

# PINK NOISE IS ALL YOU NEED: COLORED NOISE EXPLORATION IN DEEP REINFORCEMENT LEARNING

Onno Eberhard<sup>1</sup>    Jakob Hollenstein<sup>2,1</sup>    Cristina Pinneri<sup>3,1</sup>    Georg Martius<sup>1</sup>

<sup>1</sup>Max Planck Institute for Intelligent Systems    <sup>2</sup>University of Innsbruck

<sup>3</sup>Max Planck ETH Center for Learning Systems

{firstname.lastname}@tuebingen.mpg.de

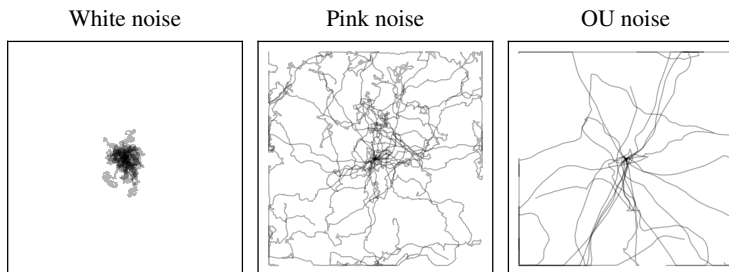
## ABSTRACT

In off-policy deep reinforcement learning with continuous action spaces, exploration is often implemented by injecting action noise into the action selection process. Popular algorithms based on stochastic policies, such as SAC or MPO, inject white noise by sampling actions from uncorrelated Gaussian distributions. In many tasks, however, white noise does not provide sufficient exploration, and temporally correlated noise is used instead. A common choice is Ornstein-Uhlenbeck (OU) noise, which is closely related to Brownian motion (red noise). Both red noise and white noise belong to the broad family of *colored noise*. In this work, we perform a comprehensive experimental evaluation on MPO and SAC to explore the effectiveness of other colors of noise as action noise. We find that pink noise, which is halfway between white and red noise, significantly outperforms white noise, OU noise, and other alternatives on a wide range of environments. Thus, we recommend it as the default choice for action noise in continuous control.

## 1 INTRODUCTION

Exploration is vitally important in reinforcement learning (RL) to find unknown high reward regions in the state space. This is especially challenging in continuous control settings, such as robotics, because it is often necessary to coordinate behavior over many steps to reach a sufficiently different state. The simplest exploration method is to use action noise, which adds small random perturbations to the policy’s actions. In off-policy algorithms, where the exploratory behavioral policy does not need to match the target policy, action noise may be drawn from any random process. If the policy is deterministic, as in DDPG (Lillicrap et al., 2016) and TD3 (Fujimoto et al., 2018), action noise is typically white noise (drawn from temporally uncorrelated Gaussian distributions) or Ornstein-Uhlenbeck (OU) noise, and is added to the policy’s actions. In algorithms where the policy is stochastic, such as SAC (Haarnoja et al., 2018) or MPO (Abdolmaleki et al., 2018), the action sampling itself introduces randomness. As the sampling noise is typically uncorrelated over time, these algorithms effectively employ a scale-modulated version of additive white noise, where the noise scale varies for different states.

Figure 1: Trajectories of pure noise agents on a bounded integrator environment (Sec. 6). White action noise (left) does not reach far in this environment, and it would not be able to collect a sparse reward at the edges: it *explores locally*.



OU noise (right) only *explores globally* and gets stuck at the edges. Pink noise (center) provides a balance of local and global exploration, and covers the state space more uniformly than the other two.

In many cases, white noise exploration is not sufficient to reach relevant states. Both MPO and SAC have severe problems with certain simple tasks like MountainCar because of inadequate exploration. As in TD3 or DDPG, the off-policy nature of these algorithms makes it possible to replace the white noise process, which is implicitly used for action sampling, by a different random process. The effectiveness of temporal correlation in the action selection has been noted before (e.g. Osband et al., 2016) and is illustrated in Figure 1, where the exploration behavior of white noise (uncorrelated) is compared to that of noises with intermediate (pink noise) and strong (OU noise) temporal correlation on a simple integrator environment (more on this in Section 6). Using highly correlated noise, such as OU noise, can yield sufficient exploration to deal with these hard cases, but it also introduces a different problem: strongly off-policy trajectories. Too much exploration is not beneficial for learning a good policy, as the on-policy state-visitation distribution must be covered during training to make statistical learning possible. Thus, a typical approach is to use white noise by default, and alternatives like OU noise only when necessary. In this work, our goal is to find a better strategy, by considering noises with intermediate temporal correlation, in the hope that these work well both on environments where white noise is enough, and on those which require increased exploration.

To this end, we investigate the effectiveness of *colored noise* as action noise in deep RL. Colored noise is a general family of temporally correlated noise processes with a parameter  $\beta$  to control the correlation strength. It generalizes white noise ( $\beta = 0$ ) and Brownian motion (red noise,  $\beta = 2$ ), which is closely related to OU noise. We find that average performance across a broad range of environments can be increased significantly by using colored action noise with intermediate temporal correlation ( $0 < \beta < 2$ ). In particular, we find *pink noise* ( $\beta = 1$ ) to be an excellent default choice. Interestingly, pink noise has also been observed in the movement of humans: the slight swaying of still-standing subjects, as well as the temporal deviations of musicians from the beat, have both been measured to exhibit temporal correlations in accord with pink noise (Duarte & Zatsiorsky, 2001; Hennig et al., 2011).

