

EVaDE : Event-Based Variational Thompson Sampling for Model-Based Reinforcement Learning

Siddharth Aravindan

SIDDHARTH.ARAVINDAN@GMAIL.COM

Dixant Mittal

DIXANT@COMP.NUS.EDU.SG

Wee Sun Lee

LEEWS@COMP.NUS.EDU.SG

Department of Computer Science, National University of Singapore

Editors: Vu Nguyen and Hsuan-Tien Lin

Abstract

Posterior Sampling for Reinforcement Learning (PSRL) is a well-known algorithm that augments model-based reinforcement learning (MBRL) algorithms with Thompson sampling. PSRL maintains posterior distributions of the environment transition dynamics and the reward function, which are intractable for tasks with high-dimensional state and action spaces. Recent works show that dropout, used in conjunction with neural networks, induces variational distributions that can approximate these posteriors. In this paper, we propose Event-based Variational Distributions for Exploration (EVaDE), which are variational distributions that are useful for MBRL, especially when the underlying domain is object-based. We leverage the general domain knowledge of object-based domains to design three types of event-based convolutional layers to direct exploration. These layers rely on Gaussian dropouts and are inserted between the layers of the deep neural network model to help facilitate variational Thompson sampling. We empirically show the effectiveness of EVaDE-equipped Simulated Policy Learning (EVaDE-SimPLe) on the 100K Atari game suite.

Keywords: Exploration; Thompson Sampling; Model-Based Reinforcement Learning

1. Introduction

Model-Based Reinforcement Learning (MBRL) has recently gained popularity for tasks that allow for a very limited number of interactions with the environment [43]. These algorithms use a model of the environment, that is learnt in addition to the policy, to improve sample efficiency in several ways; these include generating artificial training examples [43, 34], assisting with planning [22, 9, 39, 10] and guiding policy search [21, 8]. Additionally, it is easier to incorporate inductive biases derived from the domain knowledge of the task for learning the model, as the biases can be directly built into the transition and reward functions.

In this paper, we demonstrate how domain knowledge can be utilised for designing exploration strategies in MBRL. While model-free agents explore the space of policies and value functions, MBRL agents explore the space of transition dynamics and reward functions.

One method for exploring the space of transition dynamics and reward functions is Posterior Sampling for Reinforcement Learning (PSRL) [33, 25], which uses the Thompson sampling [35] method of sampling the posterior of the model to explore other plausible models. Maintaining the posterior is generally intractable and in practice, variational distributions are often used as an approximation to the posterior [3, 38, 42].



Figure 1: Rewards in Breakout, a popular Atari game. (a) shows an interaction between the ball and a brick which gives the agent a positive reward. (b) shows a state, where the paddle is unable to prevent the ball from going out of bounds. The lack of this interaction between the agent and the ball in this situation results in a penalty for the agent.

Traditionally, variational distributions are designed with two considerations in mind: inference and/or sampling should be efficient with the variational distribution, and the variational distribution should approximate the true posterior as accurately as possible. However, as the variational distribution may not fully represent the posterior, different approximations may be suitable for different purposes. In this paper, we propose to design the variational distribution to generate trajectories through parts of the state space that may potentially give high returns, for the purpose of exploration.

In MBRL, trajectories are generated in the state space by running policies that are optimised against the learned model. One way to generate useful exploratory trajectories is by perturbing the reward function in the model, so that a different part of the state space appears to contain high rewards, leading the policy to direct trajectories towards those states. Another method is to perturb the reward function, so that the parts of the state space traversed by the current policy appear sub-optimal, causing the policy to seek new trajectories.

We focus on problems where the underlying domain is object-based, meaning that the reward functions heavily depend on the locations of individual objects and their interactions, which we refer to as events. An example of such an object-based task is the popular Atari game Breakout (as shown in Figure 1). In this game, the agent receives rewards when the ball hits a brick and avoids losing a life by keeping the ball within bounds using the paddle, both of which are interactions between two objects. The rewards in this game are determined by the interactions between the ball and the bricks or the paddle.

For such domains, we introduce Event-based Variational Distributions for Exploration (EVaDE), a set of variational distributions that can help generate useful exploratory trajectories for deep convolutional neural network models. EVaDE comprises of three Gaussian dropout-based convolutional layers [32]: the noisy event interaction layer, the noisy event weighting layer, and the noisy event translation layer. The noisy event interaction layer is designed to provide perturbations to the reward function in states where multiple objects appear at the same location, randomly perturbing the value of interactions between objects. The noisy event weighting layer perturbs the output of a convolutional layer at a single location, assuming that the output of the convolutional filters captures events; this would upweight and downweight the reward associated with these events randomly. The noisy event translation layer perturbs trajectories that go through "narrow passages"; small

translations can randomly affect the returns from such trajectories, causing the policy to explore different trajectories.

These EVaDE layers can be used as standard convolutional layers and inserted between the layers of the environment network models. When included in deep convolutional networks, the noisy event interaction layers, the noisy event weighting layers, and the noisy event translation layers generate perturbations on possible object interactions, the importance of different events, and the positional importance of objects/events, respectively, through the dropout mechanism. This mechanism induces variational distributions over the model parameters [32, 12].

An interesting aspect of designing for exploration is that the variational distributions can be useful, even if they do not approximate the posterior well, as long as they assist in perturbing the policy out of local optima. After perturbing the policy, incorrect parts of the model will either be corrected by data or left unchanged if they are irrelevant to optimal behaviour.

Finally, we approximate PSRL by incorporating EVaDE layers into the reward models of Simulated Policy Learning (SimPLE) [43]. We conduct experiments to compare EVaDE-equipped SimPLE (EVaDE-SimPLE) with various popular baselines on the 100K Atari test suite. In the conducted experiments, all agents operate in the low data regime, where the number of interactions with the real environment is limited to 100K. EVaDE-SimPLE agents achieve a mean human-normalised score (HNS) of 0.682 in these games, which is 79% higher than the mean score of 0.381 achieved by a recent low data regime method, CURL [20], and 30% higher than the mean score of 0.525 achieved by vanilla SimPLE agents.

2. Background and Related Work

Posterior sampling approaches like Thompson Sampling [35] have been one of the more popular methods used to balance the exploration exploitation trade-off. Exact implementations of these algorithms have been shown to achieve near optimal regret bounds [2, 17]. These approaches, however, work by maintaining a posterior distribution over all possible environment models and/or action-value functions. This is generally intractable in practice. Approaches that work by maintaining an approximated posterior distribution [28, 4], or approaches that use bootstrap re-sampling to procure samples, [27, 24] have achieved success in recent times.

