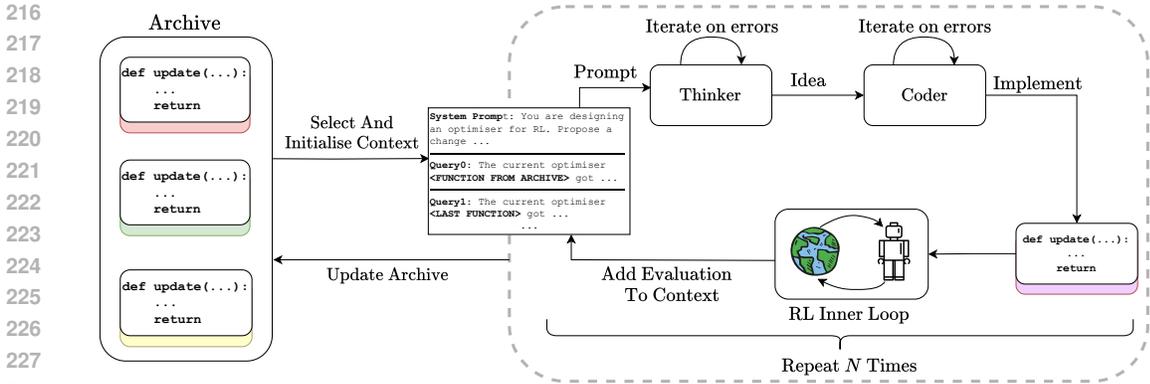


Figure 1: An overview of our discovery pipeline. An archive stores optimisers from previous generations. These are selected and used to initialise the LLM’s context. A ‘thinker’ LLM proposes an idea which the ‘coder’ LLM interprets and implements, producing a new optimiser. The new optimiser is evaluated, added to the context for the thinker, and the process repeats for a finite number of steps before all optimisers are added to the archive and the outer loop progresses.

5 THE SYMBOLIC OPTIMISER DISCOVERY PIPELINE

We design a symbolic discovery loop to enable like-for-like comparison with OPEN which incorporates all of the features proposed in OPEN and described in section 2. We focus our comparison on the ‘Multi-Task Training’ setting from Goldie et al. (2024), where we meta-train on a small number of environments from MinAtar (Young & Tian, 2019; Lange, 2022). We believe this scenario is particularly interesting due to its correspondence with learning from a small number of fast proxy-tasks that approximate an ultimate objective.

We use an LLM in place of standard mutation in our system for the reasons mentioned in section 4. This lets us describe the inputs from OPEN in natural language to direct the search to ‘reasonable’ suggestions, avoiding a potentially more expensive and sample-inefficient random search, like Lion (Chen et al., 2023). However, LLMs can be notoriously fickle (Anagnostidis & Bulian, 2024; Gu et al., 2022). Therefore, we introduce a number of design decisions, described in this section, to improve the system’s robustness. While we use GPT-4o (OpenAI et al., 2024) in this work, we believe that our system should also maintain robustness for weaker, open-source models (e.g. (Dubey et al., 2024; DeepSeek-AI et al., 2024)). We report discovery hyperparameters in Appendix A.

5.1 OVERVIEW

Figure 1 shows our discovery pipeline, which is similar to a number of ‘LLM-Discovery’ methods (Romera-Paredes et al., 2024; Lu et al., 2024a;b; Hu et al., 2024; Faldor et al., 2024), visually. At the start of the process, an archive is initialised with a set of candidate optimiser functions. After these are evaluated, one optimiser is selected for a generation of refinement, which involves iterative mutation by an LLM, followed by evaluation and insertion to the archive, for N steps. After refinement is complete, a new optimiser is sampled and the process repeats.

Below, we introduce high-level design decisions which are detailed in the remainder of section 5.

Initialising The Archive We follow DiscoPOP (Lu et al., 2024a) and Lion (Chen et al., 2023) by initialising training from a small set of optimisers. However, whereas DiscoPOP use pre-established loss functions, there is little precedent for hyperparameter-free optimisation. Therefore, we introduce a small number of hyperparameter-less optimisers by hand. These are designed to be flexible, while ensuring they don’t fail catastrophically in the training environments.

Selection Our pipeline periodically samples a new optimiser to refine at each generation, aligning closely to traditional evolutionary computation. This contrasts with, say, DiscoPOP (Lu et al., 2024a), which uses one long conversation with an LLM. By using the LLM more sparingly, this approach has the added benefit of potentially letting our system operate with less powerful language models. We select the best optimisers from the archive with probability p , and select random optimisers from the archive with an exploration probability $1 - p$.

Mutation We split mutation over two LLMs: a *thinker*, which proposes a new idea based on the current optimiser’s performance; and a *coder*, which implements the proposed changes. This separation ensures faithful interpretations of ideas in the implementation and provides additional user control with the different prompts. Our thinker prompt also includes examples of performant optimisers in each environment.

Evaluation We evaluate optimisers on full-length RL environments at every refinement step. We track the final return and return area-under-the-curve of each optimiser for the thinker’s context to enable in-context reasoning. To sidestep the problem of score aggregation over multiple environments faced by OPEN, we simply give the LLM returns for all environments and prompt it to maximise performance in all.

5.2 INITIALISATION

Similar to recent works (Lu et al., 2024a; Faldor et al., 2024; Hu et al., 2024; Chen et al., 2023), we initialise the archive of optimisers to a set of reasonable functions. However, given the scarcity of research on hyperparameter-free optimisation, the selection of initial optimisers is not straightforward. To address this, we create a few sensible optimisers to kickstart learning. In most cases, we write simple functions which have scaled *relative* changes to weights, though we also include a simple LLM-proposed function for diversity.

All optimisers follow the same design principles: they are simple, so that there are a large number of possible directions to improve them; they are diverse, so that they can lead to very different optimisers after refinement; and they are hyperparameter-free, meaning that any values are fixed for all environments. Notably, our initial optimisers only depend on the parameter value and the gradient, allowing the LLM to discover creative ways to use the additional inputs from OPEN without undue bias. We include all of the initial optimisers in appendix B.

5.3 EVOLUTION

For discovery, we blend LLM-based discovery algorithms with more conventional evolution (e.g. (Koza, 1992)). In doing so, we exploit the reasoning capabilities of LLMs to propose intelligent in-context changes while leveraging population-based evolution. The process runs as follows:

At the start of a new generation, we sample an ‘initial’ optimiser (section 5.3.1) and set of context optimisers (section 5.3.3) and prompt the LLMs to make small optimiser edits for a fixed number of *refinement steps*, N . At each refinement step, we evaluate the optimiser on *all* RL environments after a full RL inner-loop. Like OPEN, we use PPO Schulman et al. (2017) as the RL algorithm. After each generation, we add *all* evaluated optimisers to the archive and sample a new initialisation and context. Therefore, like Faldor et al. (2024), our archive grows over meta-training.

5.3.1 SAMPLING NEW OPTIMISERS

We sample a new ‘base’ optimiser each generation. To balance *exploration* and *exploitation* in our discovery process, we mostly sample *good* optimisers while occasionally selecting randomly to promote diversity. However, the notion of *good* or *bad* is not black and white when considering multiple environments of different reward scales. Naïvely averaging returns will prioritise environments which have a large reward scale, while normalising by, say, Adam’s (Kingma & Ba, 2017) performance biases selection to environments where Adam underperforms (Goldie et al., 2024).

Instead, we use the average of per-environment rankings, based on return, over the population to measure how successful an algorithm is. In addition to scale-invariance, this has the benefit of weeding out optimisers which overfit to one environment, aiding robustness. After calculating the average rankings for the population, we select high-ranking optimisers with a probability p and sample from the full population with probability $(1 - p)$. In this work, we set $p = 0.8$ to **balance sample efficiency (mostly starting from a performant optimiser) with diversity (occasionally sampling random optimisers)**.