Our work contributes a comprehensive experimental evaluation of various action noise types on MPO and SAC. We find that pink noise has not only the best average performance across our selection of environments, but that in 80% of cases it is not outperformed by any other noise type. We also find that pink noise performs on par with an oracle that tunes the noise type to an environment, while white and OU noise perform at 50% and 25% between the worst noise type selection and the oracle, respectively. To investigate whether there are even better noise strategies, we test a color-schedule that goes from globally exploring red noise to locally exploring white noise over the course of training, as well as a bandit method to automatically tune the noise color to maximize rollout returns. Both methods, though they significantly improve average performance when compared to white and OU noise, are nevertheless significantly outperformed by pink noise. In addition to the results of our experiments, we attempt to explain why pink noise works so well as a default choice, by constructing environments with simplified dynamics and analyzing the different behaviors of pink, white and OU noise. *Our recommendation is to switch from the current default of white noise to pink noise.*

## 2 BACKGROUND & RELATED WORK

Reinforcement learning (RL) has achieved impressive results, particularly in the discrete control setting, such as achieving human-level performance in Atari games with DQN (Mnih et al., 2015) or mastering the game of Go (Silver et al., 2016) by using deep networks as function approximators. In this paper, we are concerned with the continuous control setting, which is especially appropriate in robotics. In continuous action spaces, it is typically intractable to choose actions by optimizing a value function over the action space. This makes many deep RL methods designed for discrete control, such as DQN, not applicable. Instead, researchers have developed policy search methods (e.g. Williams, 1992; Silver et al., 2014), which directly parameterize a policy. These methods can be divided into on-policy algorithms, such as TRPO (Schulman et al., 2015) and PPO (Schulman et al., 2017), and off-policy algorithms such as DDPG, TD3, SAC and MPO.

All these algorithms have to address the problem of exploration, which is fundamental to RL: in order to improve policy performance, agents need to explore new behaviors while still learning to act optimally. One idea to address exploration is to add a novelty bonus to the reward (e.g. Thrun, 1992). In deep RL, this can be done by applying a bonus based on sample density (Tang et al., 2017) or prediction error (Burda et al., 2019). Another method to encourage exploration is to take inspiration from bandit methods like Thompson sampling (e.g. Russo et al., 2018), and act optimistically with

respect to the uncertainty in the Q-function (Osband et al., 2016). The simplest strategy, however, is to randomly perturb either the policy parameters (Plappert et al., 2018; Mania et al., 2018), or the actions themselves by randomly sampling a function for each episode that deterministically alters the action selection (Raffin & Stulp, 2020), by learning correlations between action dimensions and state space dimensions to induce increasing excitation in the environment (Schumacher et al., 2022), or by learning an action prior from task-agnostic data (Bagatella et al., 2022).

In this work, we consider the simplest and most common form of exploration in continuous control: action noise. Action noise can be either explicitly added to the policy, or implicitly, by randomly sampling actions from a stochastic policy. The most common form of action noise is white noise, which typically comes from sampling from independent Gaussian distributions at every time step. Apart from white action noise, Lillicrap et al. (2016) successfully used temporally correlated Ornstein-Uhlenbeck noise, and Pinneri et al. (2020) achieved improvements in model predictive control by utilizing colored noise. Inspired by these successes, in this work we investigate the effectiveness of colored action noise in the context of model-free RL, specifically on MPO and SAC.

### 3 METHOD

In this paper, we investigate exploration using action noise. In algorithms like DDPG and TD3, where the learned policy  $\mu$  is deterministic, action noise is simply added to the policy:

$$a_t = \mu(s_t) + \sigma \varepsilon_t, \tag{1}$$

where  $\varepsilon_{1:T} = (\varepsilon_1, \dots, \varepsilon_T)$  is sampled from a random process, and  $\sigma$  is a scale parameter. If  $\varepsilon_t$  is sampled independently at every time step, e.g. from a Gaussian distribution, then  $\varepsilon_{1:T}$  is called *white noise* (WN). This is the prevailing choice of action noise, though it is also common to use time-correlated Ornstein-Uhlenbeck noise ( $\varepsilon_{1:T} \sim \text{OU}_T$ ) (Uhlenbeck & Ornstein, 1930).

Algorithms which parameterize a stochastic policy, such as SAC and MPO, also use action noise. In continuous action spaces, the most common policy distribution is a diagonal Gaussian, represented by the functions  $\mu(s_t)$  and  $\sigma(s_t)$ :  $a_t \sim \mathcal{N}(\mu(s_t), \text{diag}(\sigma(s_t)))$ . This can equivalently be written as

$$a_t = \mu(s_t) + \sigma(s_t) \odot \varepsilon_t, \tag{2}$$

where  $\varepsilon_t \sim \mathcal{N}(0, I)$ . In this case, the action noise  $\varepsilon_{1:T}$  is again Gaussian white noise, which is scale-modulated by the function  $\sigma$ .

White noise is not correlated over time ( $\text{cov}[\varepsilon_t, \varepsilon_{t'}] = 0$ ). In some environments, this leads to very slow exploration, which in turn leads to inadequate state space coverage, leaving high reward regions undiscovered. Thus, it is often beneficial to use action noise with temporal correlation ( $\text{cov}[\varepsilon_t, \varepsilon_{t'}] > 0$ ), like Ornstein-Uhlenbeck (OU) noise. OU noise was recommended as the default choice for DDPG, and has been shown to lead to a significant increase in state space coverage (Hollenstein et al., 2022). OU noise is defined by the stochastic differential equation (SDE)

$$\dot{\varepsilon}_t = -\theta \varepsilon_t + \sigma \eta_t, \tag{3}$$

where  $\eta_t$  is a white noise process. If  $\theta = 0$ , then this equation defines integrated white noise, also called Brownian motion. Brownian motion is temporally correlated, but cannot be used as action noise if generated in this way, because its variance increases unboundedly over time, violating the action space limits. This problem is addressed by setting  $\theta > 0$  (a typical choice is  $\theta = 0.15$ ), which bounds the variance. More details about OU noise and Brownian motion can be found in Section A.