Variational inference procures samples from distributions that can be represented efficiently while also being easy to sample. These variational distributions are updated with observed data to approximate the true posterior as accurately as possible. Computationally cost effective methods such as dropouts have been known to induce variational distributions over the model parameters [32, 12]. Consequently, variational inference approaches that approximate the posterior distributions required by Thompson sampling have gained popularity [3, 38, 36, 40].

Model-based reinforcement learning improves sample complexity at the computational cost of maintaining and performing posterior updates to the learnt environment models. Neural networks have been successful in modelling relatively complex and diverse tasks such as Atari games [23, 14]. Over the past few years, variational inference has been used

to represent environment models, with the intention to capture environment stochasticity [15, 5, 13].

SimPLe [43] is one of the first algorithms to use MBRL to train agents to play video games from images. It is also perhaps the closest to EVaDE, as it not only employs an iterative algorithm to train its agent, but also uses an additional convolutional network assisted by an autoregressive LSTM based RNN to approximate the posterior of the hidden variables in the stochastic model. Thus, similar to existing methods [15, 5, 13], these variational distributions are used for the purpose of handling environment stochasticity rather than improving exploration. To the contrary, EVaDE-SimPLe is an approximation to PSRL, that uses a Gaussian dropout induced variational distribution over deterministic reward functions solely for the purpose of exploration. Unlike SimPLe, which uses the stochastic model to generate trajectories to train its agent, EVaDE-SimPLe agents optimize for a deterministic reward model sampled from the variational distribution and a learnt transition model. Moreover, with EVaDE, these variational distributions are carefully designed so as to explore different object interactions, importance of events and positional importance of objects/events, that we hypothesize are beneficial for learning good policies in object-based tasks.

The current state of the art scores in the Atari 100K benchmark is achieved by EfficientZero [41], which was developed concurrently with our work. Its success is a consequence of combining several improvements proposed previously in addition to integrating tree search with learning to improve the policy executed by the agent. We believe that the benefits of using the variational designs induced by the EVaDE layers proposed in this paper are complementary to such search based methods, as these layers could be used in their reward models to guide the policy search by generating useful exploratory trajectories, especially in object-based domains.

3. Event Based Variational Distributions

Event-based Variational Distributions for Exploration (EVaDE) consist of a set of variational distribution designs, each induced by a noisy convolutional layer. These convolutional layers can be inserted after any intermediate hidden layer in deep convolutional neural networks to help us construct approximate posteriors over the model parameters to produce samples from relevant parts of the model space. EVaDE convolutional layers use Gaussian multiplicative dropout to draw samples from the variational approximation of the posterior. Posterior sampling is done by multiplying each parameter, θ_{env}^i , of these EVaDE layers by a perturbation drawn from a Gaussian distribution, $\mathcal{N}(1, (\sigma_{env}^i)^2)$. These perturbations are sampled by leveraging the reparameterization trick [19, 30, 29, 11] using a noise sample from the standard Normal distribution, $\mathcal{N}(0, 1)$, as shown in Equation 1. The variance corresponding to each parameter, $(\sigma_{env}^i)^2$, is trained jointly with the model parameters θ_{env} .

$$\tilde{\theta}_{env}^i \leftarrow \theta_{env}^i (1 + \sigma_{env}^i \epsilon^i); \quad \epsilon^i \sim \mathcal{N}(0, 1) \quad (1)$$

When the number of agent-environment interactions is limited, the exploration strategy employed by the agent is critical. In object-based domains, rewards and penalties are often sparse and occur when objects interact. Hence, the agent needs to experience most of the events in order to learn a good environment model. Generating trajectories that contain

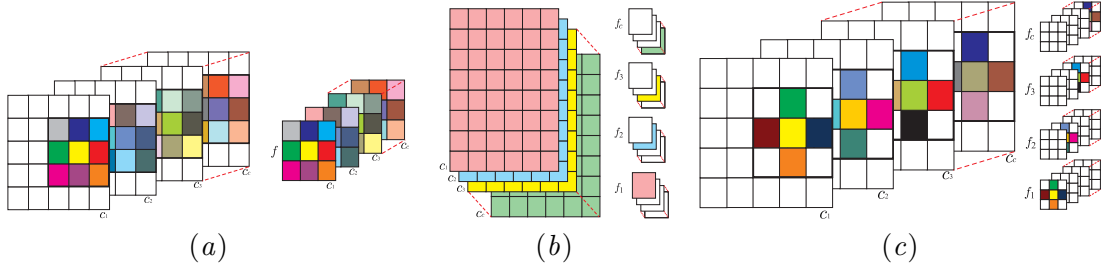


Figure 2: (a) This image shows one noisy event interaction filter acting on an input with c channels. Here f is an $m \times m$ noisy convolutional filter, which acts upon input patches at the same location across different channels, noisily altering the value of events captured at those locations. (b) This image shows how the filters of the noisy event weighting layer weight the input channels. The filters f_1, f_2, f_3 and f_c randomly upweight and downweight the events captured by the channels c_1, c_2, c_3 and c_c respectively. The white entries of the filter are entries that are set to zero, while the rest are trainable noisy model parameters. (c) The noisy event translation filter. The filters f_1, f_2, f_3 and f_c noisily translate events/objects captured by the channels c_1, c_2, c_3 and c_c respectively. The white entries of the filter are entries that are set to zero, while the rest are trainable noisy model parameters. Gaussian multiplicative dropout is applied to all the non-zero parameters of all EVaDE filters.

events is hence a reasonable exploration strategy. Additionally, a very common issue with training using a very few number of interactions is that the agent may often get stuck in a local optimum, prioritising an event, which is relatively important, but may not lead to an optimal solution. Generating potentially high return alternate trajectories that do not include that event is another reasonable exploration strategy.

With these exploration strategies in mind, we introduce three EVaDE layers, namely the noisy event interaction layer, the noisy event weighting layer and the noisy event translation layer. All the three layers are constructed with the hypothesis that the channels of the outputs of intermediate layers of deep convolutional neural networks either capture object positions, or events (interaction of multiple objects detected by multi-layer composition of the network).