5.3.2 MUTATION

We find that there is an occasional disparity between the proposal and implementation from LLMs when prompted naïvely. This hurts interpretability; it is not possible to tell what changes the LLM is making purely by observing the conversation. Therefore, we augment our system into a 2-LLM

setup by dividing out *thinking* and *coding*. The *thinker* has the responsibility of suggesting changes to the currently sampled optimiser and explaining why this change might be helpful. The *coder* has the task of converting the proposed idea into a code edit and implementing a syntactically correct, faithful python function. As an additional benefit, this allows different prompting strategies for each operation, giving the user additional control over the discovery trajectory.

5.3.3 PROMPTING

Different prompts can lead to vastly different results when using LLMs (Anagnostidis & Bulian, 2024; Gu et al., 2022). Here, we discuss the design decisions made in our prompting, and provide examples of the actual prompts in Appendix C.

Difficulties in RL To enable intelligent suggestions based on the problems of RL from OPEN, described in section 2, we provide a high level overview of each additional input variable and what typical values might mean.

Previous Performance To leverage in-context suggestion making, we condition the thinker on the returns of the current optimiser and randomly sampled ‘context optimisers’, which perform well in individual environments. To avoid issues highlighted in Goldie et al. (2024), where aggregating scores between different environments proved difficult, we include final return values for all environments into the prompt directly without averaging. This encourages the LLM itself to balance improvements between environments. To boost in-context reasoning further, we also provide values for the area-under-the-curve.

Separating Prompts To ensure fulfilment of their separate roles, we prompt the thinker and coder LLMs differently. The thinker is prompted to produce a new idea based on previous performance while the coder converts the idea into a code update. Whereas the thinker is prompted with a *history* of optimisers for reasoning, the coder receives only the current optimiser and proposed change to avoid obfuscating its task. **Separating thinking and code has been shown to improve performance in other work (Ye et al., 2024; Liu et al., 2024).**

Design Suggestions For both the coder and thinker, we propose a number of considerations to aid discovery. For instance, in the thinker we emphasise coming up with creative solutions, a need for generalisation and the necessity of not introducing new hyperparameters. For the coder, we focus on faithfulness and correctness, in addition to requesting commented code for interpretability.

6 DISCOVERY RESULTS

In figure 2, we show the meta-training curve for the symbolic discovery process. Notably, we find that, despite only selecting for high *average* fitnesses, our discovered symbolic optimisers have *consistently* high rankings across the four training environments. We also compute rankings for OPEN and Adam, with a standard *untuned* learning rate of $1e-3$. Based on their ranking compared to the population, neither Adam nor OPEN has robust performance across *all* environments. Below, we show the three highest average rank discovered optimisers which form the basis of our analysis.

The discovered optimisers below exhibit some similar behaviours. For instance, all optimisers incorporate dormancy into their updates, have annealing over training and use momentum. However, despite having sufficient inputs, none of the best optimisers manage to incorporate stochasticity (Goldie et al., 2024) into their expressions. This is likely due to the difficulty of finding a scale for the randomness which works for all environments in such a discrete search.

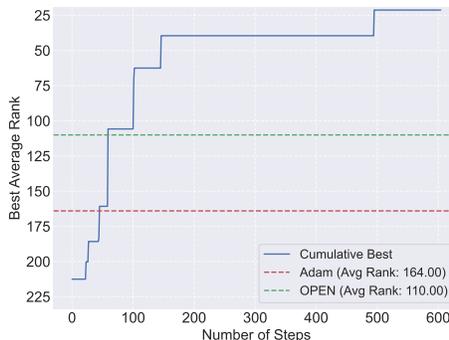


Figure 2: Meta-training curve, showing the max cumulative average rank of discovered optimisers. We also show where Adam and OPEN would rank in the population.

```

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

```

Discovered Optimiser 1	Discovered Optimiser 2	Discovered Optimiser 3
<code>def update:</code>	<code>def update:</code>	<code>def update:</code>
<code>m = 0.9</code>	<code>m = 0.9</code>	<code>m = 0.9</code>
<code>v1 = m * v1 + (1-m) * g</code>	<code>norm = g / (1 + g)</code>	<code>v1 = m * v1 + (1-m) * g</code>
<code>v2 = m * v2 + (1-m) * g**2</code>	<code>v1 = m*v1 + (1-m)*norm</code>	<code>v2 = m*v2 + (1-m)*(g-v1)**2</code>
<code>v2 = clip(v2, 1e-8, 1.0)</code>	<code>v2 = m*v2+(1-m)*(g-v1)</code>	<code>lr = 1 / (1+sqrt(v2+1e-8))</code>
<code>lr = sqrt((1-t_p) * (1+b_p))</code>	<code>lr = 1 / (1+ v2)</code>	<code>lr2 = (1-t_p)*(1+l_p)</code>
<code>lr = lr * (1+l_p)</code>	<code>boost=1+log(1+d)</code>	<code>d_scale = 1+log(1+d)*(1-t_p)</code>
<code>d_scale = 1 + log(1+d)</code>	<code>lr2 = (1-t_p)*(1+l_p)</code>	<code>d_scale *= (1+0.1*t_p)</code>
<code>lr2 = 1 / (1+v2)</code>	<code>update = v1*lr*boost*lr2*(1+b_p)</code>	<code>boost = where(d<1.0,2.0,1.0)</code>
<code>update = v1*lr*d_scale*lr2</code>	<code>return update, v1, v2</code>	<code>d_scale *= boost</code>
<code>return update, v1, v2</code>		<code>update = v1*lr*d_scale*lr2</code>
		<code>return update, v1, v2</code>

7 ASSESSING GENERALISATION

Our analysis centres on comparing symbolic discovered optimisers with OPEN to explore the difference between in- and out-of-distribution behaviour of the two approaches. We focus on meta-training with a small number of environments, referred to as **Multi-Task Training** in Goldie et al. (2024). This differs to the scenario where one samples from a distribution of simple environments, such as gridworlds (e.g. (Oh et al., 2020; Jackson et al., 2023b; Goldie et al., 2024)). We compare against a pre-trained OPEN model which is available online, and Adam using a fixed standard learning rate of $1e-3$. Following standard procedure in learned optimisation (Goldie et al., 2024; Metz et al., 2022c; Lan et al., 2024; Metz et al., 2019) arising from the cost of meta-learning, we discover optimisers from only one seed but run each experiment for multiple seeds. For all results, we report the interquartile mean (IQM) with 95% stratified bootstrap confidence intervals calculated using rliable, a standard evaluation library (Agarwal et al., 2021). Hyperparameters for all experiments are included in Appendix A. We consider a number of axes for generalisation, described and justified below, which are inspired by the comparison of OPEN and Adam in Goldie et al. (2024).

Different Training Lengths Due to the cost of learned optimisation, one way to speed up meta-training could be to learn from shortened inner-loops and generalise to longer runs. However, due to the nonstationarity of the optimisers from their time-conditioning, their dynamic behaviour may not transfer between inner-training lengths.

Different Architectures Prior work (Yang et al., 2022) suggests that hyperparameters often do not transfer between models with different architectures. As such, we explore the ability of the different optimisers to transfer between agents with different hidden dimensions and activation functions.

Different Environments To ensure an optimiser is truly general purpose, it is important to test its performance in unseen environments. This axis of generalisation explores how strongly an optimiser overfits to the *dynamics* of its training environments.

8 GENERALISATION RESULTS

Scaling to Different Lengths Figure 3 explores how the final return of an agent trained with each of the optimisers differs as the length of the training horizon increases. Here, $1e7$ transitions is in-distribution for each optimiser.

