

# Cascade - A sequential ensemble method for continuous control tasks

Schmöcker R. , Dockhorn A.

**Keywords:** Ensemble learning, reinforcement learning, continuous control, PPO.

## Summary

Though reinforcement learning has been successfully applied to a variety of domains, there is still room left for improvement, in particular, in terms of the final performance. Ensemble Reinforcement Learning (ERL) tries to enhance reinforcement learning techniques by using multiple models or algorithms. We propose a novel ERL technique, called Cascade which in the context of continuous control tasks and with PPO as the base training algorithm clearly outperforms standard PPO in terms of the final performance. To shine light on the working mechanisms of Cascade, we conduct ablation studies, showing how the different components of Cascade contribute to its overall performance. Furthermore, we demonstrate that Cascade has a robust monotonicity as the ensemble's performance increases with each additional base agent even when weak base agents are added in large numbers.

## Contribution(s)

1. The proposition of a novel Ensemble Reinforcement Learning (ERL) algorithm Cascade for continuous control tasks that outperforms its base learner when using PPO as the underlying reinforcement learning algorithm.

**Context:** To the best of our knowledge, there is no prior work where the ensemble policy uses a convex combination of its base learners and still gains a significant performance advantage.

2. By multiple ablation studies, we investigate the mechanisms contributing to Cascade's performance.

**Context:** We show that Cascade relies on all base learners being trained at all stages of the training process as well as Cascade relying on sequentially adding base learners instead of starting with the final network. Lastly, we show that Cascade can chain an arbitrary number of base learners of arbitrary strengths without a loss in performance.

# Cascade - A sequential ensemble method for continuous control tasks

**Schmöcker R. , Dockhorn A.**

`schmoecker@tnt.uni-hannover.de, dockhorn@tnt.uni-hannover.de`

<sup>1</sup>Faculty of EECS, Leibniz University Hannover

<sup>2</sup>Faculty of EECS, Leibniz University Hannover

## Abstract

Though reinforcement learning has been successfully applied to a variety of domains, there is still room left for improvement, in particular, in terms of the final performance. Ensemble Reinforcement Learning (ERL) tries to enhance reinforcement learning techniques by using multiple models or algorithms. We propose a novel ERL technique, called Cascade which in the context of continuous control tasks and with PPO as the base training algorithm clearly outperforms standard PPO in terms of the final performance. To shine light on the working mechanisms of Cascade, we conduct ablation studies, showing how the different components of Cascade contribute to its overall performance. Furthermore, we demonstrate that Cascade has a robust monotonicity as the ensemble's performance increases with each additional base agent even when weak base agents are added in large numbers.

## 1 Introduction

Ensemble Reinforcement Learning (ERL) tries to improve reinforcement learning methods by using multiple models or algorithms. Advantages of ERL methods might stem from an increase in diversity, increased representational capabilities, or avoidance of the reliance on a single, possibly inaccurate predictor/decision-maker ([Dietterich, 2000](#)).

This paper's research question was the development of a novel ERL algorithm that could perform well on continuous control tasks. Hence our contribution is the proposal and experimental investigation of an ERL algorithm named Cascade which sequentially adds agents to an ensemble as opposed to starting with the full ensemble. The latest agent simultaneously learns a policy and how to integrate itself into the ensemble consisting of all previous agents. Additionally, the ensemble is designed in such a way that all ensemble members, once they are added, can refine their policy. Furthermore, Cascade will set the ensemble up so that it can be trained end-to-end, thereby requiring the selection of an arbitrary policy-based RL algorithm (such as PPO).

We show experimentally that Cascade, when using PPO to train the ensemble, outperforms PPO on a variety of continuous control tasks in terms of the final performance. Additionally, the individual components of Cascade were investigated, discovering that the sequential nature combined with allowing every ensemble member to be trainable at all times, are the biggest performance contributors. Furthermore, we show that Cascade is robust to the number of base agents as well as their expressivity by demonstrating that even for large ensembles consisting of weak base agents, adding a base agent almost always yields a performance improvement. In particular, we show that even agents, which on their own lack the capacity to solve the considered tasks, can be effectively combined by Cascade to solve these tasks. In this paper we focused mainly on PPO as it was the training algorithm that synergizes best with Cascade.

The paper is structured as follows: In Section 2 we give a short overview of ERL methods for continuous control tasks. Afterwards, in Section 3 we will define the Cascade algorithm and explain its intuition. The subsequent Section 4 presents experiments to answer structured research questions on Cascade’s capabilities. In Section 5 we will give reasons for why Cascade works as well as it does and discuss avenues for future work. At last, in Section 6 we summarize our main findings.

## 2 Related work

A plethora of different approaches to continuous control tasks using ERL exists. Januszewski et al. (2021) use an ensemble of TD3-like (Fujimoto et al., 2018) agents to counteract the bias an individual  $Q$  function and consequently the actor  $\pi$  has by using the ensemble to average out the bias. They do this by simultaneously training an ensemble of differently initialized agents. At evaluation time, the ensemble’s action is simply the average of the individual agents’ actions.

