# How to Stay Curious while avoiding Noisy TVs using Aleatoric Uncertainty Estimation

Augustine N. Mavor-Parker [1]  Kimberly A. Young [2 3]  Caswell Barry [2 *]  Lewis D. Griffin [4 *]

## Abstract

When extrinsic rewards are sparse, artificial agents struggle to explore an environment. Curiosity, implemented as an intrinsic reward for prediction errors, can improve exploration but it is known to fail when faced with action-dependent noise sources ('noisy TVs'). In an attempt to make exploring agents robust to noisy TVs, we present a simple solution: aleatoric mapping agents (AMAs). AMAs are a novel form of curiosity that explicitly ascertain which state transitions of the environment are unpredictable, even if those dynamics are induced by the actions of the agent. This is achieved by generating separate forward predictions for the mean and aleatoric uncertainty of future states, with the aim of reducing intrinsic rewards for those transitions that are unpredictable. We demonstrate that in a range of environments AMAs are able to circumvent action-dependent stochastic traps that immobilise conventional curiosity driven agents. Furthermore, we demonstrate empirically that other common exploration approaches—previously thought to be immune to agent-induced randomness—can be trapped by stochastic dynamics. Code to reproduce our experiments is provided.

## 1. Introduction

Efficient exploration is a central problem in reinforcement learning. Agents need to be capable of finding novel information without depending on extrinsic rewards to shepherd them through the state space of a given environment (e.g.

*Equal contribution [1]Centre for Artificial Intelligence, University College London, UK [2]Department of Cell and Developmental Biology, University College London, UK [3]Boston University, Center for Systems Neuroscience, Graduate Program for Neuroscience, USA [4]Department of Computer Science, University College London, UK. Correspondence to: Augustine Mavor-Parker <a.mavorparker@cs.ucl.ac.uk>.

Sutton & Barto (2018); Pathak et al. (2017); Burda et al. (2019b), see Weng (2020) for a review). A notable exploration method that effectively deals with sparse rewards is curiosity driven learning—where agents are equipped with a self-supervised forward prediction model that employs prediction errors as intrinsic rewards (Schmidhuber, 1991b; Pathak et al., 2017; Schmidhuber, 1991a). Curiosity is built upon the intuition that in unexplored regions of the environment, the forward prediction error of the agent's internal model will be large (Schmidhuber, 1991b; Pathak et al., 2017). As a result, agents are rewarded for visiting regions of the state space that they have not previously occupied. If, however, a particular state transition is impossible to predict, it will trap a curious agent (Burda et al., 2019b; Schmidhuber, 1991a). This is referred to as the noisy TV problem (e.g. Burda et al. (2019b); Schmidhuber (1991a)), the etymology being that a naively curious agent could dwell on the unpredictability of a noisy TV screen.

Several existing curiosity-like methods (Burda et al., 2019b; Pathak et al., 2017; 2019) aim to avoid noisy TVs or "stochastic traps" (Shyam et al., 2019). Nevertheless, employing dynamics based prediction errors as intrinsic rewards is difficult as current methods either fail when stochastic traps are action-dependent, or require an ensemble of dynamics models (Pathak et al., 2017; 2019; Shyam et al., 2019; Burda et al., 2019a). Even if ensemble methods are available, we demonstrate that they cannot reliably overcome the allure of observing random observations. Additionally, we find that random network distillation—a dynamics-free exploration technique usually assumed to be robust to stochasticity—is also susceptible to noisy TVs. Fundamentally, popular intrinsic reward approaches are vulnerable to the never ending novelty of a noisy TV.

We present a simple solution to the noisy TV problem—instead of only predicting the next state, we also predict its variance (i.e it's *aleatoric uncertainty* (Kendall & Gal, 2017)). The uncertainty of a statistical model can be described as the sum of two *theoretically* distinct types of uncertainty: epistemic uncertainty and aleatoric uncertainty (e.g. Hora (1996), see Hüllermeier & Waegeman (2021) for a review). Epistemic uncertainty measures the errors of a model's prediction that can be minimised with addi-

tional experience and learning (Hüllermeier & Waegeman, 2021). As a result, an agent using epistemic uncertainties as intrinsic rewards tends to value dynamics it has not previously encountered, and hence cannot predict accurately, but could learn to predict in the future (e.g. Osband et al. (2016)). More concretely, epistemic uncertainty can be considered to be the "expected information gain" of observing the next predicted state (Mukhoti et al., 2021). On the other hand, prediction errors that are due to aleatoric uncertainties are, by definition, a result of unpredictable processes (Hüllermeier & Waegeman, 2021). Therefore, any agent that receives intrinsic rewards for aleatoric dynamics risks being trapped, as exemplified by the noisy TV problem (Schmidhuber, 1991a; Burda et al., 2019a). By predicting aleatoric uncertainties, our curious agents are able to disregard stochastic dynamics if they are consistent with the agent's predicted variance—avoiding the trap of noisy TVs. Our contributions are summarised as follows:

1. We benchmark the performance of existing exploration techniques, highlighting that they are more vulnerable to stochasticity than previously assumed

2. We present a novel form of curiosity that can operate proficiently in exploration benchmark environments in the presence of a noisy TV, while still preserving exploration performance without a noisy TV

3. We show that even in the famously deterministic domain of Atari, natural sources of randomness exists

Finally, we also highlight the connections of AMAs to experimental neuroscience. Our approach to resisting stochasticity is both inspired by and builds upon proposals developed within neuroscience (Yu & Dayan, 2005), that suggest expected uncertainties in predictions of future states are signalled by the modulation of cortical acetylcholine in the mammalian brain. The implications for neuroscience and potential animal experiments are included in the discussion (section 5.2) and Appendix B.

## 2. Background

### 2.1. Epistemic and Aleatoric Uncertainties

Estimating the epistemic uncertainty surrounding future states would be an ideal basis for a curious agent but tractable epistemic uncertainty estimation with high dimensional data is an unsolved problem (see Gal (2016) for an introduction to the field). We implicitly incentivise agents to seek epistemic uncertainties by removing the aleatoric component from the total prediction error. More specifically, fundamental to our model is the maximum likelihood approach from Kendall & Gal (2017), which we use to predict aleatoric uncertainties based on the input. We then subtract

these aleatoric uncertainties from prediction errors, which is a novel approach for efficiently estimating epistemic uncertainties.

There are similar methods that separate epistemic and aleatoric uncertainties in return predictions (Clements et al., 2019), or within a latent variable model (Depeweg et al., 2018)—allowing for the construction of policies that are rewarded for exploring their environments and punished for experiencing aleatoric uncertainty. However, as far as we are aware, we are the first to compute aleatoric uncertainties within a *scalable* curiosity framework to remove *intrinsic rewards* for those state transitions with aleatoric uncertainty, which implicitly rewards agents for experiencing epistemic uncertainties. However, we do note that a similar approach was announced shortly after ours (Jain et al., 2021), which shows that epistemic uncertainties can be estimated by subtracting aleatoric uncertainty from a predicted prediction error. Our approach is simpler, rather than trying to predict prediction errors, we use the implicitly calculated epistemic uncertainty as intrinsic rewards online.

### 2.2. Curiosity and Intrinsic Motivation in Reinforcement Learning

Curiosity-driven (Pathak et al., 2017) agents assign value to states of the environment that they deem to be "interesting" (Still & Precup, 2012; Schmidhuber, 1997). How a curiosity based method computes whether a state is "interesting" (Still & Precup, 2012; Schmidhuber, 1997) is usually its defining characteristic. The original formulation of curiosity used prediction errors directly as intrinsic rewards (Schmidhuber, 1991b). The noisy TV problem quickly emerged when using this naïve approach in stochastic environments (Schmidhuber, 1991a). In order to evade the allure of stochastic traps, the first proposed solution to the noisy TV problem implements "interesting" (Still & Precup, 2012; Schmidhuber, 1997) as prediction errors that reduce over time (Schmidhuber, 1991a; Kaplan & Oudeyer, 2007). Others consider "interesting" (Still & Precup, 2012; Schmidhuber, 1997) to mean a high dependency between present and future states and actions (i.e. "interesting" things are predictable (Still & Precup, 2012) or controllable (Mohamed & Jimenez Rezende, 2015)).

Inverse dynamics feature (IDF) curiosity (Pathak et al., 2017) rejuvenated interest in using one step prediction errors as intrinsic rewards. IDF curiosity avoids stochastic traps by computing prediction errors with features that aim to only contain information concerning stimuli the agent can affect (Pathak et al., 2017). Further experiments (Burda et al., 2019a) showed that simple one-step prediction errors also work effectively within a random representation space generated by feeding state observations through a randomly initialised network. Burda et al. (2019a) also showed

the (IDF) approach is vulnerable to action-dependent noisy TVs—demonstrated by giving the agent a 'remote control' to a noisy TV in the environment that could induce unpredictable environment transitions. This motivated random network distillation (RND) (Burda et al., 2019b), which removes dynamics from the prediction problem altogether—instructing a network to learn to predict the output of another fixed randomly initialised network at each state, using the resulting error as intrinsic rewards. RND persists as a key component in state of the art algorithms deployed in high dimensional state spaces (Badia et al., 2019).

Other exploration methods explicitly leverage uncertainty quantification for exploration. The canonical approach is "optimism under uncertainty", which in its most basic form means weighting the value of state-actions pairs inversely to the number of times they have been experienced (Sutton & Barto, 2018)(p. 36). Known as count based methods (Strehl & Littman, 2008; Bellemare et al., 2016), this approach was shown to reliably evade noise sources in minigrid environments (Raileanu & Rocktäschel, 2020). However, it is not feasible to count state visitations in many environments where there is a large number of unique states (Bellemare et al., 2016). "Pseudo-count" methods exchange tabular look up tables for density models to estimate an analogous intrinsic reward to counts in large state spaces (Bellemare et al., 2016)—related to density models, (Kim et al., 2019) use state "compressibility" as intrinsic rewards.

Attempts have been made to reward epistemic uncertainty explicitly. This typically requires a posterior distribution over model parameters, which is intractable without approximations such as ensembles or variational inference (e.g. Houthooft et al. (2016)). Osband et al. (2016) instantiated an ensemble (Lakshminarayanan et al., 2017) into the final layer of a deep Q-network—rewarding its agents for epistemic value uncertainty. Pathak et al. (2019) use the variance of ensemble predictions being used as intrinsic rewards, while Shyam et al. (2019) reward experience of epistemic uncertainty within an ensemble of environment models. Lastly, uncertainty estimation methods have recently been developed that enforce a smoothness constraint in the representation space (Mukhoti et al., 2021; van Amersfoort et al., 2021)—allowing for sensible estimations of uncertainty to be made from learned representations—but these approaches have not yet been adopted in reinforcement learning.