3.1. Noisy Event Interaction Layer

The noisy event interaction layer is designed with the motivation of increasing the variety of events experienced by the agent. This layer consists of noisy convolutional filters, each having a dimension of $m \times m \times c$, where c is the number of input channels to the layer. Every filter parameter is multiplied by a Gaussian perturbation as shown in Equation 1. The filter dimension, m , is a hyperparameter that can be set so as to capture objects within a small $m \times m$ patch of an input channel. Assuming that the input channels capture the positions of different objects, a filter that combines the c input channels locally captures the local object interaction within the $m \times m$ patch. By perturbing the filter, different combinations of interactions can be captured; if the filter is used as part of the reward function, it will correspondingly reward different interactions depending on the perturbation. The policy optimized for different perturbed reward functions is likely to generate trajectories that

contain different events. Note that convolutional filters are equivariant, so the same filter will detect the event anywhere in the image and can result in trajectories that include the event at different positions in the image.

We describe the filter in more detail. Every output pixel of the filter, $y_{i,j}^k$, representing $(i,j)^{th}$ pixel of the k^{th} output channel, can be computed as shown in Equation 2. Here x is the input to the layer with c input channels, $P_{x_{i,j}}^l$ is the $m \times m$ patch (represented as a matrix) centred around $x_{i,j}^l$, the $(i,j)^{th}$ pixel of the l^{th} input channel, $\tilde{\theta}_k^l$ is the l^{th} channel of the k^{th} noisy convolutional filter, \odot the Hadamard product operator, and $\mathbb{1}_m$ is an m dimensional column vector whose every entry is 1.

$$y_{i,j}^k = \sum_{l=0}^c \mathbb{1}_m^T \left(\tilde{\theta}_k^l \odot P_{x_{i,j}}^l \right) \mathbb{1}_m \quad (2)$$

Figure 2a shows how this filter is applied over the input channels.

3.2. Noisy Event Weighting Layer

Overemphasis on certain events is possibly one of the main causes due to which agents converge to sub-optimal policies in object based tasks. Hence, it would be useful to easily be able to increase as well as decrease the importance of an event. For this layer, we assume that each input channel is already detecting an event and design a variational distribution over model parameters that directly up-weights or down-weights the events captured by different input channels.

This layer can be implemented with the help of c 1×1 noisy convolutional filters (each having a dimension of $1 \times 1 \times c$ as shown in Figure 2b), where c is the number of input channels. We denote the l^{th} element of the k^{th} filter in the layer as θ_k^l . To implement per channel noisy weighting, we set every θ_k^k as a trainable model parameter, which has a Gaussian dropout variance parameter associated with it to facilitate noisy weighting as shown in Equation 1. All other weights, i.e., θ_k^l when $l \neq k$ are set to 0. Thus each noisy event weighting layer has c trainable model parameters and c trainable Gaussian dropout parameters. A pictorial representation of how this layer acts on its input is presented in Figure 2b.

Every output $y_{i,j}^k$, corresponding to the $(i,j)^{th}$ pixel of the k^{th} output channel, can be computed using Equation 3, where $\tilde{\theta}_k^k$ is the noisy scaling factor for the k^{th} input channel.

$$y_{i,j}^k = \tilde{\theta}_k^k x_{i,j}^k \quad (3)$$

We expect that inducing such a variational distribution that up-weights or downweights events randomly helps the agents learn from different events that are randomly emphasised by different model samples drawn from the distribution. This may eventually help them in escaping local optima caused by overemphasis of certain events.

3.3. Noisy Event Translation Layer

In object based domains, an agent often has to perform a specific sequence of actions to successfully gain some rewards and may be penalized heavily for deviation from the

sequence. We refer to the specific sequence of actions as a "narrow passage". A small translation of the positions of the environment or other objects will often cause the agent to be unsuccessful. When random translations of obstacles, events or boundaries are performed within the reward function, the optimized policy may select a different trajectory, possibly allowing it to escape from a locally optimal trajectory. We thus design the noisy event translation layer to induce a variational distribution over such model posteriors that can sample a variety of translations of relevant objects.

The noisy soft-translation on an input with c channels, is performed with the help of c convolutional filters, each having a dimension of $m \times m \times c$. These filters compute a noisy weighted sum of the corresponding input pixel and the pixels near it to effect a *noisy* translation of the channel. Similar to the noisy event weighting layer, each filter of the noisy event translation layer acts on one input channel. To achieve this, every parameter except the parameters of the k^{th} channel of the k^{th} filter, θ_k^k (which has a dimension of $m \times m$), and their corresponding dropout variances, is set to 0, for all k . Moreover in the channel θ_k^k , only the middle column and row contain trainable parameters. Figure 2c shows a detailed pictorial representation of this structure of the filters. A random translation of up to n pixels of the input can be achieved by using a $(2n + 1) \times (2n + 1)$ noisy event translation layer.

Equation 4 shows how $y_{i,j}^k$, the $(i,j)^{th}$ output pixel of the k^{th} channel, is computed. Here, $P_{x_{i,j}^k}$ is a $m \times m$ patch centred at $(i,j)^{th}$ pixel of the k^{th} input channel, $\tilde{\theta}_k^k$ is the k^{th} channel of the k^{th} noisy convolutional filter, \odot the Hadamard product operator, and $\mathbb{1}_m$ is an m dimensional column vector where all the entries are 1.

$$y_{i,j}^k = \mathbb{1}_m^T \left(\tilde{\theta}_k^k \odot P_{x_{i,j}^k} \right) \mathbb{1}_m \quad (4)$$

3.4. Representational Capabilities of EVaDE networks

Ideally, adding EVaDE layers for exploration should not hinder the network to be unable to represent the true model, even if they don't accurately approximate the posterior. Theorem 1 below states that this is indeed the case.

Theorem 1 *Let \mathfrak{m} be any neural network. For any convolutional layer l , let $m_i(l) \times n_i(l) \times c_i(l)$ and $m_o(l) \times n_o(l) \times c_o(l)$ denote the dimensions of its input and output respectively. Then, any function that can be represented by \mathfrak{m} can also be represented by any network $\tilde{\mathfrak{m}} \in \tilde{\mathcal{N}}$, where $\tilde{\mathcal{N}}$ is the set of all neural networks that can be constructed by adding any combination of EVaDE layers to \mathfrak{m} , provided that, for every EVaDE layer \tilde{l} added, \tilde{l} uses a stride of 1, $m_i(\tilde{l}) \leq m_o(l)$, $n_i(\tilde{l}) \leq n_o(l)$ and $c_i(\tilde{l}) \leq c_o(l)$.*

Proof The proof follows from the fact that every EVaDE layer \tilde{l}_i that is added is capable of representing the identity function. A detailed proof is presented in the supplementary material.