A broad family of temporally correlated noises is given by *colored noise*, which generalizes both white noise and Brownian motion (in this context called *red noise*).

**Definition 1** (Colored noise). A stochastic process is called colored noise with color parameter  $\beta$ , if signals  $\varepsilon(t)$  drawn from it have the property  $|\hat{\varepsilon}(f)|^2 \propto f^{-\beta}$ , where  $\hat{\varepsilon}(f) = \mathcal{F}[\varepsilon(t)](f)$  denotes the Fourier transform of  $\varepsilon(t)$  ( $f$  is the frequency) and  $|\hat{\varepsilon}(f)|^2$  is called the power spectral density (PSD).

The color parameter  $\beta$  controls the amount of temporal correlation in the signal. The PSDs of colored noise with different  $\beta$  are shown in Figure A.2. If  $\beta = 0$ , then the signal is uncorrelated, and the PSD is flat, meaning that all frequencies are equally represented. This noise is called white noise in analogy to light, where a signal with equal power on all visible frequencies is perceived as white.

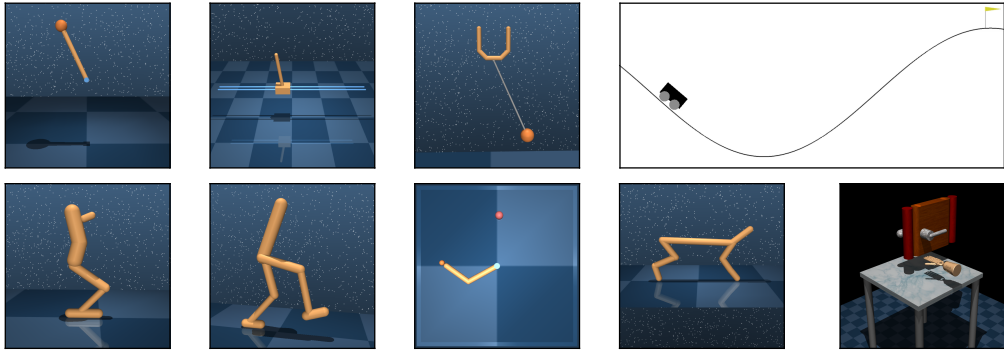


Figure 2: The environments we use: Pendulum, CartPole (balance + swingup tasks), Ball-In-Cup, MountainCar, Hopper, Walker, Reacher, Cheetah, Door. See Section C for more details. Images partly taken from Tassa et al. (2018) with permission.

Red noise ( $\beta = 2$ ) is named so, because it has more weight on lower frequencies, which in light corresponds to the red part of the spectrum. Gaussian colored noise with constant variance can be efficiently generated, and the complete noise signal for an episode can be sampled at once, to be used as action noise according to Equations (1) and (2). If generated like this, which we denote by  $\varepsilon_{1:T} \sim \text{CN}_T(\beta)$  (more details in Section A), white noise is identical to independently sampling from a Gaussian distribution at every time step. Red noise ( $\text{CN}_T(2)$ ) is very similar to OU noise with the default setting  $\theta = 0.15$ , as both are essentially Brownian motion with bounded variance (see Figure A.2).<sup>1</sup> By setting  $0 < \beta < 2$ , colored noise allows us to search for a better default action noise type with intermediate temporal correlation between white and red noise. One special case is *pink noise*, which is defined by  $\beta = 1$ .

#### 4 IS PINK NOISE ALL YOU NEED?

Fujimoto et al. (2018) found that the type of action noise (white or OU) in general does not influence the performance of TD3. In contrast to this, Hollenstein et al. (2022) found that the noise type does have an influence, but that the impact of this choice, as well as which noise type is preferable, depends entirely on the environment. We start by confirming these latter results<sup>2</sup>, and compare white and OU action noise with a selection of colored action noises ( $\beta \in [0, 2]$ ), in terms of the achieved performance. In all of our experiments we use MPO and SAC, relying on the implementations from Raffin et al. (2021) and Pardo (2020), respectively. We found that both algorithms significantly outperform TD3 across tasks, and thus only briefly discuss TD3 in Section B.1. Since the optimality of an action noise type depends on the environment, we perform experiments on a diverse set of 10 different tasks taken from the DeepMind Control Suite (Tassa et al., 2018), OpenAI Gym (Brockman et al., 2016), and the Adroit hand suite (Rajeswaran et al., 2018). These environments are shown in Figure 2 and are described in more detail in Section C.

To evaluate the performance of a training procedure (which always lasts  $10^6$  environment interactions), we run 5 evaluation rollouts every  $10^4$  interactions. We then report the performance as the mean return of all these evaluation rollouts. Since this performance is related to the area under the learning curve, it is a measure combining both the final policy performance, and the sample efficiency of an algorithm. More detailed results, including learning curves and an analysis of the final policy performance, can be found in Sections B.2 and G.

##### 4.1 DOES THE NOISE TYPE MATTER?

To assess the importance of the choice of action noise, we evaluate the performances achieved by SAC and MPO when using white noise, OU noise and colored noise as action noise (where

<sup>1</sup>Other settings of  $\theta$  for OU noise, which are less similar to red noise, are discussed in Section A. In the main text we only consider the default setting of  $\theta = 0.15$ .

<sup>2</sup>We also confirm the former results (see Section B.1), but find that colored noise (especially pink noise) outperforms both white and OU noise on TD3.