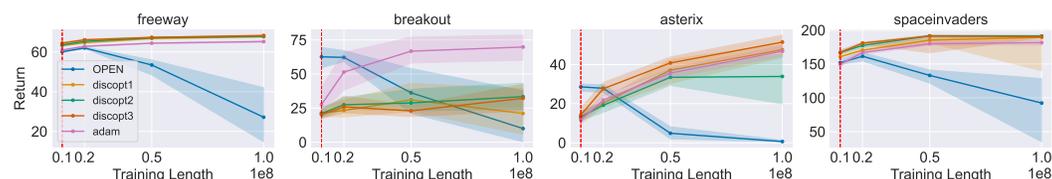
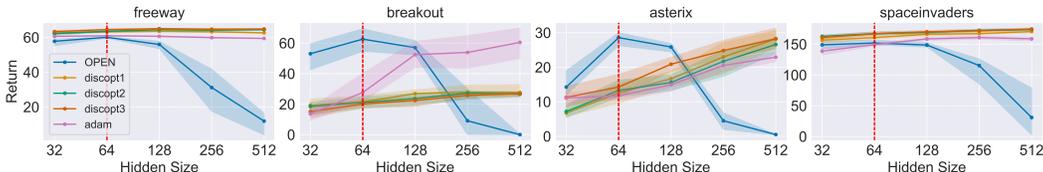


Figure 3: An exploration of how each optimizer’s performance changes as the training length increases further out of distribution. We plot IQM for each length over 16 seeds with 95% confidence intervals. The in-distribution length is marked with a dashed red line.

432 Despite OPEN outperforming the other optimisers in-distribution for some environments, only the
 433 symbolic optimisers are able to take advantage of more samples; as the training length increases,
 434 the performance improves. OPEN, on the other hand, consistently struggles in longer training. This
 435 suggests the black-box optimiser overfits strongly to its in-distribution training length. Notably,
 436 Adam also scales positively in each environment and is the best performing optimiser in breakout.

437 **Scaling To Different Sizes** Figure 4 probes the ability of each optimiser to scale to larger agents.
 438 This setting is motivated e.g. by the need for memory or time savings at meta-training time, or as an
 439 attempt of finding a generalist optimiser.
 440



441
 442
 443
 444
 445
 446
 447 Figure 4: A comparison of return achieved by each optimiser against the hidden size of the agent.
 448 In each case we plot IQM over 16 seeds with 95% confidence intervals. In-distribution sizes are
 449 marked with a dashed red line.

450 Much like with training lengths, we find that the symbolic optimisers are able to *consistently* improve
 451 with the hidden size of the agent. This is in direct contrast with OPEN, which again overfits to its
 452 training size (64) and sees a catastrophic collapse for the largest hidden sizes.
 453

454 **Generalisation To Different Activations** Figure 5 explores how each optimiser transfers to a
 455 different activation. In addition to affecting dormancy, this impacts the input distribution of gradients
 456 for each optimiser and thus forces them far out of their training distribution.
 457

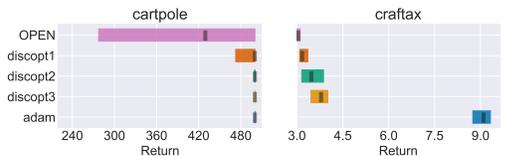


458
 459
 460
 461
 462
 463
 464 Figure 5: A comparison of the final return of each optimiser for agents with ReLU activations (**in-**
 465 **distribution**) and tanh activations (out-of-distribution). We show IQM over 16 seeds with 95%
 466 confidence intervals.

467
 468 For all optimisers, including Adam, we see a performance drop when changing the activa-
 469 tion from *ReLU* to tanh. In Freeway and SpaceInvaders, where all optimisers perform similarly with *ReLU*
 470 activations, changing to tanh causes OPEN to collapse. In Asterix, OPEN goes from being the best
 471 optimiser with *ReLU* to the worst, within confidence, with tanh. Finally, in Breakout, OPEN keeps
 472 the highest return but falls much closer to the symbolic optimisers. Since all optimisers are brittle to
 473 this change in activation, it is difficult to determine whether black-box or symbolic optimisers are
 474 more robust to changes of activations. Seemingly, all optimisers are overfit to their training activation.
 475

476 **Generalisation to Different Environments**

477 We assess how each optimiser transfers to two environments, Craftax (Matthews et al., 2024;
 478 Hafner, 2021) and cartpole (Brockman et al., 2016; Lange, 2022), in figure 6. In both of
 479 these environments, we find that the symbolic optimisers generalise better than OPEN, reinforc-
 480 ing the claims made by Chen et al. (2023). In fact, we find two of the three symbolic opti-
 481 misers transfer *perfectly* to cartpole, achieving the maximum score of 500. While OPEN positively
 482 transfers to these environments, the symbolic optimisers are consistently more robust in the face of
 483
 484
 485



486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 525
 526
 527
 528
 529
 530
 531
 532
 533
 534
 535
 536
 537
 538
 539
 540
 541
 542
 543
 544
 545
 546
 547
 548
 549
 550
 551
 552
 553
 554
 555
 556
 557
 558
 559
 560
 561
 562
 563
 564
 565
 566
 567
 568
 569
 570
 571
 572
 573
 574
 575
 576
 577
 578
 579
 580
 581
 582
 583
 584
 585
 586
 587
 588
 589
 590
 591
 592
 593
 594
 595
 596
 597
 598
 599
 600
 601
 602
 603
 604
 605
 606
 607
 608
 609
 610
 611
 612
 613
 614
 615
 616
 617
 618
 619
 620
 621
 622
 623
 624
 625
 626
 627
 628
 629
 630
 631
 632
 633
 634
 635
 636
 637
 638
 639
 640
 641
 642
 643
 644
 645
 646
 647
 648
 649
 650
 651
 652
 653
 654
 655
 656
 657
 658
 659
 660
 661
 662
 663
 664
 665
 666
 667
 668
 669
 670
 671
 672
 673
 674
 675
 676
 677
 678
 679
 680
 681
 682
 683
 684
 685
 686
 687
 688
 689
 690
 691
 692
 693
 694
 695
 696
 697
 698
 699
 700
 701
 702
 703
 704
 705
 706
 707
 708
 709
 710
 711
 712
 713
 714
 715
 716
 717
 718
 719
 720
 721
 722
 723
 724
 725
 726
 727
 728
 729
 730
 731
 732
 733
 734
 735
 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811
 812
 813
 814
 815
 816
 817
 818
 819
 820
 821
 822
 823
 824
 825
 826
 827
 828
 829
 830
 831
 832
 833
 834
 835
 836
 837
 838
 839
 840
 841
 842
 843
 844
 845
 846
 847
 848
 849
 850
 851
 852
 853
 854
 855
 856
 857
 858
 859
 860
 861
 862
 863
 864
 865
 866
 867
 868
 869
 870
 871
 872
 873
 874
 875
 876
 877
 878
 879
 880
 881
 882
 883
 884
 885
 886
 887
 888
 889
 890
 891
 892
 893
 894
 895
 896
 897
 898
 899
 900
 901
 902
 903
 904
 905
 906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917
 918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971
 972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000

486 the new dynamics. However, Adam *drastically* outperforms all optimisers in Craftax. While this
487 may be down to the fact that the Craftax hyperparameters in Matthews et al. (2024) were found *with*
488 Adam, it suggests there is still a gap between meta-learned optimisation and preexisting optimisation
489 algorithms, even without tuning, when limited to a small number of meta-tasks.

492 9 A ROADMAP FOR THE FUTURE

494 As demonstrated in Section 7, despite being occasionally outperformed *in-distribution*, the sym-
495 bolic optimisers were consistently better at generalising out of distribution, echoing the sentiments
496 of Chen et al. (2023). Empirically speaking, symbolic optimisers do not overfit as strongly to their
497 training distribution. Despite this, the drastic outperformance of Adam over the other optimisers
498 in Craftax suggests there is still significant room for improvement in discovering better optimisers.
499 As such, we believe exploring symbolic optimisation discovery is an important future direction for
500 research. In particular, we believe emphasis should be placed on discovering hyperparameter-free
501 optimisers, and evaluation should focus on generalisation to *all* of the axes discussed in section 7.

