
World Discovery Models

Mohammad Gheshlaghi Azar^{*1} Bilal Piot^{*1} Bernardo Avila Pires^{*1} Jean-Bastien Grill² Florent Althé²
Rémi Munos²

Abstract

As humans we are driven by a strong desire for seeking novelty in our world. Also upon observing a novel pattern we are capable of refining our understanding of the world based on the new information—humans can *discover* their world. The outstanding ability of the human mind for discovery has led to many breakthroughs in science, art and technology. Here we investigate the possibility of building an agent capable of discovering its world using the modern AI technology. In particular we introduce NDIGO, Neural Differential Information Gain Optimisation, a self-supervised discovery model that aims at seeking new information to construct a global view of its world from partial and noisy observations. Our experiments on some controlled 2-D navigation tasks show that NDIGO outperforms state-of-the-art information-seeking methods in terms of the quality of the learned representation. The improvement in performance is particularly significant in the presence of white or structured noise where other information-seeking methods follow the noise instead of discovering their world.

1. Introduction

Modern AI has been remarkably successful in solving complex decision-making problems such as GO (Silver et al., 2016; 2017), simulated control tasks (Schulman et al., 2015), robotics (Levine et al., 2016), poker (Moravčík et al., 2017) and Atari games (Mnih et al., 2015; Hessel et al., 2018). Despite these successes the agents developed by those methods are specialists: they perform extremely well at the tasks they were trained on but are not very successful at generalising their task-dependent skills in the form of a general domain understanding. Also, the success of the existing AI agents often depends strongly on the availability of external

feedback from their world in the form of reward signals or labelled data, for which some level of supervision is required. This is in contrast to the human mind, which is a general and self-supervised learning system that *discovers* the world around it even when no external reinforcement is available. Discovery is the ability to obtain knowledge of a phenomenon for the first time (Merriam-Webster, 2004). As discovery entails the process of learning of and about new things, it is an integral part of what makes humans capable of understanding their world in a task-independent and self-supervised fashion.

The underlying process of discovery in humans is complex and multifaceted (Hohwy, 2013). However one can identify two main mechanisms for discovery (Clark, 2017). The first mechanism is **active information seeking**. One of the primary behaviours of humans is their attraction to novelty (new information) in their world (Litman, 2005; Kidd & Hayden, 2015). The human mind is very good at distinguishing between the *novel* and the *known*, and this ability is partially due to the extensive internal reward mechanisms of *surprise*, *curiosity* and *excitement* (Schmidhuber, 2009). The second mechanism is **building a statistical world model**. Within cognitive neuroscience, the theory of statistical predictive mind states that the brain, like scientists, constructs and maintains a set of hypotheses over its representation of the world (Friston et al., 2014). Upon perceiving a novelty, our brain has the ability to validate the existing hypothesis, reinforce the ones which are compatible with the new observation and discard the incompatible ones. This self-supervised process of hypothesis building is essentially how humans consolidate their ever-growing knowledge in the form of an accurate and global model.

Inspired by these inputs from cognitive neuroscience, information-seeking algorithms have received significant attention to improve the exploration capability of artificial learning agents (Schmidhuber, 1991b; Houthoofd et al., 2016; Achiam & Sastry, 2017; Pathak et al., 2017; Burda et al., 2018). However, the scope of the existing information-seeking algorithms is often limited to the case of fully observable and deterministic environments. One of the problems with the existing novelty-seeking algorithms is that agents trained by these methods tend to become attracted to random patterns in their world and stop exploring upon en-

^{*}Equal contribution ¹DeepMind, London, UK ²DeepMind Paris, France. Correspondence to: Mohammad Gheshlaghi Azar <mazar@google.com>.

countering them, despite the fact that these random patterns contain no actual *information* on the world (Burda et al., 2018). Moreover, the performance of existing agents are often evaluated based on their ability to solve a reinforcement learning (RL) task with extrinsic reward, and not on the quality of the learned world representation, which is the actual goal of discovery. Thus, it is not clear whether the existing algorithms are capable of using the novel information to discover their world. Therefore, the problem of discovery in the general case of partially observable and stochastic environments remains open.

The main contribution of this paper is to develop a practical and end-to-end algorithm for discovery in stochastic and partially observable worlds using modern AI technology. We achieve this goal by designing a simple yet effective algorithm called NDIGO, **N**eural **D**ifferential **I**nformation **G**ain **O**ptimisation, for information seeking designed specifically for stochastic partially observable domains. NDIGO identifies novelty by measuring the increment of information provided by a new observation in predicting the future observations, compared to a baseline prediction for which this observation is withheld. We show that this measure can be estimated using the difference of prediction losses of two estimators, one of which can access the complete set of observations while the other does not receive the latest observation. We then use this measure of novelty as the intrinsic reward to train the policy using a state of the art reinforcement learning algorithm (Kapturowski et al., 2019). One of the key features of NDIGO is its robustness to noise, as the process of subtracting prediction losses cancels out errors that the algorithm cannot improve on. Moreover, NDIGO is well-suited for discovery in partially observable domains as the measure of novelty in NDIGO drives the agent to the unobserved areas of the world where new information can be gained from the observations. Our experiments show that NDIGO produces a robust performance in the presence of noise in partial observable environments: NDIGO not only finds true novelty without being distracted by the noise, but it also incorporates this information into its world representation without forgetting previous observation.

2. Related Work

It has been argued for decades in developmental psychology (White, 1959; Deci & Ryan, 1985; Csikszentmihalyi & Csikszentmihalyi, 1992), neuroscience (Dayan & Balleine, 2002; Kakade & Dayan, 2002; Horvitz, 2000) and machine learning (Oudeyer & Kaplan, 2008; Gottlieb et al., 2013; Schmidhuber, 1991a) that an agent maximising a simple intrinsic reward based on patterns that are both novel and learnable could exhibit essential aspects of intelligence such as autonomous development (Oudeyer & Smith, 2016).

More specifically, in his survey on the theory of creativity

and intrinsic motivation, Schmidhuber (2010) explains how to build the agent that could discover and understand in a self-supervised way its environment. He establishes that 4 crucial components are necessary: **i)** a world model (Ha & Schmidhuber, 2018) that encodes what is currently known. It can be a working memory component such as a Long Short Term Memory network (LSTM, Hochreiter & Schmidhuber, 1997) or a Gated Recurrent Unit network (GRU, Cho et al., 2014). **ii)** a learning algorithm that improves the world model. For instance, Guo et al. (2018) have shown that a GRU trained with a Contrastive Prediction Coding (CPC, Oord et al., 2018) loss on future frames could learn a representation of the agent’s current and past position and orientation, as well as position of objects in the environment. **iii)** An intrinsic reward generator based on the world model that produces rewards for patterns that are both novel and learnable. Different types of intrinsic rewards can be used, such as the world model’s prediction error (Stadie et al., 2015; Pathak et al., 2017), improvement of the model’s prediction error, also known as prediction gain (Achiam & Sastry, 2017; Schmidhuber, 1991a; Lopes et al., 2012), and finally information gain (Shyam et al., 2018; Itti & Baldi, 2009; Little & Sommer, 2013; Frank et al., 2014; Houthoofd et al., 2016). **iv)** the last component is an RL algorithm that finds an optimal policy with respect to the intrinsic rewards.

Recently, several implementations of intrinsically motivated agents have been attempted using modern AI technology. Most of them used the concept of prediction error as an intrinsic reward (Stadie et al., 2015; Pathak et al., 2017; Burda et al., 2018; Haber et al., 2018). However, it has been argued that agents optimising the prediction error are susceptible to being attracted to white noise (Oudeyer et al., 2007) and therefore should be avoided. To solve the white-noise problem, different types of random or learned projections (Burda et al., 2018) of the original image into a smaller feature space less susceptible to white-noise are considered. Other implementations rely on approximations of the concept of information gain (Houthoofd et al., 2016; Achiam & Sastry, 2017) via a variational lower bound argument. Indeed, as they are trying to train a probabilistic model over the set of possible dynamics, the computation of the posterior of that distribution is intractable (Houthoofd et al., 2016). Finally, models based on prediction gain are fundamentally harder to train compared to prediction error (Achiam & Sastry, 2017; Lopes et al., 2012; Pathak et al., 2017), and are less principled than information gain (Schmidhuber, 2010).

3. Setting

We consider a partially observable environment where an agent is shown an observation o_t at time t , then selects an action a_t which generates a new observation o_{t+1} at the next time step. We assume observations o_t are gen-

erated by an underlying process x_t following Markov dynamics, i.e. $x_{t+1} \sim P(\cdot|x_t, a_t)$, where P is the dynamics of the underlying process. Although we do not explicitly use the corresponding terminology, this process can be formalised in terms of Partially Observable Markov Decision Processes (POMDPs; Lovejoy, 1991; Cassandra, 1998).

The future observation o_{t+1} in a POMDP can also be seen as the output of a stochastic mapping with input the current history. Indeed, at any given time t , let the current history h_t be all past actions and observations $h_t \stackrel{\text{def}}{=} (o_0, a_0, o_1, a_1, \dots, a_{t-1}, o_t)$. Then we define $\mathbb{P}(\cdot|h_t, a_t)$ the probability distribution of o_t knowing the history and the action a_t . One can generalise this notion for k -step prediction: for any integers $t \geq 0$ and $k \geq 1$, let us denote by $t : t+k$ the integer interval $\{t, \dots, t+k-1\}$, and let $\mathcal{A}_{t:t+k} \stackrel{\text{def}}{=} (a_t, \dots, a_{t+k-1})$ and $\mathcal{O}_{t:t+k} \stackrel{\text{def}}{=} (o_t, \dots, o_{t+k-1})$ be the sequence of actions and observations from time t up to time $t+k-1$, respectively. Then o_{t+k} can be seen as a sample drawn from the probability distribution $\mathbb{P}(\cdot|h_t, \mathcal{A}_{t:t+k})$, which is the k -step open-loop prediction model of the observation o_{t+k} . We also use the short-hand notation $\mathbb{P}_{t+k|t} = \mathbb{P}(\cdot|h_t, \mathcal{A}_{t:t+k})$ as the probability distribution of o_{t+k} given the history h_t and the sequence of actions $\mathcal{A}_{t:t+k}$.