## 3. Method

Our method operates in an environment defined as a Markov decision process that consist of states $s \in \mathcal{S}$, actions $a \in \mathcal{A}$, and rewards $r \in \mathcal{R} \subset \mathbb{R}$ (Sutton & Barto, 2018). At each timestep $t$ the agent selects an action via a stochastic policy $a_t \sim \pi(\cdot|s_t)$ (Szepesvári, 2010) and then receives a reward

$r_{t+1}$ and state $s_{t+1}$ generated via the transition function $p(\mathbf{s}_{t+1}, r_{t+1}|\mathbf{s}_t, a_t)$ of the environment (Sutton & Barto, 2018). The objective of the agent is to learn a stochastic policy $\pi$, parametrised by $\xi$, which aims to maximise the expectation of the sum of discounted future rewards (e.g. Mnih et al. (2016)).

$$\max_{\pi_\xi} \mathbb{E}_{\pi_\xi}\left[\sum_{k=0}^{T} \gamma^k r_{t+k}\right] \quad (1)$$

Where $T$ is the episode length and $\gamma$ is the discount factor. Following other curiosity based methods, the total reward is the sum of the intrinsic reward provided by the intrinsic reward module of the agent and the extrinsic reward provided by the environment (e.g. Pathak et al. (2017); Badia et al. (2019); Raileanu & Rocktäschel (2020); Burda et al. (2019b)).

$$r_t = \beta r_t^i + r_t^e \quad (2)$$

Where the superscripts $i$ and $e$ indicate intrinsic and extrinsic rewards, and $\beta$ is a hyperparameter that regulates the influence of intrinsic rewards on the policy. In previous works (Burda et al., 2019a), the intrinsic reward $r_t^i$ is equal to the mean squared forward prediction error of a curiosity module. To avoid stochastic traps we subtract the aleatoric uncertainty—which is constrained to have a diagonal covariance (Kendall & Gal, 2017)—from the prediction error, so that agents are not surprised by transitions that were previously learnt to be unpredictable.

$$r_t^i = \|\mathbf{s}_{t+1} - \hat{\mu}_{t+1}\|^2 - \eta \operatorname{Tr}(\hat{\mathbf{\Sigma}}_{t+1}) \quad (3)$$

Where $\hat{\mu}_{t+1}$ is the predicted mean of the next state, $\hat{\mathbf{\Sigma}}_{t+1}$ is the predicted aleatoric uncertainty of the next state and $\eta$ is a hyperparameter that regulates by how much the predicted uncertainty of the next state effects intrinsic rewards. To learn to predict the mean of the next state $\hat{\mu}_{t+1}$ and its aleatoric uncertainty $\hat{\mathbf{\Sigma}}_{t+1}$, we follow Kendall & Gal (2017)—fitting a diagonal covariance Gaussian distribution to the elements of the next state. The predictions are made by a double-headed neural network—with a mean prediction head $\mathbf{f}$ parameterised by $\theta$ and a variance prediction head $\mathbf{g}$ parametrised by $\phi$. As employed in previous works, the separate heads of the double-headed deep network share feature extracting parameters (Kendall & Gal, 2017). The prediction network performs *heteroscedastic* aleatoric uncertainty estimation (Kendall & Gal, 2017), which in a reinforcement learning context means the prediction heads are conditioned on the current state and action,

$$p(\mathbf{s}_{1:N}|\theta, \phi) = \prod_{t=1}^{N} \mathcal{N}(\mathbf{s}_{t+1}; \mathbf{f}_\theta(\mathbf{s}_t, \mathbf{a}_t), \mathbf{g}_\phi(\mathbf{s}_t, \mathbf{a}_t)) \quad (4)$$

where $N$ is the total number of states observed during training. While Kendall & Gal (2017) use *maximum a posteriori*
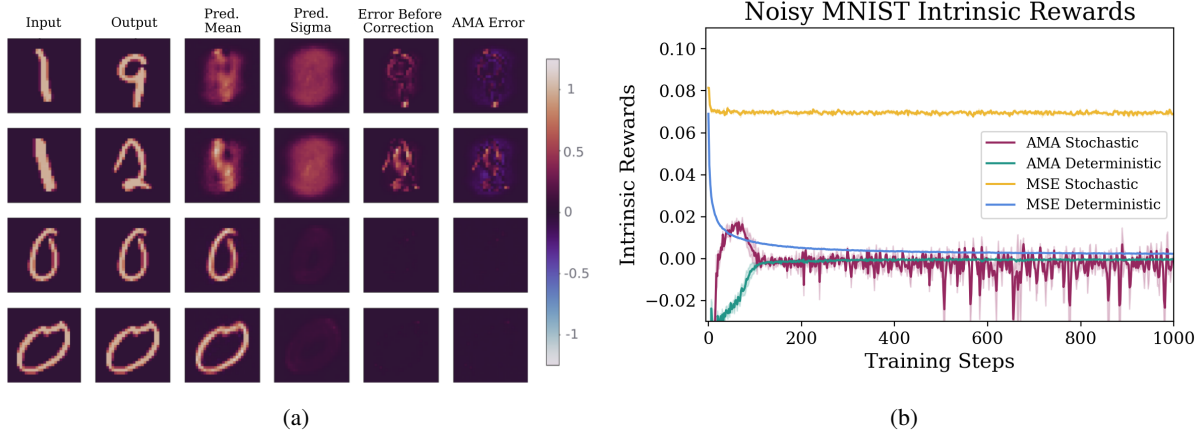
*Figure 1.* AMAs can learn to ignore stochastic transitions. (a) Example transitions from the Noisy MNIST environment along with associated predictions. The top two rows show stochastic transitions where AMA's predicted variance is high in the majority of the image allowing intrinsic reward to be small despite the stochastic transition. (b) Two reward curves for MSE and AMA are plotted where stochastic is the $1 \rightarrow \{2, ..., 9\}$ transitions and deterministic is the $0 \rightarrow 0$ transitions.

inference with a zero-mean Gaussian prior on the network parameters $\{\theta, \phi\}$, we found empirically that the resulting regularisation terms in the cost function did not improve results. Accordingly, we simply perform maximum likelihood estimation with the likelihood presented in Equation (4), resulting in the following cost function (Kendall & Gal, 2017).

$$\mathcal{L}_{t+1}(\theta, \phi) = (\mathbf{s}_{t+1} - \hat{\mu}_{t+1})^{\top} \hat{\mathbf{\Sigma}}_{t+1}^{-1}(\mathbf{s}_{t+1} - \hat{\mu}_{t+1}) + \lambda \log(\det(\hat{\mathbf{\Sigma}}_{t+1}))$$

$$(5)$$

The first term is the familiar mean squared error *divided by the uncertainty* $\hat{\mathbf{\Sigma}}_{t+1}$. The second term blocks the explosion of predicted aleatoric uncertainties (Kendall & Gal, 2017). We follow Kendall & Gal (2017)'s prescription of estimating $\log \mathbf{\Sigma}$ instead of $\mathbf{\Sigma}$ to ensure stable optimisation. Furthermore, the hyperparameter $\lambda$ was added to adjust the model's aleatoric uncertainty budget (e.g. Depeweg et al. (2018); Clements et al. (2019); Eriksson & Dimitrakakis (2019)). We use the predicted mean and aleatoric uncertianty of the next state—which are being learned online with Equation (5)—to compute intrinsic rewards according to Equation (3). Lastly, we would like to highlight that the policy network is separate to the state prediction network as in other curiosity based methods (Pathak et al., 2017).

## 4. Experiments

The purpose of this work is to improve the exploration capabilities of deep reinforcement learning algorithms in stochastic environments where current methods can fail catastrophically. As we are interested in exploration, we measure exploration directly by calculating an agent's environment coverage (when possible). Details of the noisy

TVs used—including TVs we add to environments as well as a natural noisy TV in Atari—are contained within each subsection. Extra details such as the hyperparameters and architectures used are in Appendix C. Shaded regions are standard error of the mean and we use 5 seeds for each method, except minigrid where we use 10 seeds and the MNIST experiments where only 3 were necessary.

### 4.1. Noisy MNIST

First we completed a supervised learning task, similar to the noisy MNIST environment introduced by Pathak et al. (2019). The environment does not elicit any actions from an agent. Instead, the prediction network simply needs to learn one step mappings between pairs of MNIST handwritten digits. The first images in the pairs are randomly selected 0s or 1s. When the first image is a 0 then the second image is the exact same image (these are the deterministic transitions). When the first image is a 1, then the second image is a random digit from 2-9 (these are the stochastic transitions). A prediction model capable of avoiding noisy TVs should eventually learn to compute equal intrinsic rewards for both types of transitions (Pathak et al., 2019).