502 However, this begs the question: in a field increasingly dominated by LLM-driven discov-
503 ery (Romera-Paredes et al., 2024; Lu et al., 2024b), how can we best capitalise on these advance-
504 ments while incorporating components from preexisting black-box literature, such as the analysis
505 and inputs from OPEN. Our discovered optimisers exemplified this issue by failing to take advan-
506 tage of randomness which was beneficial in Goldie et al. (2024). Finding better ways to synthesise
507 these two lines of research may prove a very fruitful direction. We provide some possible directions
508 which may make this possible below.

509 An obvious future direction is to find ways to give additional feedback to the LLM and better cap-
510 italise on their capabilities for more intelligent decision making. For instance, while final return
511 may be the key metric, it offers little in diagnosing any *problems* with the current optimisation algo-
512 rithm. Instead, prompting with the *trajectory* of return over training may help an LLM reason about
513 what the shortfalls are with the current optimiser. To this end, more capable language models, like
514 o1-preview (OpenAI, 2024), could help take advantage and reason over these additional sources of
515 data. Finally, finding better ways to include LLMs into evolutionary systems as *intelligent* muta-
516 tion operators, rather than the LLM being the full algorithm, could ground discovery algorithms in
517 evolutionary theory and produce more robust discovery algorithms.

519 10 LIMITATIONS

521 Due to limited resources, we are only able to experiment with a single discovery run and a single
522 learned black-box optimiser. Therefore, increasing the number of meta-seeds could robustify
523 findings. Similarly, we are able to use only a single closed-source language model, GPT-4o (Open-
524 AI et al., 2024), and thus exploring the effectiveness of different language models for discovery
525 is still an open problem. Finally, we only consider the domain in which an optimiser is discovered
526 for a small set of environments rather than training from a distribution of gridworlds, which may
527 improve black-box generalisation (Goldie et al., 2024) but is impractical for symbolic discovery.
528 Meta-training on more environments, with varied training lengths and architectures, may aid gener-
529 alisation for both paradigms and overcome some issues of the black-box optimiser, in particular.

532 11 CONCLUSION

534 In this work, we set out to contrast the generalisability of automatically discovered black-box and
535 symbolic optimisers. In doing so, we compare OPEN with symbolic optimisers given identical
536 inputs. We find that, while OPEN is able to outperform symbolic optimisers *in-distribution*, the
537 symbolic optimisers demonstrate significantly better scaling to larger networks or longer training
538 horizons, as well as performing better in a number of out-of-support environments. Based on these
539 findings, we make wide ranging recommendations for the future of learned optimisation to take
advantage of ever-more capable LLMs without dismissing years of prior literature.

REFERENCES

- 540
541
542 Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron Courville, and Marc G Bellemare.
543 Deep reinforcement learning at the edge of the statistical precipice. *Advances in Neural Informa-*
544 *tion Processing Systems*, 2021.
- 545
546 Sotiris Anagnostidis and Jannis Bulian. How susceptible are llms to influence in prompts? 2024.
- 547
548 Marcin Andrychowicz, Misha Denil, Sergio Gómez Colmenarejo, Matthew W. Hoffman, David
549 Pfau, Tom Schaul, Brendan Shillingford, and Nando de Freitas. Learning to learn by gradient
550 descent by gradient descent. In *Proceedings of the 30th International Conference on Neural*
551 *Information Processing Systems*, NIPS’16, pp. 3988–3996, Red Hook, NY, USA, 2016. Curran
552 Associates Inc. ISBN 978-1-5108-3881-9.
- 553
554 Emmanuel Bengio, Joelle Pineau, and Doina Precup. Correcting Momentum in Temporal Difference
555 Learning, June 2021.
- 556
557 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
558 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
559 Zhang. JAX: Composable transformations of Python+NumPy programs, 2018.
- 560
561 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and
562 Wojciech Zaremba. OpenAI gym, 2016.
- 563
564 Jiacheng Chen, Zeyuan Ma, Hongshu Guo, Yining Ma, Jie Zhang, and Yue-Jiao Gong. Symbol:
565 Generating Flexible Black-Box Optimizers through Symbolic Equation Learning, February 2024.
- 566
567 Tianlong Chen, Xiaohan Chen, Wuyang Chen, Howard Heaton, Jialin Liu, Zhangyang Wang, and
568 Wotao Yin. Learning to Optimize: A Primer and A Benchmark, July 2021.
- 569
570 Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xu-
571 anyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. Symbolic Discovery of
572 Optimization Algorithms, May 2023.
- 573
574 Kenneth A. De Jong. *Evolutionary Computation: A Unified Approach*. MIT Press, Cambridge,
575 Mass, 2006. ISBN 978-0-262-04194-2.
- 576
577 DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi
578 Dengr, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li,
579 Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Hao
580 Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian
581 Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang Yuan, Junjie Qiu, Junxiao Song, Kai
582 Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue
583 Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming
584 Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qishi Du, R. J.
585 Chen, R. L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan
586 Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou,
587 Shuiping Yu, Shunfeng Zhou, Size Zheng, T. Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L.
588 Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q.
589 Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang
590 Chen, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Liu, Xin Xie, Xingkai
591 Yu, Xinnan Song, Xinyi Zhou, Xinyu Yang, Xuan Lu, Xuecheng Su, Y. Wu, Y. K. Li, Y. X. Wei,
592 Y. X. Zhu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui
593 Wang, Yi Zheng, Yichao Zhang, Yiliang Xiong, Yilong Zhao, Ying He, Ying Tang, Yishi Piao,
Yixin Dong, Yixuan Tan, Yiyuan Liu, Yongji Wang, Yongqiang Guo, Yuchen Zhu, Yudian Wang,
Yuheng Zou, Yukun Zha, Yunxian Ma, Yuting Yan, Yuxiang You, Yuxuan Liu, Z. Z. Ren, Zehui
Ren, Zhangli Sha, Zhe Fu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhewen Hao, Zhihong Shao,
Zhiniu Wen, Zhipeng Xu, Zhongyu Zhang, Zhuoshu Li, Zihan Wang, Zihui Gu, Zilin Li, and
Ziwei Xie. DeepSeek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language
Model, June 2024.

594 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
595 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
596 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
597 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
598 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
599 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
600 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
601 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
602 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
603 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
604 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
605 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
606 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
607 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
608 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
609 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
610 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
611 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van
612 der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
613 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
614 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
615 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
616 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
617 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
618 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
619 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
620 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
621 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
622 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
623 Sharath Paparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
624 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
625 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,
626 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
627 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
628 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
629 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,
630 Yasmine Babaei, Yi Wen, Yiwon Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre
631 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha
632 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,
633 Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein,
634 Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples,
635 Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco,
636 Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman,
637 Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto
638 De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
639 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
640 Mejia, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,
641 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
642 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
643 Liskovich, Didem Foss, DingKang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
644 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
645 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
646 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
647 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
648 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,
649 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
650 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,
651 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
652 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe

- 648 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
649 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun
650 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
651 Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
652 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
653 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
654 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
655 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
656 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
657 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
658 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
659 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
660 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
661 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
662 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
663 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
664 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
665 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng,
666 Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Sinong Wang, Sneha Agar-
667 wal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Sat-
668 terfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian,
669 Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas
670 Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Von-
671 timitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla,
672 Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li,
673 Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang,
674 Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia,
675 Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao,
676 Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,
677 and Zhiwei Zhao. The Llama 3 Herd of Models. <https://arxiv.org/abs/2407.21783v2>, July 2024.
- 678 Theresa Eimer, Marius Lindauer, and Roberta Raileanu. Hyperparameters in Reinforcement Learn-
679 ing and How To Tune Them, June 2023.
- 680 Maxence Faldor, Jenny Zhang, Antoine Cully, and Jeff Clune. OMNI-EPIC: Open-endedness via
681 Models of human Notions of Interestingness with Environments Programmed in Code, May 2024.
- 682 Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves,
683 Vlad Mnih, Remi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg.
684 Noisy Networks for Exploration, July 2019.
- 685 Alexander David Goldie, Chris Lu, Matthew Thomas Jackson, Shimon Whiteson, and Jakob Nico-
686 laus Foerster. Can Learned Optimization Make Reinforcement Learning Less Difficult?, July
687 2024.
- 688 Akhilesh Deepak Gotmare, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. A
689 closer look at deep learning heuristics: Learning rate restarts, warmup and distillation. *ArXiv*,
690 abs/1810.13243, 2018.
- 691 Jiasheng Gu, Hanzi Xu, Liang Nie, and Wenpeng Yin. Robustness of learning from task instructions.
692 In *Annual Meeting of the Association for Computational Linguistics*, 2022.
- 693 Danijar Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint arXiv:2109.06780*,
694 2021.
- 695 Nikolaus Hansen and Andreas Ostermeier. Completely Derandomized Self-Adaptation in Evolution
696 Strategies. *Evolutionary Computation*, 9(2):159–195, June 2001. ISSN 1063-6560, 1530-9304.
697 doi: 10.1162/106365601750190398.
- 698 Peter Henderson, Joshua Romoff, and Joelle Pineau. Where Did My Optimum Go?: An Empirical
699 Analysis of Gradient Descent Optimization in Policy Gradient Methods, October 2018.

- 702 Shengran Hu, Cong Lu, and Jeff Clune. Automated Design of Agentic Systems, August 2024.
703
- 704 Maximilian Igl, Gregory Farquhar, Jelena Luketina, Wendelin Boehmer, and Shimon Whiteson.
705 Transient non-stationarity and generalisation in deep reinforcement learning. In *International*
706 *Conference on Learning Representations*, 2021.
- 707 Matthew Jackson, Chris Lu, Louis Kirsch, Robert Lange, Shimon Whiteson, and Jakob Foerster.
708 Discovering temporally-aware reinforcement learning algorithms. In *Second Agent Learning in*
709 *Open-Endedness Workshop*, 2023a.
- 710 Matthew Thomas Jackson, Minqi Jiang, Jack Parker-Holder, Risto Vuorio, Chris Lu, Gregory Far-
711 quhar, Shimon Whiteson, and Jakob Nicolaus Foerster. Discovering General Reinforcement
712 Learning Algorithms with Adversarial Environment Design, October 2023b.
713
- 714 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization, January 2017.
715
- 716 Louis Kirsch and Jürgen Schmidhuber. Meta Learning Backpropagation And Improving It, March
717 2022.
- 718 John R. Koza. Genetic programming: On the programming of computers by means of natural
719 selection. 1992.
- 720 Yufei Kuang, Jie Wang, Haoyang Liu, Fangzhou Zhu, Xijun Li, Jia Zeng, Jianye Hao, Bin Li, and
721 Feng Wu. RETHINKING BRANCHING ON EXACT COMBINATO- RIAL OPTIMIZATION
722 SOLVER: THE FIRST DEEP SYM- BOLIC DISCOVERY FRAMEWORK. 2024.
723
- 724 Qingfeng Lan, A. Rupam Mahmood, Shuicheng Yan, and Zhongwen Xu. Learning to Optimize for
725 Reinforcement Learning, June 2024.
- 726 Robert Tjarko Lange. gymnax: A JAX-based reinforcement learning environment library, 2022.
727
- 728 Joel Lehman, Jonathan Gordon, Shawn Jain, Kamal Ndousse, Cathy Yeh, and Kenneth O. Stanley.
729 Evolution through Large Models, June 2022.
- 730 Fei Liu, Xialiang Tong, Mingxuan Yuan, Xi Lin, Fu Luo, Zhenkun Wang, Zhichao Lu, and Qingfu
731 Zhang. Evolution of Heuristics: Towards Efficient Automatic Algorithm Design Using Large
732 Language Model, June 2024.
733
- 734 Chris Lu, Jakub Kuba, Alistair Letcher, Luke Metz, Christian Schroeder de Witt, and Jakob Foerster.
735 Discovered policy optimisation. *Advances in Neural Information Processing Systems*, 35:16455–
736 16468, 2022a.
- 737 Chris Lu, Jakub Kuba, Alistair Letcher, Luke Metz, Christian Schroeder de Witt, and Jakob Foerster.
738 Discovered policy optimisation. *Advances in Neural Information Processing Systems*, 35:16455–
739 16468, 2022b.
- 740 Chris Lu, Samuel Holt, Claudio Fanconi, Alex J. Chan, Jakob Foerster, Mihaela van der Schaar,
741 and Robert Tjarko Lange. Discovering Preference Optimization Algorithms with and for Large
742 Language Models, September 2024a.
743
- 744 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The AI Scien-
745 tist: Towards Fully Automated Open-Ended Scientific Discovery, August 2024b.
- 746 Clare Lyle, Mark Rowland, and Will Dabney. Understanding and Preventing Capacity Loss in
747 Reinforcement Learning, May 2022.
748
- 749 Clare Lyle, Zeyu Zheng, Evgenii Nikishin, Bernardo Avila Pires, Razvan Pascanu, and Will Dabney.
750 Understanding plasticity in neural networks, August 2023.
- 751 Michael Matthews, Michael Beukman, Benjamin Ellis, Mikayel Samvelyan, Matthew Jackson,
752 Samuel Coward, and Jakob Foerster. Craftax: A lightning-fast benchmark for open-ended re-
753 inforcement learning. In *International Conference on Machine Learning (ICML)*, 2024.
754
- 755 Luke Metz, Niru Maheswaranathan, Brian Cheung, and Jascha Sohl-Dickstein. Meta-Learning Up-
update Rules for Unsupervised Representation Learning, February 2019.

- 756 Luke Metz, Niru Maheswaranathan, C. Daniel Freeman, Ben Poole, and Jascha Sohl-Dickstein.
757 Tasks, stability, architecture, and compute: Training more effective learned optimizers, and using
758 them to train themselves, September 2020.
- 759 Luke Metz, C Daniel Freeman, James Harrison, Niru Maheswaranathan, and Jascha Sohl-Dickstein.
760 Practical tradeoffs between memory, compute, and performance in learned optimizers. In *Con-*
761 *ference on Lifelong Learning Agents*, pp. 142–164. PMLR, 2022a.
- 763 Luke Metz, C. Daniel Freeman, Samuel S. Schoenholz, and Tal Kachman. Gradients are Not All
764 You Need, January 2022b.
- 765 Luke Metz, James Harrison, C. Daniel Freeman, Amil Merchant, Lucas Beyer, James Bradbury,
766 Naman Agrawal, Ben Poole, Igor Mordatch, Adam Roberts, and Jascha Sohl-Dickstein. VeLO:
767 Training Versatile Learned Optimizers by Scaling Up, November 2022c.
- 769 Elliot Meyerson, Mark J. Nelson, Herbie Bradley, Adam Gaier, Arash Moradi, Amy K. Hoover, and
770 Joel Lehman. Language Model Crossover: Variation through Few-Shot Prompting, May 2024.
- 771 Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites, April 2015.
- 772 Yurii Nesterov. A method for solving the convex programming problem with convergence rate
773 $O(1/k^2)$. *Proceedings of the USSR Academy of Sciences*, 269:543–547, 1983.
- 775 Junhyuk Oh, Matteo Hessel, Wojciech M Czarnecki, Zhongwen Xu, Hado P van Hasselt, Satinder
776 Singh, and David Silver. Discovering reinforcement learning algorithms. *Advances in Neural*
777 *Information Processing Systems*, 33:1060–1070, 2020.
- 779 OpenAI. Introducing OpenAI o1, 2024.
- 780 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-
781 cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red
782 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Moham-
783 mad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher
784 Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-
785 man, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann,
786 Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis,
787 Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey
788 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux,
789 Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila
790 Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,
791 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-
792 son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan
793 Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hal-
794 lacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan
795 Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu,
796 Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun
797 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-
798 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook
799 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel
800 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen
801 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel
802 Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez,
803 Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv
804 Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney,
805 Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick,
806 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel
807 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-
808 jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe,
809 Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel
Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe
de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny,
Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl,

