

Curiosity and the exploration-exploitation dilemma

Erik J Peterson^{1,2*} and Timothy D Verstynen^{1,2,3,4}

*For correspondence:

Erik.Exists@gmail.com (EJP)

¹Department of Psychology; ²Center for the Neural Basis of Cognition; ³Carnegie Mellon Neuroscience Institute; ⁴Biomedical Engineering, Carnegie Mellon University, Pittsburgh PA

Abstract The exploration-exploitation dilemma is one of a few fundamental problems in the decision and life sciences. It has also been proven to be a mathematically intractable problem, when it is framed in terms of reward. To overcome this we have challenged the basic formulation of the problem itself, as having a single objective, namely reward. In its place we have defined independent objectives for exploration and exploitation. Through theory and numerical experiments we prove that a competition between exploration-as-curiosity and exploitation-for-reward has a tractable solution, based on the classic win-stay, lose-switch strategy from game theory. This strategy is possible because we treat information and reward as if they have equal value, and succeeds because the definition of curiosity we introduce is efficient. Besides offering a mathematical answer, this view of the problem seems more robust than the traditional approach because it succeeds in the difficult conditions where rewards are deceptive, or non-stationary.

When we observe an animal grappling with the decision to either explore or exploit, we often imagine this decision is based on reward. For example, if a bee goes in a familiar direction *to gather nectar*, it is said to be exploiting past knowledge. If it goes in an unknown direction *to gather nectar*, it is said to be exploring. Because the environment seems stochastic to the bee it cannot know which action it should be doing, exploiting or exploring *to gain more nectar*. It is this uncertainty *about rewards* that makes the choice a dilemma. It is also this uncertainty which ensures there is no tractable mathematical solution (Thrun and Möller, 1992; Dayan and Sejnowski, 1996; Ishii et al., 2002; Şimşek and Barto, 2006; Gershman, 2018). We illustrate this traditional view in Fig. 1a.

The choice between exploring and exploiting is a kind of decision that is faced routinely by learners of all kinds, including foraging bees, business organizations, humans, worms, monkeys, rodents, birds, children, and computer algorithms (Gupta et al., 2006; Sutton and Barto, 2018; Woodgate et al., 2017; Lee et al., 2011; Schulz et al., 2018; Calhoun et al., 2014; Wang and Hayden, 2019; Sumner et al., 2019; Auersperg, 2015). But is reward really fundamental to it?

In the natural world exploration already finds two explanations. If there is no reason to expect a reward, exploration is described theoretically as a search for information, what we will call curiosity (Berlyne, 1950; Schmidhuber, 1991; Kidd and Hayden, 2015; de Abril and Kanai, 2018; Jaegle et al., 2019; Friston et al., 2016). On the other hand, if reward is expected, then exploration gets re-interpreted as we described above, as a search for reward, and this specific interpretation is what leads the famous dilemma (Kelly, 1956; Berger-Tal et al., 2014; Dayan and Sejnowski, 1996; Thrun, 1992; Mehlhorn et al., 2015; Kobayashi and Hsu, 2019).

We have come to believe that considering reward value at all during exploration is a mistake. This might like a very counterintuitive suggestion but we will prove in this paper exploration is often

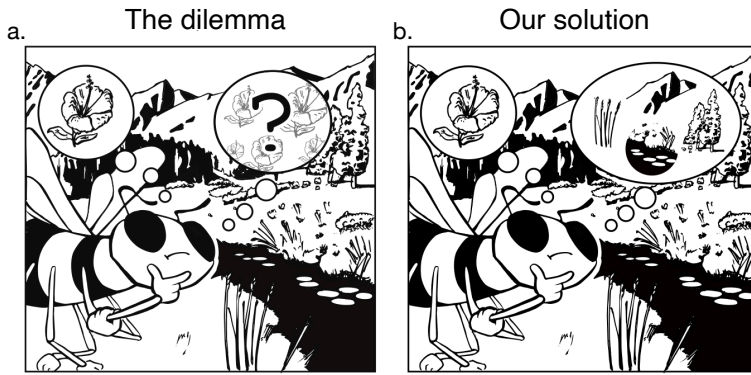


Figure 1. Two views of exploration and exploitation. **a.** The classic dilemma: either exploit an action with a known reward (e.g., return to the previous flower) or explore other actions on the chance they will return a better outcome. The central challenge here is that the outcome of exploration is uncertain, and filled with questions. **b.** An alternative view of the dilemma, with two goals: either maximize rewards or maximize information value with a curious search of the environment. *Artist credit:* Richard Grant.

43 better handled by a curiosity search. We illustrate this contrarian view in Fig. 1b.

44 The first half of this paper is devoted to developing a new view of curiosity that is simple but
 45 general. This effort, it turns out, lets us make a new contribution to information theory. It also lets
 46 us develop new proofs that define an ideal curious search. The second half uses these first results
 47 to prove curiosity is a rational solution to all explore-exploit decisions that are normally based on
 48 rewards.

49 Results

50 To understand curiosity in animals we turn to the fields of artificial life and machine learning. This
 51 is because curiosity in algorithmic form is unusually effective. In fact we are aware of no other
 52 approach to machine learning that is as flexible, or broadly effective, as curiosity. It has been
 53 shown to solve reinforcement learning problems, including sparse rewards, (*Schmidhuber, 1991*;
 54 *Pathak et al., 2017*; *Stanton and Clune, 2018*), and deceptive reward problems (*Lehman and Stan-*
 55 *ley, 2011a*), multi-objective problems where no one solution is correct (*Mouret, 2011*; *Fister et al.,*
 56 *2019*), and can lead to diverse skills (*Colas et al., 2020a*) that are robust to environment instability
 57 (*Mouret and Clune, 2015*) and robust to damage (*Cully et al., 2015*).