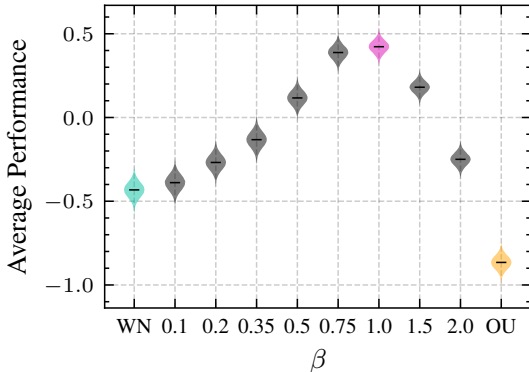


Figure 3: Bootstrap distributions for the expected average performance of MPO and SAC using different action noise types (details in Section B.2). Highlighted are white noise (WN), pink noise ( $\beta = 1$ ), and Ornstein-Uhlenbeck noise (OU).

Environment	Best noise	$p$	Pink?
Pendulum	2.0	<b>0.01</b>	<b>×</b>
Cartpole (b.)	<b>1.0 (Pink)</b>	—	<b>✓</b>
Cartpole (s.)	<b>1.0 (Pink)</b>	—	<b>✓</b>
Ball-In-Cup	0.75	0.88	<b>✓</b>
MountainCar	2.0	0.59	<b>✓</b>
Hopper	<b>1.0 (Pink)</b>	—	<b>✓</b>
Walker	0.5	0.36	<b>✓</b>
Reacher	White noise	<b>0.02</b>	<b>×</b>
Cheetah	0.75	0.62	<b>✓</b>
Door	0.75	0.65	<b>✓</b>

Table 1: A Welch  $t$ -test reveals that the performance difference between pink noise and the *best noise* is only significant in two out of ten environments. The rightmost column answers whether pink noise performs equally well as the best noise type.

$\beta \in \{0.1, 0.2, 0.35, 0.5, 0.75, 1, 1.5, 2\}$ ), on the benchmark environments shown in Figure 2. We repeat all learning runs with 20 different seeds, resulting in a total of  $20 \times 10 \times 2 \times 10 = 4000$  experiments<sup>3</sup>, each one reporting a single scalar performance. To control for the influence of the algorithm and environment on the performance of a particular noise type, we group all results by algorithm and task, and normalize them to zero mean and unit variance. We then calculate a noise type’s *average performance*: the normalized performance of all runs using this noise type, averaged across algorithms and environments. In Figure 3, bootstrap distributions for the expected average performances are shown, generated using the 20 random seeds available per task and algorithm (more details in Section B.2). It can be seen that the noise type indeed matters for performance. A clear preference for pink noise ( $\beta = 1$ ) becomes visible, which considerably outperforms white noise and OU noise across tasks. In Section B.2, this performance difference can be seen on the corresponding learning curves, where we compare white, OU and pink noise.

Achieving the best average performance across environments is not the same as being the best performing option on each individual environment. This begs the question of when pink noise (the best general option) also performs as good as an environment’s best noise type. We perform a Welch  $t$ -test for each environment to check whether the expected difference between the performances of pink noise and the task-specific best noise<sup>4</sup> is significant. The results are listed in Table 1. Although pink noise only achieves the highest mean across seeds in three of the ten tasks, the statistical analysis reveals that the difference between pink noise and the best noise type is only significant in two out of ten cases. In other words, in the tested environments, pink noise performs on par with the best choice of noise type in 80% of cases! What about the two environments, Pendulum and Reacher, where pink noise is outperformed by other noise types? On Pendulum, pink noise, on average, achieves 83% of the performance of red noise ( $\beta = 2$ ). In contrast, white noise only achieves 39% of the performance of red noise (OU performs similarly to red noise). On Reacher, pink noise achieves 99% of white noise’s performance, while OU noise achieves only 76%. So, even on the few environments where pink noise is outperformed significantly, it is clearly preferable as a default over white noise and OU noise. These results indicate that *pink noise* seems to be *all you need*.

#### 4.2 IS PINK NOISE A GOOD DEFAULT?

The best performance on a given environment is always achieved with the task-specific best action noise type. It would be nice to always use this best noise type, but it is often unpractical to run a large hyperparameter search to find it, especially when including many possible colors  $\beta$ . It is common therefore, to stick to a “default” choice, which is typically white noise. In the previous section, we saw that pink noise is a better default choice than white noise, but it is still unclear whether this is

<sup>3</sup>seeds  $\times$  tasks  $\times$  algorithms  $\times$  noise types (WN + OU + 8 colors)

<sup>4</sup>To compute the *best noise* type, we normalize out the contribution of the algorithm (similarly to the *average performance*) and take the mean performance over random seeds on each environment.

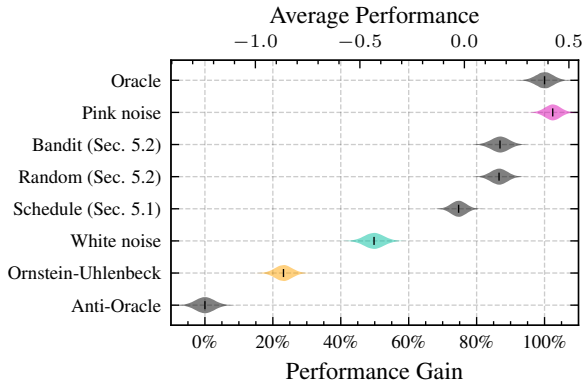


Figure 4: Bootstrap distributions for the expected average performances of all methods we discuss in this paper. Highlighted are again white noise and OU noise (the popular options), as well as pink noise (our suggestion). While OU noise and white noise only achieve about 23% and 50% of the possible performance gain of an oracle method, pink noise performs equally to the oracle! Pink noise is also not outperformed by a color-scheduling method (Section 5.1), nor by a bandit algorithm (Section 5.2).

enough, or if a hyperparameter search might be needed for good performance. In this section, we will analyze how much performance is lost by sticking to a default value of white noise, pink noise, or OU noise, compared to using the task-specific best noise type.