If the added EVaDE layers induce distributions that poorly approximate the posterior, performance can indeed be poorer. But with enough data, the correct model should still be learnable since it is representable, as long as the optimization does not get trapped in a poor local optimum.

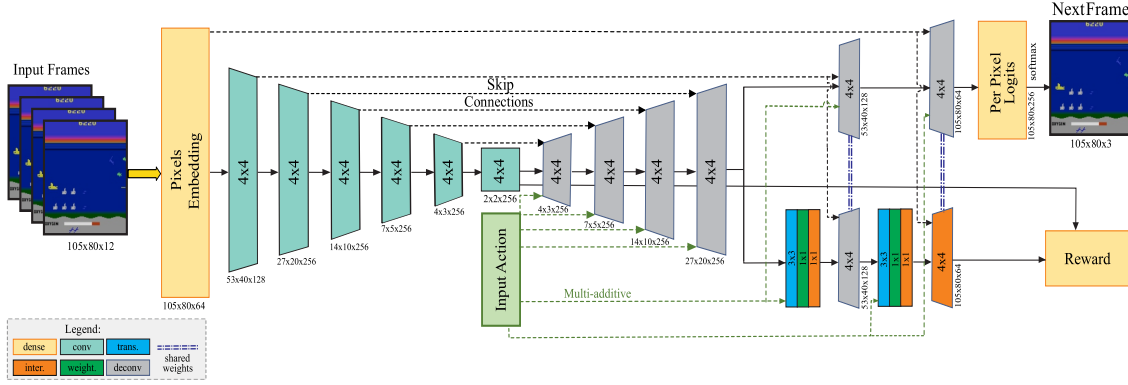


Figure 3: The network architecture of the environment model used to train EVaDE-SimPLE.

3.5. Approximate PSRL with EVaDE equipped Simulated Policy Learning

Simulated Policy Learning (SimPLE) [43] is an iterative model based reinforcement learning algorithm, wherein the environment model learnt is used to generate artificial episodes to train the agent policy. In every iteration, the SimPLE agent first interacts with the real environment using its current policy. After being trained on the set of all collected interactions, the models of the transition and reward functions are then used as a substitute to the real environment to train the policy of the agent to be followed by it in its next interactions with the real environment. PSRL [33, 25], which augments MBRL with Thompson sampling, has a very similar iterative structure as that of SimPLE. However, instead of maintaining a single environment model, PSRL maintains a posterior distribution over all possible environment models given the interactions experienced by the agent with the real environment. The agent then optimizes a policy for an environment model sampled from this posterior distribution. This policy is used in its real environment interactions of the next iteration. EVaDE equipped SimPLE approximates PSRL, by maintaining an approximate posterior distribution of the reward function with the help of the variational distributions induced by the three EVaDE layers.

Being an approximation of PSRL, an EVaDE-SimPLE agent has the same iterative training structure where it acts in the real environment using its latest policy to collect training interactions, learns a transition model and an approximate posterior over the reward model by jointly optimizing the environment model parameters, θ_{env} , and the Gaussian dropout parameters of the reward model, σ_{rew} , with the help of supervised learning. It then optimizes its policy with respect to an environment characterized by the learnt transition function and a reward model sample that is procured from the posterior with the help of Gaussian dropout as shown in Equation 1. This policy is then used by the agent to procure more training interactions in the next iteration.

4. Experiments

We conduct our experiments on the 100K Atari game test suite first introduced by [43]. This test suite consists of 26 Atari games where the number of agent-environment interactions

is limited to 100K. Due to its diverse range of easy and hard exploration games [6], this test-suite has become a popular test-bed for evaluating reinforcement algorithms.

4.1. Network Architecture

In our experiments we use the network architecture of the deterministic world model introduced by [43] to train the environment models of the SimPLe agents, but do not augment it with the convolutional inference network and the autoregressive LSTM unit. Readers are referred to [43] for more details.

The architecture of the environment model used by EVaDE-SimPLe agents is shown in Figure 3. This model is very similar to the one used by SimPLe agents, except that it has a combination of a 3×3 noisy event translation layer, a noisy event weighting layer and a 1×1 noisy event interaction layer inserted before the fifth and sixth de-convolutional layers. The final de-convolutional layer acts as a noisy event interaction filter when computing the reward, while it acts as a normal de-convolutional layer while predicting the next observation. Sharing weights between layers allows us to achieve this. We insert EVaDE layers in a way that it perturbs only the reward function and not the transition dynamics.

We reuse the network architecture of [43] to train the policies in both the SimPLe and EVaDE-SimPLe agents using Proximal Policy Optimization (PPO) [31]. All the hyperparameters used for training the policy network and environment are the same as the ones used in [43].

4.2. Experimental Details

The training regimen that we use to train all the agents is the same and is structured similarly to the setup used by [43]. The agents, initialized with a random policy and collect 6400 real environment interactions before starting the first training iteration. In every subsequent iteration, every agent trains its environment model with its collection of real world interactions, refines its policy by interacting with the environment model, if it is a vanilla-SimPLe agent, or a transition model and a reward model sampled from the approximate posterior, if it is an EVaDE-SimPLe agent, and then collects more interactions with this refined policy.

PSRL regret bounds scale linearly with the length of an episode experienced by the agent in every iteration [26]. Shorter horizons, however, run the risk of the agent not experiencing anything relevant before episode termination. To balance these factors, we set the total number of iterations to 30, instead of 15. We allocate an equal number of environment interactions to each iteration, resulting in 3200 agent-environment interactions per iteration. The total number of interactions that each SimPLe and EVaDE-SimPLe agent procures (about 102K) is similar to SimPLe agents trained in [43], which allocates double the number of interactions per iteration, but trains for only 15 iterations. To disambiguate between the different SimPLe agents referred to in this paper, we refer to the SimPLe agents trained in our paper and [43] as SimPLe(30) and SimPLe respectively from here on.

We try to keep the training schedule of EVaDE-SimPLe and SimPLe(30) similar to the training schedule of the deterministic model in [43] so as to keep the comparisons fair. We train the environment model for 45K steps in the first iteration and 15K steps in all subsequent iterations. In every iteration of simulated policy training, 16 parallel PPO agents

Table 1: Comparison of the performances achieved by popular baselines and five independent training runs of EVaDE-SimPLe and SimPLe(30) agents with 100K agent-environment interactions in the 26 game Atari 100K test suite.