Li et al. (2023) propose a modification, called TEEN, which also trains an ensemble of actor-critics but such that their state-action visit distributions are as diverse as possible to encourage exploration. Additionally, they directly decrease the bias of the  $Q$ -functions targets by using a shared target that contains the average of a random subset of the  $Q$ -function ensemble.

Saphal et al. (2021) developed SEERL which creates an ensemble of agents within just one training cycle by regularly saving the agent’s current parameters. To promote diversity between the ensemble members, they use a cyclical learning rate which ensures convergence before parameters are saved and enough disturbance to get out of any local minimum. After training, they pick a subset of the saved agents and combine them without any extra training phase by averaging or binning the individual action outputs.

Chen & Huang (2023) introduce a hierarchical ensemble method named HED that trains an ensemble of actor-critics on two different levels. First, on a lower level, where the agents are trained independently of each other, and secondly on a higher level where the ensemble’s policy, which is the average of the individual actor outputs, is updated in a fashion that promotes cooperation among the base learners.

Liu et al. (2023) took a slightly different approach that aims at increasing the stability of the ensemble algorithm. During training, they maintain an ensemble of actor-critics and compute their standard policy gradients which are used to construct a single so-called robust policy gradient which is then applied to all agents. The robust policy gradient is chosen in such a fashion that when applied to all agents at once, maximizes the average of the expected returns of all ensemble agents.

Yet another approach, SUNRISE, is demonstrated by Lee et al. (2021) who also use an actor-critic ensemble. They use the standard deviation of the ensemble’s  $Q$ -functions as an uncertainty measure which is used to weigh samples and to construct an upper-confidence bound (UCB) during inference. More specifically, each ensemble agent samples an action during inference. The action that maximizes the UCB (which is a weighted sum of the  $Q$ -functions mean and standard deviation) is chosen as the ensemble’s action.

Jaderberg et al. (2017) showed that an ensemble of agents can be used to optimize a hyperparameter schedule. They keep a fixed size population of agents which are trained individually in parallel. Regularly, the lowest scored agents are discarded and replaced by the better performing ones. To promote further hyperparameter exploration, once the training state and hyperparameters have been copied, the hyperparameters are slightly perturbed. Their method is applicable to continuous control tasks even though it was not specifically designed for it.

In contrast to these works, Cascade does not necessarily produce an ensemble of agents that would perform well on their own. The idea is that not every agent needs to know everything about the task. They just need to be able to account for the weaknesses of other agents. Hence, Cascade solely focuses on the performance of the entire ensemble and how the agents can optimally cooperate and complement each other.

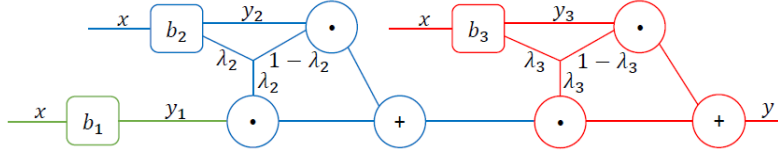


Figure 1: Cascade net of size 3. Except for the bottom net  $b_1$ , each net outputs a fallback action  $\lambda_i \in [0, 1]$  which is used for the convex combination.

More generally, [Song et al. \(2024\)](#) published a comprehensive survey on ERL thereby giving an overview of the entire field. Their scope is broader than the restriction to continuous control as done here.

### 3 Method

The main idea behind Cascade is to sequentially train base agents in such a fashion that a new base agent always tries to optimally complement the current ensemble. This means that the base agent should only learn to perform well or learn to adjust the ensemble in a subspace where the ensemble performs poorly and otherwise delegate the task of choosing an action to the ensemble. The new base agent can then be integrated into the ensemble to obtain a new, improved ensemble.

A naive approach to this is to train a base agent on a surrogate environment that extends the action space of the original environment by an action we call the *fallback action*. The values of this fallback action lie in the range  $[0, 1]$ . With a probability equal to the fallback action, instead of executing the given action on the original environment, the action of the current ensemble is executed. The newly trained base agent can then be integrated into the ensemble. The action of the new ensemble is the action of the latest base agent with probability equal to its fallback action and otherwise the action of the previous ensemble. Since the base agent could learn to fall back with a probability that is arbitrarily close to 1 in all states, it is ensured that theoretically an arbitrary amount of base agents can be chained without a loss in performance.

However, by sampling and then taking either the action of the base agent or the ensemble at each step as well as treating the ensemble as a blackbox (surrogate environment), non-differentiability is introduced. Additionally, the weaknesses of older base agents cannot be corrected, as only the latest base agent is trained. To overcome this, a convex-combination of the current ensemble-policy  $\pi_e$  and the new base agent  $\pi_b$  is used to obtain a policy  $\pi$ :

$$\pi(x) = \lambda_w(x) \cdot \pi_e(x) + (1 - \lambda_w(x)) \cdot \pi_b(x) \quad (1)$$

where  $x$  is a state and  $\lambda_w$  is the fallback action of the base agent. This policy  $\pi$  is then used for training. Both the parameters of  $\pi_e$  and  $\pi_b$  remain trainable.