We choose the task-specific best noise type via an “oracle” method, which can be thought of as a very extensive grid search: an environment’s best noise type is selected by looking at the results of all noise types, and choosing the best performing option across 10 seeds. By also doing an “anti-oracle” experiment, which selects the worst noise type on each environment (i.e. the most unlucky pick of noise types possible), we can define a new “performance gain” measure, which might be easier to interpret than the average performance in the previous section.<sup>5</sup> The performance gain of a noise type specifies where its average performance falls between the anti-oracle’s performance (0%) and the oracle’s performance (100%). Figure 4 presents the performance gains of using white noise, Ornstein-Uhlenbeck noise and pink noise as a default for all environments.

By always sticking to Ornstein-Uhlenbeck noise, only about 25% of the highest possible performance gain is achieved, and the resulting performance would be closer to using the anti-oracle. By using white noise instead, we already achieve a performance gain of over 50%. However, picking pink noise does not appear to sacrifice any performance compared to the oracle!<sup>6</sup> The gain achieved by switching from white noise as the default to pink noise is both considerable and significant, and we recommend switching to *pink noise as the default action noise*.

## 5 ALL THE COLORS OF THE RAINBOW

We have found that pink noise is the best default action noise over a broad range of environments. There are still some environments, however, where pink noise is outperformed by other noise types, specifically by white and red noise on the Reacher and Pendulum tasks, respectively. This indicates that there may not exist a single noise type which performs best on all environments. However, this consideration is only valid if the noise is kept constant over the course of training. If we instead try a different approach, and choose the noise type separately for each rollout, we may find a strategy that is outperformed nowhere. In this section we discuss two such non-constant methods, which differ in the way a rollout’s noise type is selected: a color-schedule going from  $\beta = 2$  to  $\beta = 0$ , and a bandit approach with the intention of finding the optimal color for an environment.

### 5.1 IS COLOR-SCHEDULING BETTER THAN PINK NOISE?

To find a more effective exploration method than pink noise, it is helpful to understand why learning in the Pendulum and Reacher environments works better with other noise types. The Pendulum environment is underactuated and requires a gradual build-up of momentum by slowly swinging back and forth. Strongly correlated action noise, such as red noise, makes this behavior much more likely. The Reacher task, on the other hand, has neither a particularly large state space, nor does it exhibit

<sup>5</sup>The (anti-)oracle is evaluated on the 10 seeds not used for noise type selection to avoid sampling bias. We repeat this by selecting (evaluating) once on the first (latter) 10 seeds, and once on the latter (first) 10 seeds.

<sup>6</sup>It looks as if pink noise is even exceeding the oracle’s performance. This difference is not statistically significant and is due to the oracle only having access to the 10 seeds not used for evaluation.

Environment	P > S	P < S	Pink $\geq$ Schedule?	P > B	P < B	Pink $\geq$ Bandit?
Pendulum		<0.01	✗		0.73	✓
Cartpole (b.)	<0.01		✓	0.10		✓
Cartpole (s.)	0.01		✓	<0.01		✓
Ball-In-Cup	<0.01		✓	0.47		✓
MountainCar		0.15	✓	0.29		✓
Hopper	0.05		✓	0.15		✓
Walker	<0.01		✓		0.28	✓
Reacher	<0.01		✓		0.50	✓
Cheetah	0.05		✓	0.64		✓
Door	0.15		✓	0.27		✓

Table 2: How does pink noise compare to a schedule and a bandit method? For each environment, we perform a Welch  $t$ -test to test for inequality of the performances of pink noise vs. the schedule/bandit method. The  $p$ -values are arranged to show which performance is larger. Pink noise performs significantly better than the schedule on most environments. Compared to the bandit algorithm, pink noise performs better overall, but on most environments the difference is not significant. In the “Pink  $\geq$  ...?” columns, (✓) means that pink noise does not perform significantly worse than the alternative.

under-actuation, such that white noise is well suited to explore the space. Here, temporally correlated noise will only lead to off-policy trajectory data, thereby inhibiting learning. In general, strongly correlated noise leads to more global exploration, while uncorrelated noise explores more locally. We return to these ideas in Section 6, where we analyze the effects of high and low correlation on two simple environments.

Our method should work well on both of these environments, and on environments which require a mix of local and global exploration. Thus, we are looking for a strategy which balances local and global exploration. A simple idea to do this is a color-schedule: start with highly correlated red noise ( $\beta = 2$ ) and then slowly decrease  $\beta$  to white noise ( $\beta = 0$ ) over the course of training. The rationale behind this strategy is that high-reward regions can be quickly discovered at the beginning of training when  $\beta$  is large, while the trajectories get more on-policy over time, helping with environments like Reacher. Indeed, a similar approach that schedules the action noise *scale*, has been shown to work quite well (Hollenstein et al., 2022).

We implement a  $\beta$ -schedule, which linearly goes from  $\beta = 2$  to  $\beta = 0$ , on MPO and SAC and repeat the experiment with 20 random seeds on all environments. Bootstrap distributions for the expected average performance across environments are shown in Figure 4 (denoted by *Schedule*). The results indicate that the schedule is generally better than OU and white noise, but does not outperform pink noise. Indeed, pink noise is significantly better, as the confidence intervals do not overlap. If we take a more detailed look at the individual environments (Table 2), we see that thanks to the additional highly correlated noise, the schedule does outperform pink noise on the Pendulum environment, as expected. However, in all other environments pink noise either significantly outperforms the schedule or they perform on par, so our recommendation to use pink noise as a default remains.