Game	SimPLe	SimPLe(30)	CURL	OTRainbow	Eff. Rainbow	EVaDE-SimPLe
Mean HNS	0.443	0.525	0.381	0.264	0.285	0.682
Median HNS	0.144	0.151	0.175	0.204	0.161	0.267
Vs EVaDE (W/L)	7W,19L	3W,23L	9W,17L	6W,20L	9W,17L	-
Best Performing	5	2	4	1	3	11

collect $z * 1000$ batches of 50 environment interactions each, where $z = 1$ in all iterations except iterations 8, 12, 23 and 27 where $z = 2$ and in iteration 30, where $z = 3$. The policy is also trained when the agent interacts with the real environment. However, the effect of these interactions (numbering 102K) on the policy is minuscule when compared to the 28.8M transitions experienced by the agent when interacting with the learnt environment model. Additional experimental details as well as the anonymized code for our agents are provided in the supplementary, which is available at <https://tinyurl.com/3zb8nywx>.

4.3. Results

We report the mean and median Human Normalized Scores (HNS) achieved by SimPLe(30), EVaDE-SimPLe and popular baselines SimPLe [43], CURL [20], OverTrained Rainbow [18] and Data-Efficient Rainbow [37] in Table 1. For each baseline, we report the number of games in which it is the best performing, among all compared methods, as well as the number of games in which it scores more (or less) than EVaDE-SimPLe, which are counted as its wins (or losses).

EVaDE-SimPLe agents achieve the highest score in 11 of the 26 games in the test suite, outperforming every other method on at least 17 games. Moreover, the effectiveness of the noisy layers to improve exploration can be empirically verified as EVaDE-SimPLe manages to attain higher mean scores than SimPLe(30) in 23 of the 26 games, even though both methods follow the same training routine. EVaDE-SimPLe also scores a mean HNS of 0.682, which is 79% higher than the score of 0.381 achieved by a popular baseline, CURL, and 30% more than the mean HNS of 0.525 achieved by SimPLe(30). Additionally, EVaDE-SimPLe agents also surpass the human performances [7] in 5 games, namely Asterix, Boxing, CrazyClimber, Freeway and Krull.

We also compute the Inter-Quartile Means (IQM), ¹ a metric resilient to outlier games and runs, of SimPLe(30) and EVaDE-SimPLe agents. EVaDE-SimPLe agents achieve an IQM of 0.339, which is 68% higher than the IQM of 0.202 achieved by Simple(30) agents. This affirms that the improvements obtained due to the addition of the EVaDE layers are robust to outlier games and runs. In the supplementary material, we provide the scores achieved by all five independent runs of SimPLe(30) and EVaDE-SimPLe agents, which were used to compute the IQM.

Furthermore, a paired t-test on the mean HNS achieved by EVaDE-SimPLe and SimPLe(30) agents on each of the 26 games yields a single-tailed p-value of 3×10^{-3} confirming

1. IQM is well regarded in the reinforcement learning community, advocated by [1], which won the Outstanding Paper Award at NeurIPS 2021

Table 2: Scores (mean \pm 1 standard error) achieved by SimPLe agents when equipped with all three EVaDE filters individually and when equipped with all filters simultaneously. All scores are averaged over five independent training runs.

Game	SimPLe (30)	Inter. Layer	Weight. Layer	Trans. Layer	EVaDE-SimPLe
BankHeist	78.6 \pm 31.7	107.5 \pm 29.2	168.4 \pm 19.9	180.7 \pm 16.7	224.2 \pm 35.4
BattleZone	4544 \pm 803	6688 \pm 1617	7525 \pm 2164	7750 \pm 1355	11094 \pm 572
Breakout	18.9 \pm 1.7	19.8 \pm 3.6	22.4 \pm 5.5	19.5 \pm 1.5	24 \pm 3.4
CrazyClimber	43458 \pm 7709	59546 \pm 3164	64191\pm5196	59006 \pm 3282	60716 \pm 4082
DemonAttack	120.7 \pm 18.2	136.3 \pm 24.4	132 \pm 14.7	133.7 \pm 26	141.8 \pm 12.5
Frostbite	260.3 \pm 2.5	254.6 \pm 5.5	254.4 \pm 3.6	263.2 \pm 2.1	274.2 \pm 11
JamesBond	245.6\pm11.2	202.2 \pm 65.2	182.5 \pm 56.2	160.3 \pm 68.4	235.6 \pm 50.2
Kangaroo	576 \pm 330	2201\pm993	837.5 \pm 345	1297 \pm 321	1186.5 \pm 168
Krull	4532 \pm 883	3117 \pm 781	4818 \pm 440	5185 \pm 991	5335\pm455
Qbert	2583 \pm 746	1953 \pm 674	932 \pm 148	3333 \pm 575	2764 \pm 783
RoadRunner	2385 \pm 888	7178 \pm 1227	4853 \pm 1322	6070 \pm 1834	7799 \pm 1296
Seaquest	321.6 \pm 52	644.4 \pm 91.6	608.5 \pm 144.9	644.2 \pm 56.9	617.5 \pm 118.1
HNS	0.52	0.56	0.65	0.69	0.77
IQM	0.22	0.29	0.26	0.29	0.4
Vs SimPLe(30) (W/L)	-	8W,4L	9W,3L	11W,1L	11W,1L

that the performance improvements over SimPLe(30) of EVaDE-SimPLe agents are statistically significant as an algorithm when applied to multiple games.

4.4. Ablation Studies

We also conduct ablation studies by equipping SimPLe(30) with just one of the three EVaDE layers to ascertain whether all of them aid in exploration. We do this by just removing the other two layers from the EVaDE environment network model (see Figure 3). Note that reward models that do not equip the noisy event interaction filter, also do not apply the Gaussian multiplicative dropout to the sixth de-convolutional layer.

We use a randomly selected suite of 12 Atari games in our ablation study. The games were chosen by arranging the 26 games of the suite in the alphabetical order, and then using the numpy [16] random function to sample 12 numbers from 0 to 25 without replacement. The corresponding games were then picked. Coincidentally, the chosen test suite contains easy exploration games such as Kangaroo, RoadRunner and Seaquest as well as BankHeist, Frostbite and Qbert, which are hard exploration games [6].

The mean scores, HNS and IQM achieved when SimPLe(30) is equipped with only one type of noisy convolutional layer and those of SimPLe(30) and EVaDE-SimPLe are presented in Table 2.