4. Learning the World Model

The world model should capture what the agent currently knows about the world so that he could make predictions based on what it knows. We thus build a model of the world by predicting future observations given the past (see e. g., Schmidhuber, 1991a; Guo et al., 2018). More precisely, we build an internal representation b_t by making predictions of futures frames o_{t+k} conditioned on a sequence of actions $\mathcal{A}_{t:t+k}$ and given the past h_t . This is similar to the approach of Predictive State Representations (Littman et al., 2002), from which we know that if the learnt representation b_t is able to predict the probability of any future observation conditioned on any sequence of actions and history, then this representation b_t contains all information about the belief state (i.e., distribution over the ground truth state x_t).

4.1. Architecture

We propose to learn the world model by using a recurrent neural network (RNN) f_θ fed with the concatenation of observation features z_t and the action a_t (encoded as a one-hot vector). The observation features z_t are obtained by applying a convolutional neural network (CNN) f_ϕ to the observation o_t . The RNN is a Gated Recurrent Unit (GRU) and the internal representation is the hidden state of the GRU, that is, $b_t = f_\theta(z_t, a_{t-1}, b_{t-1})$, as shown in Figure 1. We initialise this GRU by setting its hidden state

to the null vector 0, and using $b_0 = f_\theta(z_0, a, 0)$ where a is a fixed, arbitrary action and z_0 are the features corresponding to the original observation o_0 . We train this representation b_t with some future-frame prediction tasks conditioned on sequences of actions and the representation b_t . These frame prediction tasks consist in estimating the probability distribution, for various $K \geq k \geq 1$ (with $K \in \mathbb{N}^*$ to be specified later), of future observation o_{t+k} conditioned on the internal representation b_t and the sequence of actions $\mathcal{A}_{t:t+k}$. We denote these estimates by $\hat{p}_{t+k|t}(\cdot|b_t, \mathcal{A}_{t:t+k})$ or simply by $\hat{p}_{t+k|t}$ for conciseness and when no confusion is possible. As the notation suggests, we will use $\hat{p}_{t+k|t}$ as an estimate of $\mathbb{P}_{t+k|t}$. The neural architecture consists in K different neural nets $\{f_{\psi_k}\}_{k=1}^K$. Each neural net f_{ψ_k} receives as input the concatenation of the internal representation b_t and the sequence of actions $\mathcal{A}_{t:t+k}$, and outputs the distributions over observations: $\hat{p}_{t+k|t} = f_{\psi_k}(b_t, \mathcal{A}_{t:t+k})$. For a fixed $t \geq 0$ and a fixed $K \geq k \geq 1$, the loss function $L(o_{t+k}, \hat{p}_{t+k|t})$ at time step $t+k-1$ associated with the network f_{ψ_k} is a cross entropy loss: $L(o_{t+k}, \hat{p}_{t+k|t}) = -\ln(\hat{p}_{t+k|t}(o_{t+k}))$. We finally define for any given sequence of actions and observations the *representation loss* function L_{repr} as the sum of these cross entropy losses: $L_{\text{repr}} = \sum_{t \geq 0, K \geq k \geq 1} L(o_{t+k}, \hat{p}_{t+k|t})$. +

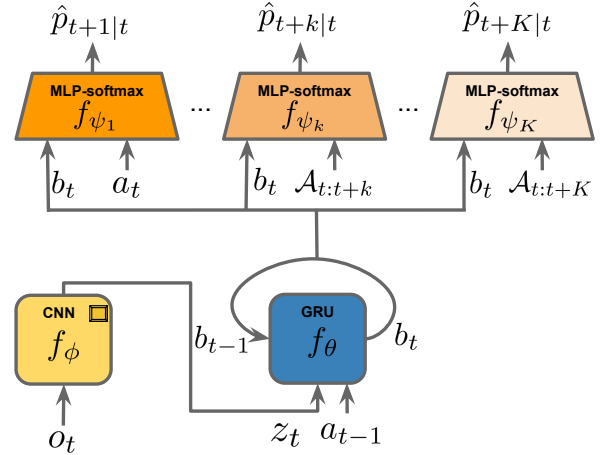


Figure 1. World Model: a CNN and a GRU encode the history h_t into an internal representation b_t . Then, K frame predictions tasks are trained in order to shape the representation b_t .

4.2. Evaluation of the learnt representation

In the POMDP setting, the real state x_t represents all there is to know about the world at time t . By constructing a belief state, which is a distribution $P_b(\cdot|h_t)$ over the possible states conditioned on the history h_t , the agent can assess its uncertainty about the real state x_t given the history h_t . Therefore, in order to assess the quality of the learnt representation b_t , we use the glass-box approach described in Figure 12 to build a belief state of the world. It consists

simply in training a neural network f_τ fed by the internal representation b_t to predict a distribution $\hat{p}_b(\cdot|b_t)$ over the possible real state x_t . This kind of approach is only possible in artificial or controlled environments where the real state is available to the experimenter but yet not given to the agent. We also make sure that no gradient from f_τ is being back-propagated to the internal representation b_t such that the evaluation does not influence the learning of the representation and the behaviour of the agent. For a fixed $t \geq 0$, the loss used to train f_τ is a cross entropy loss (For a more detailed description of the approach see Guo et al., 2018): $L_{\text{discovery}}(x_t, \hat{p}_b(\cdot|b_t)) \stackrel{\text{def}}{=} -\ln(\hat{p}_b(x_t|b_t))$. We call this loss *discovery loss*, and use it as a measure of how much information about the whole world the agent is able to encode in its internal representation b_t , i.e., how much of the world has been discovered by the agent.

5. NDIGO Agent

Our NDIGO agent is a discovery agent that learns to seek new information in its environment and then incorporate this information into a world representation. Inspired by the intrinsic motivation literature (Schmidhuber, 2010), the NDIGO agent achieves this information-seeking behaviour as a result of optimising an intrinsic reward. Therefore, the agent’s exploratory skills depend critically on designing an appropriate reward signal that encourages discovering the world. Ideally, we want this reward signal to be high when the agent gets an observation containing new information about the real state x_t . As we cannot access x_t at training time, we rely on the accuracy of our future observations predictions to estimate the information we have about x_t .

Intuitively, for a fixed horizon $H \in \mathbb{N}^*$, the *prediction error loss* $L(o_{t+H}, \hat{p}_{t+H|t}) = -\ln(\hat{p}_{t+H|t}(o_{t+H}))$ is a good measure on how much information b_t is lacking about the future observation o_{t+H} . The higher the loss, the more uncertain the agent is about the future observation o_{t+H} so the less information it has about this observation. Therefore, one could define an intrinsic reward directly as the prediction error loss, thus encouraging the agent to move towards states for which it is the less capable of predicting future observations. The hope is that the less information we have in a certain belief state, the easier it is to gain new information. Although this approach may have good results in deterministic environments, it is however not suitable in certain stochastic environments. For instance, consider the extreme case in which the agent is offered to observe white noise such as a TV displaying static. An agent motivated with prediction error loss would continually receive a high intrinsic reward simply by staying in front of this TV, as it cannot improve its predictions of future observations, and would effectively remain fascinated by this noise.

5.1. The NDIGO intrinsic reward

The reason why the naive prediction error reward fails in such a simple example is that the agent identifies that a lot of information is lacking, but does not acknowledge that no progress is made towards acquiring this lacking information. To overcome this issue, we introduce the NDIGO reward, for a fixed $K \geq H \geq 1$, as follows:

$$r_{t+H-1}^{\text{NDIGO}} \stackrel{\text{def}}{=} L(o_{t+H}, \hat{p}_{t+H|t-1}) - L(o_{t+H}, \hat{p}_{t+H|t}), \quad (1)$$

where o_{t+H} represents the future observation considered and H is the horizon of NDIGO. The two terms in the right-hand side of Equation (1) measure how much information the agent lacks about the future observation o_{t+H} knowing all past observations prior to o_t with o_t either excluded (left term) or included (right term). Intuitively, we take the difference between the information we have at time t with the information we have at time $t-1$. This way we get an estimate of how much information the agent gained about o_{t+H} by observing o_t . Note that the reward r_{t+H-1}^{NDIGO} is attributed at time $t+H-1$ in order to make it dependent on h_{t+H-1} and a_{t+H-1} only (and not on the policy), once the prediction model \hat{p} has been learnt. If the reward had been assigned at time t instead (time of prediction) it would have depended on the policy used to generate the action sequence $\mathcal{A}_{t:t+H-1}$, which would have violated the Markovian assumption required to train the RL algorithm. Coming back to our broken TV example, the white noise in o_t does not help in predicting the future observation o_{t+H} . The NDIGO reward is then the difference of two large terms of similar amplitude, leading to a small reward: while acknowledging that a lot of information is missing (large prediction error loss) NDIGO also realises that no more of it can be extracted (small difference of prediction error loss). Our experiments show that using NDIGO allows the agent to avoid being stuck in the presence of noise, as presented in Section 6, thus confirming these theoretical considerations.

5.2. Algorithm

Given the intrinsic reward r_{t+H-1}^{NDIGO} , we use the state-of-the-art RL algorithm R2D2 (Kapturowski et al., 2019) to optimise the policy. The NDIGO agent interacts with its world using the NDIGO policy to obtain new observation o_{t+k} , which is used to train the world model by minimising the future prediction loss $L_{t+k|t} = L(o_{t+k}, \hat{p}_{t+k|t})$. The losses $L_{t+k|t}$ are then used to obtain the intrinsic reward at the next time step, and the process is then repeated. An in-depth description of the complete NDIGO algorithm can be found in Appendix B.5.