## 5.2 IS BANDIT COLOR SELECTION BETTER THAN PINK NOISE?

The results in the previous section indicate that, while changing the noise type over the training process can improve performance, simply moving from globally exploring red noise to more locally exploring white noise does not outperform pink noise. Instead of trying to find a different schedule to fit all environments, in this section we consider an adaptive approach. By using a bandit algorithm to select the action noise color for each rollout on the basis of past rollout returns, it might be possible to find not only the general best noise for a given environment, but even to automatically adapt the noise to different stages of training. The bandit algorithm we use is based on Thompson sampling, the details are explained in Section D.

We use the rollout return itself as the bandit reward signal. The reasoning for this is that in environments where strong exploration is necessary (such as Pendulum and MountainCar), high return will only be achieved by strongly correlated actions. On the other hand, if environments do not require correlated actions, or a capable policy has been learned, the highest return should be achieved by the action noise which least disturbs the policy, i.e. noise with a low correlation.

As an additional baseline, we also perform an experiment where a color ( $\beta$ ) is randomly selected for each rollout.<sup>7</sup> For both methods we use the same list of  $\beta$  values as in Section 4.1 (incl.  $\beta = 0$ ), and repeat the experiments with 20 random seeds. The results on MPO and SAC are shown in Figure 4 (marked with *Bandit* and *Random*) and Table 2. As visible in Figure 4, the bandit method is again outperformed by pink noise: a bootstrapping test yields a highly significant difference in expected average performance across environments ( $p = 0.005$ ). Looking at the results on the individual tasks (Table 2), it seems like the bandit method indeed outperforms pink noise on the two problematic environments (Pendulum and Reacher), however, this difference is not significant. A detailed comparison of the bandit and its random baseline can be found in Table D.1, which shows that neither of the two methods significantly outperforms the other on any environment. This indicates that, while there may be merit in changing the noise type over the training process, the rollout return appears to contain too little information to effectively guide the noise type selection. Thus, our recommendation to use pink noise as a default remains unchanged.

## 6 HOW DO ACTION NOISE AND ENVIRONMENT DYNAMICS INTERACT?

Why is pink noise such a good default noise type? In Section 5.1, we briefly discussed the concepts of local and global exploration, and hypothesized that the best exploration behavior provides a balance of the two, such that high reward regions will be found, while trajectories are not too off-policy. To analyze how different noise types behave, we will look at a simplified *bounded integrator* environment: a velocity-controlled 2D particle moving in a box. If we control this particle purely by noise, we can analyze the exploration behavior in isolation of a policy. As a first test, we run 20 episodes of 1000 steps in an environment of size  $250 \times 250$  with white noise, pink noise, and OU noise (all with scale  $\sigma = 1$ ,  $x$ - and  $y$ -velocity controlled independently). The resulting trajectories are shown in Figure 1. It can be seen that pink noise provides the best combination of local and global exploration: it reaches the edges (unlike white noise), but does not get stuck there (unlike OU noise).

A good mix of local and global exploration gives rise to a more uniform state space coverage, as can be seen in Figure 1. Thus, how well a noise type explores depends highly on the size of the environment: if the environment was much smaller, white noise would be enough to cover the space and pink noise trajectories would look similar to the OU trajectories shown here. On the other hand, if the environment were bigger, then pink noise would not reach the edges and OU noise would explore better. The uniformity of the state space coverage is measured by the entropy of the state-visitation distribution. We estimate the entropy induced by a noise type using a histogram density approximation: we partition the state space into a number of boxes ( $50 \times 50 = 2500$  boxes), sample  $10^4$  trajectories, and count the number of sampled points in each box.

Figure 5 shows the entropy achieved by white noise, OU noise and pink noise as a function of the environment size. The sizes are chosen to reflect the complete sensible range for episode lengths of 1000 steps and noise scales of  $\sigma = 1$ : from very small ( $50 \times 50$ ) to very large ( $2000 \times 2000$ ). Pink noise is not “special” in the sense that it performs best on all environments, as we already saw in the previous sections. However, it performs best on “medium scales”, as determined by the episode length, and does not suffer from severe degradation in performance over the whole spectrum of sensible environments. If we do not know where on this spectrum a given environment lies, then pink noise is clearly a better default choice than white noise or OU noise!

Besides integrating actions, another common aspect of environment dynamics is oscillation. Oscillation dynamics are dominant in the Pendulum and MountainCar environments<sup>8</sup>, but also relevant in other domains, like Ball-In-Cup, Cartpole, and Walker. To model these dynamics, we construct a second environment: a simple harmonic oscillator. This system is a frictionless 1-dimensional physical setup, in which a mass  $m$  is attached to an ideal spring of stiffness  $k$ . The state space consists of the mass’s position and velocity, and the action describes a force that is applied to the mass. The goal is to maximize the energy in the oscillator system (which is equivalent to maximizing the amplitude), similar to the MountainCar and Pendulum tasks, where this is necessary to collect the sparse reward.

The oscillator environment is parameterized by the resonant frequency  $f$  of the system, which is fixed by setting the stiffness  $k = 4\pi^2$  and the mass  $m = 1/f^2$  (more details in Section F). Figure 5

<sup>7</sup>This method would be roughly equivalent to the bandit method if we provided no bandit reward signal.

<sup>8</sup>See Section E for a simple method exploiting this property to solve MountainCar.