We trained two different neural networks on this task (adapted from Liao (2020)), one with a mean squared error (MSE) loss function—as a baseline—and the other with the AMA loss function (Kendall & Gal, 2017). The networks are equivalent except that the AMA network has two prediction heads. Both networks contain skip connection from the input layer to the output layer and were optimised with Adam (Kingma & Ba, 2015) at a learning rate of 0.001 and a batch size of 32. The uncertainty budget hyperparameter $\lambda$ and the uncertainty weighting hyperparameter $\eta$ were set to 1 for the AMA network. The MSE prediction network is un-
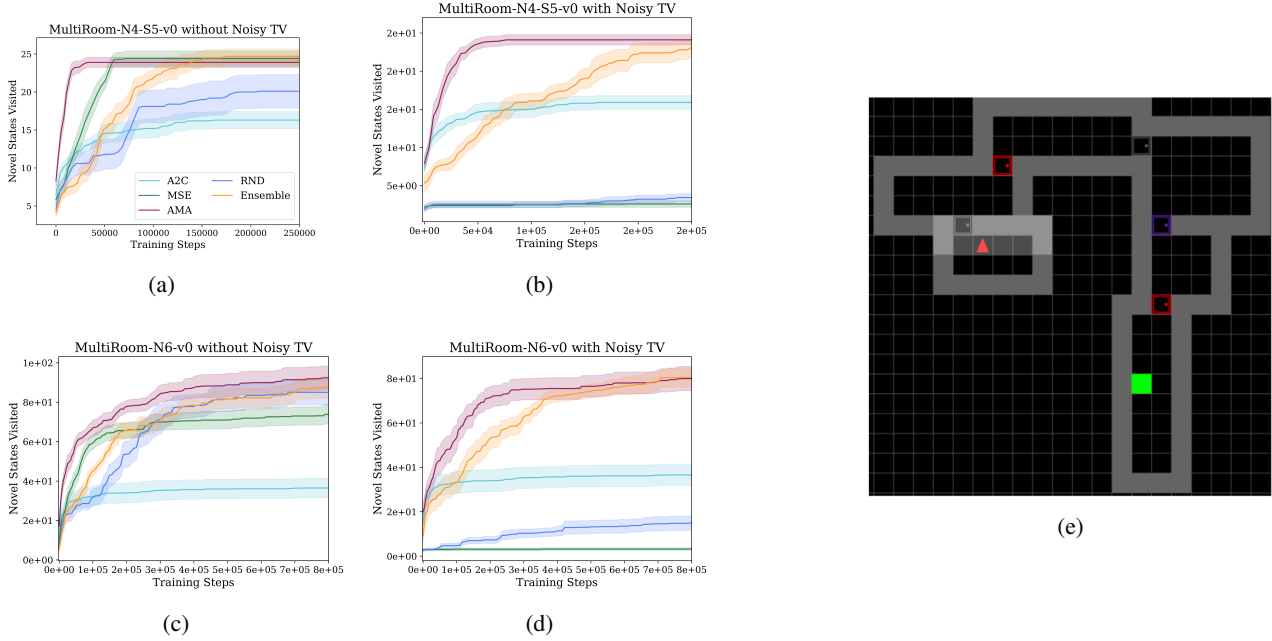
*Figure 2.* AMA agents effectively explore sparse reward minigrid environments that contain action dependent stochastic traps. (a) and (b) panel show performance on the easiest environment, containing four rooms, while the (c) and (d) show performance on a more challenging environment with six rooms. AMA and MSE have similar exploration performance when the noisy TV is absent, outperforming a no-intrinsic-reward baseline—but when a noisy TV is present only the AMA curiosity approach is able to significantly explore the environment. Ensemble methods are robust to noisy TVs in this case, but not random network distillation. Panel (e) shows an example six room environment. Standard error represents seed variation. See appendix C.3 for plots of extrinsic rewards.

able to reduce prediction errors for the stochastic transitions, causing it to produce much larger intrinsic rewards than the deterministic transitions, consistent with Pathak et al. (2019). On the other hand, the AMA prediction network is able to cut its losses by attributing high variance to the stochastic transitions, making them just as rewarding as the deterministic transitions.

### 4.2. Minigrid

Next we test AMAs on the Gym MiniGrid environment (Chevalier-Boisvert et al., 2018), which allows for resource limited deep reinforcement learning. The agent receives tensor observations describing its receptive field at each timestep. The channels of the observations represent semantic features (e.g. blue door, grey wall, empty, etc.) of each grid tile. The action space is discrete (containing actions: turn left, turn right, move forward, pick up, drop, toggle objects and done) allowing the agent to move around the environment as wells well as open and close doors. We used singleton environments but with different seeds for each run—resulting in different environment configurations for each seed. We measure exploration by counting the number of unique states visited throughout training by the zeroth actor in a distributed RL agent. An action-dependent noisy TV was added, inspired by other minigrid experiments with

noisy TVs from (Raileanu & Rocktäschel, 2020), by setting approximately half of the state observation to uniformly sampled integers within the range of possible minigrid values. When the agent selects the 'done' action the noisy TV is activated in the next observation. This is the only effect of the 'done' action. We perform policy optimisation with a synchronous advantage actor critic (A2C) implementation recommended by the gym minigrid README (Mnih et al., 2016; Willems, 2020). For the minigrid experiments we train on intrinsic and extrinsic rewards with their relative weighting being equal. The actor critic weights were optimised with RMSProp (Tieleman & Hinton, 2014) at a learning rate of 0.001, while the intrinsic reward module was optimised with the Adam (Kingma & Ba, 2015) optimizer at a learning rate of 0.001 for the AMA agent and 0.0001 for the MSE agent. All methods used the same A2C base implementation (Willems, 2020) with the same default hyperparameters. Intrinsic reward specific hyperparameters were optimised for each baseline individually with and without the noisy TV, see Appendix C.2 for details.

The uncertainty budget $\lambda$ of the AMA network was set to 0.1 as the environment representations from minigrid are very sparse, which we found empirically reduces the prediction networks willingness to predict uncertainties. The uncertainty weighting $\eta$ was set to 1 (its natural value). We found
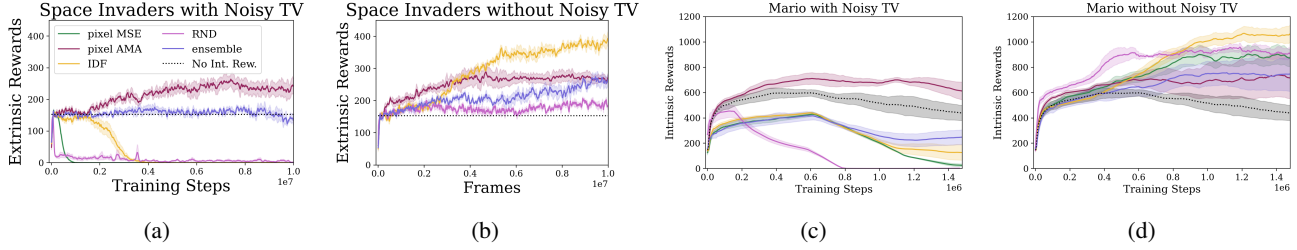
*Figure 3.* Pixel AMA performs significantly better than all baselines with a noisy TV (a) and (c) and without the distracting noisy TV AMA nearly matches its most directly comparable method Pixel MSE (b) and (d). No extrinsic rewards were used for policy optimisation. In Mario distance covered and extrinsic reward are equivalent. The y-axis plots extrinsic rewards per episode. The no intrinsic reward baseline for Space Invaders is the mean performance of the random agent data provided by Burda et al. (2019a). For Mario we were required to compute the no intrinsic reward baseline from scratch, which is the PPO agent from Burda et al. (2019b) with intrinsic rewards turned off *without* a noisy TV (but repeated in all panels for easy comparison).

that clipping intrinsic rewards to the range $[0, \infty]$ compensated for possible over predictions of uncertainty, which we implemented for AMA. The robustness of performance to hyperparameters is analysed in Section 4.6.

The forward prediction module of AMA works in the observation space as opposed to a learned feature space as is implemented in other curiosity driven methods (Pathak et al., 2017; Burda et al., 2019b). Pixel based curiosity was chosen due to its simplicity. The forward prediction model is a double headed CNN, which builds upon a previous intrinsic motivation implementation on minigrid (Raileanu & Rocktäschel, 2020). The ensembles' intrinsic rewards were implemented based on Pathak et al. (2019) using the same forward prediction network as AMA. The RND approach was adapted from Raileanu & Rocktäschel (2020)'s implementation.

We perform experiments in four and six room configurations of the minigrid (see Figure 2(e) for an example six room environment). Without a noisy TV both AMA and MSE reward functions generate visit more states compared to the no intrinsic reward baseline. On the other hand, the presence of a noisy TV profoundly affects the performance of the MSE curiosity agent, greatly reducing the number of states visited. In contrast, AMA agents are almost unaffected by the presence of an action dependent noisy TV. RND (Burda et al., 2019b) shows good performance without a noisy TV but performs poorly when a noisy TV is present. Ensemble disagreement (Pathak et al., 2019) demonstrate robustness to this version of a noisy TV in minigrid—preserving its exploration boost over A2C with and without the TV.

### 4.3. Mario and Space Invaders

We have shown AMAs can learn to ignore stochasticity in the supervised setting of MNIST and also explore proficiently in the sparse reward multiroom environments of minigrid. In this section we test whether AMAs can learn

to explore in stochastic versions of high dimensional retro video game environments. We isolated two games that have been used as curiosity benchmarks in the relevant exploration literature: Space Invaders (from Bellemare et al. (2013) used in Burda et al. (2019a); Pathak et al. (2019)) and Mario (from Nichol et al. (2018) used in Pathak et al. (2017); Burda et al. (2019b); Pathak et al. (2019)).

While these environments are useful for testing curiosity, they are mostly deterministic. Sticky action environments (Machado et al., 2018) were previously developed to make atari RL algorithms more robust to stochasticity. However, the stochasticity in sticky action environments is independent of the agents policy, meaning an agent cannot trap itself by selecting actions that generate random dynamics. This is unlike random dynamics that are likely to be found in real world applications—for example a curious warehouse robot could become distracted with watching another robots actions whose goal it does not understand. Consequently, we developed a noisy TV wrapper for atari game environments, where the action space is extended with an action that induces random grayscale tiled CIFAR-10 (Krizhevsky et al., 2014) images in place of game frames for the next observation (using[1], see Appendix C.6 for examples). When the noisy TV action is selected, the zeroth action of the action space is sent to the game emulator (the choice of the zeroth action is arbitrary).

We adapt the proximal policy optimisation (PPO) (Schulman et al., 2017) curiosity implementation from Burda et al. (2019a) into an AMA curiosity system. When using the pixel feature space, we extend Burda et al. (2019a)'s U-Net (Ronneberger et al., 2015) to use two output heads to predict the mean and variance of future states. We leave all PPO hyperparameters equal to their values from Burda et al. (2019a). We set the uncertainty budget hyperparameter
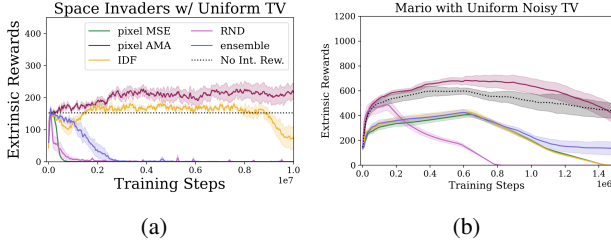
---

[1] https://github.com/snatch59/load-cifar-10

(a)                                          (b)



(a)                                          (b)

*Figure 4.* AMA is robust to Noisy TVs of a very random uniform noise distribution (random pixels from 0-255) while other baselines are also trapped by this additional noisy TV. We verify in C.8 the eventual collapse in Mario is a property of the curiosity methods tested and not a feature of AMA.

