# Benchmarking Bonus-Based Exploration Methods on the Arcade Learning Environment

Adrien Ali Taïga 12 William Fedus 12 Marlos C. Machado 2 Aaron Courville 13 Marc G. Bellemare 23

### **Abstract**

This paper provides an empirical evaluation of recently developed exploration algorithms within the Arcade Learning Environment (ALE). We study the use of different reward bonuses that incentives exploration in reinforcement learning. We do so by fixing the learning algorithm used and focusing only on the impact of the different exploration bonuses in the agent's performance. We use Rainbow, the state-of-the-art algorithm for value-based agents, and focus on some of the bonuses proposed in the last few years. We consider the impact these algorithms have on performance within the popular game MON-TEZUMA'S REVENGE which has gathered a lot of interest from the exploration community, across the the set of seven games identified by Bellemare et al. (2016) as challenging for exploration, and easier games where exploration is not an issue. We find that, in our setting, recently developed bonuses do not provide significantly improved performance on MONTEZUMA'S REVENGE or hard exploration games. We also find that existing bonus-based methods may negatively impact performance on games in which exploration is not an issue and may even perform worse than  $\epsilon$ -greedy exploration.

## 1. Introduction

Despite recent entreaties for better practices to yield reproducible research (Henderson et al., 2018; Machado et al., 2018b), the literature on exploration in reinforcement learning still lacks a *systematic* comparison between existing methods. In the context of the the Arcade Learning Environment

Proceedings of the  $2^{nd}$  Exploration in Reinforcement Learning Workshop at the  $36^{th}$  International Conference on Machine Learning, Long Beach, California, PMLR 97, 2019. Copyright 2019 by the author(s).

ronment (ALE; Bellemare et al., 2013), we observe comparisons of agents trained under different regimes: with or without reset, using varying number of training frames, with and without sticky actions (Machado et al., 2018b) and evaluating only on a small subset of the available games. This makes it nearly impossible to assess the field's progress towards efficient exploration.

Our goal here is to revisit some of the recent bonus based exploration methods using a common evaluation regime. We do so by

- Comparing all methods on the same set of Atari 2600 games;
- Applying these bonuses on the same value-based agent architecture, Rainbow (Hessel et al., 2018);
- Fixing the number of samples each algorithm uses during training to 200 million game frames.

As an additional point of comparison, we also evaluate in the same setting NoisyNets (Fortunato et al., 2018), part of the original Rainbow algorithm and  $\epsilon$ -greedy exploration.

We study three questions relevant to exploration in the ALE:

- How well do different methods perform on MON-TEZUMA'S REVENGE?
- Do these methods generalize to Bellemare et al.'s set of "hard exploration games", when their hyperparameters are tuned only on MONTEZUMA'S REVENGE?
- Do they generalize to other Atari 2600 games?

We find that, despite frequent claims of state-of-the-art results in MONTEZUMA'S REVENGE, when the learning algorithm and sample complexity are kept fixed across the different methods, little to no performance gain can be observed over older methods. Furthermore, our results suggest that performance on MONTEZUMA'S REVENGE is not indicative of performance on other hard exploration games. In fact, on 5 out of 6 hard exploration games performance of considered bonus-based methods is on-par with an  $\epsilon$ -greedy algorithm, and significantly lower than human-level performance. Finally, we find that, while exploration bonuses improve performance on hard exploration games, they typically hurt performance on the easier Atari 2600 games.

<sup>&</sup>lt;sup>1</sup>MILA, Université de Montréal <sup>2</sup>Google Research, Brain Team <sup>3</sup>CIFAR Fellow. Correspondence to: Adrien Ali Taïga <adrien.ali.taiga@umontreal.ca>.

Taken together, our results suggests that more research is needed to make bonus-based exploration robust and reliable, and serve as a reminder of the pitfalls of developing and evaluating methods primarily on a single domain.

#### 2. Related Work

Exploration methods may encourage agents toward unexplored parts of the state space in different ways. Count-based methods generalize previous work that was limited to tabular methods (Strehl & Littman, 2008) to estimate counts in high dimension (Bellemare et al., 2016; Ostrovski et al., 2017; Tang et al., 2017; Choshen et al., 2018; Machado et al., 2018a). Prediction error has also been used as a novelty signal to compute an exploration bonus (Stadie et al., 2015; Pathak et al., 2017; Burda et al., 2019). Another class of exploration methods apply the Thompson sampling heuristic to reinforcement learning (Osband et al., 2016; O'Donoghue et al., 2017; Touati et al., 2018).