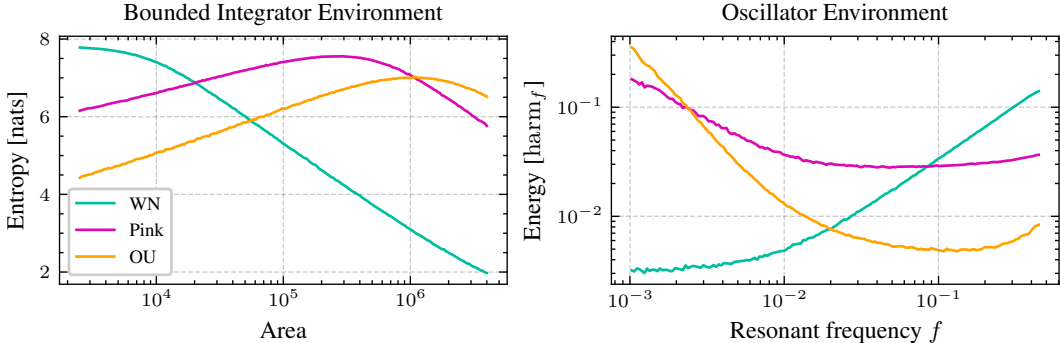


Figure 5: Pink noise strikes a favorable middle ground between white noise and Ornstein-Uhlenbeck noise on a wide range of environments. On both a bounded integrator environment parameterized by its size (left), and on a simple harmonic oscillator environment parameterized by its resonant frequency (right), it is much more general in terms of the range of parameters which yield good results, and performs well on the complete range of reasonable parameterizations. We argue that this quality is what makes it a good default.

shows the average energy in the oscillator system (over episodes of 1000 steps) as a function of the resonant frequency  $f$ , which we vary from very low ( $f = \frac{1}{1000}$ , episode length = 1 period) to very high ( $f = \frac{1}{2}$ , Nyquist frequency), when driven by white noise, pink noise, and OU noise. The energy is measured relative to the average energy achieved by optimal actions (sinusoidal excitation at the resonant frequency), denoted  $\text{harm}_f$ . Even though this is a completely different setup to the bounded integrator, and we are measuring a completely different quantity, the two plots look remarkably similar. Again, this shows the power of pink noise as a default action noise: if we do not know the resonant frequency of the given environment, pink noise is the best choice.

These two environments (bounded integrator and oscillator) are rather simplistic. However, the dynamics of many real systems undoubtedly contain parts which resemble oscillations (when a spring or pendulum is present), single or double integration (when velocities/steps or forces/torques are translated into positions) or contact dynamics (such as the box in the bounded integrator). If an environment’s dynamics are very complex, i.e. they contain many such individual parts, then the ideal action noise should score highly on each of these “sub-tasks”. However, if all these individual parts have different parameters (like the environment size or resonant frequency above), it stands to reason that the best single action noise would be the one which is general enough to play well with all parameterizations, i.e. *pink noise*. On the flip side, the average performance in Figure 3 over all environments may be interpreted as the performance over a very complicated environment, with the sub-tasks being the “actual” environments. This might explain why we see this curve: all sub-tasks have very different parameters, and require different action noises (as seen in Table 1), but pink noise is general enough to work well on all sub-tasks, and thus easily outperforms noise types like white noise or OU noise, which are only good on very specific environments (see Figure 5).

## 7 CONCLUSION

In this work we performed a comprehensive experimental evaluation of colored noise as action noise in deep reinforcement learning for continuous control. We compared a variety of colored noises with the standard choices of white noise and Ornstein-Uhlenbeck noise, and found that pink noise outperformed all other noise types when averaged across a selection of standard benchmarks. Pink noise is only significantly outperformed by other noise types on two out of ten environments, and overall performs equally well to an oracle selection of the noise type. Additionally, we compared pink noise to more sophisticated methods that change the noise type over the course of training: a color-schedule, a bandit method, and a random selection scheme. No method outperforms pink noise, and our recommendation is to *use pink noise as the default action noise*. Finally, we studied the behaviors of pure noise agents on two simplified environments: a bounded integrator and a harmonic oscillator. The results showed that pink noise is much more general with respect to the environment parameterization than white noise and OU noise, which sheds some light on why it performs so well as the default choice.

## ACKNOWLEDGMENTS

We want to thank Marco Bagatella, Sebastian Blaes, and Pierre Schumacher for helpful feedback on earlier revisions of this text, and the Max Planck ETH Center for Learning Systems for supporting Cristina Pinneri. Funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy – EXC number 2064/1 – Project number 390727645.

## REFERENCES

- Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Rémi Munos, Nicolas Heess, and Martin A. Riedmiller. Maximum a posteriori policy optimisation. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=SlANxQW0b>. 1
- Marco Bagatella, Sammy Joe Christen, and Otmar Hilliges. SFP: State-free priors for exploration in off-policy reinforcement learning. *Transactions on Machine Learning Research*, 2022. URL <https://openreview.net/forum?id=qYNfwFCX9a>. 3
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *CoRR*, abs/1606.01540, 2016. URL <http://arxiv.org/abs/1606.01540>. 4
- Yuri Burda, Harrison Edwards, Amos J. Storkey, and Oleg Klimov. Exploration by random network distillation. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=H1lJJnR5Ym>. 2
- Marcos Duarte and Vladimir M. Zatsiorsky. Long-range correlations in human standing. *Physics Letters A*, 283(1):124–128, 2001. ISSN 0375-9601. doi: [https://doi.org/10.1016/S0375-9601\(01\)00188-8](https://doi.org/10.1016/S0375-9601(01)00188-8). URL <https://www.sciencedirect.com/science/article/pii/S0375960101001888>. 2
- Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1582–1591. PMLR, 2018. URL <http://proceedings.mlr.press/v80/fujimoto18a.html>. 1, 4
- 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.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1856–1865. PMLR, 2018. URL <http://proceedings.mlr.press/v80/haarnoja18b.html>. 1
- Holger Hennig, Ragnar Fleischmann, Anneke Fredebohm, York Hagmayer, Jan Nagler, Annette Witt, Fabian J. Theis, and Theo Geisel. The nature and perception of fluctuations in human musical rhythms. *PLOS ONE*, 6(10):1–7, 10 2011. doi: 10.1371/journal.pone.0026457. URL <https://doi.org/10.1371/journal.pone.0026457>. 2
- Jakob Hollenstein, Sayantan Auddy, Matteo Saveriano, Erwan Renaudo, and Justus Piater. Action noise in off-policy deep reinforcement learning: Impact on exploration and performance. *CoRR*, abs/2206.03787, 2022. doi: 10.48550/arXiv.2206.03787. URL <https://doi.org/10.48550/arXiv.2206.03787>. 3, 4, 7
- 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. In Yoshua Bengio and Yann LeCun (eds.), *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL <http://arxiv.org/abs/1509.02971>. 1, 3