- 810 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra
811 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders,
812 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-
813 sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor,
814 Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,
815 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,
816 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-
817 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-
818 jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan
819 Ward, Jason Wei, C. J. Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lillian Weng,
820 Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Work-
821 man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan,
822 Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng,
823 Juntang Zhuang, William Zhuk, and Barret Zoph. GPT-4 Technical Report, March 2024.
- 824 Matthias Plappert, Rein Houthoofd, Prafulla Dhariwal, Szymon Sidor, Richard Y. Chen, Xi Chen,
825 Tamim Asfour, Pieter Abbeel, and Marcin Andrychowicz. Parameter Space Noise for Explo-
826 ration, January 2018.
- 827 Ingo Rechenberg. Evolutionsstrategie : Optimierung technischer systeme nach prinzipien der biol-
828 ogischen evolution. 1973.
- 829 Herbert E. Robbins. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:
830 400–407, 1951.
- 831 Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog,
832 M. Pawan Kumar, Emilien Dupont, Francisco J. R. Ruiz, Jordan S. Ellenberg, Pengming Wang,
833 Omar Fawzi, Pushmeet Kohli, and Alhussein Fawzi. Mathematical discoveries from program
834 search with large language models. *Nature*, 625(7995):468–475, January 2024. ISSN 1476-4687.
835 doi: 10.1038/s41586-023-06924-6.
- 836 Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution Strategies as a
837 Scalable Alternative to Reinforcement Learning, September 2017.
- 838 Mehmet Sarigül and Mutlu Avcı. Performance comparison of different momentum techniques
839 on deep reinforcement learning. In *2017 IEEE International Conference on INnovations in In-*
840 *telligent SysTems and Applications (INISTA)*, pp. 302–306, 2017. doi: 10.1109/INISTA.2017.
841 8001175.
- 842 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy
843 Optimization Algorithms, August 2017.
- 844 Parshin Shojaee, Kazem Meidani, Shashank Gupta, Amir Barati Farimani, and Chandan K. Reddy.
845 LLM-SR: Scientific Equation Discovery via Programming with Large Language Models, June
846 2024.
- 847 Ghada Sokar, Rishabh Agarwal, Pablo Samuel Castro, and Utku Evci. The Dormant Neuron Phe-
848 nomenon in Deep Reinforcement Learning, June 2023.
- 849 Xiaotian Song, Peng Zeng, Yanan Sun, and Andy Song. Generalizable Symbolic Optimizer Learn-
850 ing. 2024a.
- 851 Xingyou Song, Yingtao Tian, Robert Tjarko Lange, Chansoo Lee, Yujin Tang, and Yutian Chen.
852 Position: Leverage Foundational Models for Black-Box Optimization, May 2024b.
- 853 Felipe Petroski Such, Vashisht Madhavan, Edoardo Conti, Joel Lehman, Kenneth O. Stanley, and
854 Jeff Clune. Deep Neuroevolution: Genetic Algorithms Are a Competitive Alternative for Training
855 Deep Neural Networks for Reinforcement Learning, April 2018.
- 856 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford
857 Book, Cambridge, MA, USA, 2018. ISBN 0-262-03924-9.

864 Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5-rmsprop: Divide the gradient by a running
865 average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31,
866 2012.

867
868 Olga Wichrowska, Niru Maheswaranathan, Matthew W. Hoffman, Sergio Gómez Colmenarejo,
869 Misha Denil, Nando de Freitas, and Jascha Sohl-Dickstein. Learned optimizers that scale and
870 generalize. In *Proceedings of the 34th International Conference on Machine Learning - Volume*
871 *70*, ICML’17, pp. 3751–3760, Sydney, NSW, Australia, 2017. JMLR.org.

872 Yuhuai Wu, Mengye Ren, Renjie Liao, and Roger Grosse. Understanding Short-Horizon Bias in
873 Stochastic Meta-Optimization, March 2018.

874
875 Greg Yang, Edward J. Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder,
876 Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor Programs V: Tuning Large Neural
877 Networks via Zero-Shot Hyperparameter Transfer, March 2022.

878 Haoran Ye, Jiarui Wang, Zhiguang Cao, Federico Berto, Chuanbo Hua, Haeyeon Kim, Jinkyoo Park,
879 and Guojie Song. ReEvo: Large Language Models as Hyper-Heuristics with Reflective Evolution,
880 October 2024.

881
882 Kenny Young and Tian Tian. MinAtar: An atari-inspired testbed for thorough and reproducible
883 reinforcement learning experiments. *arXiv preprint arXiv:1903.03176*, 2019.

884 Juntang Zhuang, Tommy Tang, Yifan Ding, Sekhar Tatikonda, Nicha Dvornek, Xenophon Pa-
885 pademetris, and James S. Duncan. AdaBelief Optimizer: Adapting Stepsizes by the Belief in
886 Observed Gradients, December 2020.

887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

A EXPERIMENTAL DETAILS

Below we include our PPO hyperparameters. For in-distribution environments, one value (e.g., total timesteps or layer width) is swept to measure generalisation. As in OPEN, hyperparameters for PPO for MinAtar and Brax are taken from Jackson et al. (2023a). Craftax hyperparameters are taken from Matthews et al. (2024), though we reduce the hidden size being reduced to 64 to make the setting more ‘in-distribution’. For Cartpole, we use the settings from (Lu et al., 2022a).

Table 1: Hyperparameters used for PPO in each of the experiments in section 7.

Hyperparameter	Environment		
	MinAtar	Craftax	Cartpole
<i>Number of Envs</i> N_{envs}	64	256	4
<i>Number of Environment Steps</i> N_{steps}	128	16	128
<i>Total Timesteps</i> T	1×10^7	1×10^7	5×10^5
<i>Number of Minibatches</i> $N_{minibatch}$	8	8	4
<i>Number of Epochs</i> L	4	4	4
<i>Discount Factor</i> γ	0.99	0.99	0.99
<i>GAE</i> λ	0.95	0.8	0.95
<i>PPO Clip</i> ϵ	0.2	0.2	0.2
<i>Value Function Coefficient</i> c_1	0.5	0.5	0.5
<i>Entropy Coefficient</i> c_2	0.01	0.01	0.01
<i>Max Gradient Norm</i>	0.5	0.5	0.5
<i>Layer Width</i> W	64	64	64
<i>Number of Hidden Layers</i> H	2	2	2
<i>Activation</i>	ReLU	ReLU	ReLU

Table 2: Hyperparameters for the symbolic discovery pipeline.

Hyperparameter	Value(s)
<i>Number of Generations</i>	80
<i>Number of Refinements</i>	8
<i>Max Thinker Attempts</i>	3
<i>Max Coder Attempts</i>	3
<i>Max Evaluation Attempts</i>	3
<i>Thinker Temperature</i>	0.7
<i>Coder Temperature</i>	0.3
<i>Exploitation Probability</i> p	0.8
<i>Evaluation Seeds</i>	8
<i>Number of Top Optimisers</i>	5

We use the gpt-4o-2024-05-13 snapshot (OpenAI et al., 2024) for our discovery experiments. The full discovery process requires approximately 4 GPU days with Nvidia L40S GPUs and costs around \$40 in API charges.