*Figure 5.* IDF AMA (a) and Pixel AMA (b) avoids the natural trap in Bank Heist and so explores more of the maze on average than the IDF MSE. (b) shows the different seeds from Pixel AMA plotted individually due to the very long episodes of pixel AMA. An implementation detail means the values of the x-axis above are a close approximation to their true values see C.7 for details.

$\lambda$ to 1 and the uncertainty weighting hyperparameter $\eta$ to 2—doubling the punishment given to our agents for experiencing aleatoric uncertainty. For the these experiments we do not clip intrinsic rewards. The hyperparameter optimisation process is described in Appendix C.4. We also test the robustness of AMA to alternative noise distributions (Section 4.5) as well as different settings of hyperparameters (Section 4.6). We compare AMA to four alternative intrinsic reward methods: random network distillation (RND) (Burda et al., 2019b), inverse dynamics feature (IDF) curiosity (Pathak et al., 2017), MSE pixel based curiosity (Burda et al., 2019a) and ensemble disagreement (Pathak et al., 2019). It is important to stress *policy optimisation is done with intrinsic rewards only* following Burda et al. (2019a).

Compared to the relatively weak baseline of pixel based curiosity, the Space Invaders and Mario experiments show similar results to the minigrid experiments—MSE and AMA pixel based curiosity have comparable performance when no noisy TV is present (Figure 3(b) and 3(d)), while with the noisy TV AMA greatly outperforms MSE pixel based curiosity (Figure 3(a) and Figure 3(c)). Unsurprisingly, RND and IDF curiosity maintain their superiority over pixel based methods without a noisy TV (Figure 3(b) and 3(d)). However, unlike AMA curiosity, both these baselines are vulnerable to action dependent noisy TVs (Figure 3(a) and Figure 3(c)). While Ensembles and RND have previously been shown to be able to evade stochastic traps (Burda et al., 2019b), it seems that if the number of novel states in the trap is large then their approximation of epistemic uncertainty breaks down.

### 4.4. Bank Heist: An Atari Game with a Natural Noisy TV

Although Atari games are famously deterministic (Machado et al., 2018), we identified a naturally ocurring stochastic trap in the Bank Heist gameplay videos of the original IDF curiosity paper (Pathak et al., 2017). The objective of Bank Heist is to simultaneously avoid police cars and navigate to
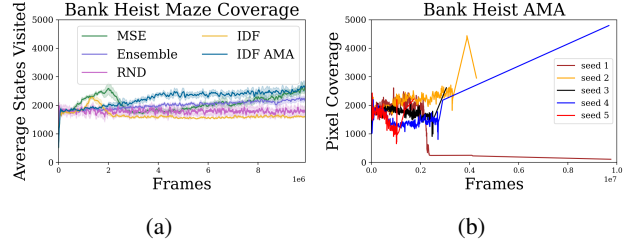
banks distributed across four 2D mazes—which can be entered and exited through the sides of the screen. Importantly, with each enter/exit the bank locations reset randomly.

When trained on purely intrinsic rewards, IDF curiosity will perpetually enter and exit the maze while also dropping dynamite. This can create high prediction error as it is impossible to predict when the dynamite will explode and where the banks will regenerate. An example video of this behaviour orignally from (Pathak et al., 2017) is provided[2]. To measure the effect of this pathological behaviour on exploration, we count the number of pixels covered by the car on average in an episode.

As expected, the IDF agent initially begins to explore the environment before coverage decreases (Figure 5(a)). On the other hand, most of the other baselines seem to be relatively unaffected by this trap—although they do not seem to explore particularly proficiently. Pixel AMA seems to develop an interesting strategy which makes its performance difficult to plot—so we plot all of its seeds individually in Figure 5(b). Pixel AMA discovers does not rob any banks so that it is not be chased by the police, allowing it to explore unimpeded with very long episodes (so long that they exceed the $10^7$ timesteps used for training hence why the line is cut off at different points for different seeds). Overall, pixel AMA explores the maze well, with 2 seeds getting good performance, 2 seeds getting extremely good performance and 1 seed becoming bored and deciding to give up on exploration. Lastly, we integrated the AMA prediction paradigm into the IDF approach, predicting the mean and aleatoric variance of future state representations and computing rewards with Equation (3) with AMA hyperparameters $\lambda$ and $\eta$ set to 1. This ablation was not susceptible to the natural Bank Heist trap (Figure 5). Overall, this section highlights that noisy TVs can be an inherent issue even in environments with

---

[2]https://www.youtube.com/watch?v=S4YdZe70XMQ

simple dynamics.

### 4.5. Different Noise Distributions

A key desiderata of AMA is that exploration performance should be robust to noisy TVs no matter the distribution they follow. In addition to the CIFAR Noisy TV experiments, we performed further experiments to see how AMA and the baselines would react to uniform "salt and pepper" noise in Mario. Although uniform noise is a simpler distribution than natural images, it is even more distracting than CIFAR images due to its high entropy. Consequently, AMA is again the only method that does not catastrophically fail when a noisy TV is present as can be seen in Figure 4. Interestingly, ensembles perform considerably worse with uniform noise than with CIFAR noise. AMA's extrinsic reward decreases slightly in Mario—we verified in appendix C.8 that this is a result of the well known boredom behaviour of curiosity systems discovered by (Pathak et al., 2017), rather than being distracted by a Noisy TV.

### 4.6. Hyperparameter Ablations

The hyperparameters $\lambda$ and $\eta$ have simple interpretations—$\lambda$ is a prior on the width of the noise distribution, while $\eta$ controls how much AMA is punished for experiencing random dynamics. Although $\lambda$ and $\eta$ are interpretable, it is important to understand the effect on exploration performance when hyperparameters are perturbed from their optimal values. Table 1 shows that Mario performance in relatively robust to hyperparameter selection but as expected $\eta$ and $\lambda$ can be tuned to be more susceptible to noisy TVs. Table 2 shows that it is not crucial to tune $\eta$ minigrid but due to the nature of minigrid observations it is important to have a wide prior on the variance of states. An interesting direction for further work would be to try to estimate $\lambda$ online from experience.

## 5. Discussion

### 5.1. Limitations

The AMA reward function implicitly rewards epistemic uncertainty by assuming the total uncertainty can be decomposed into epistemic and aleatoric uncertainties. While theoretically true (Kendall & Gal, 2017), there is no guarantee that AMAs are able to surgically subtract those errors due to aleatoric dynamics from the total prediction error. Additionally, aleatoric uncertainty estimates are not guaranteed to be reliable for out of distribution data, meaning intrinsic rewards could become less reliable the further the agent travels into novel territory (Mukhoti et al., 2021). In practice we find that a stochastic policy—and clipping in the case of minigrid—offsets potentially deceptive intrinsic rewards. Furthermore, trust region methods in the

| $\eta$ | $\lambda$ | Max. X Distance | Max. X Distance w/ TV |
|---|---|---|---|
| 2 | 0.1 | 598 | 641 |
| 2 | 1 | 616 | 715 |
| 2 | 10 | 1092 | 451 |
| 1 | 1 | 725 | 814 |
| 0.5 | 1 | 568 | 516 |

*Table 1.* Mario results with different hyperparameters. Hyperparameter tuning is not crucial in Mario. However, increasing $\lambda$ or decreasing $\eta$ can make AMA susceptible to noisy TVs.

| $\eta$ | $\lambda$ | Novel States | Novel States w/ TV |
|---|---|---|---|
| 1 | 0.1 | 100 | 92 |
| 1 | 1 | 85 | 5 |
| 1 | 10 | 44 | 3 |
| 2 | 0.1 | 115 | 107 |
| 0.5 | 0.1 | 100 | 70 |

*Table 2.* Minigrid results with different hyperparameters. Due to the sparse observations of minigrid, it is important to place a prior of being more eager to predict uncertainty by decreasing $\lambda$. For minigrid, tuning $\eta$ is not very important.

retro games (i.e. using PPO (Schulman et al., 2017) instead of A2C (Mnih et al., 2016)) may also compensate for occasionally deceptive rewards—suggested by the fact that intrinsic reward clipping was not necessary for the retro game experiments. Finally, we note that like all curiosity approaches (e.g. Burda et al. (2019a)), our method generates non-stationary rewards, which is known to make learning difficult for RL agents.

### 5.2. Acetylcholine

AMAs build upon (Yu & Dayan, 2003)'s work on Acetylcholine—a neurotransmitter associated with "expected" uncertainty signalling in the brain (we point the interested reader to Appendix A for a literature review). The exact nature of the uncertainty signalled by acetylcholine—whether it is epistemic or aleatoric—is an open question in theoretical neuroscience posed by (Yu & Dayan, 2005). Previously, it was obvious why a biological agent would find epistemic uncertainty predictions useful as they can be used to maximise information gain (Mukhoti et al., 2021). In this work, we present a use case for aleatoric uncertainty predictions—rejuvenating interest in (Yu & Dayan, 2005)'s call for experiments to analyse the kinds of uncertainty signalled in brain. To codify the relevance of AMAs to theoretical neuroscience, we propose an animal experiment in B as well as theoretical predictions of cholinergic activity in the proposed task.

# 6. Conclusion

We have shown AMAs are able to avoid action-dependent stochastic traps that destroy the exploration capabilities of conventional curiosity driven agents in environments with high entropy noisy TVs (Burda et al., 2019a; Pathak et al., 2019; Burda et al., 2019b). AMAs tractably avoid stochastic traps by decreasing intrinsic rewards in regions with high estimated aleatory. Future RL research should aim to integrate the AMA approach into curiosity methods that operate in feature spaces besides pixels or even within those methods that circumvent dynamics altogether (e.g. Burda et al. (2019b)), with the aim of achieving state of the exploration even when noisy TVs are present. Lastly, more work should be done to understand the cause of the failure cases of popular intrinsic reward methods in stochastic environments.

# Acknowledgements

# References

Acquas, E., Wilson, C., and Fibiger, H. C. Conditioned and unconditioned stimuli increase frontal cortical and hippocampal acetylcholine release: effects of novelty, habituation, and fear. *Journal of Neuroscience*, 16(9): 3089–3096, 1996.

Badia, A. P., Sprechmann, P., Vitvitskyi, A., Guo, D., Piot, B., Kapturowski, S., Tieleman, O., Arjovsky, M., Pritzel, A., Bolt, A., et al. Never give up: Learning directed exploration strategies. In *International Conference on Learning Representations*, 2019.