58 What is missing to extend these insights to understand animal behavior is a mathematical un-
 59 derstanding of curiosity. We need an understanding that can span fields, in other words.

60 Here we define curiosity as an open-ended search to learn any information (*Kidd and Hayden,*
 61 *2015*). We assume it should be efficient and motivated by value. What we would like is a way
 62 to describe curiosity so it does not depend on any one kind of learning, because curiosity itself
 63 depends on many different kinds of learning. In the most extreme case curiosity would then allow
 64 for any kind of learning. So this is the setting we have chosen.

65 Memory and learning

66 Any notion of information value requires that there must be a memory, Θ , and a way to learn, f .
 67 For our general mathematical theory we will define memory as,

- 68 • A set of real valued numbers
- 69 • Embedded in a finite space, with N dimensions
- 70 • Closed and bounded

71 Which lets us define learning in a matching way. We say that a learning function *can be any*
 72 *function* that maps observations x into memory. Using the function-arrow notation this kind of

Box 1. Norms.

A **norm** is a mathematical way to measure the overall size or extent of an “object” in a space. The object we are concerned with is a mathematical object with a form given by, $\Delta\Theta = f_{\Theta}(x) - \Theta$. But to define the norm let us as the first step consider a generic vector v . Keeping with convention will have its norm denoted by $\|v\|$, and defined by three properties:

1. $\|v\| > 0$ when $v \neq 0$ and $\|v\| = 0$ only if $v = 0$ (point separating)
2. $\|kv\| = |k| \|v\|$ (absolute scalability)
3. $\|v + w\| \leq \|v\| + \|w\|$ (the triangle inequality)

What does this mean for $\Delta\Theta$? A norm over a changing memory provides most the properties we require for our definition of E . It depends only on learning (and forgetting), that dependence must be monotonic, a norm is zero only when nothing in the space/memory changes and is positive definite. These satisfy the first four properties. For the fifth, see Box 2.

learning is denoted by, $f(x; \Theta) \rightarrow \Theta'$. Likewise a forgetting function *can be any function* that inverts the memory, $f^{-1}(x; \Theta') \rightarrow \Theta$.

So what are observations? It is convenient to assume that all observations are also embedded in a finite real space, $x \sim X \in \mathbb{R}^m$, as are actions $a \sim A \in \mathbb{R}^k$, and that rewards are positive real numbers, $R \in (\mathbb{R}^+)$.

In calling them “observations” we do not mean to imply they must arise from the senses. Observations could be internally generated from thought, or externally generated from the senses, or both in combination. Observations may contain action information, or they may not. Observations may contain reward information, or they may not. We have no preference because we wish to accommodate all the different kinds of cognition.

Information value

Let us say we have a memory Θ , which has been learned by f over a history of observations, $[f(x_0; \Theta_0), f(x_1; \Theta_1), f(x_2; \Theta_2), \dots]$. Can we measure how much value the next observation $f(x_t)$ should have? If there is such a measure E , we think it is reasonable to require that it have certain properties.

1. E should depend only on a memory, and what can be immediately learned. (x, Θ , and f)
2. An observation that is known completely, or cannot be learned, should have a value of 0. (If $\Delta\Theta = 0$, then $E = 0$)
3. E should be non-negative because learning is never itself harmful. ($E \geq 0$)
4. E should increase monotonically with the change in memory. ($E \propto \Delta\Theta$)
5. E for the same observation should become decelerating in finite time. ($E \rightarrow 0$ as $t \rightarrow T^*$)

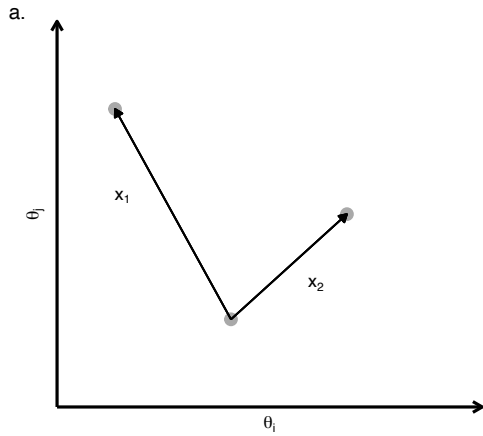
These properties are in part inspired by Shannon’s axioms for information theory, and in part found in the theories of information maximization, and information foraging (*Inglis et al., 2001; Reddy et al., 2016*). What we have done is remove Bayes rule as the learning rule. This is important not for its own sake but because anything we can prove in this more general setting must be true then for all practical learning algorithms (*MacKay, 2003*).

We can satisfy all these properties by taking a geometric view of learning, as given in the definitions below. They are also shown graphically in Fig. 2 and described in detail in Boxes 1-2.

Definition 1 Let E be the norm of $\|f_{\Theta}(x) - \Theta\|$.

Definition 2 Let $\|f_{\Theta}(x)\|$ be constrained such that $\nabla^2\Theta < 0$ for all $t \geq T^*$.