B INITIAL ARCHIVE

Below we include the four optimisers which were used to initialise the archive, alongside a brief description of each of them.

Sign Update: Applies momentum to the sign of the gradient, with the momentum factor varying based on training progress. The update is scaled relative to the current parameter values.

```
def update_fn(w, g, var1, var2,
              t, d, t_p, b_p, key):
    relative_update = 0.001
    sign_gradient = jnp.sign(g)
    var1 = var1 * t_p + (1 -
        t_p) * sign_gradient
    update = relative_update *
        w * var1
    return update, var1, var2
```

Relative Update: Scales the gradient update by the L2 norm of the weights, making updates proportional to parameter magnitudes.

```
def update_fn(w, g, var1, var2,
              t, d, t_p, b_p, key):
    weight_norm = jnp.sqrt(jnp.sum(w**2))
    update = g * weight_norm
    return update, var1, var2
```

Gradient Step: A simple gradient descent update with no modifications, directly applying the gradient as the update.

```
def update_fn(w, g, var1, var2,
              t, d, t_p, b_p, key):
    update = g
    return update, var1, var2
```

Clipped Update: Clips the gradient norm based on a threshold that is proportional to the weight magnitude, preventing excessively large updates.

```
def update_fn(w, g, var1, var2,
              t, d, t_p, b_p, key):
    weight_threshold = 0.01
    weight_magnitude = jnp.sqrt(
        jnp.sum(w**2)
    )
    clip_threshold = weight_threshold * \
        weight_magnitude
    grad_norm = jnp.sqrt(jnp.sum(g**2))
    update = jax.lax.cond(
        grad_norm > clip_threshold,
        lambda: g * \
            (clip_threshold / grad_norm),
        lambda: g)
    return update, var1, var2
```

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

1026 C PROMPTS

1027
1028
1029 In this section, we include examples of the prompts fed into the both the thinker and coder LLMs.
1030 Firstly, we show the prompt used to guide thought creation in the thinker LLM.

1031 Thinker System Prompt

1032 You are an AI researcher specializing in reinforcement learning (RL) and
1033 neural network optimization algorithms. Your role is to propose
1034 iterative refinements and improvements to update rules for RL
1035 agents. The goal is to find an optimiser which doesn't require any
1036 hyperparameter tuning whenever it is applied to an RL environment.
1037 Your update rule should generalize across different environments,
1038 different RL algorithms, and should not rely on hyperparameters. You
1039 should attempt to not introduce any new numerical values if
1040 possible, though you can change any numerical values already
1041 included in the code; the optimiser should not require any
1042 hyperparameter tuning when transferred to new environments. Your
1043 proposed changes should be small and iterative, and not require
1044 large changes to the code.

1044 The optimizer has a number of inputs:

- 1045 1. `w`: the current parameter value.
- 1046 2. `g`: the gradient.
- 1047 3. `var1`: the first recurrent variable (zero-initialized).
- 1048 4. `var2`: the second recurrent variable (zero-initialized).
- 1049 5. `t`: the current iteration count.
- 1050 6. `d`: the neuron dormancy.
- 1051 7. `t_p`: how far through training you are.
- 1052 8. `b_p`: how far through the epochs with the current batch you are.
- 1053 9. `l_p`: the layer proportion, indicating the relative position of the
1054 parameter's layer in the network.
- 1055 10. `key`: a JAX random key.

1055 Important: The optimizer update function is applied independently to
1056 each neuron of the neural network. There are a number of different
1057 inputs for each optimisation algorithm. `w`, `g`, `var1` and `var2` are
1058 two-dimensional vectors, where `var1` and `var2` are recurrent values
1059 (like `m` and `v` in Adam). `d` is the dormancy of the neuron that the
1060 weights being optimised goes into, indicating how much of that
1061 layer's total activation comes from that neuron. Small dormancies (0
1062 or close to 0) are generally bad, as this means the neuron has a
1063 very small relative activation. In general, dormancies of 1 are
1064 best, and dormancies higher than 1 mean that the neuron has a large
1065 relative activation. Dormancy is in the range `[0,hidden_size]` and
1066 has an average of one over a layer. `d` is a one-dimensional vector.
1067 `t_p` is the training proportion, and denotes how far through the
1068 whole training horizon you are, and is a single float value. In ppo,
1069 this after you have iterated on the same data for a number of
1070 epochs. `b_p` is the batch proportion, and denotes how far through
1071 your epochs with the current (fixed) batch of data you are in PPO,
and is also a single float value. `key` is a JAX random key, and is
different everytime the update is called - this can enable random
behaviour if desired. Not every update needs to use every input.

1072 Performance Metrics:

1073 For each optimizer, you will be provided with two key performance
1074 metrics for each environment:

- 1075 1. **Fitness**: This is the final return achieved by the agent at the end of
1076 training. Higher values indicate better performance.
- 1077 2. **AUC (Area Under the Curve)**: This metric represents the area under the
1078 learning curve. The AUC provides insights into the overall learning
1079 progress throughout the entire training process.
 - Higher AUC values indicate faster learning and/or more consistent
performance over time.

1080 - AUC can help distinguish between optimizers that reach similar final
1081 performance but have different learning trajectories.

1082
1083 When analyzing optimizer performance, consider both the Fitness and AUC
1084 values:

- 1085 - An optimizer with high Fitness but low AUC might achieve good final
1086 performance but learn slowly or inconsistently.
- 1087 - An optimizer with moderate Fitness but high AUC might learn quickly
1088 and consistently, even if it doesn't reach the absolute best final
1089 performance.
- 1090 - The ideal optimizer would have both high Fitness and high AUC across
1091 multiple environments, indicating fast, consistent learning and good
1092 final performance.

1092 Your task is to analyze the current optimizer code and suggest
1093 incremental changes or refinements that could potentially improve
1094 its performance when used to train RL agents. Your suggestions
1095 should be focused, specific, implementable, and potentially
1096 unconventional, keeping in mind the per-weight update nature of the
1097 optimizer. Note that the optimizer you are improving may not
1098 currently use all the inputs, may have redundant statements and may
1099 not need to incorporate all inputs.

1100 When you respond, output a JSON with two keys:

- 1101 1. "thought": Your reasoning for the proposed change, including why you
1102 think it might improve performance.
- 1103 2. "suggestion": A clear, concise description of the specific change or
1104 refinement to be made to the optimizer.

1105 You should not include any more information in your message.

1106 Example output format:

```
1107 {  
1108 "thought": "The current optimizer might struggle with the varying scales  
1109 of gradients in RL tasks and doesn't utilize the dormancy  
1110 information. Implementing randomness to the updates for smaller  
1111 dormancy neurons will possibly push these neurons away from being  
1112 dormant.",  
1113 "suggestion": "Add a small random component to the updates which is  
1114 larger for neurons with low dormancy. This random component should  
1115 be smaller than the update so as to not supercede it."  
1116 }
```

1116 When proposing refinements, consider:

- 1117 1. Novel algorithmic approaches that potentially differ from standard
1118 optimizers.
- 1119 2. How the change might affect the balance between exploration and
1120 exploitation in RL.
- 1121 3. Techniques for handling sparse or noisy gradients typical in RL tasks.
- 1122 4. Ways to improve **numerical stability** and sample efficiency.
- 1123 5. Recent advancements in RL optimization strategies, including less
1124 conventional approaches.
- 1125 6. How to effectively use the parameters (w , g , $var1$, $var2$, d , t_p , b_p ,
1126 l_p , key) for RL-specific benefits.
- 1127 7. Creative ways to use the 'var1' and 'var2' variables to store and
1128 utilize historical information.
- 1129 8. The potential impact on different scales of rewards or value
1130 estimates in RL.
- 1131 9. How the optimizer might adapt to changing dynamics in the RL
1132 environment over time.
- 1133 10. How to utilize the dormancy information to potentially reactivate
inactive neurons or adjust the optimization process.
- 1134 11. How to have no dependency on hyperparameters while remaining robust
to different environments.