Barry, C., Heys, J. G., and Hasselmo, M. E. Possible role of acetylcholine in regulating spatial novelty effects on theta rhythm and grid cells. *Frontiers in neural circuits*, 6:5, 2012.

Bellemare, M., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. Unifying count-based exploration and intrinsic motivation. *Advances in neural information processing systems*, 29:1471–1479, 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.

Burda, Y., Edwards, H., Pathak, D., Storkey, A., Darrell, T., and Efros, A. A. Large-scale study of curiosity-driven learning. In *International Conference on Learning Representations*, 2019a.

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

Chevalier-Boisvert, M., Willems, L., and Pal, S. Minimalistic gridworld environment for openai gym. https://github.com/maximecb/gym-minigrid, 2018.

Clements, W. R., Robaglia, B.-M., Van Delft, B., Slaoui, R. B., and Toth, S. Estimating risk and uncertainty in deep reinforcement learning. *arXiv preprint arXiv:1905.09638*, 2019.

Depeweg, S., Hernandez-Lobato, J.-M., Doshi-Velez, F., and Udluft, S. Decomposition of uncertainty in bayesian deep learning for efficient and risk-sensitive learning. In *International Conference on Machine Learning*, pp. 1184–1193. PMLR, 2018.

Eriksson, H. and Dimitrakakis, C. Epistemic risk-sensitive reinforcement learning. *arXiv preprint arXiv:1906.06273*, 2019.

Gal, Y. Uncertainty in deep learning. *University of Cambridge*, 1(3), 2016.

Giovannini, M., Rakovska, A., Benton, R., Pazzagli, M., Bianchi, L., and Pepeu, G. Effects of novelty and habituation on acetylcholine, gaba, and glutamate release from the frontal cortex and hippocampus of freely moving rats. *Neuroscience*, 106(1):43–53, 2001.

Hasselmo, M. E. The role of acetylcholine in learning and memory. *Current opinion in neurobiology*, 16(6):710–715, 2006.

Hora, S. C. Aleatory and epistemic uncertainty in probability elicitation with an example from hazardous waste management. *Reliability Engineering & System Safety*, 54(2-3):217–223, 1996.

Houthooft, R., Chen, X., Duan, Y., Schulman, J., De Turck, F., and Abbeel, P. Vime: Variational information maximizing exploration. *Advances in Neural Information Processing Systems*, 29:1109–1117, 2016.

Hüllermeier, E. and Waegeman, W. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning*, 110(3):457–506, 2021.

Jain, M., Lahlou, S., Nekoei, H., Butoi, V., Bertin, P., Rector-Brooks, J., Korablyov, M., and Bengio, Y. Deup:

Direct epistemic uncertainty prediction. *arXiv preprint arXiv:2102.08501*, 2021.

Kaplan, F. and Oudeyer, P.-Y. In search of the neural circuits of intrinsic motivation. *Frontiers in neuroscience*, 1:17, 2007.

Kendall, A. and Gal, Y. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pp. 5574–5584, 2017.

Kim, Y., Nam, W., Kim, H., Kim, J.-H., and Kim, G. Curiosity-bottleneck: Exploration by distilling task-specific novelty. In *International Conference on Machine Learning*, pp. 3379–3388. PMLR, 2019.

Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *ICLR (Poster)*, 2015.

Krizhevsky, A., Nair, V., and Hinton, G. The cifar-10 dataset. *online: http://www. cs. toronto. edu/kriz/cifar. html*, 55 (5), 2014.

Lakshminarayanan, B., Pritzel, A., and Blundell, C. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in Neural Information Processing Systems*, 30, 2017.

Liao, X. pytorch simple autoencoder. `https://github.com/L1aoXingyu/pytorch-beginner/blob/master/08-AutoEncoder/simple_autoencoder.py`, 2020.

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, 2018.

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 *International conference on machine learning*, pp. 1928–1937, 2016.

Mohamed, S. and Jimenez Rezende, D. Variational information maximisation for intrinsically motivated reinforcement learning. *Advances in Neural Information Processing Systems*, 28:2125–2133, 2015.

Mukhoti, J., Kirsch, A., van Amersfoort, J., Torr, P. H., and Gal, Y. Deterministic neural networks with appropriate inductive biases capture epistemic and aleatoric uncertainty. *arXiv preprint arXiv:2102.11582*, 2021.

Nichol, A., Pfau, V., Hesse, C., Klimov, O., and Schulman, J. Gotta learn fast: A new benchmark for generalization in rl. *arXiv preprint arXiv:1804.03720*, 2018.

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.

Parikh, V., Kozak, R., Martinez, V., and Sarter, M. Prefrontal acetylcholine release controls cue detection on multiple timescales. *Neuron*, 56(1):141–154, 2007.

Parr, T. and Friston, K. J. Uncertainty, epistemics and active inference. *Journal of The Royal Society Interface*, 14 (136):20170376, 2017.

Pathak, D., Agrawal, P., Efros, A. A., and Darrell, T. Curiosity-driven exploration by self-supervised prediction. In *International Conference on Machine Learning*, pp. 2778–2787. PMLR, 2017.

Pathak, D., Gandhi, D., and Gupta, A. Self-supervised exploration via disagreement. In *International conference on machine learning*, pp. 5062–5071. PMLR, 2019.

Pepeu, G. and Giovannini, M. G. Changes in acetylcholine extracellular levels during cognitive processes. *Learning & memory*, 11(1):21–27, 2004.

Raileanu, R. and Rocktäschel, T. Ride: Rewarding impact-driven exploration for procedurally-generated environments. In *International Conference on Learning Representations*, 2020. URL `https://openreview.net/forum?id=rkg-TJBFPB`.

Ranganath, C. and Rainer, G. Neural mechanisms for detecting and remembering novel events. *Nature Reviews Neuroscience*, 4(3):193–202, 2003.

Ronneberger, O., Fischer, P., and Brox, T. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.

Schmidhuber, J. Adaptive confidence and adaptive curiosity. In *Institut fur Informatik, Technische Universitat Munchen, Arcisstr. 21, 800 Munchen 2*. Citeseer, 1991a.

Schmidhuber, J. A possibility for implementing curiosity and boredom in model-building neural controllers. In *Proc. of the international conference on simulation of adaptive behavior: From animals to animats*, pp. 222–227, 1991b.

Schmidhuber, J. What's interesting?, 1997.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Shyam, P., Jaśkowski, W., and Gomez, F. Model-based active exploration. In *International Conference on Machine Learning*, pp. 5779–5788, 2019.

Still, S. and Precup, D. An information-theoretic approach to curiosity-driven reinforcement learning. *Theory in Biosciences*, 131(3):139–148, 2012.

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.

Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.

Szepesvári, C. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103, 2010.

Tieleman, T. and Hinton, G. Rmsprop gradient optimization. *URL http://www. cs. toronto. edu/tijmen/csc321/slides/lecture_slides_lec6. pdf*, 2014.

van Amersfoort, J., Smith, L., Jesson, A., Key, O., and Gal, Y. On feature collapse and deep kernel learning for single forward pass uncertainty. *arXiv preprint arXiv:2102.11409*, 2021.

Weng, L. Exploration strategies in deep reinforcement learning. *lilianweng.github.io/lil-log*, 2020. URL https://lilianweng.github.io/lil-log/2020/06/07/exploration-strategies-in-deep-reinforcement-learning.html.

Willems, L. RL Starter Files. https://github.com/lcswillems/rl-starter-files, 2020.

Yu, A. and Dayan, P. Expected and unexpected uncertainty: Ach and ne in the neocortex. *Advances in neural information processing systems*, pp. 173–180, 2003.

Yu, A. and Dayan, P. Uncertainty, neuromodulation, and attention. *Neuron*, 46(4):681–692, 2005.

# Appendix

## A. Acetylcholine

In the mammalian brain acetylcholine is implicated in a range of processes including learning and memory, fear, novelty detection, and attention (Ranganath & Rainer, 2003; Pepeu & Giovannini, 2004; Acquas et al., 1996; Barry et al., 2012; Yu & Dayan, 2005; Hasselmo, 2006; Giovannini et al., 2001; Parikh et al., 2007). Traditional views—supported by the rapid increase in cholinergic tone in response to environmental novelty and demonstrable effects on neural plasticity—emphasised its role as a learning signal, generating physiological changes that favour encoding of new information over retrieval (Hasselmo, 2006).

Notably, (Yu & Dayan, 2003) proposed an alternative perspective, suggesting that acetylcholine signals the *expected uncertainty* of top down *predictions*, while modulation of norepinephrine is a result of *unexpected uncertainties*. More concretelty, (Yu & Dayan, 2003)'s model can be seen as favouring bottom up sensory input over top down predictions if predictions are believed to be inaccurate—consistent with evidence that shows acetylcholine inhibits feedback connections and strengthens sensory inputs (Hasselmo, 2006). However this approach does not explicitly separate epistemic and aleatoric uncertainties (Yu & Dayan, 2005). In contrast, the utility of quantifying epistemic uncertainties for exploration has been widely recognised in the RL literature (e.g. (Osband et al., 2016; Pathak et al., 2019)). Here we demonstrate a potential use of aleatoric uncertainties in exploring agents both biological and artificial. Namely, aleatoric uncertainties can be used to divert attention away from unpredictable dynamics when using prediction errors as intrinsic rewards. This is similar to a model proposed by (Parr & Friston, 2017), suggesting acetylcholine may indicate expected uncertainties in top down predictions within an MDP.

In this context we propose an extension to (Yu & Dayan, 2005)'s dichotomy. Specifically, we suggest that in the mammalian brain, cortical acetylcholine signals expected aleatoric uncertainties, while norepinephrine is modulated by epistemic uncertainties both expected and unexpected. This formulation is attractive in an ML framework, providing a means to avoid stochastic traps, while also being consistent with empirical biological data (Hasselmo, 2006; Yu & Dayan, 2003).

## B. A Proposed Test for the Aleatoric Model of Acetylcholine

Inspired by (Yu & Dayan, 2005), we propose that in the mammalian brain acetylcholine signals aleatoric uncertainty surrounding future states. However, we are not aware of any experimental neuroscience data that elucidates the specific nature of the uncertainty signalled by acetylcholine. As a result, this section proposes a 1D rodent VR task designed to test the specific nature of cholinergic uncertainty signalling in the mammalian brain which we hope will be picked up by experimental neuroscientists. To supplement our experimental proposal, we compute theoretical predictions of cholinergic activity within either an aleatoric or epistemic acetylcholine model—two competing interpretations of (Yu & Dayan, 2005)'s work. The aleatoric model uses aleatoric uncertainties as a theoretical acetylcholine signal (Kendall & Gal, 2017), while the epistemic model uses ensemble variance as an acetylcholine signal (Pathak et al., 2019).