5.3. Relation to information gain

Information gain has been widely used as the novelty signal in the literature (Houthoofd et al., 2016; Little

& Sommer, 2013). A very broad definition of the information gain (Schmidhuber, 2010) is the distance (or divergence) between distributions on any random event of interest ω before and after a new sequence of observations. Choosing the random event to be the future observations or actions and the divergence to be the Kullback-Leiber divergence then the k -step predictive information gain $IG(o_{t+k}, \mathcal{O}_{t:t+k} | h_t, \mathcal{A}_{t:t+k})$ of the future event o_{t+k} with respect to the sequence of observations $\mathcal{O}_{t:t+k}$ is defined as: $IG(o_{t+k}, \mathcal{O}_{t:t+k} | h_t, \mathcal{A}_{t:t+k}) \stackrel{\text{def}}{=} \text{KL}(\mathbb{P}_{t+k|t+k-1} || \mathbb{P}_{t+k|t-1})$, and measures how much information can be gained about the future observation o_{t+k} from the sequence of past observations $\mathcal{O}_{t:t+k}$ given the whole history h_t up to time step t and the sequence of actions $\mathcal{A}_{t:t+k}$ from t up to $t + H - 1$. In the case of $k = 1$ we recover the 1-step information gain on the next observation o_{t+1} due to o_t . We also use the following short-hand notation for the information gain $IG_{t+k|t} = IG(o_{t+k}, \mathcal{O}_{t:t+k} | h_t, \mathcal{A}_{t:t+k})$ for every $k \geq 1$ and $t \geq 0$. Also by convention we define $IG_{t|t} = 0$.

We now show that the NDIGO intrinsic reward r_{t+H-1}^{NDIGO} can be expressed as the difference of information gain due to $\mathcal{O}_{t:t+H}$ and $\mathcal{O}_{t+1:t+H}$. For a given horizon $H \geq 1$ and $t \geq 0$, the intrinsic reward for time step $t + H - 1$ is:

$$r_{t+H-1}^{\text{NDIGO}} \stackrel{\text{def}}{=} L(o_{t+H}, \hat{p}_{t+H|t-1}) - L(o_{t+H}, \hat{p}_{t+H|t}) \\ = \ln \left(\frac{\hat{p}_{t+H|t}(o_{t+H})}{\hat{p}_{t+H|t-1}(o_{t+H})} \right).$$

Given that $\hat{p}_{t+H|t}$ and $\hat{p}_{t+H|t-1}$ are respectively an estimate of $\mathbb{P}_{t+H|t}$ and $\mathbb{P}_{t+H|t-1}$, and based on the fact that these estimates become more accurate as the number of samples increases, we have:

$$\mathbb{E} [r_{t+H-1}^{\text{NDIGO}}] = \mathbb{E}_{o_{t+H} \sim \mathbb{P}_{t+H|t+H-1}} \left[\ln \left(\frac{\hat{p}_{t+H|t}(o_{t+H})}{\hat{p}_{t+H|t-1}(o_{t+H})} \right) \right] \\ \approx \mathbb{E}_{o_{t+H} \sim \mathbb{P}_{t+H|t+H-1}} \left[\ln \left(\frac{\mathbb{P}_{t+H|t}(o_{t+H})}{\mathbb{P}_{t+H|t-1}(o_{t+H})} \right) \right] \\ = \text{KL}(\mathbb{P}_{t+H|t+H-1} || \mathbb{P}_{t+H|t-1}) \\ - \text{KL}(\mathbb{P}_{t+H|t+H-1} || \mathbb{P}_{t+H|t}) \\ = IG_{t+H|t} - IG_{t+H|t-1}. \quad (2)$$

The first term $IG_{t+H|t}$ in Equation (2) measures how much information can be gained about o_{t+H} from the sequence of past observations $\mathcal{O}_{t:t+H}$ whereas the second term $IG_{t+H|t-1}$ measures how much information can be gained about o_{t+H} from the sequence of past observations $\mathcal{O}_{t+1:t+H}$. Therefore, as $\mathcal{O}_{t+1:t+H} = \mathcal{O}_{t:t+H} \setminus \{o_t\}$, the expected value of the NDIGO reward at step $t + H - 1$ is equal to the amount of additional information that can be gained by the observation o_t when trying to predict o_{t+H} .

6. Experiments

We evaluate the performance of NDIGO qualitatively and quantitatively on five experiments, where we demonstrate different aspects of discovery with NDIGO. In all experiments there are some hidden objects which the agent seeks to discover. However the underlying dynamics of the objects are different. In the simplest case, the location of objects only changes at the beginning of every episode, whereas in the most complex the objects are changing their locations throughout the episode according to some random walk strategy. We investigate (i) whether the agent can efficiently search for novelty, i.e., finding the location of objects; (ii) whether the agent can encode the information of object location in its representation of the world such that the discovery loss of predicting the objects is as small as possible.

6.1. Baselines

We compare our algorithm NDIGO- H , with H being the horizon and taking values in $\{1, 2, 4\}$, to different information seeking and exploration baselines considered to be state of the art in the intrinsic motivation literature. **Prediction Error (PE)** (Haber et al., 2018; Achiam & Sastry, 2017): The PE model uses the same architecture and the same losses than NDIGO. The only difference is that the intrinsic reward is the predictor error: $r_t^{\text{PE}} = L(\hat{p}_{t+1|t}, o_{t+1})$. **Prediction Gain (PG)** (Achiam & Sastry, 2017; Ostrovski et al., 2017): Our version of PG uses the same architecture and the same losses than NDIGO. In addition, at every $n = 2$ learner steps we save a copy of the prediction network into a fixed target network. The intrinsic reward is the difference in prediction error, between the up-to-date network and the target network predictions: $r_t^{\text{PG}} = L(\hat{p}_{t+1|t}^{\text{target}}, o_{t+1}) - L(\hat{p}_{t+1|t}, o_{t+1})$, where $\hat{p}_{t+1|t}^{\text{target}}$ is the distribution computed with the weights of the fixed target network. **Intrinsic Curiosity Module (ICM)** (Pathak et al., 2017; Burda et al., 2018): The method consists in training the internal representation b_t to be less sensitive to noise using a self-supervised inverse dynamics model. Then a forward model is used to predict the future internal representation \hat{b}_{t+1} from the actual representation b_t and the action a_t (more details on this model are in Appendix D). The intrinsic reward $r_t^{\text{FPE}} = \|\hat{b}_{t+1} - b_{t+1}\|_2^2$.

6.2. Test environments

The 5 rooms environment. The 5 rooms environment (see Figure 2) is a local-view 2D environment composed of 5 rooms implemented using the `pycolab` library¹. In `pycolab`, the environment is composed of cells that contain features such as walls, objects or agents. In the 5 rooms

¹<https://github.com/deepmind/pycolab>

environment, there is one central 5×5 room and four peripheral rooms (composed of 48 cells) that we will refer to as upper, lower, left and right rooms. Each of the four peripheral rooms may contain different types of “objects” that occupy a cell exclusively. At every episode, the agent starts in the middle of the central room and the starting position of each object is randomised. The objects may or may not move, but as a general rule in any episode they never leave the room they started in. Finally, we only place objects in the peripheral rooms, and in each room there is never more than one object.

The maze environment. The `maze` environment (see Figure 3) is also a `pycolab` local-view 2D environment. It is set-up as a maze composed of six different rooms connected by corridors. The agent starts at a fixed position in the environment in an otherwise empty room 0; rooms are numbered from 0 to 5 based on the order in which they can be reached, i.e. the agent cannot reach room number 3 without going through rooms 1 and 2 in this order. A `white noise` object is always present in room 1, and there is single `fixed` in rooms 2, 3 and 4. Room 5 contains a special `movable`, which should attract the agent even when the environment is completely learned.

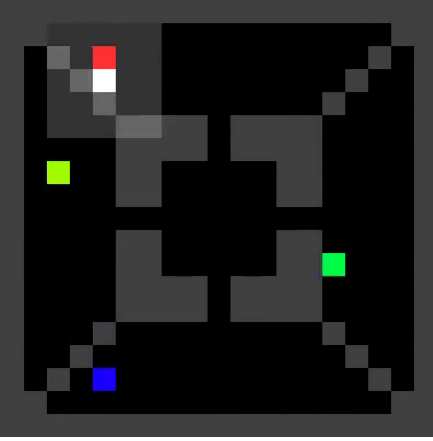


Figure 2. The 5 rooms environment: in this instance, we can see in white the agent, 4 fixed objects in each of the 4 peripheral rooms and in grey the impenetrable walls. The shaded area around the agent represents its 5×5 region-cell local view.

Objects. We consider five different types of objects: `fixed`, `bouncing`, `Brownian`, `white noise` and `movable`. `fixed` objects are fixed during episodes, but change position from episode to episode. They provide information gain about their position when it is not already encoded in the agent’s representation. `bouncing` objects bounce in a straight line from wall to wall inside a room. In addition to providing information gain similar to `fixed` objects, they allow us to test the capacity of the representation to encode predictable object after the object is no longer in

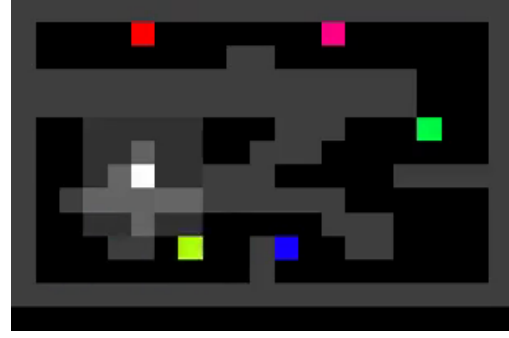


Figure 3. The `maze` environment: in this instance, we can see in white the agent, 4 `fixed` objects in blue, green, pink and red. `white noise` is the closest object to the agent location also in green.

the agent’s view. `Brownian` objects follow a Brownian motion within a room, by moving uniformly at random in one of the four directions. `white noise` objects change location instantly to any position inside the same room, uniformly at random, at each time step, and are therefore unpredictable. Finally, `movable` objects do not move by themselves, but the agent can cause them to move to a random location by attempting to move into their cells. Interacting with these objects allows more information gain to be generated.

Agent’s observations and actions. The observation o_t at time t consists in a concatenation of images (called channels) of 25 pixels representing the different features of the 5×5 local view of the agent. This can be represented by multidimensional array $(5, 5, c)$ where c is the number of channels. The first channel represents the walls in the local view: 1 indicates the presence of a wall and 0 the absence of a wall. Then, each of the remaining channels represents the position of an object with a one-hot array if the object is present in the local view or with a null array otherwise. The possible actions a_t are stay, up, down, right, left and are encoded with a one-hot vector of size 5.