Burda et al. (2018) benchmarks various exploration methods based on prediction error within a set of simulated environment including some Atari 2600 games. However their study differs from ours as their setting ignore the environment reward and instead learns exclusively from the intrinsic reward signal.

# 3. Exploration methods

We focus on bonus-based methods that encourage exploration through a reward signal. At each time-step the agent is trained with the reward  $r_t = e_t + \beta \cdot i_t$  where  $e_t$  is the extrinsic reward provided by the environment,  $i_t$  the intrinsic reward computed by agent and  $\beta>0$  a scaling parameter. We now summarize different ways to compute the intrinsic reward i.

#### 3.1. Pseudo-counts

Pseudo-counts (Bellemare et al., 2016; Ostrovski et al., 2017) were proposed as way to estimate counts in high dimension states spaces using a density model. The agent is then encouraged to visit states with a low visit count. Let  $\rho$  be a density model over the state space and  $\rho_t(s)$  the density assigned to s after being trained on a sequence of states  $s_1, ..., s_t$ . We will write  $\rho'_t(s)$  the density assigned to s if  $\rho$  were to be updated with s. We require  $\rho$  to be learning positive (i.e  $\rho'_t(s) \geq \rho_t(s)$ ) and define the prediction gain as  $\operatorname{PG}_t(s) = \log \rho'_t(s) - \log \rho_t(s)$ . The pseudo-count  $\hat{N}_t(s_t) \approx \left(e^{\operatorname{PG}_t(s_t)} - 1\right)^{-1}$  can then be used to compute the intrinsic reward

$$i^{\text{PSC}}(s_t) := (\hat{N}_t(s_t))^{-1/2}.$$
 (1)

CTS (Bellemare et al., 2014) and PixelCNN (Van den Oord et al., 2016) have been both used as density models. We will disambiguate these agent by the name of their density model.

#### 3.2. Intrinsic Curiosity Module

Intrinsic Curiosity Module (ICM, Pathak et al., 2017) promotes exploration via curiosity. Pathak et al. formulates curiosity as the agent's ability to predict the consequence of its own actions in a learned feature space. ICM includes a learned embedding, a forward and an inverse model. The embedding is trained through the inverse model, which in turn, has to predict the agent's action between two states  $s_t$  and  $s_{t+1}$  using their embedding  $\phi(s_t)$  and  $\phi(s_{t+1})$ . Given a transition  $(s_t, a_t, s_{t+1})$  the intrinsic reward is then given by the error of the forward model in the embedding space between  $\phi(s_{t+1})$  and the predicted estimate  $\hat{\phi}(s_{t+1})$ 

$$i^{\text{ICM}}(s_t) = \|\hat{\phi}(s_{t+1}) - \phi(s_{t+1})\|_2^2. \tag{2}$$

#### 3.3. Random Network Distillation

Random Network Distillation (RND, Burda et al., 2019) derives a bonus from the prediction error of a random network. The intuition is that the prediction error will be low on states that are similar to those previously visited and high on newly visited states?. A neural network  $\hat{f}$  with parameters  $\theta$  is trained to predict the output of a fixed randomly initialized neural network f:

$$i^{\text{RND}}(s_t) = \|\hat{f}(s_t; \theta) - f(s_t)\|_2^2$$
 (3)

#### 3.4. NoisyNets

Though is does not generate an exploration bonus, we also evaluate NoisyNets (Fortunato et al., 2018) as it was chosen as the exploration strategy of the original Rainbow implementation (Hessel et al., 2018). NoisyNets add noise in parameter space and propose to replace standard fully-connected layers y = Ws + b by a noisy version that combines a deterministic and a noisy stream:

$$y = (W + W_{noisy} \odot \epsilon^{W})s + (b + b_{noisy} \odot \epsilon^{b}), \quad (4)$$

where  $\epsilon^W$  and  $\epsilon^b$  are random variable and  $\odot$  denotes elementwise multiplication.

#### 4. Evaluation protocol

We evaluate two key properties of exploration methods in the ALE:

- Sample efficiency: obtaining a decent policy quickly.
- **Robustness**: performing well across different games of the ALE with the same set of hyperparameters.