We present the learning curves of the trained EVaDE and SimPLe(30) agents in Figure 4. We omit the error bars here for clarity. Looking at the learning curves presented, it can possibly be said that an increase in scores of SimPLe(30) equipped with one of the EVaDE layers at a particular iteration would mean an increase in scores of EVaDE-SimPLe, albeit in later iterations. This pattern can clearly be seen in the games of BankHeist, Frostbite, Kangaroo, Krull and Qbert. This delay in learning could possibly be attributed to the agent draining its interaction budget by exploring areas suggested by one of the layers that are ineffective for that particular game. However, we hypothesise that since all the layers provide different types of exploration, their combination is more often helpful than wasteful.

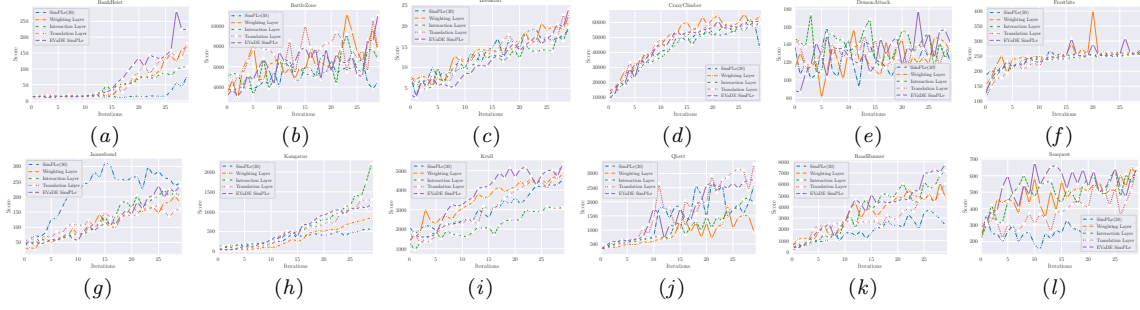


Figure 4: Learning curves of EVaDE-SimPLE agents, SimPLE(30) agents and agents which only add one of the EVaDE layers trained on the 12 game subset of the Atari 100K test suite.

382 This is validated by the fact that EVaDE-SimPLE achieves higher mean HNS and IQM than
 383 any other agent in this study.

384 We also look at the policies learnt by these agents in RoadRunner, where every EVaDE
 385 layer seems to improve the performance of the SimPLE agent even when added individually.
 386 We find that the vanilla-SimPLE RoadRunner agent learns a suboptimal policy where it
 387 frequently either collides with a moving obstacle or gets caught by the coyote. When the
 388 agent adds noisy interaction layers into its reward models, the policy learns to trick the
 389 coyote into colliding with the obstacle, an interaction beneficial to the roadrunner. On the
 390 other hand, when the reward models use only the noisy translation layers, the agent seems
 391 to learn a less *risky* policy, as it aggressively keeps away from both the obstacle as well as the
 392 coyote. The agent that adds only the noisy weighting layers seems to prioritize collecting
 393 points. This allows the coyote to get near the roadrunner, which could be undesirable.
 394 The policy learnt by the EVaDE-SimPLE agent combines the properties of the agents that
 395 add only the interaction and translation layers, as it tricks the coyote to colliding with
 396 the obstacle, while keeping a safe distance from both. The behaviour of the agent in this
 397 game provides some evidence that EVaDE layers can allow us to design different types of
 398 exploration based on our prior knowledge. We include videos of one episode run of each
 399 agent in the supplementary.

400 From Table 2, it can be seen that individually, each filter achieves a higher HNS than
 401 SimPLE(30), thus indicating that, on average, all the filters help in aiding exploration.
 402 Moreover, we see that with the exception of the noisy event interaction layer, the increase in
 403 HNS when the other two EVaDE layers are added individually is considerable. Additionally,
 404 all agents that add any of the noisy convolutional layers achieve higher IQM scores than
 405 baseline SimPLE agents. Furthermore, we hypothesise that all the layers provide different
 406 types of exploration since their combination is more often helpful than wasteful. This is
 407 validated by the fact that EVaDE-SimPLE achieves a higher mean HNS and IQM than any
 408 other agent in this study.

409 While having more parameters enlarges the class of functions representable by the model
 410 used by EVaDE agents, we emphasize that it is the design of the layers and not the higher
 411 number of parameters that is the reason for the performance gains. The agents that include
 412 only the noisy translation layers outperform SimPLE(30) on all metrics but add just 4K

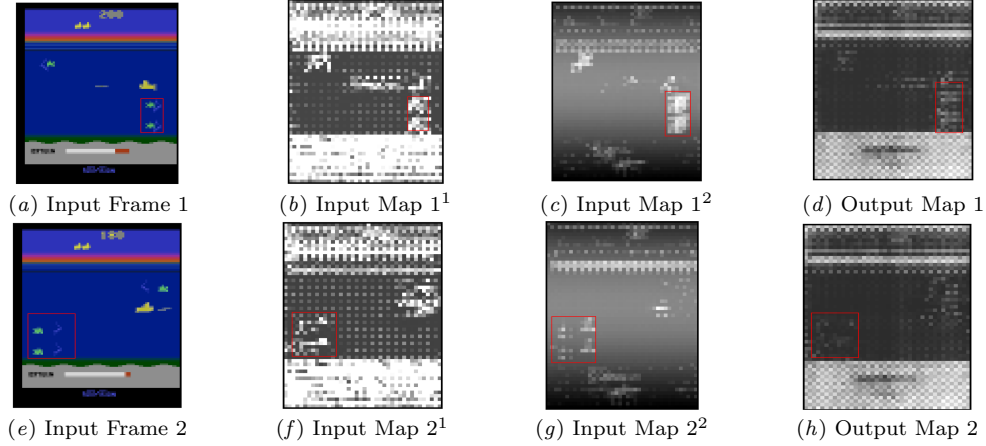


Figure 5: This figure shows the output map that captures interactions between two input maps when passed through the noisy event interaction layer.

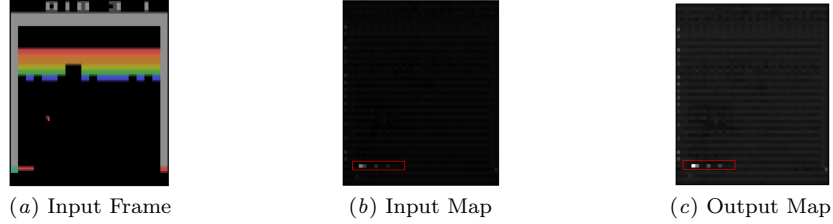


Figure 6: This figure shows an output map (channel) that up-weights the corresponding input feature map when passed through the noisy event weighting layer.

parameters to the reward model which contains about 10M parameters. Similarly, EVaDE-SimPLe adds only 0.43M parameters for its improvements.