Definition 1 – distance



Definition 2 – deceleration

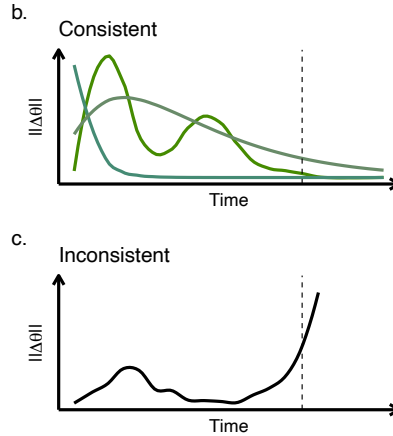


Figure 2. A diagram of our definitions. **a.** This panel illustrates a two dimensional memory. The information value of two observations x_1 and x_2 depends on the norm of memory with learning. Here we show this distance as a euclidean norm, denoted as a black arrow. **b-c** This panel illustrates learning dynamics with time (over a series of observations that are not shown). If information value becomes decelerating in finite time bound, then we say that learning is consistent with Def. 2. This is shown in panel b. The bound is depicted as a dotted line. If learning does not decelerate, then it is said to be inconsistent with Def. 2 (Panel c). *It is important to note:* our account of curiosity does not apply when learning is inconsistent.

Box 2. The Laplacian.

Here we use ∇^2 to represent the **Laplacian**, which provides the second derivative of vector valued functions, giving us the remaining property we need. (See Box 1 for the others). To explain why, suppose we are working with a 1- dimensional memory. In this case the Laplacian reduces to the second derivative on our memory, $\frac{d^2\theta}{dx^2} < 0$. A negative second derivative will force the changes in memory to be decelerating, by definition. In higher dimensions of memory then, a reader can think of the Laplacian as the “second derivative” of vector-valued functions like f . That is, negative Laplacian imply decelerating by definition.

Box 3. Optimal substructure

To use Bellman's principle of optimality (Box 4) the problem we want to solve needs to have **optimal substructure**. This opaque term can be understood by looking in turn at another theoretical construct, Markov spaces.

In Markov spaces there are a set of states or observations (x_0, x_1, \dots, x_T) , where the transition to the next x_t depends only on the previous state x_{t-1} . This limit means that if we were to optimize over these states, as we do in reinforcement learning, we know that we can treat each transition as its own “subproblem” and therefore the overall situation has “optimal substructure”.

In reinforcement learning the problem of reward collection is defined by fiat as happening in a Markov space (**Sutton and Barto, 2018**). The problem for curiosity is it relies on memory which is necessarily composed of many past observations, in many orders, and so it cannot be a Markov space. So if we wish to use dynamic programming to maximize E , we need to find another way to establish optimal substructure. This is the focus of Theorem 1.

Exploration as a dynamic programming problem

Curiosity is a search to maximize information value, again in the broadest sense possible. We decided to seek a dynamic programming solution to maximizing curiosity (**Bellmann, 1954**). This is useful because it would guarantee value is maximized, which, because of the way we have made our definitions, also guarantees that learning is complete. It is also convenient because dynamic programming is the basis for reinforcement learning (**Sutton and Barto, 2018**).

The normal route to find a dynamic programming solution is to first prove your problem has “optimal substructure”, then to derive the final result using Bellman optimality, and this is the route we will follow. For a more complete introduction to optimal substructure, see Box 3. For more on dynamic programming and Bellman optimality, see Box 4.

In Theorem 1 (mathematical appendix) we prove our definition of memory has optimal substructure. This proof relies on the forgetting function f^{-1} to succeed. Even though we have said little about it until now, this is why we include forgetting in our early definitions.

The Bellman derivation leads to an especially practical and simple result (Eq. 13). Given an arbitrary starting value $E_0 > 0$, the dynamic programming solution to maximizing E is given by Eq. 2. A full derivation is provided in the mathematical appendix. We denote our policy that follows this solution as π_E .

$$V_E^\pi(x) = \operatorname{argmax}_A [E_0 + V_E(x_1)] \quad (2)$$

However as long learning about the different states x is independent, and as long as learning continues until steady-state, $E_t = 0$, then there are many equally good solutions to the search. Under these conditions we can simplify Eq. 2 further, giving Eq. 3.

$$V_E^\pi(x) = \operatorname{argmax}_A [E_0 + E(x_1)] \quad (3)$$

Optimal exploration, and the importance of boredom

What would it mean for a curiosity search to be efficient? Does our Bellman solution ensure it?

Exploration is limited in practice by available energy, hunger, greed, thirst and other biological priorities. More importantly search should be limited to avoid the central failure of curiosity: it can become distracted by what we call minutia. Those details which have no use. A good example of this was given by **Pathak et al. (2017)** who asks us to, “Imagine a scenario where the agent is observing the movement of tree leaves in a breeze. Since it is inherently hard to model breeze, it

Box 4. Dynamic programming and Bellman optimality.