Any model that is used to represent base agents has to output the fallback action  $\lambda(x) \in [0, 1]$  in addition to its usual action outputs. In the case of artificial neural networks (ANNs), the policy  $\pi$  can be approximated by taking a convex combination of the network’s action outputs directly (instead of the distributions they represent. For example, the convex combination of two normal distributions is in general not the same as the distribution induced by the convex combination of their parameters). If several ANNs are stacked this way, then this results in an architecture called *Cascade net*. Fig. 1 shows how 3 base agents are chained which results in a Cascade net of size 3. For reference, the pseudocode of the proposed technique to sequentially train a Cascade net is presented in the supplementary materials. This algorithm will be referred to as *Cascade*.



## 4 Experiments

First, it will be demonstrated that Cascade is able to clearly outperform the baseline on most MuJoCo (Todorov et al., 2012) environments using PPO. The subsequent experiments will then shift their focus towards studying individual components of Cascade to gain insights into the reasons for this performance increase.

### 4.1 Experiment setup

Cascade as described in Section 3 was run on various MuJoCo tasks with a focus on the Ant-v4 and Walker2d-v4 environments because these two differ greatly in terms of dynamics and complexity which gives a good impression of how Cascade behaves in different settings. When not mentioned otherwise, PPO was used as the algorithm to train the Cascade, 6 million environment steps took place, and the Cascade net was extended every 1 million steps. For details and how performances (and other metrics) were measured see supplementary materials Section B. The code for our experiments is publicly available at <https://github.com/codebro634/Cascade>.

### 4.2 Standard Cascade

Fig. 2 shows the result when Cascade without any modifications as described in the experiment setup was run on Ant-v4, Walker2d-v4, Humanoid-v4, Hopper-v4, and HalfCheetah-v4. We refer to this setup as *standard Cascade*. For comparison, each graph also shows the PPO baseline (see supplementary materials Section B.4) trained for 6 million steps. However, as mentioned in the study design, different baseline performances emerge for different training durations since learning rate annealing is used (see supplementary materials Fig. 9). Except for the Hopper-v4 and Walker environment where Cascade performs slightly worse or only slightly better than the baseline, Cascade decisively outperforms the baseline in all other environments after 6M steps. While in Ant-v4, Humanoid-v4, and HalfCheetah-v4 the baseline performances are about equal when Cascade consists only of one base net, Cascade takes the lead in performance shortly after a second base net is added and grows the lead from there.

Next, we compared Cascade to a slightly different PPO agent that mimics a base agent of Cascade as closely as possible. Therefore, a single base net (see supplementary materials Section B) will be trained with PPO for 6 million steps. To mimic the conditions under which Cascade trained, the learning rate was cyclical, linearly decreasing from its initial value to 0, every 1 million timesteps. Other than changing the network architecture and learning rate schedule, everything else is identical to the baseline PPO. Fig. 3 shows the performance results compared with standard Cascade.

Surprisingly, in Ant-v4, a single base agent reached a significantly higher performance than the PPO baseline which just differs in the network size and learning rate schedule, which might suggest that the byproduct of having a cyclical learning rate (remember, a new PPO instance is used every time the Cascade net expands) adds to the performance of Cascade. This is not the case however, as Cascade is very robust to different learning rate schedules (for details, see supplementary materials F). Nonetheless, there is still a clear gap between the base agent and Cascade both in Ant-v4 and Walker2d-v4.

Additionally, we evaluated the Cascade algorithm with SAC and DDPG as the base training algorithm, as detailed in the supplementary materials in Section H. Using these training algorithms Cascade did not outperform PPO, therefore we chose PPO as the default training algorithm for subsequent experiments. It remains to be researched why Cascade synergizes so well with PPO in particular.