6.3. Performance evaluation

The agent’s performance is measured by its capacity to estimate the underlying state of the world from its internal representation (discovery loss, see Section 4.2). In `pycolab`, it is possible to compute a discovery loss for each aspect of the world state (location of each object for instance). So that it is easy to understand which aspects of the world the agent can understand and keep in its internal representation. Once again we stress the fact that no gradient is back-propagated from that evaluation procedure to the internal representation. In addition, we provide other statistics such as average values of first-visit time and visit counts of a given object to describe the behavior of the agent. The first-visit time is the number of episode time steps the agent

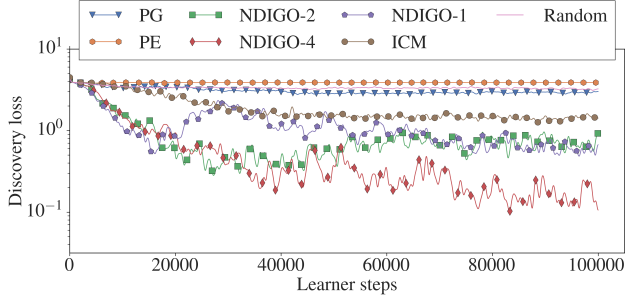


Figure 4. Experiment 1: Average discovery loss of the fixed object. The results are averaged over 10 seeds.

needs before first observing a given object; the visit count is the total number of time steps where the agent observes the object. Finally, we also provide more qualitative results with videos of the agent discovering the worlds (see <https://www.youtube.com/channel/UC5OPHK7pvsZE-jVclZMvbmQ>).

6.4. Experimental results

In this section we evaluate the performance of NDIGO on some controlled navigation task (for the implementation details and the specification of the prediction and policy networks and the training algorithms see Appendix B).

Experiment 1. We evaluate the discovery skills of NDIGO by testing how effectively it can ignore the white noise, from which there is nothing to learn, and discover the location of the fixed object. Here, we use a 5 rooms setting with a fixed object in the upper room, and a white noise object in the lower room.

	Visit count		First visit time	
	fixed	w. noise	fixed	w. noise
Random	14.1 \pm 14.3	24.6 \pm 12.6	339.0 \pm 40.5	225.6 \pm 50.4
PE	0.1 \pm 0.2	158.3 \pm 3.7	392.6 \pm 18.1	15.5 \pm 4.0
PG	27.3 \pm 22.0	22.5 \pm 10.3	306.4 \pm 49.4	233.7 \pm 56.6
ICM	144.8 \pm 37.2	23.8 \pm 12.4	132.4 \pm 41.2	238.3 \pm 55.0
NDIGO-1	120.9 \pm 43.4	19.1 \pm 9.3	78.4 \pm 28.5	279.4 \pm 42.9
NDIGO-2	154.0 \pm 45.5	7.4 \pm 6.7	112.6 \pm 46.2	345.8 \pm 36.5
NDIGO-4	300.4 \pm 22.2	1.4 \pm 1.2	40.8 \pm 9.7	330.7 \pm 47.4

Table 1. Experiment 1: Average values of the visit counts and first visit time of the trained agent for the fixed and white noise objects in one episode.

We report in Figure 4 the learning curves for the discovery loss of the fixed object. This result shows the quality of the learned representation in terms of encoding the location of fixed object. We observe that the long-horizon variant of NDIGO (NDIGO-4) outperforms the best baseline (ICM) by more than an order of magnitude. Also the asymptotic performance of NDIGO-4 is significantly better than NDIGO-1 and NDIGO-2.

In Table 1 we also report the average value and standard deviation of visit count and first visit time of the trained agents for the fixed object and the white noise object in an episode². We observe that different variants of NDIGO are driven towards the fixed object and manage to find it faster than the baselines while avoiding the white noise object. While ICM is also attracted by the fixed object, it is not doing it as fast as NDIGO. PE, as expected, is only attracted by the white noise object where its reward is the highest. We also observe that the performance of NDIGO improves as we increase the prediction horizon. From now on, in the tables, we report only the ICM results as it is the only competitive baseline. Exhaustive results are reported in Appendix E.1.

Experiment 2. To demonstrate better the information-seeking behaviour of our algorithm, we place randomly a fixed object in either the upper, left or right room and a white noise object in the lower room. Thus, to discover the object, the agent must actively look for it in all but the lower room.

Similar to Experiment 1, We report the average discovery loss of the fixed object in Figure 5. We observe that all variants of NDIGO perform better than the baselines by a clear margin. Though ICM performance is not far behind NDIGO (less than two times worse than NDIGO-4). We also observe no significant difference between the different variants of NDIGO in this case. We also report in Table 2 the first visit and visit counts for the fixed object and the white noise object in an episode. NDIGO again demonstrates a superior performance to the baselines. We also observe that NDIGO in most case is not attracted towards the white noise object. An interesting observation is that, as we increase the horizon of prediction in NDIGO, it takes more time for the agent to find the fixed object but at the same time the visit counts increases as well, i.e, the agent stay close to the object for longer time after the first visit.

As a qualitative result, we also report top-down-view snapshots of the behavior of NDIGO-2 up to the time of discovery of fixed in the right room in Figure 6. We also depicts the predicted view of the world from the agent’s representation in Figure 6. As the location of object is unknown to the agent, we observe that the agent searches the top-side, left-side and right-side rooms until it discovers the fixed object in the right-side room. It also successfully avoids the bottom-side room containing the white noise object. Also as soon as the agent finds the fixed object the uncertainty about the location of fixed object completely vanishes (as the agent has learned there is only one fixed object exists in the world).

² Each episode is set to end after 400 time steps; if an agent does not find the object by the end of the episode, the first visit time is set to 400.

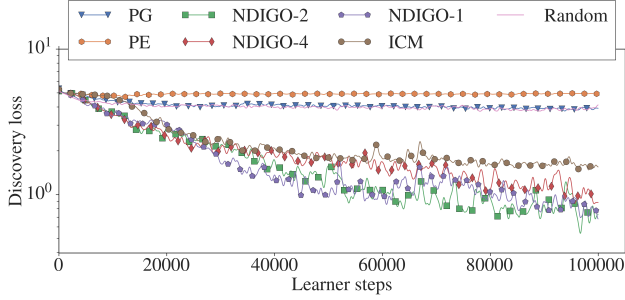


Figure 5. Experiment 2: Average discovery loss of the fixed object. The results are averaged over 10 seeds.

	Visit count		First visit time	
	fixed	w. noise	fixed	w. noise
ICM	151.7 \pm 33.0	15.6 \pm 9.0	142.1 \pm 40.8	198.7 \pm 55.1
NDIGO-1	180.2 \pm 42.7	12.8 \pm 6.9	101.1 \pm 31.1	237.2 \pm 49.4
NDIGO-2	209.3 \pm 34.9	3.5 \pm 2.3	121.1 \pm 36.5	306.4 \pm 43.4
NDIGO-4	233.7 \pm 41.6	5.3 \pm 3.7	126.7 \pm 43.3	268.2 \pm 53.1

Table 2. Average values of the visit counts and first visit time of the trained agent for the fixed and white noise objects in Experiment 2.

Experiment 3. We investigate whether NDIGO is able to discover and retain the dynamics of moving (but still predictable) objects even when not being in its field of view. For this, we used a 5 rooms setting with two bouncing objects in upper and lower rooms and a white noise object in the right room.

	Visit count		First visit time	
	upper obj.	lower obj.	upper obj.	lower obj.
ICM	80.5 \pm 28.3	89.1 \pm 28.6	174.8 \pm 53.4	127.8 \pm 51.4
NDIGO-1	41.0 \pm 8.5	45.2 \pm 11.6	34.4 \pm 18.7	38.8 \pm 16.1
NDIGO-2	108.5 \pm 25.1	31.3 \pm 20.9	118.3 \pm 50.4	312.6 \pm 50.6
NDIGO-4	198.7 \pm 33.4	44.2 \pm 28.8	64.5 \pm 38.8	320.8 \pm 47.5

Table 3. Average values of the visit counts and first visit time of the trained agent for the bouncing objects in Experiment 3.

We report the discovery loss in Figure 7. We observe that all variants of NDIGO outperforms the baselines by a large margin in terms of the discovery loss of the bouncing object. As the discovery loss for both bouncing objects is small, this indicates that NDIGO can encode the dynamics of bouncing objects in its representation. We report the first-visit and visit counts for the bouncing objects in Table 3. NDIGO has a superior performance than the baselines both in terms of visit counts and visit time to the fixed objects except for the visit count of the lower object in which ICM produces the best performance. Finally, as a qualitative result, we also report top-down-view snapshots of the behavior of NDIGO-1 after the discovery of each bouncing object in Figure 8. We observe that the agent can estimate the location of both bouncings in the first

visit. Also after departing from the green bouncing object and moving towards the red bouncing object, still it can track the dynamics of the green bouncing object with some small error. This is despite the fact that the green bouncing object is not anymore observed by the agent.

Experiment 4. We investigate if the horizon H affects the performance of the agents in terms of its sensitivity to structured noise. For this we evaluated which objects the agent seeks in a 5 rooms setting with a Brownian object in the upper room and a fixed object in the lower room. In the upper room, the Brownian moves at every time step. For the Brownian, unlike white noise, it is not guaranteed that the reward of NDIGO is zero. However by increasing the horizon, one may expect that the intrinsic reward due to the Brownian object becomes negligible because it becomes harder to predict with higher horizons.

	Visit count		First visit time	
	Brownian	fixed	Brownian	fixed
ICM	358.3 \pm 9.4	0.5 \pm 0.9	34.0 \pm 8.3	385.1 \pm 24.6
NDIGO-1	356.1 \pm 6.9	0.0 \pm 0.0	23.4 \pm 6.4	398.9 \pm 8.9
NDIGO-2	350.7 \pm 5.4	0.1 \pm 0.3	21.1 \pm 4.8	383.9 \pm 25.6
NDIGO-4	0.4 \pm 1.0	290.5 \pm 31.4	395.5 \pm 12.4	68.4 \pm 29.8

Table 4. Average values of the visit counts and first visit time of the trained agent for the Brownian and fixed objects in Experiment 4, with all baselines.