The proposed task places an animal in a VR corridor containing a series of spatial landmarks and two reward zones in which it must respond in order to have a chance of receiving a reward. Responding in reward zone A causes the animal to teleport to a random position along the track. Responding in reward zone B causes the animal to teleport to a fixed position on the track.

To compute how both models predict the cholinergic signal should respond in the proposed rodent VR experiment, we simulate the task with a simple multi-armed bandit environment. In our bandit model of the task an agent predicts a 1D function by sampling minibatches from different regions of the input. In one region of sample space the function takes a simple sinuisoidal form, analogous to zone A of the VR track, in a second region the function consists of points randomly sampled from a standard normal distribution at each timestep, analogous to zone B (Figure 6(c)). As described previously, we applied two models to this task, in the first acetycholine was identified with aleatoric uncertainty, while in the second—as a control—acetylcholine tracks epistemic uncertainty.

We trained an action value based multi-armed bandit to maximise intrinsic rewards for two kinds of forward prediction models: a double headed network trained to optimise the AMA objective and an ensemble of networks where each member is minimising their own MSE (e.g. (Pathak et al., 2019)). The aleatoric model uses the AMA reward function whereas the epistemic model is (intrinsically) rewarded for variance in ensemble predictions. We plot both models' uncertainties in each reward zone over time—recovering a clear prediction of cholinergic activity in both cases. The aleatoric uncertainty of AMA
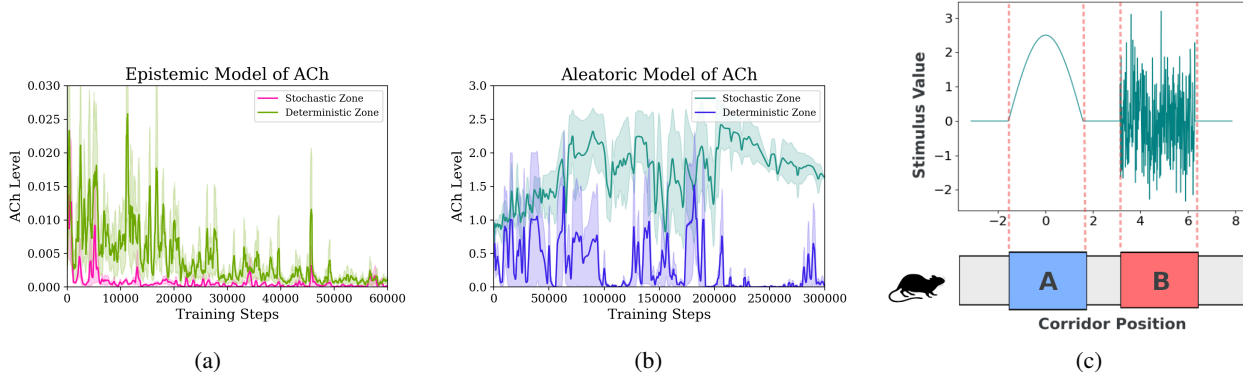
(a)    (b)    (c)

*Figure 6.* Predictions on a theoretical experiment to illuminate the epistemic or aleatoric nature of cholinergic signalling in the brain. (a) and (b) show that the epistemic model predicts acetycholine will eventually decrease in both zones while the aleatoric model predicts acetylcholine decreasing to zero in the stable zone but remaining high in the noisy zone. Panel (c) shows our 1D model of the proposed animal experiment where the bandit samples different 'corridor positions' and receives scalar stimuli, which it is trying to predict—using the resulting prediction errors as intrinsic rewards. Plots are smoothed with the same method as the Bank Heist plots C.7.

remains high in reward zone B but decreases in reward zone A. On the other hand, the epistemic model shows a decrease in acetylcholine in both reward zones over time. We hope these clear and distinct predictions on the nature of cholinergic uncertainty signalling will be tested by the experimental neuroscience community in a task similar to the one we propose.

## C. Implementation Details

### C.1. Noisy MNIST

We use three random seeds for the repeats of the MNIST experiments. The results in the graph show test set performance. The hyperparameters used are listed below. The learning rate was manually tuned so that the identity transformation was learned for the deterministic transitions (hence very low loss for the MSE and AMA network) and the AMA network produced sensible uncertainty estimates for the stochastic transitions.

| Hyperparameter | Value |
|---|---|
| MSE Learning Rate | 0.001 |
| AMA Learning Rate | 0.0001 |
| Batch Size | 32 |
| AMA uncertainty budget $\lambda$ | 1 |
| AMA uncertainty coefficient $\eta$ | 1 |

*Table 3.* Noisy MNIST hyperparameters

### C.2. Minigrid

We used 5 seeds for the minigrid experiments when tuning the learning rate for respective curiosity modules via grid search: $\in [0.01, 0.001, 0.0001]$, we selected the learning rate for each method that achieved both good exploration (as measured by novel states visited) and extrinsic rewards (if any were achieved) for both the noisy TV and no noisy TV settings on the 6 room environment. We give equal weighting to intrinisc and extrinsic rewards. We used different seeds for the grid search and the final results. For the final results we used 10 seeds. We also adapt an implementation of the Welford algorithm from stack overflow for normalising rewards [3]. We found to return to goal states it is important to normalise the sum of extrinsic and intrinsic rewards, which we now implement for the novel states visited results in the main text (the previous manuscript only normalised intrinsic rewards in the novel states visited results shown in the main text). The architecture for forward prediction is adapted from the implementation from (Raileanu & Rocktäschel, 2020) but in preliminary experimentation we

---
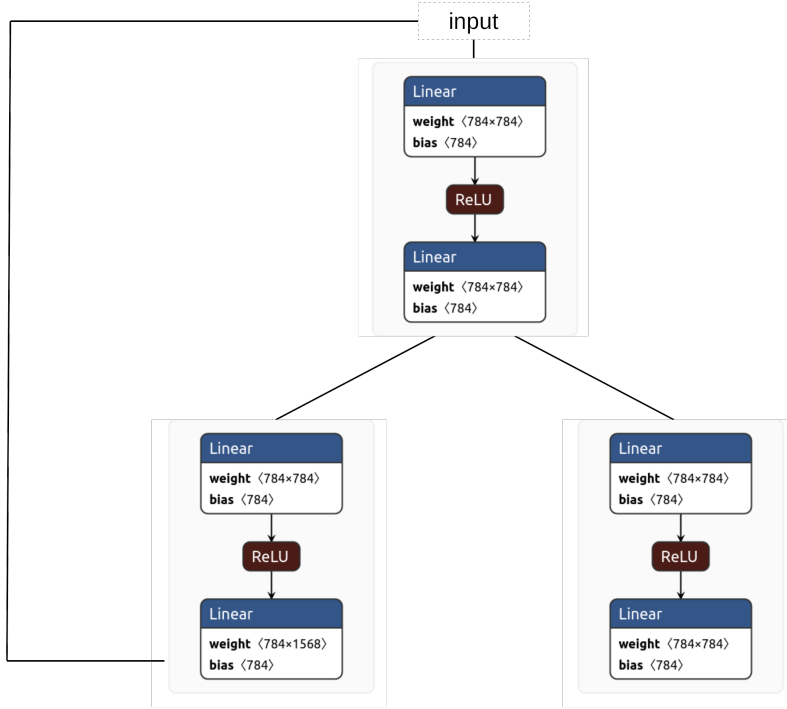[3] https://stackoverflow.com/a/5544108/13216535

*Figure 7.* MNIST image to image architecture. For the MSE network the second log variance branch was discarded. A skip connection was provided from the input to the final layer for both AMA and MSE.

ended up changing their prediction architecture dramatically.

### C.3. Extrinsic Rewards Experiments in Minigrid

Additional results showing the extrinsic rewards achieved by different agents in gym minigrid is shown in figure 8.

### C.4. Atari

For each run we used five seeds. We use the official implementations (Burda et al., 2019b;a) for the baselines we compare to (with their default hyperparameters). For AMA and Pixel MSE we adapt from (Burda et al., 2019a). The hyperparameters used for our AMA experiments can be found in Table 5. We did not change the PPO/vanilla curiosity hyperparameters from the original implementation we adapted and only changed the AMA hyperparameters. The hyperparameters were chosen by first exploring different configurations on smaller minigrid environments and evaluating promising configurations on the Atari environments.

The UNet architecture used for the forward predictions is described below. We duplicate the decoder head to create a two headed output but we only describe the encoder and decoder here. For the Pixel MSE baselines we use identical architectures but don't use the uncertainty predictions and train on MSE only. There are UNet style residual connections between the corresponding encoder and decoder layers. Leaky ReLU activations are used in the encoder layers and Tanh activations are used in the decoder layers. Batch normalisation is used throughout the hidden layers. Action information is concatenated at each layer. See supplementary code for further details.

To integrate AMA into the IDF approach, we did not share any representations between mean and variance prediction heads, instead we used two prediction MLPs for the mean and variance. Leaky ReLU is used throughout hidden layers and action information is concatenated at each layer. We used five layers with 512 units each and UNet style residual connections.
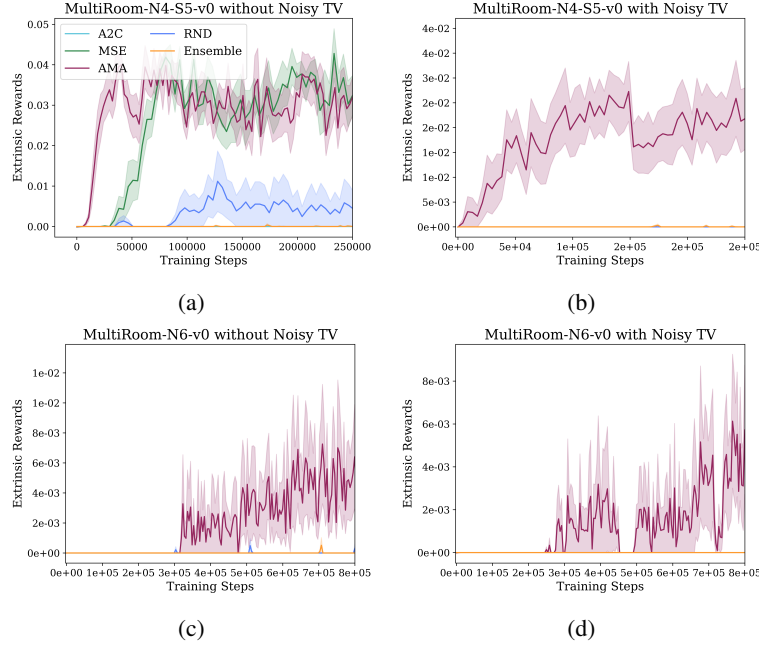
Figure 8. AMA is the only approach that can achieve significant extrinsic rewards both with and without a noisy TV present in minigrid. Mean squared error intrinsic rewards are also able to achieve significant extrinsic rewards on the 4 room environment without the noisy TV present. Interestingly, aleatoric mapping agents are the only method that is able to receive significant extrinsic rewards on the 6 room environment even when no noisy TV is present, suggesting aleatoric mapping agents may improve exploration behaviours even when a prominent noisy TV is not present—this is consistent with the novel states visited results. One reason could be that the partial observability of the minigrid environments mean that one step predictions are sometimes unpredictable. For example, an empty observation that has no objects or walls may be presented to the agent in many different positions making the prediction of the next observation impossible.