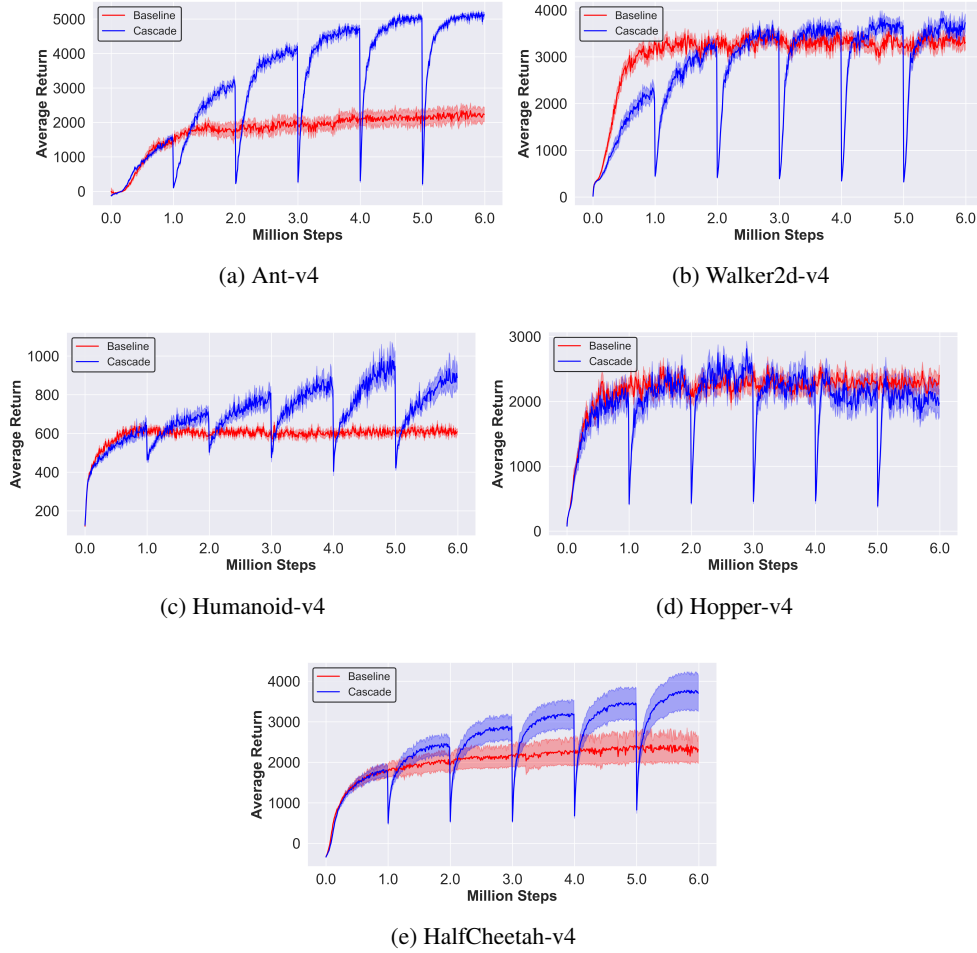


Figure 2: Performance graphs of standard Cascade (blue) vs PPO (red) on Ant-v4 (top left), Walker2d-v4 (top right), Humanoid-v4 (bottom left), Hopper-v4 (bottom middle) and HalfCheetah-v4 (bottom right) for 6 million environment steps.

**Question 1:** Can Cascade outperform any RL-algorithm while using the same RL-algorithm to train its base agents?

Yes, Cascade using PPO as the training algorithm outperforms PPO in most of the considered environments. In the cases Cascade outperformed PPO it managed to do so rather early, requiring only two base agents. And in the environments Cascade does not or barely outperform the baseline, the peak performance throughout training at least draws even with peak PPO performance.

Though this paper’s focus was continuous action space environments, we also applied a slightly modified version of Cascade to discretized versions of the MuJoCo environments (see supplementary materials Section I) and found that though Cascade also outperforms PPO in most of the discrete settings, the margins were very slim and Cascade took far longer to converge.

The final Cascade nets, obtained after running standard Cascade on Ant-v4 and Walker2d-v4, were investigated in more detail to see how many changes were applied to the initial base agent after its initial training phase. *Ant*: In this case, the entire Cascade mostly relies on the bottom base net (i.e. the net at the bottom of the Cascade net with no fallback action): While the entire Cascade net achieves an average performance of around 5.1k the bottom net alone reaches an average return

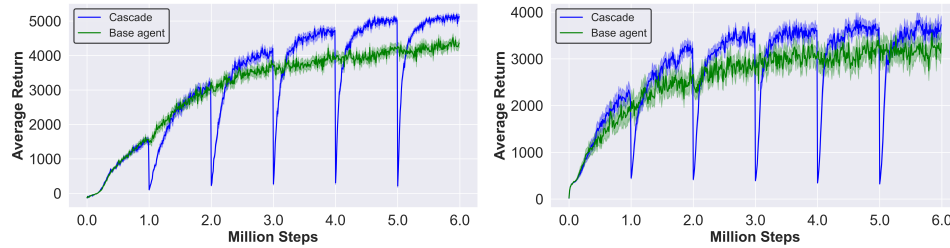


Figure 3: Performance graph of Cascade (blue) and PPO with cyclical learning rate on a base net of Cascade (green) for 6 million environment steps on Ant-v4 (left) and Walker2d-v4 (right).

of approximately 3.2k. All other policies further up the chain have a negative average return. As hypothesized in the algorithm-description this base policy improved even past 1 million steps when other agents were starting to be stacked on top. At 1 million steps the base policy performance was at around 1.5k while it was approximately 3.2k by the end of training. *Walker2d*: These characteristics differ for Walker2d-v4, as instead of improving past 1 million steps, the base policy degraded, going from 2.5k to approximately 600. However, all other policies further up the chain still perform poorly on their own, each with a performance of less than 100. Nonetheless, the entire chain still performs well, performing slightly better than the PPO baseline with a performance of roughly 3.9k at the end of training.

**Question 2:** Does Cascade still apply big changes to base agents once they are no longer at the top of the cascade?

Yes, judging by the performance of the bottom base agent. However, depending on the environment its performance may degrade or improve. Either way, the Cascade improves.