4.5. Visualizations of the EVaDE layers

We also present some visualizations that help us understand the functionality of the EVaDE layers. All these visualizations were obtained after the completion of training, with the learned weights and variances of the final trained model.

In Figure 5 we show illustrations of output map of a noisy event interaction layer detecting interactions between the right facing green-coloured enemy ships and the right facing blue-coloured divers given different input images from the game of Seaquest. We also show two input feature maps, which seem to capture the positions of these objects at the same locations. We observe that the pixels in the output feature map in Figure 5d are brighter at the locations where the two objects are close to each other, whereas in the same feature map these pixels are dimmer when the two objects are separated by some distance (Figure 5h).

In Figures 6 and 7, we show two feature maps before and after passing them through a noisy event weighting layer. The input images for these visualizations were taken from the game of Breakout. The input feature maps to the noisy event weighting layer seem

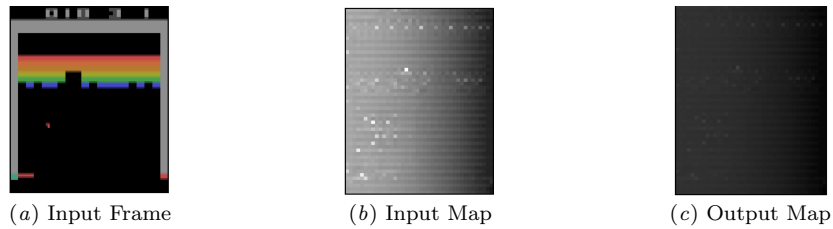


Figure 7: This figure shows an output map (channel) that down-weights the corresponding input feature map when passed through the noisy event weighting layer.

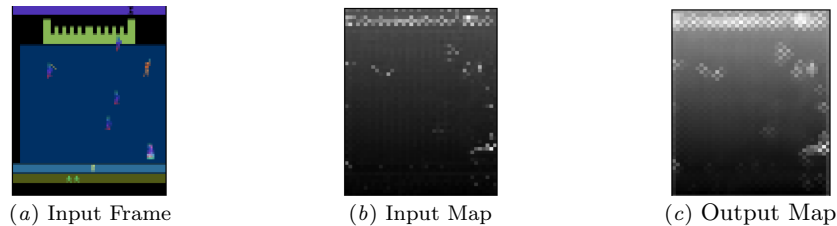


Figure 8: This figure shows the function of the noisy translation layer. The output map translates the input pixels to its top, bottom, left and right to different degrees

to capture the objects from the input image. The output feature map in Figure 6 is an upweighted version of its input, as the pixels seem to be brighter. On the other hand, the output feature map in Figure 7 seems to down-weight its input feature map, as the pixels seem a lot dimmer. The weighting factors for the input-output pairs shown in Figures 6 and 7 are 1.93 and 0.57 respectively. In Figure 8, we show the input and output feature maps of a game state from Krull, before and after passing it through a noisy event translation layer. The input feature map seems to capture different objects from the input image. The translation effect in output feature map can be seen clearly as every light pixel in the input seems to have lightened up the pixels to its top, bottom, left and right to different degrees.

5. Conclusion

In this paper, we introduce Event-based Variational Distributions for Exploration (EVaDE), a set of variational distributions over reward functions. EVaDE consists of three noisy convolutional layers: the noisy event interaction layer, the noisy event weighting layer, and the noisy event translation layer. These layers are designed to generate trajectories through parts of the state space that may potentially give high rewards, especially in object-based domains. We can insert these layers between the layers of the reward network models, inducing variational distributions over the model parameters. The dropout mechanism then generates perturbations on object interactions, importance of events, and positional importance of objects/events. We draw samples from these variational distributions and generate simulations to train the policy of a Simulated Policy Learning (SimPLe) agent.

We conduct experiments on the Atari 100K test suite, a test suite comprising 26 games where the agents are only allowed 100K interactions with the real environment. EVaDE-SimPLe agents outperform vanilla-SimPLe(30) on this test suite by achieving a mean HNS

of 0.682, which is 30% more than the mean HNS of 0.525 achieved by vanilla-SimPLe(30) agents. EVaDE-SimPLe agents also surpass human performances in five games of this suite. Additionally, these agents also outperform SimPLe(30) on median HNS as well as IQM scores, which is a metric that is resilient to outlier runs and games. Furthermore, a paired-t test also confirms the statistical significance of the improvements in HNS on the test suite ($p = 3 \times 10^{-3}$). We also find, through an ablation study, that each noisy convolutional layer, when added individually to SimPLe results in a higher mean HNS and IQM. Additionally, the three noisy layers complement each other well, as EVaDE-SimPLe agents, which include all three EVaDE layers, achieve higher mean HNS and IQM than agents which add only one noisy convolutional layer. Finally, we also present visualizations that help us understand the functionality of the EVaDE layers.