- Horia Mania, Aurelia Guy, and Benjamin Recht. Simple random search of static linear policies is competitive for reinforcement learning. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 1805–1814, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/7634ea65a4e6d9041cfd3f7de18e334a-Abstract.html>. 3
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin A. Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nat.*, 518(7540):529–533, 2015. doi: 10.1038/nature14236. URL <https://doi.org/10.1038/nature14236>. 2
- Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped DQN. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 4026–4034, 2016. URL <https://proceedings.neurips.cc/paper/2016/hash/8d8818c8e140c64c743113f563cf750f-Abstract.html>. 2, 3
- Fabio Pardo. Tonic: A deep reinforcement learning library for fast prototyping and benchmarking. *CoRR*, abs/2011.07537, 2020. URL <https://arxiv.org/abs/2011.07537>. 4
- Cristina Pinneri, Shambhuraj Sawant, Sebastian Blaes, Jan Achterhold, Joerg Stueckler, Michal Rolínek, and Georg Martius. Sample-efficient cross-entropy method for real-time planning. In Jens Kober, Fabio Ramos, and Claire J. Tomlin (eds.), *4th Conference on Robot Learning, CoRL 2020, 16-18 November 2020, Virtual Event / Cambridge, MA, USA*, volume 155 of *Proceedings of Machine Learning Research*, pp. 1049–1065. PMLR, 2020. URL <https://proceedings.mlr.press/v155/pinneri21a.html>. 3
- Matthias Plappert, Rein Houthoofd, Prafulla Dhariwal, Szymon Sidor, Richard Y. Chen, Xi Chen, Tamim Asfour, Pieter Abbeel, and Marcin Andrychowicz. Parameter space noise for exploration. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=ByBA12eAZ>. 3
- Antonin Raffin and Freek Stulp. Generalized state-dependent exploration for deep reinforcement learning in robotics. *CoRR*, abs/2005.05719, 2020. URL <https://arxiv.org/abs/2005.05719>. 3
- Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *J. Mach. Learn. Res.*, 22:268:1–268:8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>. 4
- Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. In Hadas Kress-Gazit, Siddhartha S. Srinivasa, Tom Howard, and Nikolay Atanasov (eds.), *Robotics: Science and Systems XIV, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA, June 26-30, 2018*, 2018. doi: 10.15607/RSS.2018.XIV.049. URL <http://www.roboticsproceedings.org/rss14/p49.html>. 4
- Daniel Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen. A tutorial on thompson sampling. *Found. Trends Mach. Learn.*, 11(1):1–96, 2018. doi: 10.1561/22000000070. URL <https://doi.org/10.1561/22000000070>. 2
- John Schulman, Sergey Levine, Pieter Abbeel, Michael I. Jordan, and Philipp Moritz. Trust region policy optimization. In Francis R. Bach and David M. Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 1889–1897. JMLR.org, 2015. URL <http://proceedings.mlr.press/v37/schulman15.html>. 2

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://arxiv.org/abs/1707.06347>. 2
- Pierre Schumacher, Daniel Häufle, Dieter Büchler, Syn Schmitt, and Georg Martius. Dep-rl: Embodied exploration for reinforcement learning in overactuated and musculoskeletal systems. arXiv: 2206.00484, 2022. URL <https://arxiv.org/abs/2206.00484>. 3
- David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin A. Riedmiller. Deterministic policy gradient algorithms. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, volume 32 of *JMLR Workshop and Conference Proceedings*, pp. 387–395. JMLR.org, 2014. URL <http://proceedings.mlr.press/v32/silver14.html>. 2
- David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Vedavyas Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy P. Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nat.*, 529(7587):484–489, 2016. doi: 10.1038/nature16961. URL <https://doi.org/10.1038/nature16961>. 2
- Haoran Tang, Rein Houthoofd, Davis Foote, Adam Stooke, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. #exploration: A study of count-based exploration for deep reinforcement learning. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 2753–2762, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/3a20f62a0af1aa152670bab3c602feed-Abstract.html>. 2
- Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy P. Lillicrap, and Martin A. Riedmiller. Deepmind control suite. *CoRR*, abs/1801.00690, 2018. URL <http://arxiv.org/abs/1801.00690>. 4
- Sebastian Thrun. Efficient exploration in reinforcement learning. Technical Report CMU-CS-92-102, Carnegie Mellon University, Pittsburgh, PA, January 1992. 2
- G. E. Uhlenbeck and L. S. Ornstein. On the theory of the brownian motion. *Phys. Rev.*, 36:823–841, Sep 1930. doi: 10.1103/PhysRev.36.823. URL <https://link.aps.org/doi/10.1103/PhysRev.36.823>. 3
- Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Mach. Learn.*, 8:229–256, 1992. doi: 10.1007/BF00992696. URL <https://doi.org/10.1007/BF00992696>. 2