The goal of dynamic programming is to find a series of actions (a_1, a_2, \dots, a_T) , drawn from a set A , that maximize each payout (r_1, r_2, \dots, r_T) so the total payout received is as large as possible. If there is a policy $a = \pi(x)$ to take actions, based on a series of observations (x_0, x_1, \dots, x_T) , then an optimal policy π^* will always find the maximum total value $V^* = \operatorname{argmax}_A \sum_T r_i$. The question that Bellman and dynamic programming solve then is how to make the search for finding π^* tractable. In the form above one would need to reconsider the entire sequence of actions for any one change to that sequence, leading to a combinatorial explosion. Bellman's insight was a way to make the problem simpler, and break it down into a small set of problems that we can solve in a tractable way without an explosion in complexity. This is his principle of optimality, which reads:

An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision. (*Bellmann, 1954*)

Mathematically this allows us to translate the full problem, $V^* = \operatorname{argmax}_\pi \sum_T r_1, r_2, \dots, r_T$ to a recursive one, which is given below.

$$V^* = \operatorname{argmax}_\pi [r_0 + V(x_1)] \quad (1)$$

Using this new form we can find an optimal solution by moving through the sequence of action iteratively, and solving each step independently. The question which remains though is how can we ensure our problem can be broken down this way? This is addressed in Box 3.

is even harder to predict the location of each leaf". This will imply, they note, that a curious agent will always remain curious about the leaves, to its own detriment.

To limit curiosity we will draw on its opposite, boredom, which we consider to be an adaptive trait. Others have considered boredom this way, arguing it is a useful way to motivate aversive tasks (*Bench and Lench, 2013*). We take a slightly different view and consider boredom as a means to ignore the marginal. We treat it as a tunable free parameter, $\eta \geq 0$, and then say all curious exploration should cease once $E \leq \eta$. In other words, once learning becomes marginal the action policy must change, as given by Eq. 4.

$$E \leq \eta : \pi_E \rightarrow \pi_\emptyset \quad (4)$$

We use π_\emptyset to denote any other action policy that is not π_E . Informally Eq. 4 means that once E is below the boredom threshold, E is marginal, and the learner should change its behavior from exploration to some other policy. This change in behavior could be towards an aversive policy (as in (*Bench and Lench, 2013*)) or it could be towards reward seeking, as we will assume.

We think it is reasonable to call a curious search optimal if it meets the three criteria below.

1. Exploration should visit all available states of the environment at least once.
2. Exploration should cease when learning has plateaued.
3. Exploration should take as few steps as possible to achieve 1 and 2.

In Theorem ?? we prove that π_E satisfies these, when $\eta > 0$.

Information collection

The work so far has built up the idea that the most valuable, and most efficient, curious search will come from a deterministic algorithm. That is, every step strictly maximizes E . It is this determinism

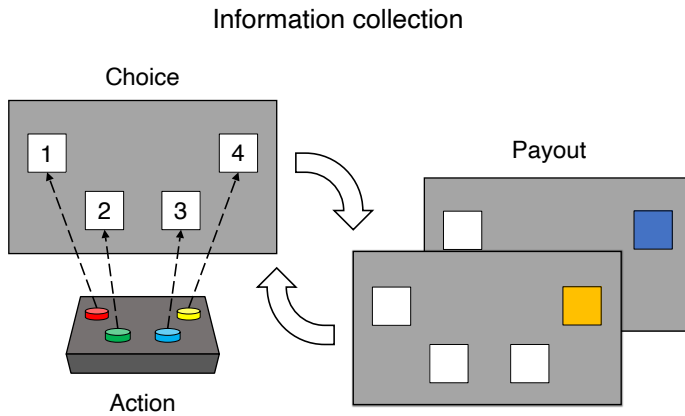


Figure 3. A simple four choice information foraging task. The information in this task is a yellow or blue stimulus, which can change from trial to trial. A good learner in this task is one who tries to learn the probabilities of all the symbols in each of the four choices. The more random the stimuli are in a choice, the more potential information/entropy there is to learn. *Note:* the information payout structure is shown below (Fig. 6a).

which will let us resolve the dilemma, later on. A deterministic view of exploration seems at odds with how the problem is generally viewed today, which involves searching with some amount of randomness. Whereas if our analysis is correct, randomness is not needed or desirable, as it must lead to less value and a longer search.

We confirm and illustrate this result using an example of a simple information foraging task (Task 1; Fig.). This variation of the bandit task (Sutton and Barto, 2018) replaces rewards with information, in this case colors. On each selection, the agent sees one of two colors according to a specific probability shown in Fig. 6a. When the relative color probabilities are more similar that arm has more entropy and so more to learn and more information value. Arms that are certain to present only one color lose their informative value quickly.

The results of this simple information seeking task are shown in Figure 4. Deterministic curiosity in this task generated more information value, in less time, and with zero regret when compared against a stochastic agent using more directed random search. As our work here predicts, noise only made the search happen slower and lead to regret. Besides offering a concrete example, performance in Task 1 is a first step in showing that even though π_E 's optimality is proven for a deterministic environment, it can perform well in other settings.

Two conjectures

How can we justify using curiosity for the dilemma when the goal of exploitation is reward collection? Intuition suggests that this is wrong-headed. Searching in an open-ended way for information must be less efficient than a direct search. In answer to this we make two conjectures.

- **Conjecture 1** Reward value and information value are equally important.

Rewards are a primary motivator for behavior in exploration (Sutton and Barto, 2018). Information is a primary motivator of behavior in exploration (Inglis et al., 2001). So which view is right? We believe that trying to assess the overall worth of reward or information is a mistake. We cannot, and should not, weigh these two perspectives on the whole. They are equals. To decide between them, it is their values in the moment which matter most.

If both reward and information are equally important, then time is not wasted in a curious search and so the process is not *necessarily* inefficient. The question becomes instead, is curiosity practical enough? This leads us to the next conjecture.