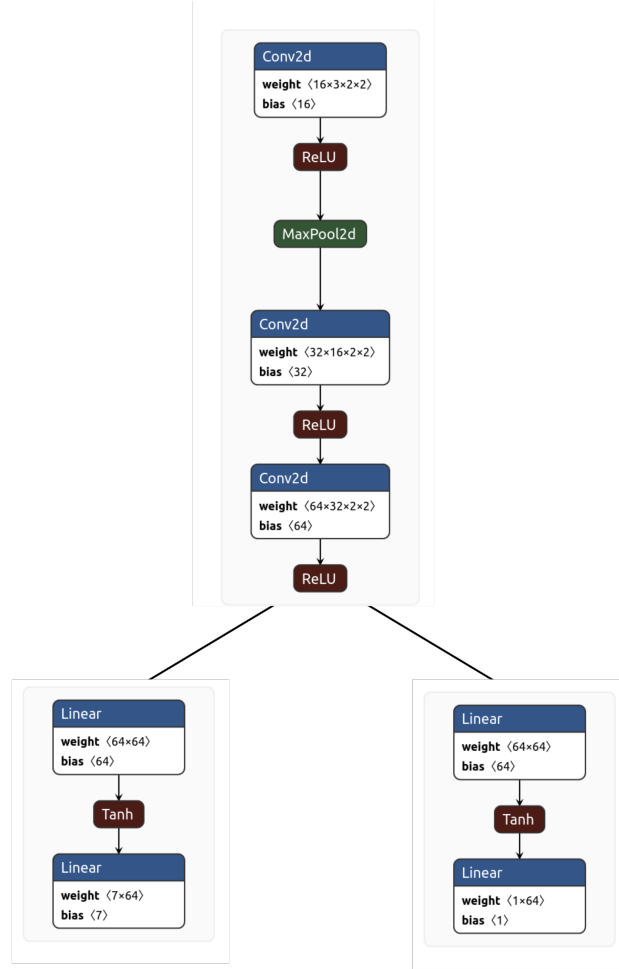| Hyperparameter | Value |
|---|---|
| AMA learning rate | 0.001 |
| MSE learning rate | 0.0001 |
| Disagreement learning rate | 0.001 |
| Random Network Distillation learning rate | 0.001 |
| RMS Prop $\alpha$ | 0.99 |
| RMS Prop $\epsilon$ | 1.000e-8 |
| number of actors | 16 |
| unroll length | 5 |
| discount factor $\gamma$ | 0.99 |
| policy learning rate | 0.001 |
| GAE $\lambda$ | 0.95 |
| entropy coefficient | 0.01 |
| value loss coefficient | 0.5 |
| max grad norm | 0.5 |
| AMA uncertainty budget $\lambda$ | 0.1 |
| AMA uncertainty coefficient $\eta$ | 1 |

*Figure 9.* Actor critic architecture for the policy network in the minigrid experiments.

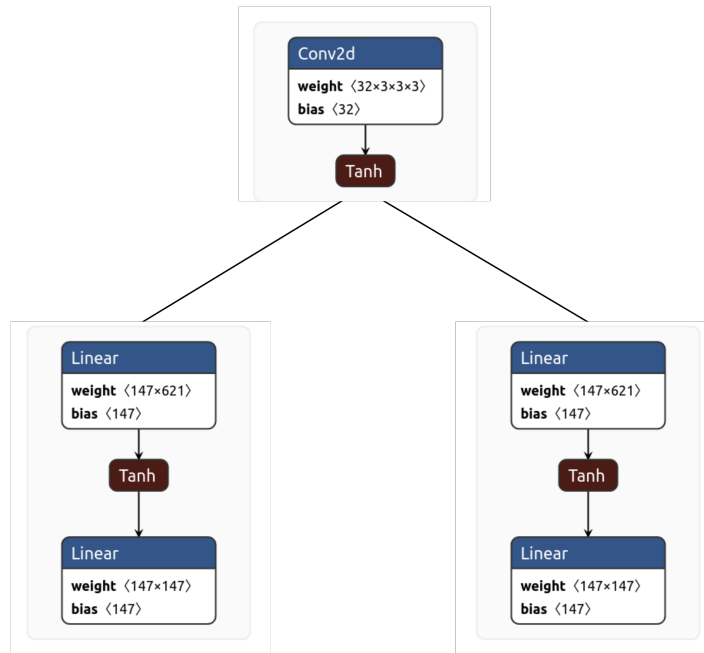| Hyperparameter | Value |
|---|---|
| global learning rate | 0.0001 |
| normalise rewards | True |
| number of PPO epochs | 3 |
| number of actors | 128 |
| unroll length | 128 |
| discount factor $\gamma$ | 0.99 |
| GAE $\lambda$ | 0.95 |
| entropy coefficient | 0.001 |
| value loss coefficient | 0.5 |
| policy gradient clip range | [-1.1, 1.1] |
| Pixel AMA uncertainty budget $\lambda$ | 1 |
| Pixel AMA uncertainty coefficient $\eta$ | 2 |
| IDF AMA uncertainty budget $\lambda$ | 1 |
| IDF AMA uncertainty coefficient $\eta$ | 1 |

*Table 4.* Retro game policy hyperparameters

*Figure 10.* Curiosity forward prediction architecture for the minigrid experiments. For the MSE baseline the variance predictions are not used and loss is computed via a standard MSE.

| Layer Type | Filters | Kernel Size | Stride |
|---|---|---|---|
| Conv2d | 32 | 8 | 3 |
| Conv2d | 64 | 8 | 2 |
| Conv2d | 64 | 4 | 2 |
| Dense (512 Units) | N/A | N/A | N/A |
| Conv2d Transpose | 64 | 4 | 2 |
| Conv2d Transpose | 32 | 8 | 2 |
| Conv2d Transpose | 4 | 8 | 2 |

*Table 5.* Retro game forward prediction hyperparameters

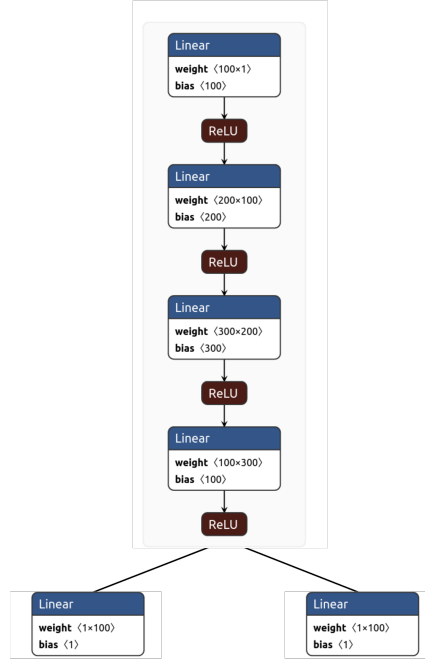| Hyperparameter | Value |
|---|---|
| epistemic learning rate | 0.0001 |
| epistemic batch size | 32 |
| aleatoric learning rate | 0.001 |
| aleatoric batch size | 1000 |
| aleatoric uncertainty budget $\lambda$ | 1 |
| aleatoric uncertainty coefficient $\eta$ | 1 |
| $\epsilon$ greedy $\epsilon$ | 0.1 |

*Table 6.* Bandit hyperparameters.



*Figure 11.* AMA prediction network for bandit tasks.

## C.5. Bandit

We performed 3 repeats to produce the standard error regions show in the graph. Learning rate was tuned by hand, observing how well the network performed in making predictions as the bandit sampled different regions of the environment. The intrinsic reward method for the epistemic bandit is based on (Pathak et al., 2019; Lakshminarayanan et al., 2017). We use an action value based bandit algorithm with $\epsilon$-greedy exploration (Sutton & Barto, 2018)(p. 31).

## C.6. Different Noise Distribution for Retro Games

For the CIFAR noisy TV we tiled a random CIFAR image (from the training set) for each frame observed on the noisy TV. This required around around 2.5 tiles to fill the $84 \times 84$ pixels of the retro game frames. An example frame can be seen in Figure 12.