#### 4.3 Keeping base nets frozen

Next, it was investigated how allowing every base net to be trainable at every stage of the training process impacted the performance. For this, only the parameters of the latest base agent were trainable, the rest were kept frozen. Fig. 4 shows the performance graphs compared to standard Cascade. It can be observed that even though each training iteration leads to a minor improvement in performance for the frozen version, the overall performance is drastically worse compared to the non-frozen version for Ant-v4 and slightly worse for Walker2d-v4. This suggests that it is in fact vital for all base agents to remain trainable to correct for imperfections, adapt to new normalization states (observation/reward normalization is used by default, for details see supplementary materials Section B.4) and in general make the entire network less rigid. The difference could in fact be partially explained by the ability to adapt to new normalization states as both versions only slightly differ in their final performance when no normalization is used (see supplementary materials Fig. 13). However, in the normed setting, older base agents are modified beyond just adapting to new normalization states as could be seen in the previous section where the policy of the oldest base agent changed drastically which was indicated by the drastic change in its average performance.

**Question 3:** What is the impact on the performance when keeping the parameters of the current ensemble frozen when training a new base agent?

Keeping the parameters frozen severely worsens performance.

#### 4.4 Starting with the final Cascade net

Next, the effect of sequentially adding base agents to the Cascade was tested. For this, the Cascade net obtained at the end of training was used right from the start. To have the same training conditions, the learning rate was cyclical, linearly decreasing from its initial value to 0, every 1 million environment

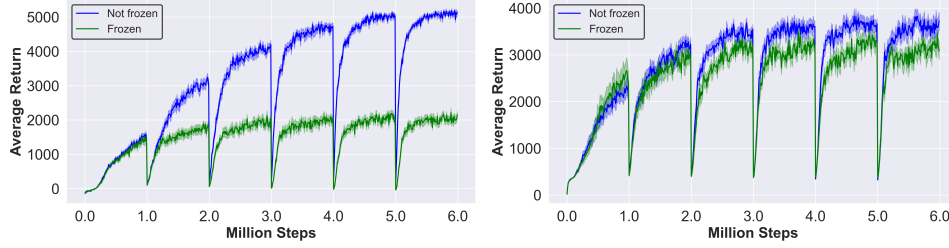


Figure 4: Performance graphs of standard Cascade (blue) vs Cascade with the weights of all but the latest base-net frozen (green) on Ant-v4 (top) and Walker2d-v4 for 6 million environment steps.

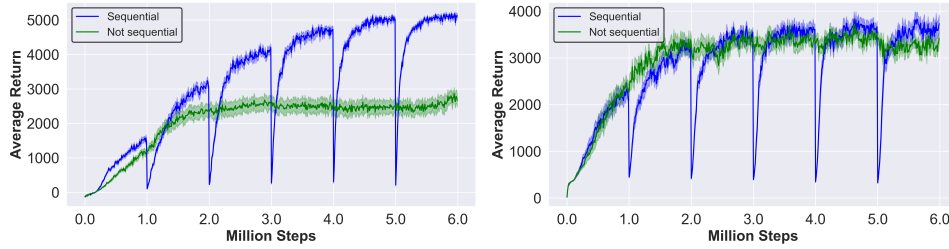


Figure 5: Performance graphs of standard Cascade (blue) versus Cascade where the final net of standard Cascade is trained right from the start (green) on Ant-v4 (top) and Walker2d-v4 (bottom) for 6 million environment steps.

steps (as PPO used to train the base agents in standard Cascade uses learning rate annealing. See supplementary materials Section B). Fig. 5 shows the performance results for this experiment on Ant-v4 and Walker2d-v4.

In both environments, the sequential variant clearly outperforms its non-sequential counterpart. While in Ant-v4, the sequential variant outperforms right from the start and ends with a final performance almost twice that of the non-sequential version, in Walker2d-v4 the sequential variant only takes the lead in performance during the third iteration and ends only with a small advantage.

This shows that the architecture itself is not the deciding factor for performance. In fact, there seems to be nothing inherently special to this architecture, as it performs worse than the standard feedforward architecture used for the baseline experiments.

**Question 4:** How does the performance of training the final Cascade net right from the start compare to adding the base agents sequentially?

Training the final Cascade net right from the start performs significantly worse than adding the base agents sequentially, thus proving that it is in fact vital to slowly build up the Cascade.

#### 4.5 Fallback action characteristics

Now, the tendency of base agents to fallback will be looked into. This was done by measuring the average fallback actions of different base nets over the course of training. This was done for three different fallback action initializations, namely  $\{0.05, 0.5 \text{ (Default)}, 0.9\}$  (To see how this is done, see supplementary materials Section E).