- **Conjecture 2** Curiosity is a sufficient solution to exploration (when learning is the goal)

Deterministic versus stochastic curiosity in a random world

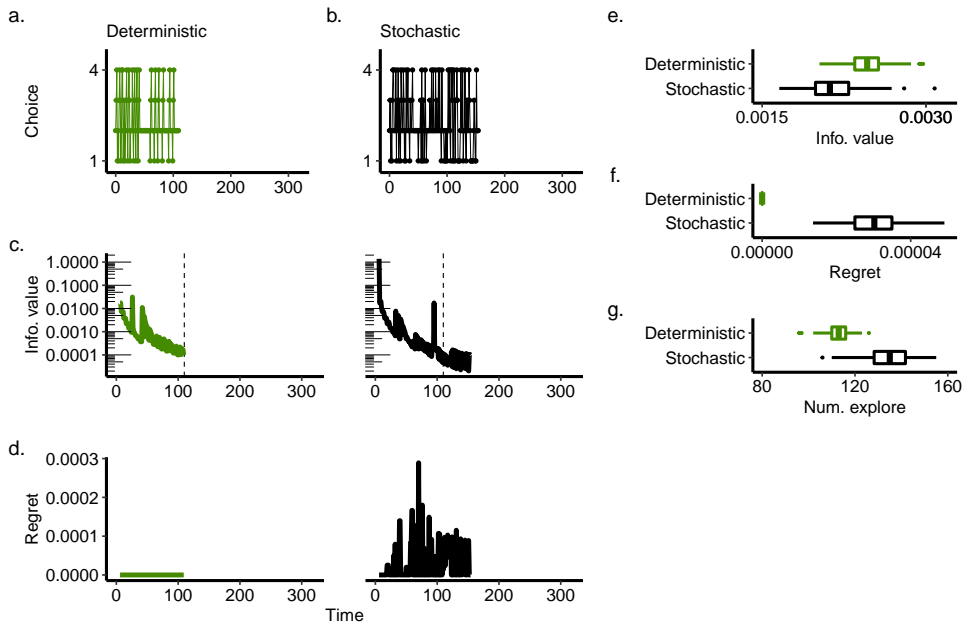


Figure 4. Comparing deterministic versus stochastic variations of the same curiosity algorithm, in a simple information foraging task (Task 1). Deterministic results are shown in the left column, and stochastic are shown in the right. The hyperparameters for both models were found by random search, which is described in the Methods. **a-b.** Examples of choice behavior. **c-d.** Information value plotted with time for the behavior shown in a-b. **e-f.** Regret plotted with time for the behavior shown in a-b. Note how our only deterministic curiosity generates zero regret, inline with theoretical predictions. **g.** Average information value for 100 independent simulations. Large values mean a more efficient search. **h.** Average regret for 100 independent simulations. Ideal exploration should have no regret. **i.** Number of steps it took to reach the boredom threshold ϵ . Smaller values imply a faster search.

241 We detail why this is the case in the next few sections.

242 **A win-stay, lose-switch solution**

243 If reward and information are equally important, how can an animal go about balancing them?
244 If you accept our conjectures, then answering this last question is the same as answering the
245 dilemma itself.

246 We posit that, based on our framing of E , this boils down to a scheduling problem. When to
247 shift from maximizing information to maximizing rewards, and vice versa?

248 Here we propose that a win-stay, lose-switch strategy is in fact an optimal and tractable solution
249 to this scheduling problem.

250 Win-stay lose-shift is a strategy famous from game theory (*Nowak and Sigmund, 1993*), where a
251 player will keep its current strategy under winning conditions, and shift otherwise. We see the same
252 kind of patterning as useful here, except the “players” are two different policies inside the animal.
253 They are competing for behavioral control. Imagine that in the previous action our bee from Fig.
254 1 observed 0.7 units of information value, E , and no reward (i.e. 0). Using our rule on the next
255 action, the bee should then choose its exploration policy. Let us say that this time E decreases to
256 0.6, and a reward value of 1.0 is observed. Now the bee should change its strategy for the next
257 round, to exploit instead. In general then if there was more reward value than information value
258 last round ($R_{t-1} > E_{t-1}$), our rule dictates an animal should choose the reward policy. If information
259 value dominated, $E_{t-1} \geq R_{t-1}$, then it should choose the information gathering policy.

260 To prove the optimality of our rule we will use regret minimization, which is a standard metric
261 of optimality in reinforcement learning (*Sutton and Barto, 2018*). Regret is a numerical quantity
262 of the difference between the most valuable choice and the value of the choice that was made.
263 To approximate regret G we use Eq. 5, where V is the value of the chosen action, and V^* is the
264 maximum value.

$$G = V^* - V \quad (5)$$

265 An optimal value algorithm always makes the most valuable choice. It therefore will have zero
266 regret, by Eq. 5. This is impossible to achieve, of course, using stochastic search. To find the al-
267 gorithm that maximizes both information and reward value with zero regret, we therefore need a
268 deterministic method.

269 To prove things about our solution we will move from game theory, which gave us our inspi-
270 ration, to the field of optimal scheduling (*Bellmann, 1954; Roughgarden, 2019*). Put another way,
271 we will now answer the question, when does win-stay, lose-switch strategy work as an optimal
272 scheduler?