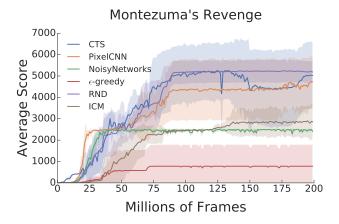


Figure 1. A comparison of different exploration methods on MON-TEZUMA'S REVENGE.

Sample efficiency is a key objective for exploration methods, yet, because published agents are often trained under different regimes it is often not possible to directly compare their performance. They often employ different reinforcement learning algorithms, varying quantity of training frames or inconsistent hyperparameter tuning. As a remedy, we fix our training protocol and train bonus-based methods with a common agent, the Rainbow implementation provided by the Dopamine framework (Castro et al., 2018) which includes Rainbow's three most important component: n-step updates (Mnih et al., 2016), prioritized experience replay (Schaul et al., 2015) and distributional reinforcement learning (Bellemare et al., 2017). To avoid introducing bias in favor of a particular method we also kept the original hyperparameters fixed. Our agents are trained for 200 million frames following Mnih et al.'s original setting. Nevertheless we also acknowledge the emerging trend of training agents an order of magnitude longer in order to produce a highscoring policy, irrespective of the sample cost (Espeholt et al., 2018; Burda et al., 2019; Kapturowski et al., 2019).

The ALE was designed with the assumption that few games would be used for training and the remaining ones for evaluation. Nonetheless it has become common to do hyperparameter tuning on MONTEZUMA'S REVENGE and only evaluate on other ALE's hard exploration games with sparse rewards: FREEWAY, GRAVITAR, SOLARIS, VENTURE, PRIVATE EYE. While this may be due to limited computational resources doing so however may come to a price on easier exploration problems as we will see later on. For this reason we chose to also evaluate performance on the original Atari training set<sup>1</sup>. Except for FREEWAY these are all considered easy exploration problems (Bellemare et al., 2016).

## 5. Empirical Results

In this section we present an experimental study of exploration methods using the protocol described previously.

#### 5.1. MONTEZUMA'S REVENGE

We begin by establishing a benchmark of bonus-based methods on MONTEZUMA'S REVENGE when each method is tuned on the same game. Details regarding implementation and hyperparameter tuning may be found in Appendix B. Figure 1 shows training curves (averaged over 5 random seeds) for Rainbow augmented with different exploration bonuses.

As anticipated,  $\epsilon$ -greedy exploration performs poorly. Other strategies are able to consistently reach 2500 points and often make further progress. We find pseudo-count with CTS matches recent bonuses and reaches a score of 5000 points within 200 millions frames. Of note, the performance we report for each method improves on the performance originally reported by the authors. This is mostly due to the fact these methods are based on weaker Deep Q-Network (Mnih et al., 2015) variants. This emphasize again the importance of the agent architecture to evaluate exploration methods.

Regarding RND performance, we note that our implementation only uses Eq. (3) bonus and does not appeal to other techniques presented in the same paper that were shown to be critical to the final performance of the algorithm. Though, we might expect that such techniques would also benefit other bonus based methods and leave it to future work.

#### 5.2. Hard exploration games

We now turn our attention to the set of games categorized as hard exploration games by Bellemare et al. (2016) that is often used as an evaluation set for exploration methods. Training curves for few games are shown in Figure 2, the remaining ones are in Appendix A. We find that performance of each method on MONTEZUMA'S REVENGE does not correlate with performance on other hard exploration problems and the gap between different methods is not as large as it was on MONTEZUMA'S REVENGE. Surprisingly, in our setting, there is also no visible difference between  $\epsilon$ -greedy exploration and more sophisticated exploration strategies.  $\epsilon$ -greedy exploration remains competitive and even outperforms other methods by a significant margin on GRAVITAR. Similar results have been reported previously (Machado et al., 2018a; Burda et al., 2019). These games were originally classified as hard exploration problems because DQN with  $\epsilon$ -greedy exploration was unable to reach a high scoring policy; however, these conclusions must be revisited with stronger base agents. Progress in these games may be due to better credit assignment methods and not to the underlying exploration bonus.