We report the results in Figure 9. We observe that the ICM baseline as well as the variants of NDIGO with the short horizon are being attracted to the structured randomness generated by the Brownian object. Only NDIGO-4 can ignore the Brownian object and discover the fixed object. As a result NDIGO-4 is the only algorithm capable of minimising the discovery loss of the fixed object.

Experiment 5. We now compare discovery ability of the agents in a complex maze environment (see Figure 3) with no extrinsic reward. Here, the agent starts in a fixed position in the maze environment, and is given no incentive to explore but its intrinsic reward. This setting is challenging for discovery and exploration, since to go the end of the maze the agents need to take a very long and specific sequence of actions. This highlights the importance of intrinsic rewards that encourage discovery. We report the learning curves of NDIGO as well as the baselines in Figure 10. We observe that in this case different variants of NDIGO outperform the baselines by a wide margin in terms of discovery loss, while NDIGO-1 and NDIGO-2 outperforming NDIGO-5. Note that due to the presence of movable object, which is unpredictable upon re-spawning, the average loss in this experiment is higher than the prior fixed object experiments. We also evaluate the discovery performance of the agent as the number of rooms it is capable of exploring within the duration of the episode. We present the average visit

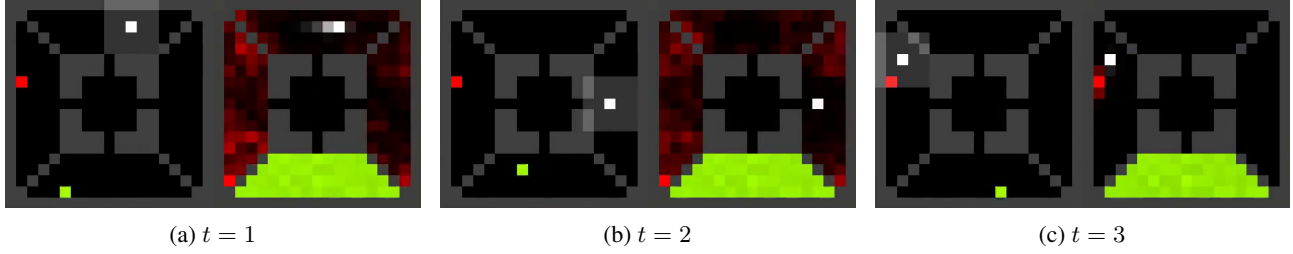


Figure 6. Experiment 2: top-down-view snapshots of the behavior of the NDIGO-4 agent. (a) after entering the top-side room (b) after entering the right-side room (c) after discovering the `fixed` object in the left-side room. In each subpanel the left-side image depicts the ground-truth top-down-view of the world and the right-side image depicts the predicted view from the agent’s representation. All times are in seconds.

	Room 1 white noise	Room 2 fixed	Visit frequency Room 3 fixed	Room 4 fixed	Room 5 movable
ICM	100.0% \pm 0.0%	26.8% \pm 25.7%	13.8% \pm 20.0%	6.5% \pm 14.3%	—
NDIGO-1	94.7% \pm 12.9%	66.4% \pm 27.4%	71.7% \pm 26.1%	70.4% \pm 26.4%	67.8% \pm 27.1%
NDIGO-2	100.0% \pm 0.0%	78.3% \pm 23.9%	84.8% \pm 20.9%	83.7% \pm 21.4%	81.5% \pm 22.5%
NDIGO-5	100.0% \pm 0.0%	49.6% \pm 29.0%	47.4% \pm 28.9%	18.8% \pm 22.6%	—
NDIGO-10	100.0% \pm 0.0%	84.1% \pm 21.4%	95.5% \pm 12.2%	45.5% \pm 29.1%	—

Table 5. Average frequency of visits to each room for the trained agents.

	Room 1 white noise	Room 2 fixed	First visit time Room 3 fixed	Room 4 fixed	Room 5 movable
ICM	4.4 \pm 3.0	324.7 \pm 79.5	375.0 \pm 44.0	391.8 \pm 24.3	—
NDIGO-1	40.6 \pm 57.2	203.0 \pm 90.6	190.5 \pm 86.0	199.9 \pm 85.2	212.7 \pm 83.2
NDIGO-2	12.9 \pm 10.7	171.5 \pm 79.5	159.4 \pm 68.8	174.5 \pm 68.9	192.8 \pm 68.9
NDIGO-5	6.8 \pm 11.5	245.1 \pm 94.1	255.9 \pm 91.4	344.9 \pm 68.7	—
NDIGO-10	8.6 \pm 5.9	128.0 \pm 75.8	119.1 \pm 53.4	283.1 \pm 81.4	—

Table 6. Average time of first visit to each room for the trained agents.

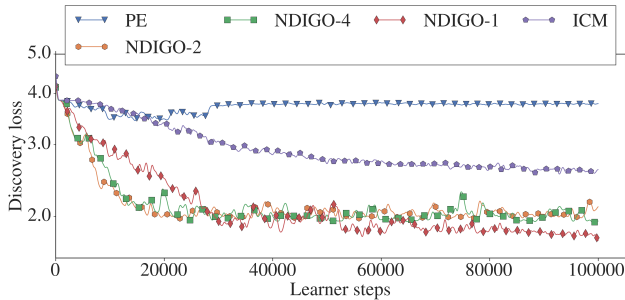


Figure 7. Experiment 3: Average discovery loss of bouncing objects. The results are averaged over 10 seeds.

frequency and first visit time of each room for the trained agents (see Tables 5 and 6). NDIGO-1 and NDIGO-2 appear as the only agents capable of reaching the final room, whereas NDIGO-4 explores 4 out of 5. The rest can not go beyond the `white noise` object.

As a qualitative result, we also report top-down-view snapshots of the behavior of NDIGO-1 up to the time of discovery of the last `fixed` in room 2 in Figure 11. We also depicts the predicted view of the world from the agent’s representation in Figure 6. We observe the agent drives across the maze all the way from room 1 to room 5 and in the process discovers the `fixed` objects in rooms 3-4 (see Figure 6a) and the movable object in room 5 (see Figure 6c). It then chases `movable` object until `movable` object gets fixated on the top-left corner of the world. The agent then moves back to room 2 (see Figure 6c) and discovers the last

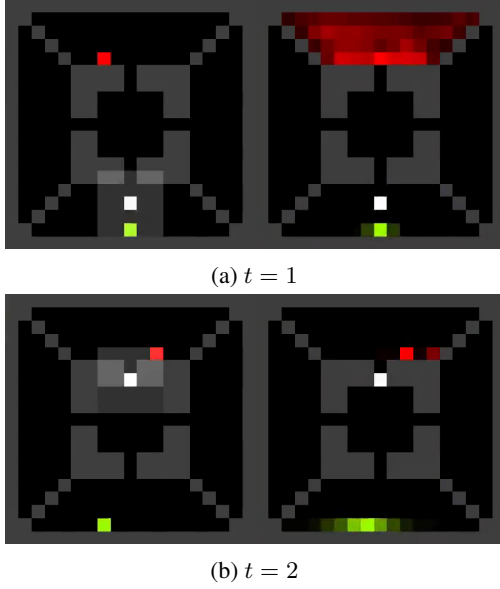


Figure 8. Experiment 3: top-down-view snapshots of the behavior of the NDIGO-1 agent. (a) after discovering the green bouncing object in the bottom-side room (b) after discovering the red bouncing object in the top-side room. In each subpanel the left-side image depicts the ground-truth top-down-view of the world and the right-side image depicts the predicted view from the agent’s representation. All times are in seconds.

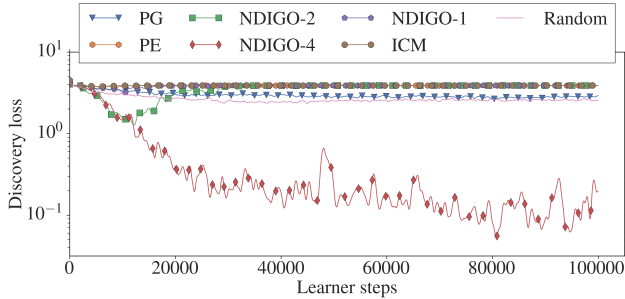


Figure 9. Experiment 4: Average discovery loss of the fixed object. The results are averaged over 10 seeds.

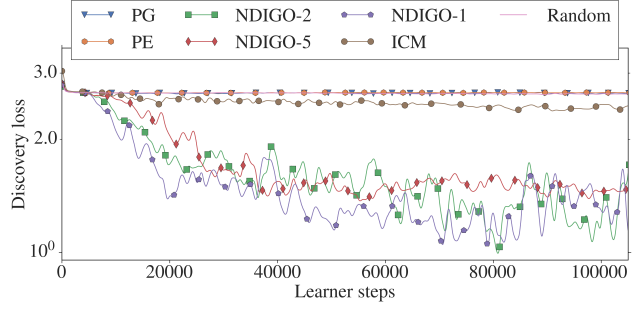


Figure 10. Experiment 5: Average discovery loss of the fixed and movable objects. The results are averaged over 10 seeds.

blue fixed object there, while maintaining its knowledge of the other objects. The reason for ignoring the blue fixed object in room 2, in the first place, might be due to the fact that the agent can obtain more intrinsic rewards by chasing the movable. So it tries to reach to room 5 as fast as possible at the expense of ignoring the blue fixed object in room 2.

7. Conclusion

We aimed at building a proof of concept for a world discovery model by developing the NDIGO agent and comparing its performance with the state-of-the-art information-seeking algorithms in terms of its ability to discover the world. Specifically, we considered a variety of simple local-view 2D navigation tasks with some hidden randomly-placed objects and looked at whether the agent can discover its environment and the location of objects. We evaluate the ability of our agent for discovery through the glass-box approach which measures how accurate location of objects can be predicted from the internal representation. Our results showed that in all these tasks NDIGO produces an effective information seeking strategy capable of discovering the hidden objects without being distracted by the white noise, whereas the baseline information seeking methods in most cases failed to discover the objects due to the presence of noise.