First, simply the product of all fallback actions of all base agents within the Cascade net was measured. This is the contribution of the bottom base net to the final output. Fig. 6 shows this product over the course of training on Ant-v4 and Walker2d-v4. In all settings, the fallback product slowly decreases as the Cascade grows. However, during every iteration, the product mostly increases (except for the first two iterations with the 0.9 initialization). This shows that the base agents have a clear

Env	Init	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$
Ant	0.05	0.97	0.99	0.99	0.97	<b>0.55</b>
Ant	0.5	0.96	0.96	0.96	0.96	<b>0.88</b>
Ant	0.9	<b>0.95</b>	0.97	0.97	0.97	0.96
Walker2d	0.05	0.66	0.82	0.79	0.80	<b>0.56</b>
Walker2d	0.5	0.82	0.85	0.92	0.90	<b>0.81</b>
Walker2d	0.9	<b>0.84</b>	0.89	0.91	0.95	0.94

Table 1: The final average fallback action of each base net in the Cascade net when standard Cascade was run on Ant-v4, Walker2d-v4 and with different fallback initializations for 6 million environment steps.  $\lambda_2$  is the average fallback of the second oldest base net and  $\lambda_6$  is the fallback of the one added last. The first base net isn't listed because it does not have a fallback action.

tendency towards fallbacking. This tendency is slightly stronger in Ant-v4 than in Walker2d-v4. For reference, if each base net would assign 0.5 to its fallback action then the bottom base net would only account for  $0.5^5 = 0.03125$  of the final output. Since in each setting (except the 0.05 initialization in Walker2d-v4) the fallback product is  $\geq 0.5$  means that it is the bottom base net that contributes most to the final output.

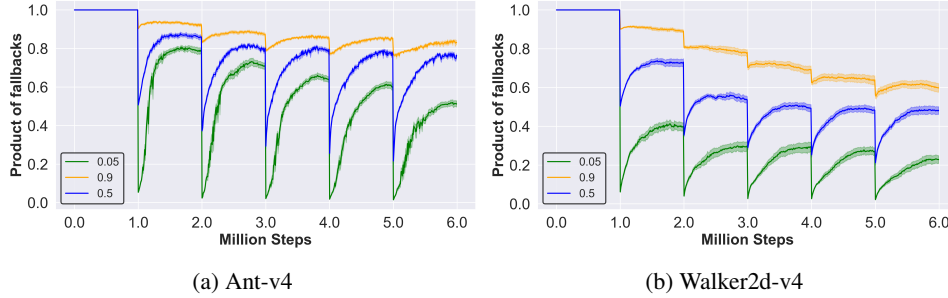


Figure 6: Product of all base agents' fallback actions over the course of training for 6 million environment steps with fallback initializations of 0.9 (yellow), 0.5 (blue) and (0.05) green.

To see how the product is made up and how different base nets contribute to it, the fallback action of the individual base nets was tracked. Tab. 1 lists the final average fallback action of each base net for all settings. In the case of the 0.05 or 0.5 initialization, it is the 6th base net that contributes either the most or second most to the final output of the Cascade net. In the case of the 0.9 initialization, the base nets (except the bottom one) contribute roughly equally to the final output. Note that the contribution of the  $i$ -th net is  $(1 - \lambda_i) \cdot \lambda_{i+1} \cdot \dots \cdot \lambda_6$  for  $i > 1$  and  $\lambda_2 \cdot \dots \cdot \lambda_6$  for  $i = 1$ . Still, in all cases (except for 0.9 initialization) and no matter the position in the Cascade, the final base agent's fallback is always significantly higher than what it started with, in most cases ending up well above 0.8.

Why for the cases 0.05 and 0.5 the 6th base net is the biggest or second biggest output-contributor can be best understood by looking at the fallback graph over the course of the training of an individual base agent. Fig. 7 shows this graph for the second to bottom base net for Ant-v4 and Walker2d-v4.

It becomes evident, that the base nets only reach their highest fallback value once they are no longer at the top of the Cascade. So naturally it is the last added base net (in this case the 6th) that has the lowest fallback value. On top of that, its final output contribution is directly given by  $1 - \lambda_6$  (i.e. no other nets can lower the contribution other than itself). This explains why it is either the biggest or second biggest contributor to the final output.

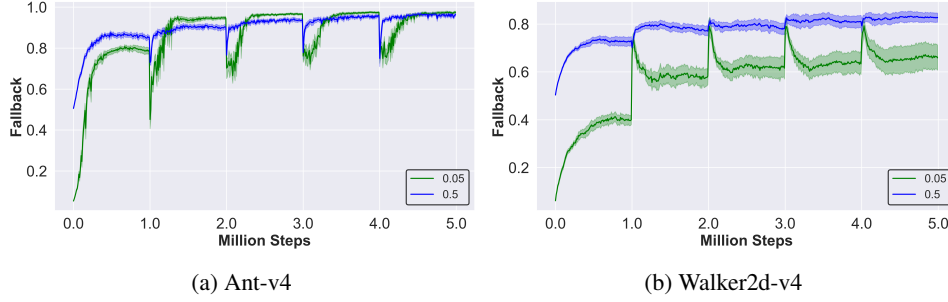


Figure 7: The average fallback action value for the second oldest base agent (i.e. the first one that has a fallback action) over the course of 5 million environment steps after this base agent has been added. This value has been plotted for Cascade with 0.5 (blue) and 0.05 (green) as the fallback initialization. Note that the spiky fallback change that occurs whenever a base net is added to the Cascade comes from the sudden change in the state distribution that goes along with adding a new base net.