## C.7. Bank Heist Plotting

As briefly mentioned in the main text, the lines plotted for Bank Heist have approximate x points as the exact frames were not recorded directly with the average number of pixels covered in an episode. However, the logging step was recorded with the average pixel coverage. In the main text we use the frame count from the nearest recorded step to the step used for pixel coverage. To show the trends of these graphs show similar results we plot pixel coverage against training steps below. Note
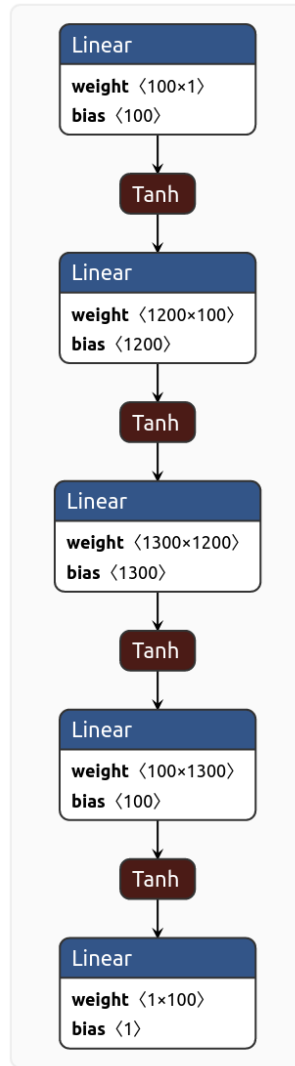
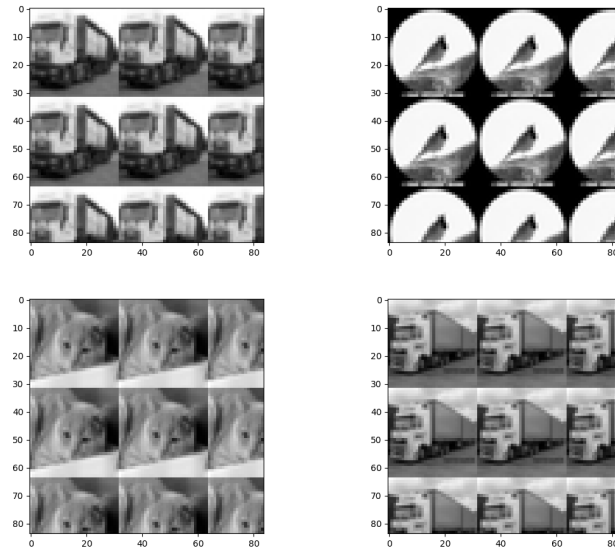*Figure 12.* Epistemic prediction network for the bandit task.

*Figure 13.* Four examples of frames that an agent might see when interacting with the CIFAR noisy TV. This is a complex noise distribution picked to test the limits of the heteroscedastic aleatoric uncertainty estimation.

for all methods besides pixel AMA the Bank Heist plots are smoothed with following code snippet[4] so that the crowded plot is readable.

### C.8. Eventual Decrease in Exploration During Mario Training

In initial experiments we noticed the AMA agent lost motivation to explore its environment after reaching it peak extrinsic reward. To ensure that this was not an inherent problem with AMA, we ran the other curiosity methods for further frames and found similar eventual decreases in extrinsic reward (Figure 14). Presumably the cause of this is that once a significant portion of the environment has been explored the agent is no longer motivated to return to those regions (as prediction error decreases)—this issue has been noted by previous authors (Pathak et al., 2017).

### C.9. Potential Negative Social Impacts

The work presented is very far from any real world deployment. If it were to be deployed in any real world context then extensive testing would need to be done to understand how the curiosity agents would behave in novel environments as erratic behaviours could be dangerous in, for example, a robotic control context. The AMA objective is overarching (like other curiosity methods) and so care should be taken if deploying in something like a recommender system as the agent could find certain behaviours intrinsically rewarding that you might not have intended it to (like the noisy TV problem). Lastly, although the AMA system contains notions of uncertainty quantification, that does not mean it is able to completely understand the limits of its predictions and so one should not be overconfident in its abilities to do so.

### C.10. Hardware

The experiments were performed on three different machines depending on their availability: A 32 core CPU with one GeForce GTX TITAN X, a 12 core CPU with two GeForce GTX TITAN Xs and one 8 core CPU with two GeForce 2080Ti GPUs.
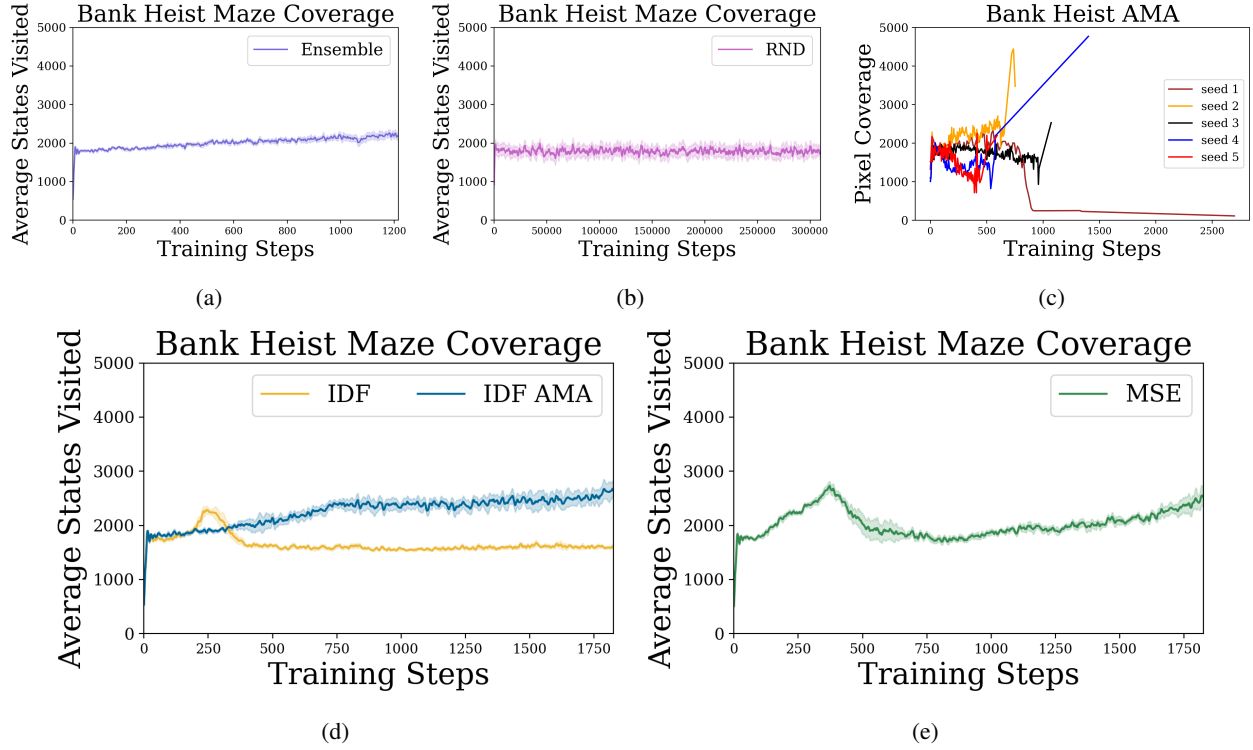
---

[4] https://stackoverflow.com/a/11352216

*Figure 14.* Trends of Bank Heist exploration plotted exactly against logging steps, showing similar trends to the approximate x-axis used in the main text. The training steps are not comparable between different methods. The pixel MSE results cannot be plotted together meaningfully against training steps as slightly different logging was used with different seeds.
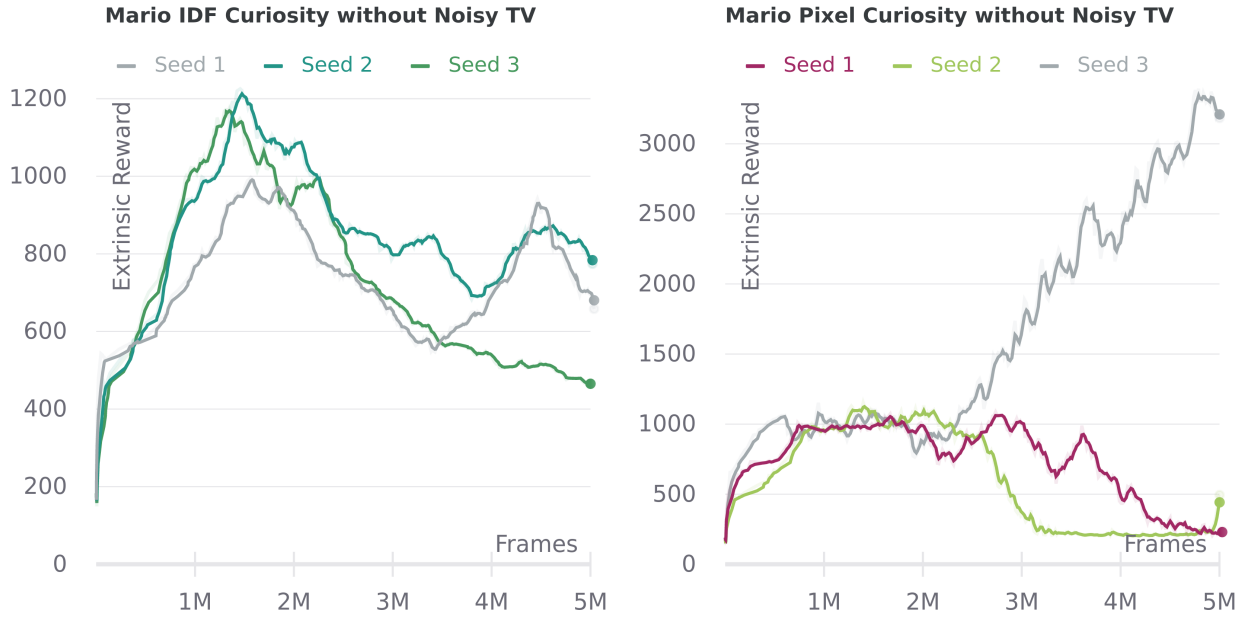


*Figure 15.* Repeats ran for more frames on Mario for MSE Pixel curiosity and MSE IDF curiosity. 3/3 of seeds see a decrease in performance for IDF curiosity and 2/3 in the pixel curiosity case.

We list times here on the 2080Ti machine, the other machines were as much as $2\times$ slower. The minigrid experiments took around 40 minutes per run, the Space Invader experiments took around 12 hours per run, the Mario experiments took around 1 hour and 20 minutes per run.

### C.11. Licensing

The repository from (Willems, 2020) has an MIT license. The code used from (Chevalier-Boisvert et al., 2018) has an Apache License 2.0. (Raileanu & Rocktäschel, 2020) has a creative commons license. Besides those listed we are not aware of any further code licensing. We adapted a copyright free rat silhouette image for the bandit figure [5].

### C.12. Further Code Acknowledgements

Although we did not directly use their code we would like to acknowledge the following open source contributions that provided a useful reference when implementing Kendall and Gal's (Kendall & Gal, 2017) aleatoric uncertainty estimation algorithms:

```
https://github.com/ShellingFord221/My-implementation-of-What-Uncertainties-
Do-We-Need-in-Bayesian-Deep-Learning-for-Computer-Vision
```

```
https://github.com/pmorerio/dl-uncertainty
```

```
https://github.com/hmi88/what
```

---

[5]`https://pixabay.com/vectors/rat-rodent-silhouette-gold-chinese-5184465/`