There remains much interesting future work to pursue. The ability of our agent to discover its world can be very useful in improving performance in multi-task and transfer settings as the NDIGO model can be used to discover the new features of new tasks. Also in this paper we focused on visually simple tasks. To scale up our model to more complex visual tasks we need to consider more powerful prediction models such as Pixel-CNN (van den Oord et al., 2016), VAE (Kingma & Welling, 2013), Info-GAN (Chen et al., 2016) and Draw (Gregor et al., 2015) capable of providing high accuracy predictions for high-dimensional visual scenes. We also can go beyond predicting only visual

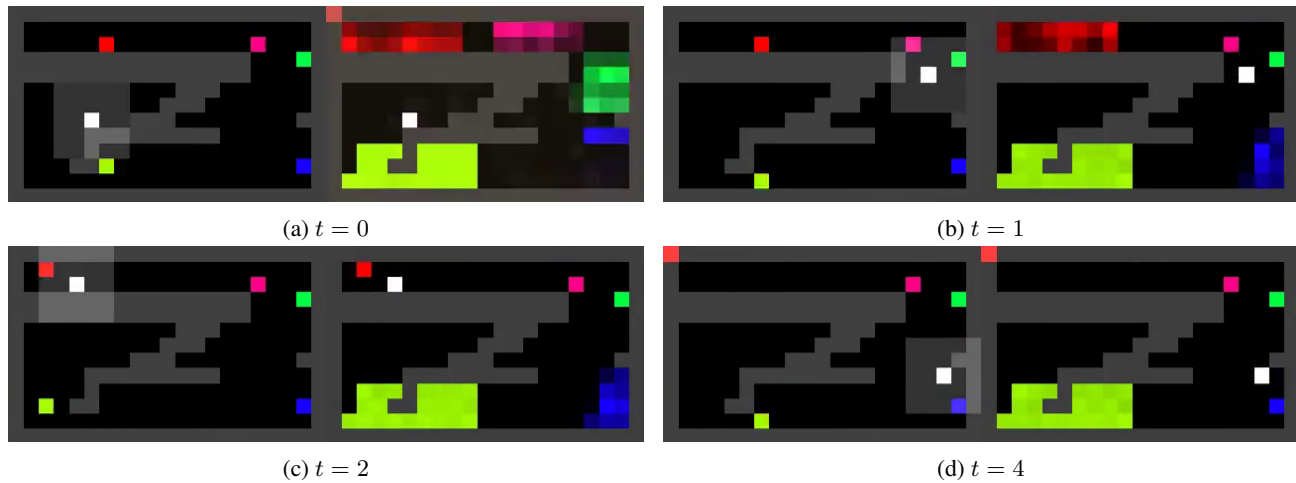


Figure 11. Experiment 5: top-down-view snapshots of the behavior of the NDIGO-1 agent in the maze problem: (a) at the beginning of the episode (b) after discovering the `fixed` objects in room 3 and 4 (c) after discovering the `movable` object in room 5 (d) after discovering the `fixed` object in room 2. In each subpanel the left-side image depicts the ground-truth top-down-view of the world and the right-side image depicts the predicted view from the agent’s representation. All times are in seconds.

observations to other modalities of sensory inputs, such as proprioception and touch sensors (Amos et al., 2018).

Acknowledgements

We would like to thank Daniel Guo, Theophane Webber, Caglar Gulcehre, Toby Pohlen, Steven Kapturovski and Tom Stepleton for insightful discussions, comments and feedback on this work.

References

- Achiam, J. and Sastry, S. Surprise-based intrinsic motivation for deep reinforcement learning. *arXiv preprint arXiv:1703.01732*, 2017.
- Amos, B., Dinh, L., Cabi, S., Rothörl, T., Colmenarejo, S. G., Muldal, A., Erez, T., Tassa, Y., de Freitas, N., and Denil, M. Learning awareness models. *Sixth International Conference on Learning Representations (ICLR 2018)*, 2018. arXiv preprint arXiv:1804.06318.
- 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.
- Cassandra, A. R. *Exact and approximate algorithms for partially observable Markov decision processes*. Brown University, 1998.
- Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., and Abbeel, P. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, pp. 2172–2180, 2016.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Clark, A. A nice surprise? predictive processing and the active pursuit of novelty. *Phenomenology and the Cognitive Sciences*, pp. 1–14, 2017.
- Csikszentmihalyi, M. and Csikszentmihalyi, I. S. *Optimal experience: Psychological studies of flow in consciousness*. Cambridge university press, 1992.
- Dayan, P. and Balleine, B. W. Reward, motivation, and reinforcement learning. *Neuron*, 36(2):285–298, 2002.
- Deci, E. and Ryan, R. M. *Intrinsic motivation and self-determination in human behavior*. Springer Science & Business Media, 1985.
- Frank, M., Leitner, J., Stollenga, M., Förster, A., and Schmidhuber, J. Curiosity driven reinforcement learning for motion planning on humanoids. *Frontiers in neuro-robotics*, 7:25, 2014.
- Friston, K. J., Stephan, K. E., Montague, R., and Dolan, R. J. Computational psychiatry: the brain as a phantastic organ. *The Lancet Psychiatry*, 1(2):148–158, 2014.
- Gottlieb, J., Oudeyer, P.-Y., Lopes, M., and Baranes, A. Information-seeking, curiosity, and attention: computational and neural mechanisms. *Trends in cognitive sciences*, 17(11):585–593, 2013.

- Gregor, K., Danihelka, I., Graves, A., Rezende, D. J., and Wierstra, D. Draw: A recurrent neural network for image generation. *arXiv preprint arXiv:1502.04623*, 2015.
- Guo, Z. D., Azar, M. G., Piot, B., Pires, B. A., Pohlen, T., and Munos, R. Neural predictive belief representations. *arXiv preprint arXiv:1811.06407*, 2018.
- Ha, D. and Schmidhuber, J. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- Haber, N., Mrowca, D., Fei-Fei, L., and Yamins, D. L. Learning to play with intrinsically-motivated self-aware agents. *arXiv preprint arXiv:1802.07442*, 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.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Hohwy, J. *The predictive mind*. Oxford University Press, 2013.
- Horvitz, J. C. Mesolimbocortical and nigrostriatal dopamine responses to salient non-reward events. *Neuroscience*, 96(4):651–656, 2000.
- Houthoofd, R., Chen, X., Duan, Y., Schulman, J., De Turck, F., and Abbeel, P. Vime: Variational information maximizing exploration. In *Advances in Neural Information Processing Systems*, pp. 1109–1117, 2016.
- Itti, L. and Baldi, P. Bayesian surprise attracts human attention. *Vision research*, 49(10):1295–1306, 2009.
- Kakade, S. and Dayan, P. Dopamine: generalization and bonuses. *Neural Networks*, 15(4-6):549–559, 2002.
- Kapturowski, S., Ostrovski, G., Dabney, W., Quan, J., and Munos, R. Recurrent experience replay in distributed reinforcement learning. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rllYtJqYX>.
- Kidd, C. and Hayden, B. Y. The psychology and neuroscience of curiosity. *Neuron*, 88(3):449–460, 2015.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Levine, S., Finn, C., Darrell, T., and Abbeel, P. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- Litman, J. Curiosity and the pleasures of learning: Wanting and liking new information. *Cognition & emotion*, 19(6):793–814, 2005.
- Little, D. Y.-J. and Sommer, F. T. Learning and exploration in action-perception loops. *Frontiers in neural circuits*, 7:37, 2013.
- Littman, M. L., Sutton, R. S., and Singh, S. Predictive representations of state. In *Advances in neural information processing systems*, pp. 1555–1561, 2002.
- Lopes, M., Lang, T., Toussaint, M., and Oudeyer, P.-Y. Exploration in model-based reinforcement learning by empirically estimating learning progress. In *Advances in Neural Information Processing Systems*, pp. 206–214, 2012.
- Lovejoy, W. S. A survey of algorithmic methods for partially observed markov decision processes. *Annals of Operations Research*, 28(1):47–65, 1991.
- Merriam-Webster, I. *Merriam-Webster's collegiate dictionary*. Merriam-Webster, 2004.
- 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.
- Moravčík, M., Schmid, M., Burch, N., Lisý, V., Morrill, D., Bard, N., Davis, T., Waugh, K., Johanson, M., and Bowling, M. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017.
- Munos, R., Stepleton, T., Harutyunyan, A., and Bellemare, M. Safe and efficient off-policy reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 1054–1062, 2016.
- Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Ostrovski, G., Bellemare, M. G., Oord, A. v. d., and Munos, R. Count-based exploration with neural density models. *arXiv preprint arXiv:1703.01310*, 2017.
- Oudeyer, P.-Y. and Kaplan, F. How can we define intrinsic motivation? In *Proceedings of the 8th International Conference on Epigenetic Robotics: Modeling Cognitive Development in Robotic Systems, Lund University Cognitive Studies, Lund: LUCS, Brighton*. Lund University Cognitive Studies, Lund: LUCS, Brighton, 2008.