273 We imagine the policies for exploration and exploitation acting as two possible “jobs” compet-
274 ing for control of behavior, a fixed resource. We know by definition that each of these jobs pro-
275 duces non-negative values which an optimal job scheduler could use: E for information or R for
276 reward/reinforcement learning. We can also ensure, again by definition, that each of the jobs takes
277 a constant amount of time and that each policy can only take one action at a time. These two prop-
278 erties, and one further assumption, are sufficient.

279 We believe it is also reasonable to assume that there will be no model of the environment avail-
280 able to the learner at the start of a problem, and that its environment is nonlinear but deterministic.

281 We arrived at our solution by taking the simplest possible path available. We just wrote down
282 the simplest equation that could have zero regret answer, and began to study that. This is Eq. 7.
283 We refer to it as a meta-greedy algorithm, as it decides in greedy fashion which of our independent
284 policies to use.

$$\pi^\pi = \left[\pi_E, \pi_R \right] \quad (6)$$

Box 5. Reward homeostasis.

We have been studying a kind of reinforcement learning where the value of a reward is fixed. We'll call this the economic style. It is the most common kind of model found in machine learning, psychology, neuroeconomics, and neuroscience (*Sutton and Barto, 2018*). In natural behavior though it is common that the motivational value of reward declines as the animal gets "full", or reaches satiety. We'll call this the homeostatic style (*Keramati and Gutkin, 2014; Juechems and Summerfield, 2019; Münch et al., 2020*). Fortunately it has already been proven that a reward collection policy designed for one style will work for the other (at least this is so an infinite time-horizon (*Keramati and Gutkin, 2014*)). However to use Eq. 7 for reward homeostasis requires that we make one modification. Ties between values should be broken in favor of exploration, and information seeking. This is in contrast to the economic reinforcement learning approach where ties break in favor of exploiting rewards. This modification leads to Eq. 8.

$$\operatorname{argmax}_{\pi^{\pi}} \left[E_{t-1}, R_{t-1} \right]_{(1, 2)} \quad (8)$$

$$\operatorname{argmax}_{\pi^{\pi}} \left[E_{t-1} - \eta, R_{t-1} \right]_{(2, 1)} \quad (7)$$

We have modified traditional argmax notation in Eq. 7. It now includes a subscript vector, (2, 1) to denote a ranking for which value/policy should be preferred in the event of a tie. Smaller numbers are preferred. We did this because an exact rule for tie breaking is critical for our proofs. In Eq. 7 we also limit the probability of having zero reward to be, itself, non-zero. That is, $p(R_t = 0) > 0$. We also limit E_0 , the first set of values, to be positive and non-zero when used with boredom, $(E_0 - \eta) \geq 0$. Both are necessary to ensure exploration.

In Theorem 5 we prove Eq. 7 is Bellman optimal, and so it will always find a maximum value sequence of E and R values. Given, that is, that the environment is deterministic. In Theorem 5 we also prove that with a judicious choice in boredom η , the reward policy π_R will also converge to its optimal value, assuming that it has one. Full details for these proofs are found in the mathematical appendix.

Reward collection

We will now focus on practical performance in reward collection. The general form of the 7 tasks we will study is depicted in Fig 5. As with our first task (Fig.) they are all variations of the classic multi-armed bandit (*Sutton and Barto, 2018*). The payouts for each of the tasks are shown separately in Fig. 6. Each was designed to either test exploration performance in a way that matches recent experimental study(s), or to test the limits of curiosity. Every trial has a set of n choices. Each choice returns a "payout", according to a predetermined probability. Payouts are information, a reward, or both. Note that, as in Task 1, information was symbolic, denoted by a color code, "yellow" and "blue" and as is custom reward was a positive real number. To evaluate the performance of our curiosity algorithm, we compare it against a range of other exploration strategies found in the literature. For a brief description of each, see Table 1. We feel this set is a broad and fair representation of today's state-of-the-art approaches.

They all have in common that their central goal is to maximize reward value, though this is often supplemented by some other goal. The results in full for all tasks and exploration strategies are shown in Fig. ???. All learners, including ours, used the same exploitation policy, based on the temporal difference learning rule (*Sutton and Barto, 2018*).

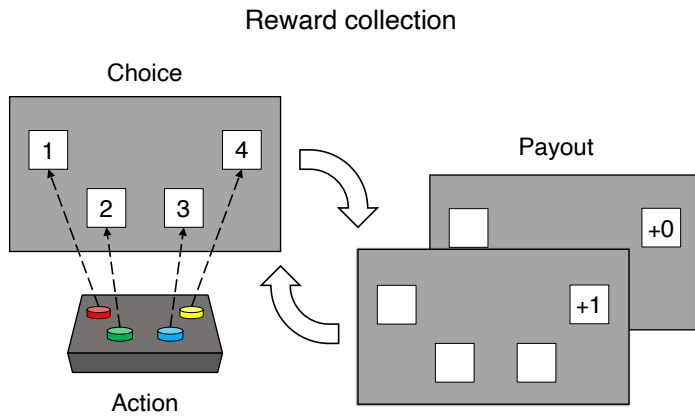


Figure 5. A 4 choice reward collection task. The reward is numerical value, here a 1 or 0. A good learner in this task is one who collects the most rewards. The task depicted here matches that in Fig 6b. Every other reward collection task we study has the same basic form, only the number of choices increases and the reward spaces are more complex.

Table 1. Exploration strategies.