<sup>&</sup>lt;sup>1</sup>Freeway, Asterix, Beam Rider, Seaquest, Space Invaders

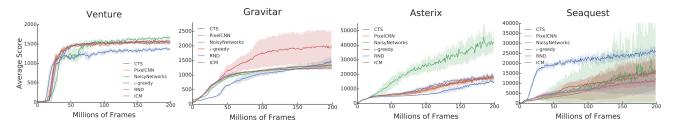


Figure 2. Evaluation of different bonus-based exploration methods on the ALE. VENTURE and GRAVITAR are hard exploration games whereas ASTERIX and SEAQUEST are easy ones.

#### 5.3. ALE training set

While the benefit of exploration bonuses has been shown on a few games they can also have a negative impact by skewing the reward landscape. To get a more complete picture, we also evaluated our agents on the original Atari training set which includes many easy exploration games. Figure 2 shows training curves for ASTERIX and SEAQUEST, the remaining games can be found in Appendix A. In this setting we noticed a reversed trend than the one observed on MONTEZUMA'S REVENGE. The pseudo-count method ends up performing worse on every game except SEAQUEST. RND and ICM are able to consistently match the level of  $\epsilon$ -greedy exploration, but not exceed it. The earlier benefits conferred by pseudo-counts result in a considerable detriment when the exploration problem is not difficult. Finally, since NoisyNets optimizes the true environment reward, and not a proxy reward, it consistently matches  $\epsilon$ -greedy and occasionally outperforms. Overall we found that bonusbased methods are generally detrimental in the context of easy exploration problems. Despite its limited performance on MONTEZUMA'S REVENGE NoisyNets gave the most consistent results across our evaluation despite its limited performance on MONTEZUMA'S REVENGE.

#### 6. Conclusion

Many exploration methods in reinforcement learning are introduced with confounding factors – longer training duration, different model architecture and new hyper parameters. This obscures the underlying signal of the exploration method. Therefore, following a growing trend in the reinforcement learning community, we advocate for better practices on empirical evaluation for exploration to fairly assess the contribution of newly proposed methods. In a standardized training environment and context, we found that  $\epsilon$ -greedy exploration can often compete with more elaborate methods on the ALE. This shows that more work is still needed to address the exploration problem in complex environments.

# 7. Acknowledgements

The authors would like to thank Hugo Larochelle, Benjamin Eysenbach, Danijar Hafner and Ahmed Touati for insightful discussions as well as Sylvain Gelly for careful reading and comments on an earlier draft of this paper.

#### References

Bellemare, M., Veness, J., and Talvitie, E. Skip context tree switching. In *International Conference on Machine Learning*, pp. 1458–1466, 2014.

Bellemare, M., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems*, pp. 1471–1479, 2016.

Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.

Bellemare, M. G., Dabney, W., and Munos, R. A distributional perspective on reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 449–458. JMLR. org, 2017.

Burda, Y., Edwards, H., Pathak, D., Storkey, A., Darrell, T., and Efros, A. A. Large-scale study of curiosity-driven learning. *arXiv preprint arXiv:1808.04355*, 2018.

Burda, Y., Edwards, H., Storkey, A., and Klimov, O. Exploration by random network distillation. In *Proceedings* of the International Conference on Learning Representations, 2019.

Castro, P. S., Moitra, S., Gelada, C., Kumar, S., and Bellemare, M. G. Dopamine: A research framework for deep reinforcement learning. arXiv preprint arXiv:1812.06110, 2018.

Choshen, L., Fox, L., and Loewenstein, Y. Dora the explorer: Directed outreaching reinforcement action-selection. *arXiv* preprint arXiv:1804.04012, 2018.

- Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T., Dunning, I., Legg, S., and Kavukcuoglu, K. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *ICML*, 2018.
- Fortunato, M., Azar, M. G., Piot, B., Menick, J., Osband, I., Graves, A., Mnih, V., Munos, R., Hassabis, D., Pietquin, O., et al. Noisy networks for exploration. In *Proceedings* of the International Conference on Learning Representations, 2018.
- Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., and Meger, D. Deep reinforcement learning that matters. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Hessel, M., Modayil, J., Van Hasselt, H., Schaul, T., Ostrovski, G., Dabney, W., Horgan, D., Piot, B., Azar, M., and Silver, D. Rainbow: Combining improvements in deep reinforcement learning. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Kapturowski, S., Ostrovski, G., Quan, J., Munos, R., and Dabney, W. Recurrent experience replay in distributed reinforcement learning. In *Proceedings of the International Conference on Learning Representations*, 2019.
- Machado, M. C., Bellemare, M. G., and Bowling, M. Count-Based Exploration with the Successor Representation. *CoRR*, abs/1807.11622, 2018a.
- Machado, M. C., Bellemare, M. G., Talvitie, E., Veness, J., Hausknecht, M., and Bowling, M. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018b.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529, 2015.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In *Proceedings of the International Conference on Machine Learning*, pp. 1928–1937, 2016.
- O'Donoghue, B., Osband, I., Munos, R., and Mnih, V. The uncertainty bellman equation and exploration. *arXiv* preprint arXiv:1709.05380, 2017.
- Osband, I., Blundell, C., Pritzel, A., and Van Roy, B. Deep exploration via bootstrapped dqn. In *Advances in neural information processing systems*, pp. 4026–4034, 2016.