- Oudeyer, P.-Y. and Smith, L. B. How evolution may work through curiosity-driven developmental process. *Topics in Cognitive Science*, 8(2):492–502, 2016.
- Oudeyer, P.-Y., Kaplan, F., and Hafner, V. V. Intrinsic motivation systems for autonomous mental development. *IEEE transactions on evolutionary computation*, 11(2): 265–286, 2007.
- Pathak, D., Agrawal, P., Efros, A., and Darrell, T. Curiosity-driven exploration by self-supervised prediction. In *International Conference on Machine Learning (ICML)*, volume 2017, 2017.
- Pohlen, T., Piot, B., Hester, T., Azar, M. G., Horgan, D., Budden, D., Barth-Maron, G., van Hasselt, H., Quan, J., Večerík, M., et al. Observe and look further: Achieving consistent performance on atari. *arXiv preprint arXiv:1805.11593*, 2018.
- Schmidhuber, J. Curious model-building control systems. In *Neural Networks, 1991. 1991 IEEE International Joint Conference on*, pp. 1458–1463. IEEE, 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. Simple algorithmic theory of subjective beauty, novelty, surprise, interestingness, attention, curiosity, creativity, art, science, music, jokes. *Journal of SICE*, 48(1), 2009.
- Schmidhuber, J. Formal theory of creativity, fun, and intrinsic motivation (1990–2010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230–247, 2010.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust region policy optimization. In *International Conference on Machine Learning*, pp. 1889–1897, 2015.
- Shyam, P., Jaśkowski, W., and Gomez, F. Model-based active exploration. *arXiv preprint arXiv:1810.12162*, 2018.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484, 2016.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.
- Stadie, B. C., Levine, S., and Abbeel, P. Incentivizing exploration in reinforcement learning with deep predictive models. *arXiv preprint arXiv:1507.00814*, 2015.
- van den Oord, A., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al. Conditional image generation with pixelcnn decoders. In *Advances in Neural Information Processing Systems*, pp. 4790–4798, 2016.
- Wang, Z., Schaul, T., Hessel, M., Van Hasselt, H., Lanctot, M., and De Freitas, N. Dueling network architectures for deep reinforcement learning. *arXiv preprint arXiv:1511.06581*, 2015.
- White, R. W. Motivation reconsidered: The concept of competence. *Psychological review*, 66(5):297, 1959.

A. NDIGO Global Network Architecture

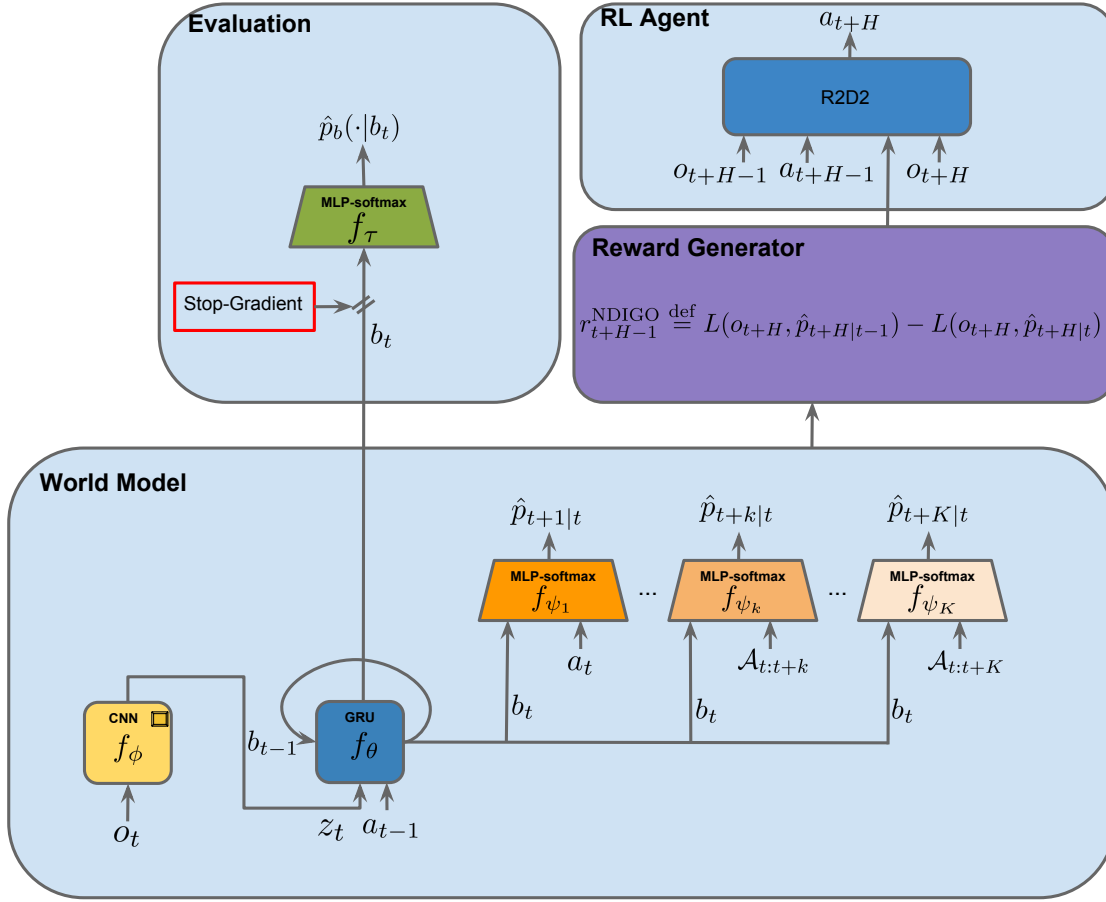


Figure 12. Global Architecture of the NDIGO agent

B. NDIGO Agent Implementation Details

B.1. World Model

- Convolution Neural Network (CNN) f_ϕ : Observations o_t are fed through a two-layer CNN (16 3×3 filters with 1×1 stride, then 16 3×3 filters with 2×2 stride; edges are padded if needed), then through a fully connected single layer perceptron with 256 units, then through a ReLU activation, resulting in a transformed observation z_t .
- Gated Recurrent Unit (GRU) f_θ : 128 units GRU.
- Frame predictors $\{f_{\psi_k}\}_{k=1}^K$ are MultiLayer Perceptrons (MLPs): one hidden-layer of 64 units followed by a ReLU activation and the output layer of 25 units (5×5 is the size of the local view which the size of the observation o_t) followed by a ReLU activation.
- $K = 10$
- Optimiser for frame predictions: Adam optimiser (Kingma & Ba, 2014) with batch size 256 and learning rate 5×10^{-4} .

B.2. Reward Generator

- The horizon $H = \{1, 2, 4\}$ can take one of these values in our experiments.

B.3. Evaluation

- The MLP f_τ : one hidden-layer of 64 units followed by a ReLU activation and the output layer of 361 units (19×19 is the size of the global view of the 5 rooms environment which is also the size of the real state x_t) followed by a ReLU activation.
- Optimiser for evaluation: Adam optimiser (Kingma & Ba, 2014) with batch size 256 and learning rate 5×10^{-4} .

B.4. RL Agent

We use the Recurrent Replay Distributed DQN (R2D2) (Kapturowski et al., 2019) with the following parameters:

- Replay: replay period is 100, replay trace length is 100, replay size is 1×10^6 and we use uniform prioritisation.
- Network architecture: R2D2 uses a two-layer CNN, followed by a GRU which feeds into the advantage and value heads of a dueling network (Wang et al., 2015), each with a hidden layer of size 128 units. The CNN and GRU of the RL agent have the same architecture and parameters as the one described for the World Model (see Sec.B.1) but do not share the same weights.
- Algorithm: Retrace Learning update (Munos et al., 2016) with discount factor $\gamma = 0.99$ and eligibility traces coefficient $\lambda = 0.97$, target network with update period 1024, no reward clipping and signed-hyperbolic re-scaling (Pohlen et al., 2018).
- Distributed training: 100 actors and 1 learner, actor update period every 100 learner steps.
- Optimiser for RL: Adam optimiser (Kingma & Ba, 2014) with batch size 256 and learning rate 1×10^{-4} .
- The intrinsic rewards are provided directly to the RL agent without any scaling.

B.5. Training loop pseudocode

Algorithm 1 NDIGO training loop.

Input: Policy π , history h_T , $K \geq H \geq 1$, weights \mathbf{W}

```

1:  $b_0 \leftarrow 0$ 
2: for  $t = 1 \dots T - K$  do
3:    $z_t = f_\phi(o_t; \mathbf{W})$  ▷ Observation CNN
4:    $b_t \leftarrow f_\theta(z_t, a_{t-1}, b_{t-1}; \mathbf{W})$  ▷ Belief GRU
5:    $\mathcal{A} \leftarrow [a_{t-1}]$ 
6:   for  $k = 1 \dots K$  do
7:      $\hat{p}_{t+k|t} \leftarrow f_{\psi_k}(b_t, \mathcal{A}; \mathbf{W})$  ▷ Prediction MLP
8:      $L_{t+k|t} \leftarrow -\ln \hat{p}_{t+k|t}(o_{t+k})$ 
9:      $\mathcal{A} \leftarrow \mathcal{A} + [a_{t+k-1}]$ 
10:  end for
11:   $r_{t+H-1}^{\text{NDIGO}} \leftarrow L_{t+H|t-1} - L_{t+H|t}$ 
12: end for
13: Update  $\mathbf{W}$  to minimise  $\sum_{t=1}^{T-K} \sum_{k=1}^K L_{t+k|t}$ 
14: Update  $\pi$  using the set of rewards  $\{r_{t+H-1}^{\text{NDIGO}}\}_{t=1:T-K+1}$  with the RL Algo.
```

C. NDIGO Alternative Architecture

An alternative architecture for NDIGO consists in encoding the sequence of actions $\mathcal{A}_{t:t+k}$ into a representation using a GRU f_ξ . The hidden state of this GRU is $a_{t,k} = f_\psi(a_{t+k}, a_{t,k-1})$, with the initialisation $a_{t,1} = f_\xi(a_t, 0)$. Then we use a single neural network f_ψ to output, for any k , the probability distribution $\hat{p}_{t,k}(\cdot | b_t, a_{t,k})$ when given the input $a_{t,k}$ concatenated with b_t . The loss function for the network f_ψ at time step $t + k - 1$ is a cross entropy loss:

$$L(o_{t+k}, \hat{p}_{t,k}(\cdot|b_t, a_{t,k})) = -\ln(\hat{p}_{t,k}(o_{t+k}|b_t, a_{t,k})).$$

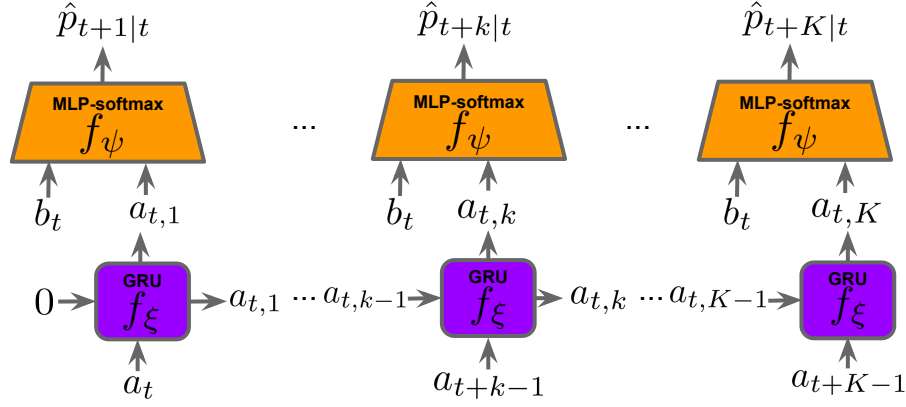


Figure 13. Alternative architecture of NDIGO for the frame prediction tasks.