Name	Class	Exploration strategy
Curiosity	Deterministic	Maximize information value
Random/Greedy	Random	Alternates between random exploration and greedy with probability ϵ .
Decay/Greedy	Random	The ϵ parameter decays with a half-life τ
Random	Random	Pure random exploration
Reward	Extrinsic	Softmax sampling of reward value
Bayesian	Extrinsic + Intrinsic	Sampling of reward value + information value
Novelty	Extrinsic + Intrinsic	Sampling of reward value + novelty signal
Entropy	Extrinsic + Intrinsic	Sampling of reward value + action entropy
Count (EB)	Extrinsic + Intrinsic	Sampling of reward value + visit counts
Count (UCB)	Extrinsic + Intrinsic	Sampling of reward value + visit counts

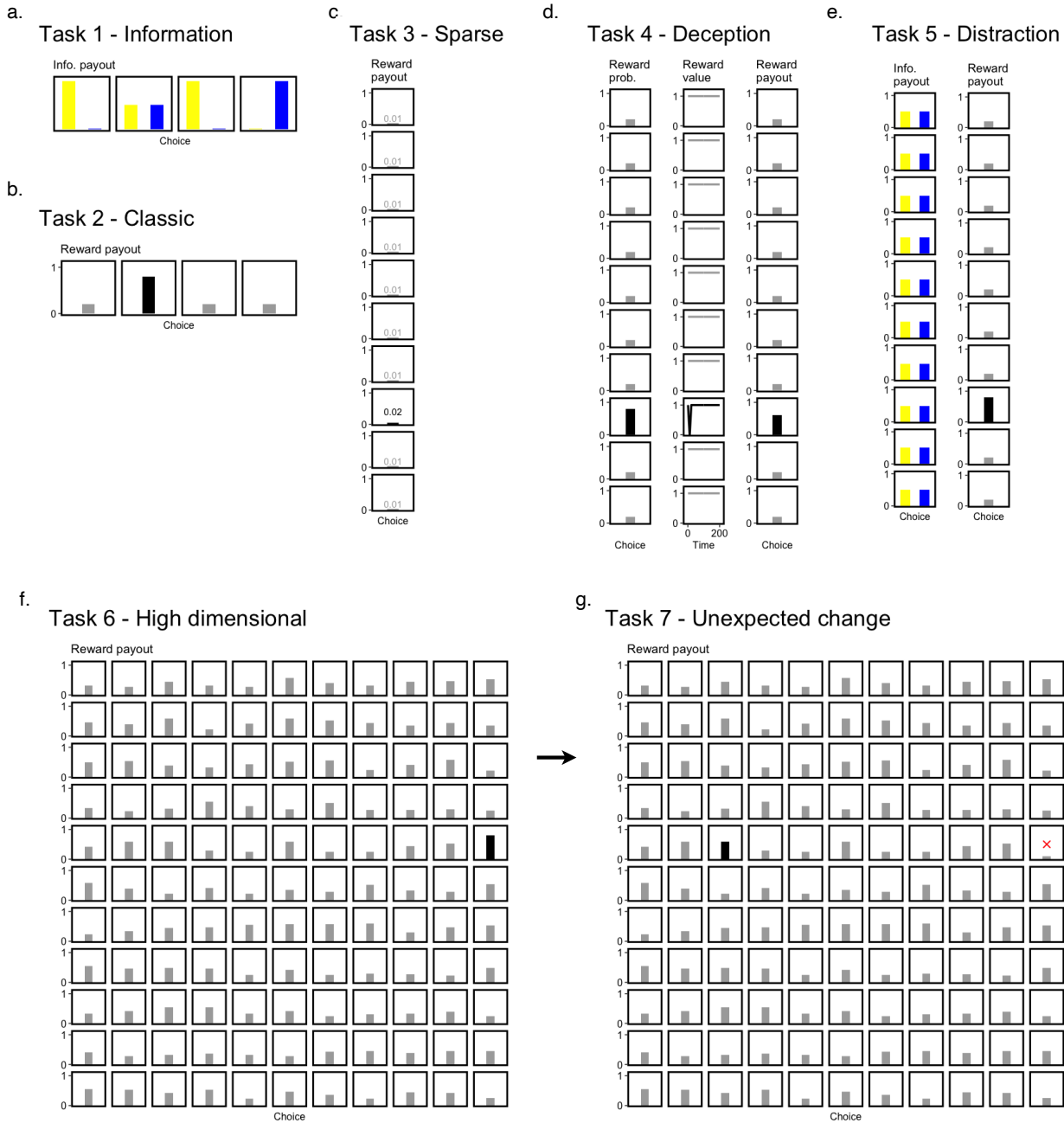
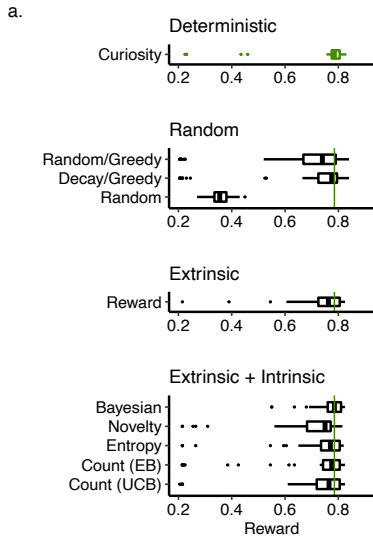
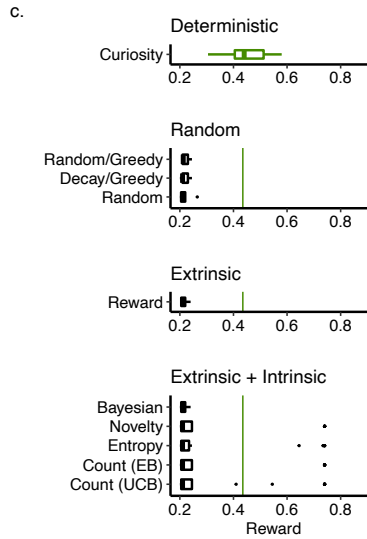