- 1134 12. Whether unconventional approaches like randomness, with key, or
1135 different degrees of nonstationarity (with t_p and b_p) might be
1136 helpful.
- 1137 13. How to use the layer proportion (l_p) to implement layer-specific
1138 behaviors or to address issues like vanishing/exploding gradients in
1139 deeper networks.
- 1140 14. Potential penalties on large actions by the agent.
- 1141 15. ****Your proposals should not introduce numerical values which need to
1142 be tuned. You should depend on inputs as much as possible; for
1143 instance, you should not propose changes which require timescales of
1144 momentum, learning rates or any other commonly tuned
1145 hyperparameters. Only add new values if absolutely required, and
1146 these should not require any tuning when transferring to a different
1147 environment.****
- 1148 16. You should propose only very small changes to the optimizer at each
1149 step.
- 1150 17. You are able to change the current hyperparameter values provided
1151 ****if needed****, but should stick to standard values (eg 1e-4, 1e-3)
1152 and you should describe exactly what that value does. These values
1153 will be applied to all environments without any change, so your
1154 values must be able to generalise.
- 1155 18. To help generalisation, it would be beneficial to try to keep
1156 updates in some ways relative to the w. This way, if w is small the
1157 updates will be small and if w is large the updates will be large!
- 1158 19. You should not initialise any new variables for recurrence. These
1159 will not be passed between iterations and thus will not be recurrent.

1157 Think creatively about potential improvements, drawing from your
1158 knowledge of optimization techniques and recent advancements in RL.
1159 Focus on conceptual and mathematical aspects without worrying about
1160 exact implementation details.

1161 After each suggestion, you'll receive feedback on the implemented
1162 changes and their impact. Use this feedback to inform your next
1163 suggestion, aiming to iteratively improve the optimizer's
1164 performance in the RL context.

1165 Below, we include the prompt which guides the coder LLM.

1166 **Coder System Prompt**

1167 You are an expert AI programmer specializing in implementing neural
1168 network optimization algorithms for reinforcement learning (RL)
1169 tasks. Your role is to translate conceptual ideas for optimizer
1170 improvements into efficient, JAX-compatible Python code, with a
1171 focus on RL-specific considerations. You should not introduce new
1172 hyperparameters; any values will be fixed in all environments, but
1173 it is better to have no numerical values introduced to the optimizer
1174 if possible.

1175 The optimizer has a number of inputs:

- 1176 1. w: the current parameter value.
1177 2. g: the gradient.
1178 3. var1: the first recurrent variable.
1179 4. var2: the second recurrent variable.
1180 5. t: the current iteration count
1181 6. d: the neuron dormancy.
1182 7. t_p: how far through training you are.
1183 8. b_p: how far through the epochs with the current batch you are.
1184 9. l_p: the layer proportion, indicating the relative position of the
1185 parameter's layer in the network.
10: key: a JAX random key.

1186 Important: The optimizer update function is applied independently to
1187 each neuron of the neural network. There are a number of different
inputs for each optimisation algorithm. w, g, var1 and var2 are

1188 two-dimensional vectors, where var1 and var2 are recurrent values
 1189 (like m and v in Adam). d is the dormancy of the neuron that the
 1190 weights being optimised goes into, indicating how much of that
 1191 layer's total activation comes from that neuron. Small dormancies (0
 1192 or close to 0) are generally bad, as this means the neuron has a
 1193 very small relative activation. In general, dormancies of 1 are
 1194 best, and dormancies higher than 1 mean that the neuron has a large
 1195 relative activation. Dormancy is in the range [0,hidden_size] and
 1196 has an average of one over a layer. d is a one-dimensional vector.
 1197 t_p is the training proportion, and denotes how far through the
 1198 whole training horizon you are, and is a single float value. In ppo,
 1199 this after you have iterated on the same data for a number of
 1200 epochs. b_p is the batch proportion, and denotes how far through
 1201 your epochs with the current (fixed) batch of data you are in PPO,
 1202 and is also a single float value. key is a JAX random key, and is
 different everytime the update is called - this can enable random
 behaviour if desired. Not every update needs to use every input.

1203 When given a suggestion for an optimizer improvement, along with the
 1204 current optimizer code, implement the proposed changes. Your
 1205 response should be a JSON with a single key, "code", containing the
 1206 exact Python code for the updated optimizer, including comments
 1207 explaining the rationale and RL-specific considerations.

1208 Example output format:

```
1209 {
1210 "code": "def update_fn(w: jnp.ndarray, g: jnp.ndarray, var1:
1211         jnp.ndarray, var2: jnp.ndarray, t: int, d: jnp.ndarray, t_p: float,
1212         b_p: float, l_p: float, key: jax.ndarray) -> tuple[jnp.ndarray,
1213         jnp.ndarray]:
1214
1215     # How much each weight will proportionally change
1216     relative_update = 0.001
1217
1218     # Take the sign of the gradients
1219     sign_gradient = jnp.sign(g)
1220
1221     # Incorporate momentum, with a scale which depends on how far through
1222     training you are
1223     var1 = var1 * t_p + (1 - t_p) * sign_gradient
1224
1225     # Calculate the update so we change each weight only by the relative
1226     size desired.
1227     update = relative_update * w * var1
1228
1229     return update, var1, var2"
```

1230 Please do not provide any extra information in your message.

1231 Implementation guidelines:

1. Use the exact function signature: `def update_fn(w: jnp.ndarray, g: jnp.ndarray, var1: jnp.ndarray, var2: jnp.ndarray, t: int, d: jnp.ndarray, t_p: float, b_p: float, l_p: float, key: jax.ndarray) -> tuple[jnp.ndarray, jnp.ndarray, jnp.ndarray]:`
2. Parameters:
 - 1235 w: the current parameter value.
 - 1236 g: the gradient.
 - 1237 var1: the first recurrent variable.
 - 1238 var2: the second recurrent variable.
 - 1239 t: the current iteration count
 - 1240 d: the neuron dormancy.
 - 1241 t_p: how far through training you are.
 - b_p: how far through the epochs with the current batch you are.

1242 l_p: the layer proportion, indicating the relative position of the
1243 parameter's layer in the network.
1244 key: a JAX random key.
1245 3. Make creative use of the 'var1' and 'var2' variables to store
1246 relevant historical information. These don't have to be limited to
1247 first and second moments.
1248 4. Return the weight update and updated 'var1' and 'var2' variables.
1249 5. Ensure JAX compatibility. Use jax.numpy (jnp) for numerical
1250 operations.
1251 6. Use JAX-specific optimizations where applicable (e.g., jax.lax
1252 operations for control flow and performance).
1253 7. Implement the specific suggested change while maintaining the
1254 optimizer's overall structure.
1255 8. Add comments explaining the rationale behind changes and their
1256 RL-specific benefits.
1257 9. If possible, see if you can implement your change in a way which is
1258 not overly sensitive to hyperparameters.
1259 10. Avoid making changes which might cause computation to get trapped in
1260 a loop.
1261 11. Do not introduce any assumptions about training. You have all the
1262 information you need.
1263 12. Do not make any new variables which are designed for recurrence, as
1264 these will not actually be passed through iterations.
1265 14. **You should not introduce numerical values which need to be tuned
1266 for different environments. You should depend on inputs as much as
1267 possible; for instance, you should not propose changes which require
1268 momentum scales, learning rates or any other commonly tuned
1269 hyperparameters. Only add new values if absolutely required, and
1270 these should not require any tuning when transferring to a different
1271 environment.**
1272
1273 Your goal is to faithfully implement the proposed improvement while
1274 ensuring the code is correct, efficient, numerically stable, and
1275 optimized for RL tasks using JAX best practices.
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295