- Ostrovski, G., Bellemare, M. G., van den Oord, A., and Munos, R. Count-based exploration with Neural Density Models. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2721–2730. PMLR, 2017.
- Pathak, D., Agrawal, P., Efros, A. A., and Darrell, T. Curiosity-driven exploration by self-supervised prediction. In *Proceedings of the International Conference on Machine Learning*, 2017.
- Schaul, T., Quan, J., Antonoglou, I., and Silver, D. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.
- Stadie, B. C., Levine, S., and Abbeel, P. Incentivizing exploration in reinforcement learning with deep predictive models. *arXiv preprint arXiv:1507.00814*, 2015.
- Strehl, A. L. and Littman, M. L. An analysis of model-based interval estimation for Markov decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.
- Tang, H., Houthooft, R., Foote, D., Stooke, A., Chen, O. X.,
  Duan, Y., Schulman, J., DeTurck, F., and Abbeel, P. #
  Exploration: A Study of Count-Based Exploration for
  Deep Reinforcement Learning. In *Advances in Neural Information Processing Systems*, pp. 2750–2759, 2017.
- Touati, A., Satija, H., Romoff, J., Pineau, J., and Vincent, P. Randomized value functions via multiplicative normalizing flows. *arXiv* preprint arXiv:1806.02315, 2018.
- Van den Oord, A., Kalchbrenner, N., and Kavukcuoglu, K. Pixel recurrent neural networks. arXiv preprint arXiv:1601.06759, 2016.

# A. Additional figures

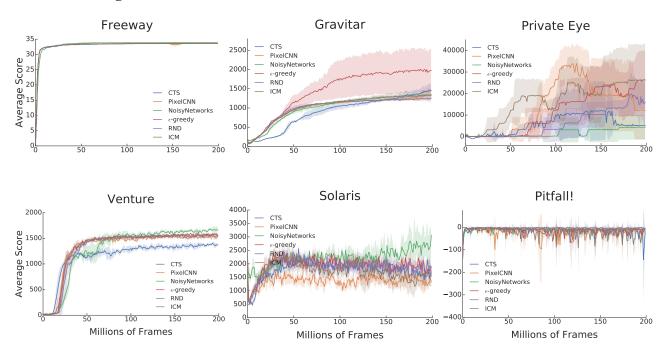


Figure 3. Results of different bonus-based exploration methods on hard exploration games. The relative ranking of methods differs from the one observed on MONTEZUMA'S REVENGE. We find that  $\epsilon$ -greedy also performs competitively. This suggests that previous claims of progress in these games has been driven by more advanced reinforcement learning algorithms, not necessarily better exploration strategies.

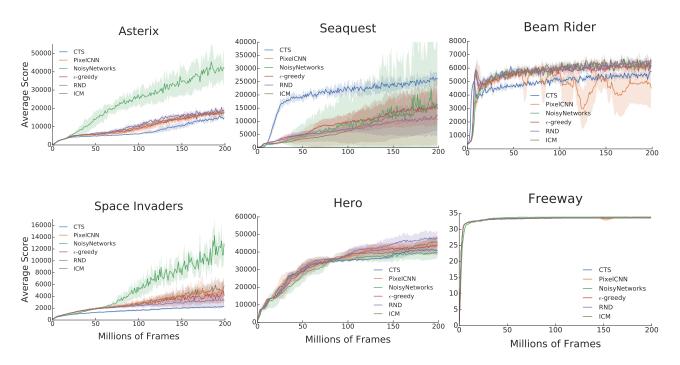


Figure 4. Evaluation of different bonus-based exploration methods on the Atari training set, except for FREEWAY all these games were classified as easy exploration problems. Rainbow with  $\epsilon$ -greedy exploration performs as well as other more complex exploration method.