D. Pathak et al. (2017)’s ICM Model for Partially Observable Environments

The method consists in training the internal representation b_t to be less sensitive to noise using a self-supervised inverse dynamics model. To do so, one inverse dynamics model f_β fed by (b_t, z_{t+1}) (concatenation of the internal representation and the transformed observation z_{t+1}) outputs a distribution \hat{p}_A over actions that predicts the action a_t . This network is trained by the loss: $L(\hat{p}_A, a_t) = -\ln(\hat{p}_A(a_t))$. Then a forward model f_α fed by (b_t, a_t) (concatenation of the internal representation and the action) outputs a vector \hat{b}_{t+1} that directly predict the future internal representation b_{t+1} . The forward model f_α is trained with a regression loss: $L_2(\hat{b}_{t+1}, b_{t+1}) = \|\hat{b}_{t+1} - b_{t+1}\|_2^2$. The neural architecture is shown in Fig. 14. Finally, the intrinsic reward is defined as:

$$r_t^{\text{FPE}} = L_2(\hat{b}_{t+1}, b_{t+1}).$$

This is slightly different from the architecture proposed by Pathak et al. (2017) in order to be compatible with partially observable environments.

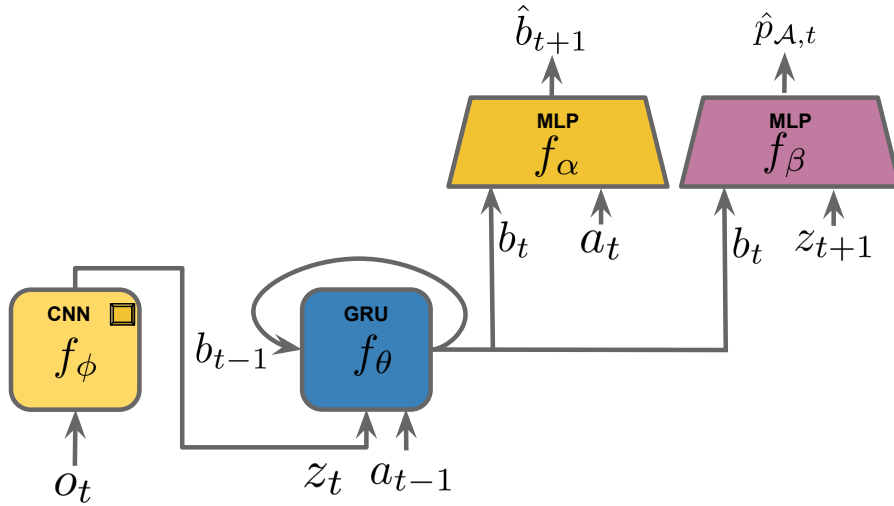


Figure 14. Pathak et al. (2017)’s ICM Model for Partially Observable Environments

D.1. Details of the ICM Model’s Architecture

The ICM agent shares exactly the same architecture than the NDIGO agent except that the forward predictors $\{f_{\psi_k}\}_{k=1}^K$ are replaced by an inverse model f_β and a forward model f_α .

- The inverse model f_β is an MLP: one hidden layer of 256 units and the output layer of 5 units (one-hot action size).
- The forward model f_α is an MLP: one hidden layer of 256 units and the output layer of 128 units (size of the GRU).

E. Additional results

E.1. Additional results for Experiment 2-4

Tables 7 to 9 contain the results (including baselines) for experiments 2 to 4.

	Visit count		First visit time	
	fixed	white noise	fixed	white noise
Random	20.8 ± 16.8	51.0 ± 15.8	332.3 ± 42.4	148.5 ± 56.9
PE	0.3 ± 0.8	161.7 ± 3.4	388.8 ± 21.0	11.2 ± 2.5
PG	32.6 ± 26.6	11.7 ± 7.7	309.9 ± 49.9	293.2 ± 50.9
ICM	151.7 ± 33.0	15.6 ± 9.0	142.1 ± 40.8	198.7 ± 55.1
NDIGO-1	180.2 ± 42.7	12.8 ± 6.9	101.1 ± 31.1	237.2 ± 49.4
NDIGO-2	209.3 ± 34.9	3.5 ± 2.3	121.1 ± 36.5	306.4 ± 43.4
NDIGO-4	233.7 ± 41.6	5.3 ± 3.7	126.7 ± 43.3	268.2 ± 53.1

Table 7. Average values of the visit counts and first visit time of the trained agent for the fixed and white noise objects in Experiment 2, with all baselines.

	Visit count		First visit time	
	upper obj.	lower obj.	upper obj.	lower obj.
PE	0.6 ± 0.5	0.1 ± 0.2	343.8 ± 42.6	390.3 ± 18.2
ICM	80.5 ± 28.3	89.1 ± 28.6	174.8 ± 53.4	127.8 ± 51.4
NDIGO-1	41.0 ± 8.5	45.2 ± 11.6	34.4 ± 18.7	38.8 ± 16.1
NDIGO-2	108.5 ± 25.1	31.3 ± 20.9	118.3 ± 50.4	312.6 ± 50.6
NDIGO-4	198.7 ± 33.4	44.2 ± 28.8	64.5 ± 38.8	320.8 ± 47.5

Table 8. Average values of the visit counts and first visit time of the trained agent for the bouncing objects in Experiment 3, with the PE and ICML baselines.

	Visit count		First visit time	
	Brownian	fixed	Brownian	fixed
Random	16.7 ± 11.6	41.1 ± 24.7	309.3 ± 46.0	244.3 ± 56.1
PE	357.3 ± 4.6	0.1 ± 0.3	15.6 ± 3.5	399.2 ± 8.5
PG	23.0 ± 14.0	38.2 ± 25.7	281.7 ± 52.4	268.4 ± 57.2
ICM	358.3 ± 9.4	0.5 ± 0.9	34.0 ± 8.3	385.1 ± 24.6
NDIGO-1	356.1 ± 6.9	0.0 ± 0.0	23.4 ± 6.4	398.9 ± 8.9
NDIGO-2	350.7 ± 5.4	0.1 ± 0.3	21.1 ± 4.8	383.9 ± 25.6
NDIGO-4	0.4 ± 1.0	290.5 ± 31.4	395.5 ± 12.4	68.4 ± 29.8

Table 9. Average values of the visit counts and first visit time of the trained agent for the Brownian and fixed objects in Experiment 4, with all baselines.

E.2. Additional results for Experiment 5

Tables 10 and 11 present the complete results of Experiment 5; note that a room is considered as visited when the agent has actually seen the object inside that room. As the object in Room 2 can be missed by the agent if it appears in the lower-right corner of the maze, the reported frequency of visits to Room 2 can be lower than that of Rooms 3 and beyond, as this is the case for the reported figures of the NDIGO-1 and NDIGO-2 agents. A dash symbol indicates that the agent never visits the corresponding room.

	Visit frequency				
	Room 1 white noise	Room 2 fixed	Room 3 fixed	Room 4 fixed	Room 5 movable
Random	100.0% \pm 0.0%	0.9% \pm 5.5%	—	—	—
PE	100.0% \pm 0.0%	—	—	—	—
PG	93.6% \pm 14.3%	—	—	—	—
ICM	100.0% \pm 0.0%	26.8% \pm 25.7%	13.8% \pm 20.0%	6.5% \pm 14.3%	—
NDIGO-1	94.7% \pm 12.9%	66.4% \pm 27.4%	71.7% \pm 26.1%	70.4% \pm 26.4%	67.8% \pm 27.1%
NDIGO-2	100.0% \pm 0.0%	78.3% \pm 23.9%	84.8% \pm 20.9%	83.7% \pm 21.4%	81.5% \pm 22.5%
NDIGO-5	100.0% \pm 0.0%	49.6% \pm 29.0%	47.4% \pm 28.9%	18.8% \pm 22.6%	—
NDIGO-10	100.0% \pm 0.0%	84.1% \pm 21.4%	95.5% \pm 12.2%	45.5% \pm 29.1%	—

Table 10. Average frequency of visits to each room for the trained agents.

	First visit time				
	Room 1 white noise	Room 2 fixed	Room 3 fixed	Room 4 fixed	Room 5 movable
Random	15.1 \pm 18.6	399.9 \pm 6.9	-	-	-
PE	8.6 \pm 4.3	-	-	-	-
PG	33.9 \pm 56.3	-	-	-	-
ICM	4.4 \pm 3.0	324.7 \pm 79.5	375.0 \pm 44.0	391.8 \pm 24.3	-
NDIGO-1	40.6 \pm 57.2	203.0 \pm 90.6	190.5 \pm 86.0	199.9 \pm 85.2	212.7 \pm 83.2
NDIGO-2	12.9 \pm 10.7	171.5 \pm 79.5	159.4 \pm 68.8	174.5 \pm 68.9	192.8 \pm 68.9
NDIGO-5	6.8 \pm 11.5	245.1 \pm 94.1	255.9 \pm 91.4	344.9 \pm 68.7	-
NDIGO-10	8.6 \pm 5.9	128.0 \pm 75.8	119.1 \pm 53.4	283.1 \pm 81.4	-

Table 11. Average time of first visit to each room for the trained agents.