**Question 5:** Do the base agents have a tendency towards fallbacking or not and which base agents contribute most to the final output?

Given different fallback action initializations of different magnitudes, all base agents always learned to assign a large value to the fallback action. Hence there is a clear tendency to fallback. In general, it is the bottom base net that contributes the most and the last added base net that contributes the second most to the final output.

#### 4.6 Tiny Cascade nets

To push the hypothesis to its limits, Cascade will be applied to extremely small base nets with only one hidden layer of size 1. A lower training duration of only half a million steps will also be considered to make it hard for Cascade to effectively chain base agents. 16 of these tiny base nets will be combined with Cascade. Fig. 8 shows the results of this experiment.

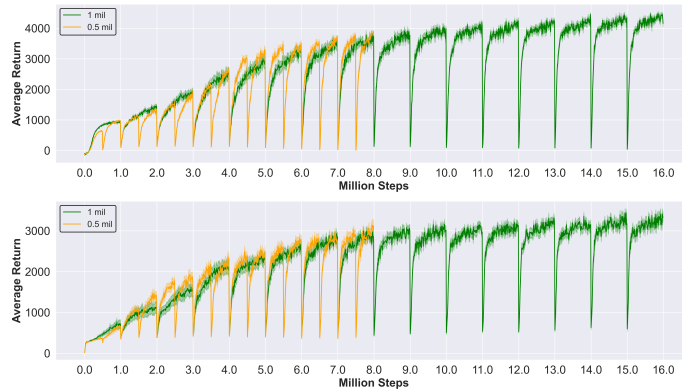


Figure 8: Performance graphs of Cascade with minimal base agents (one hidden layer of size 1) that are added every 1 million (green) or 0.5 million (yellow) steps for 16 million environment steps on Ant-v4 (top) and Walker2d-v4 (bottom).

As expected, the performances are much worse than those of standard Cascade. However, even though 16 base agents were chained, there was always a monotone performance increase from one iteration to another, even near the end when the Cascade consisted already of more than 10 base agents. At worst, the performance after a base agent was added only drew even to that of the previous iteration. This even held true for the case with only half a million training steps per iteration.

**Question 6:** Can a high number of base agents be chained without a loss in performance?

Yes. In all considered settings (Standard Cascade and Cascade with tiny base agents) adding additional agents to the Cascade has always been beneficial regardless of the training duration per iteration or the expressive capabilities of the added base agents.

## 5 Discussion

### 5.1 Reasons for the performance gain

Though the conducted experiments showed that the performance increase of Cascade over the standard PPO baseline is nontrivial and coming mainly from the fact that the base agents are added sequentially and that even older base nets are trainable at all times, the reasons why this is so effective remain unclear. Likely, there are many factors contributing to it. Two of these might be the following: First, regularly adding a net base agent helps with exploration: Directly after adding a base net the output of the entire cascade is greatly perturbed such that it initially becomes essentially random, as can be observed by the sharp drops in the performance graphs, thereby causing a huge shift in the state-distribution. Secondly, adding fresh networks prevents the Cascade net from losing its plasticity. Neural plasticity loss (Lyle et al., 2023) is a phenomenon where neural networks slowly lose their ability to train and adapt to new optimization targets which is present in reinforcement learning. In Section J in the supplementary materials, we measured the plasticity gain caused by adding a new ensemble member to the cascade.

### 5.2 Future work

First, pinpointing the exact causes of performance gain remains an open question that should be addressed in future work. In particular, it can be investigated why Cascade synergizes so well PPO yet yields only very little improvement in conjunction with other RL-algorithms. Even when considering only PPO, there are performance differences depending on the environment. It is worth investigating why in Hopper specifically, Cascade did not outperform PPO.

Future work might also include generalizing the strictly sequential nature of the resulting chain. It could be interesting to see what happens if the base nets choose between  $k$  nets (instead of  $k = 1$ ) to fall back to. One possibility to implement this is to train groups of  $k$  nets (perhaps as diverse as possible), each of which has fallback access to each of the  $k$  nets of the previous group. Or more extremely, the entire output is directly determined by the latest base agent in the cascade, which assigns a fallback weight to all other base nets. In standard Cascade, all agents remain trainable and it showed that freezing all base nets except the top one mostly hampers performance. However, it could be interesting to see if this holds true if some sort of attention mechanism is used to dynamically control which base nets are frozen and which remain trainable. Perhaps at a certain size of the cascade, no new agents should be added and it is beneficial to train inner agents and leave the top agent frozen.

Application to different domains: The behavior of Cascade could be studied on a larger suite of environments. For example, surveying the Atari benchmark could lead to new insights. Additionally, multi-agent problems with a combinatorial action space might be worth investigating as they come with a natural way of delegating the task between the action nets (one action net for each agent).

Hyperparameter optimization: The RL-baselines used for comparison were highly optimized version of PPO, SAC and DDPG. However, little to no hyperparameter optimization took place for Cascade. And the hyperparameters that were tuned, were high-level ones such as the number of base agent's training steps in Cascade. However, Cascade might be sensitive to the hyperparameters of the employed training algorithm.