References

- [1] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare. Deep Reinforcement Learning at the Edge of the Statistical Precipice. *Advances in Neural Information Processing Systems*, 34, 2021.
- [2] Shipra Agrawal and Randy Jia. Optimistic Posterior Sampling for Reinforcement Learning: Worst-Case Regret Bounds. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/3621f1454cacf995530ea53652ddf8fb-Paper.pdf>.
- [3] Siddharth Aravindan and Wee Sun Lee. State-Aware Variational Thompson Sampling for Deep Q-Networks. In *20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2021)*, pages 124–132, 2021.
- [4] Kamyar Azizzadenesheli, Emma Brunskill, and Animashree Anandkumar. Efficient Exploration through Bayesian Deep Q-Networks. In *2018 Information Theory and Applications Workshop (ITA)*, pages 1–9. IEEE, 2018.
- [5] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H. Campbell, and Sergey Levine. Stochastic Variational Video Prediction. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rk49Mg-CW>.
- [6] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying Count-Based Exploration and Intrinsic Motivation. In *Advances in Neural Information Processing Systems*, pages 1471–1479, 2016.
- [7] Marc Brittain, Joshua R. Bertram, Xuxi Yang, and Peng Wei. Prioritized Sequence Experience Replay. *CoRR*, abs/1905.12726, 2019. URL <http://arxiv.org/abs/1905.12726>.
- [8] Yevgen Chebotar, Mrinal Kalakrishnan, Ali Yahya, Adrian Li, Stefan Schaal, and Sergey Levine. Path Integral Guided Policy Search. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 3381–3388. IEEE, 2017.
- [9] Rémi Coulom. Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search. In *International conference on computers and games*, pages 72–83. Springer, 2006.
- [10] Sebastian Curi, Felix Berkenkamp, and Andreas Krause. Efficient Model-Based Reinforcement Learning through Optimistic Policy Search and Planning. *Advances in Neural Information Processing Systems*, 33, 2020.
- [11] Meire Fortunato, Mohammad Azar, Bilal Piot, Jacob Menick, Matteo Hessel, Ian Osband, Alex Graves, Volodymyr Mnih, Remi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg. Noisy Networks for Exploration. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rywHCPkAW>.
- [12] Yarín Gal and Zoubin Ghahramani. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [13] Karol Gregor, George Papamakarios, Frederic Besse, Lars Buesing, and Theophane Weber. Temporal Difference Variational Auto-Encoder. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=S1x4ghC9tQ>.
- [14] David Ha and Jürgen Schmidhuber. Recurrent World Models Facilitate Policy Evolution. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 2455–2467, 2018.
- [15] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning Latent Dynamics for Planning from Pixels. In *International Conference on Machine Learning*, pages 2555–2565. PMLR, 2019.
- [16] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020. doi: 10.1038/s41586-020-2649-2. URL <https://doi.org/10.1038/s41586-020-2649-2>.
- [17] Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-Optimal Regret Bounds for Reinforcement Learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.

- [18] Kacper Kielak. Importance of Using Appropriate Baselines for Evaluation of Data-Efficiency in Deep Reinforcement Learning for Atari. *CoRR*, abs/2003.10181, 2020. URL <https://arxiv.org/abs/2003.10181>.
- [19] Diederik P Kingma, Tim Salimans, and Max Welling. Variational Dropout and the Local Reparameterization Trick. In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*, pages 2575–2583, 2015.
- [20] Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: Contrastive Unsupervised Representations for Reinforcement Learning. In *International Conference on Machine Learning*, pages 5639–5650. PMLR, 2020.
- [21] Sergey Levine and Vladlen Koltun. Guided Policy Search. In *International conference on machine learning*, pages 1–9. PMLR, 2013.
- [22] Anusha Nagabandi, Gregory Kahn, Ronald S. Fearing, and Sergey Levine. Neural Network Dynamics for Model-Based Deep Reinforcement Learning with Model-Free Fine-Tuning. *CoRR*, abs/1708.02596, 2017. URL <http://arxiv.org/abs/1708.02596>.
- [23] Junhyuk Oh, Xiaoxiao Guo, Honglak Lee, Richard Lewis, and Satinder Singh. Action-Conditional Video Prediction Using Deep Networks in Atari Games. In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*, pages 2863–2871, 2015.
- [24] Ian Osband and Benjamin Van Roy. Bootstrapped Thompson Sampling and Deep Exploration. *arXiv preprint arXiv:1507.00300*, 2015.
- [25] Ian Osband and Benjamin Van Roy. Why is Posterior Sampling Better than Optimism for Reinforcement Learning? In *International Conference on Machine Learning*, pages 2701–2710, 2017.
- [26] Ian Osband, Daniel Russo, and Benjamin Van Roy. (More) Efficient Reinforcement Learning via Posterior Sampling. In *Advances in Neural Information Processing Systems*, pages 3003–3011, 2013.
- [27] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep Exploration via Bootstrapped DQN. In *Advances in Neural Information Processing Systems*, pages 4026–4034, 2016.
- [28] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and Exploration via Randomized Value Functions. In *Proceedings of the 33rd International Conference on Machine Learning-Volume 48*, pages 2377–2386. JMLR. org, 2016.
- [29] 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 *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=ByBA12eA2>.
- [30] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution Strategies as a Scalable Alternative to Reinforcement Learning. *arXiv preprint arXiv:1703.03864*, 2017.
- [31] 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>.
- [32] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014. URL <http://jmlr.org/papers/v15/srivastava14a.html>.
- [33] Malcolm Strens. A Bayesian Framework for Reinforcement Learning. In *International Conference on Machine Learning*, volume 2000, pages 943–950, 2000.
- [34] Richard S Sutton. Dyna, an Integrated Architecture for Learning, Planning, and Reacting. *ACM Sigart Bulletin*, 2(4):160–163, 1991.
- [35] William R Thompson. On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples. *Biometrika*, 25(3/4):285–294, 1933.
- [36] Iñigo Urteaga and Chris Wiggins. Variational Inference for the Multi-Armed Contextual Bandit. In *International Conference on Artificial Intelligence and Statistics*, pages 698–706. PMLR, 2018.
- [37] Hado van Hasselt, Matteo Hessel, and John Aslanides. When to Use Parametric Models in Reinforcement Learning? In *NeurIPS*, pages 14322–14333, 2019.
- [38] Zhendong Wang and Mingyuan Zhou. Thompson Sampling via Local Uncertainty. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 10115–10125. PMLR, 13–18 Jul 2020. URL <http://proceedings.mlr.press/v119/wang20ab.html>.
- [39] Grady Williams, Andrew Aldrich, and Evangelos Theodorou. Model Predictive Path Integral control Using Covariance Variable Importance Sampling. *arXiv preprint arXiv:1509.01149*, 2015.
- [40] Sirui Xie, Junnang Huang, Lanxin Lei, Chunxiao Liu, Zheng Ma, Wei Zhang, and Liang Lin. NADPEX: An On-Policy Temporally Consistent Exploration Method for Deep Reinforcement Learning. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rkxciiC9tm>.
- [41] Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games with limited data. *Advances in Neural Information Processing Systems*, 34, 2021.
- [42] Ruiyi Zhang, Zheng Wen, Changyou Chen, Chen Fang, Tong Yu, and Lawrence Carin. Scalable Thompson Sampling via Optimal Transport. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 87–96. PMLR, 2019.
- [43] Kaiser Lukasz, Babaeizadeh Mohammad, Milos Piotr, Osiński Błażej, Campbell Roy, H, Czechowski Konrad, Erhan Dumitru, Finn Chelsea, Kozakowski Piotr, Levine Sergey, Mohiuddin Afroz, Sepassi Ryan, Tucker George, and Michalewski Henryk. Model Based Reinforcement Learning for Atari. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=S1xCPJHtDB>.