The variance of the return on MONTEZUMA'S REVENGE is high because the reward is a step function, for clarity we also

provide all the training curves in Figure 5

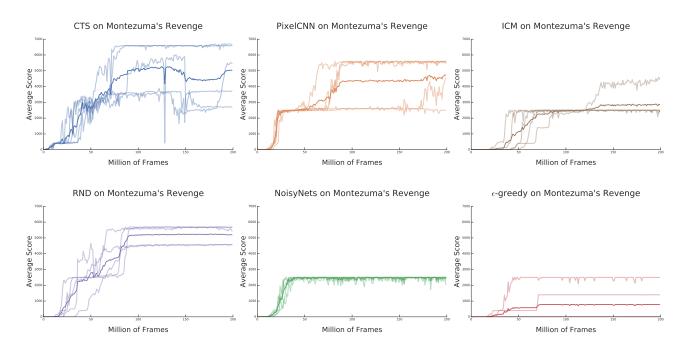


Figure 5. Training curves on MONTEZUMA'S REVENGE

# **B.** Hyperparameter tuning

Except for NoisyNets, all other methods are tuned with respect to their final performance on MONTEZUMA'S REVENGE after training on 200 million frames on five runs.

#### **B.1.** Rainbow and Atari preprocessing

We used the standard architecture and Atari preprocessing from Mnih et al. (2015). Following Machado et al. (2018b) recommendations we enable sticky actions and deactivated the termination on life loss heuristic. The remaining hyperparameters were chosen to match Hessel et al. (2018) implementation.

Hyperparameter	Value
Discount factor $\gamma$	0.99
Min history to start learning	80K frames
Target network update period	32K frames
Adam learning rate	$6.25 \times 10^{-5}$
Adam $\epsilon$	$1.5 \times 10^{-4}$
Multi-step returns n	3
Distributional atoms	51
Distributional min/max values	[-10, 10]

Every method except NoisyNets is trained with  $\epsilon$ -greedy following the scheduled used in Rainbow with  $\epsilon$  decaying from 1 to 0.01 over 1M framces.

## **B.2.** NoisyNets

We kept the original hyperparameter  $\sigma_0 = 0.5$  used in Fortunato et al. (2018) and Hessel et al. (2018).

#### Benchmarking Bonus-Based Exploration Methods on the Arcade Learning Environment

#### **B.3. Pseudo-counts**

We followed Bellemare et al.'s preprocessing, inputs are  $42 \times 42$  greyscale images, with pixel values quantized to 8 bins.

#### B.3.1. CTS

We tuned the scaling for  $\beta \in \{0.5, 0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001\}$  and found that  $\beta = 0.0005$  worked best.

#### B.3.2. PIXELCNN

We tuned the scaling factor and the prediction gain decay constant c. We ran a sweep with the following values:  $\beta \in \{5.0, 1.0, 0.5, 0.1, 0.05\}$ ,  $c \in \{5.0, 1.0, 0.5, 0.1, 0.05\}$  and found  $\beta = 0.1$  and c = 1.0 to work best.

#### **B.4. ICM**

We tuned the scaling factor and the scalar  $\alpha$  that weighs the inverse model loss against the forward model. We ran a sweep with  $\alpha = \{0.4, 0.2, 0.1, 0.05, 0.01, 0.005\}$  and  $\beta = \{2.0, 1.0, 0.5, 0.1, 0.05, 0.01, 0.005, 0.001, 0.0005\}$ . We chose  $\alpha = 0.005$  and  $\beta = 0.005$  to work best.

#### **B.5. RND**

Following Burda et al. (2019) we did not clip the intrinsic reward while the extrinsic reward was clipped (we also found in our initial experiments that clipping the intrinsic reward led to worse performance). We tuned the reward scaling factor and Adam learning rate used by RND optimizer. We ran a sweep with  $\beta = \{0.5, 0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001, 0.0005, 0.0001, 0.0005\}$  and  $lr = \{0.001, 0.0005, 0.0002, 0.0001, 0.00005\}$ . We found that  $\beta = 0.0001$  and lr = 0.0002 worked best.