## 6 Summary

In this work, an ERL method, called Cascade has been proposed and experimentally verified on some MuJoCo environments. With some minor high-level hyperparameter tweaking, Cascade managed to decisively outperform the PPO baseline on 4 out of 5 MuJoCo environments. Also, standard Cascade proved to be an effective tool for combining a large number of base agents no matter their training time or expressive capabilities. The characteristics of Cascade were investigated and it was found that the base agents have a natural tendency to assign a large value to their fallback action which led to the bottom base agent contributing most to the Cascade’s output. By observing the bottom base agent’s performance, it was discovered that Cascade makes heavy use of the freedom to change older agents in the cascade. In fact, keeping older base agents frozen proved detrimental to performance in a normalized environment. It also proved vital that the Cascade is built sequentially since starting with the final Cascade net also hampered performance thus showing that the architecture is not special but rather its sequential buildup.

## References

- Gang Chen and Victoria Huang. Ensemble reinforcement learning in continuous spaces – a hierarchical multi-step approach for policy training. In Edith Elkind (ed.), *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23*, pp. 3514–3522. International Joint Conferences on Artificial Intelligence Organization, 8 2023. DOI: 10.24963/ijcai.2023/391. URL <https://doi.org/10.24963/ijcai.2023/391>. Main Track.
- Thomas G. Dietterich. Ensemble Methods in Machine Learning. In *Multiple Classifier Systems*, pp. 1–15, Berlin, Heidelberg, 2000. Springer Berlin Heidelberg. ISBN 978-3-540-45014-6.
- Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1587–1596. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/fujimoto18a.html>.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In Jennifer G. Dy and Andreas Krause (eds.), *International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1856–1865. PMLR, 2018.
- Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Kinal Mehta, and João G.M. Araújo. Cleanrl: High-quality single-file implementations of deep reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022. URL <http://jmlr.org/papers/v23/21-1342.html>. Code last accessed on 9th July 2023.
- Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M. Czarnecki, Jeff Donahue, Ali Razavi, Oriol Vinyals, Tim Green, Iain Dunning, Karen Simonyan, Chrisantha Fernando, and Koray Kavukcuoglu. Population based training of neural networks. *CoRR*, abs/1711.09846, 2017. URL <http://arxiv.org/abs/1711.09846>.
- Piotr Januszewski, Mateusz Olko, Michał Królikowski, Jakub Swiatkowski, Marcin Andrychowicz, Łukasz Kuciński, and Piotr Miłoś. Continuous control with ensemble deep deterministic policy gradients. In *Deep RL Workshop NeurIPS 2021*, 2021. URL <https://openreview.net/forum?id=TIUfoXsnxB>.
- Kimin Lee, Michael Laskin, Aravind Srinivas, and Pieter Abbeel. SUNRISE: A simple unified framework for ensemble learning in deep reinforcement learning. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 6131–6141. PMLR, 2021. URL <http://proceedings.mlr.press/v139/lee21g.html>.

- Chao Li, Chen Gong, Qiang He, and Xinwen Hou. Keep various trajectories: Promoting exploration of ensemble policies in continuous control. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning, 2019.
- Guoqiang Liu, Gang Chen, and Victoria Huang. Policy ensemble gradient for continuous control problems in deep reinforcement learning. *Neurocomputing*, 548:126381, 2023. ISSN 0925-2312. DOI: <https://doi.org/10.1016/j.neucom.2023.126381>. URL <https://www.sciencedirect.com/science/article/pii/S0925231223005040>.
- Clare Lyle, Zeyu Zheng, Evgenii Nikishin, Bernardo Ávila Pires, Razvan Pascanu, and Will Dabney. Understanding plasticity in neural networks. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 23190–23211. PMLR, 2023. URL <https://proceedings.mlr.press/v202/lyle23b.html>.
- Rohan Saphal, Balaraman Ravindran, Dheevatsa Mudigere, Sasikanth Avancha, and Bharat Kaul. SEERL: sample efficient ensemble reinforcement learning. In Frank Dignum, Alessio Lomuscio, Ulle Endriss, and Ann Nowé (eds.), *AAMAS '21: 20th International Conference on Autonomous Agents and Multiagent Systems, Virtual Event, United Kingdom, May 3-7, 2021*, pp. 1100–1108. ACM, 2021. DOI: 10.5555/3463952.3464080. URL <https://www.ifaamas.org/Proceedings/aamas2021/pdfs/p1100.pdf>.
- Yanjie Song, Ponnuthurai Nagaratnam Suganthan, Witold Pedrycz, Junwei Ou, Yongming He, Ying-Wu Chen, and Yutong Wu. Ensemble reinforcement learning: A survey. *Appl. Soft Comput.*, 149 (Part A):110975, 2024. DOI: 10.1016/J.ASOC.2023.110975. URL <https://doi.org/10.1016/j.asoc.2023.110975>.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE, 2012. DOI: 10.1109/IROS.2012.6386109.