Figure 6. Payouts for Tasks 1 - 7. Payouts can be information, reward, or both. For comments on general task design, see Fig 5. **a.** A classic four-choice design for information collection. A good learner should visit each arm, but quickly discover that only arm two is information bearing. **b.** A classic four-choice design for testing exploration under reward collection. The learner is presented with four choices and it must discover which choice yields the highest average reward. In this task that is Choice 2. **c.** A ten choice sparse reward task. The learner is presented with four choices and it must discover which choice yields the highest average reward. In this task that is Choice 8 but the very low overall rate of rewards makes this difficult to discover. Solving this task with consistency means consistent exploration. **d.** A ten choice deceptive reward task. The learner is presented with 10 choices, but the choice which is the best on the long-term (>30 trials) has a lower value in the short term. This value first declines, then rises (see column 2). **e.** A ten choice information distraction task. The learner is presented with both information and rewards. A good learner though will realize the information does not predict reward collection, and so will ignore it. **f.** A 121 choice task with a complex payout structure. This task is thought to be at the limit of human performance. A good learner will eventually discover choice number 57 has the highest payout. **g.** This task is identical to **a.**, except for the high payout choice being changed to be the lowest possible payout. This task tests how well different exploration strategies adjust to simple but sudden change in the environment.

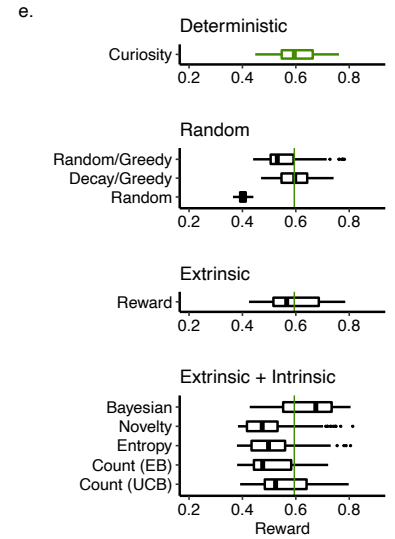
Task 2 – Classic



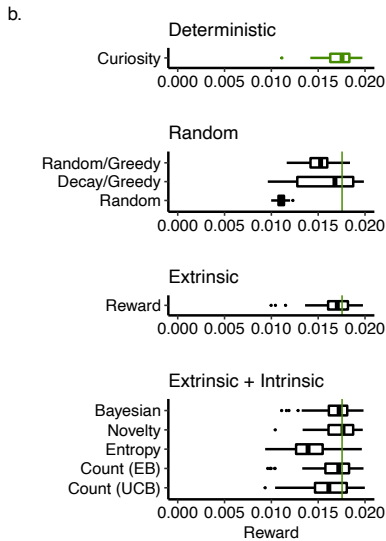
Task 4 – Deception



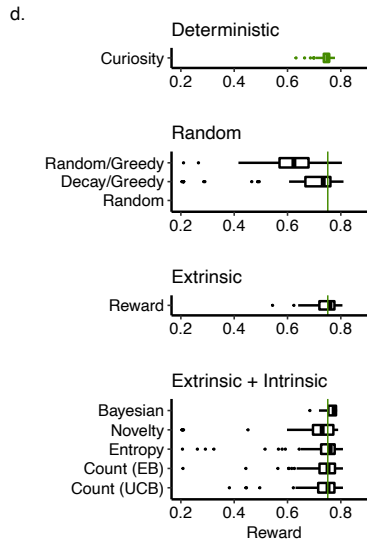
Task 6 – High dimensional



Task 3 – Sparse



Task 5 – Distraction



Task 7 – Unexpected change

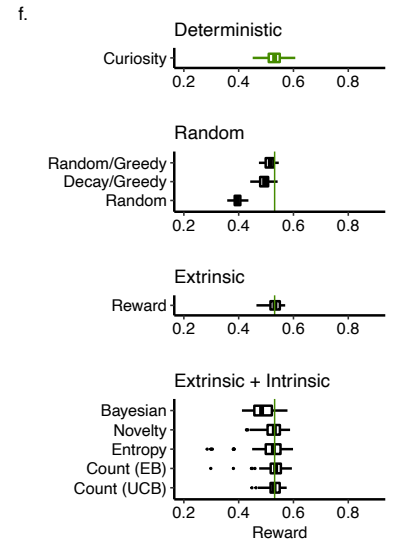


Figure 7. Summary of reward collection (*Tasks 2-7*). The strategies in each panel are grouped according to the class of search they employed (Curiosity, Random, Extrinsic reward or Extrinsic + Intrinsic rewards). **a.** Results for Task 2, which has four choices and one clear best choice. **b.** Results for Task 3, which has 10 choices and very sparse positive returns. **c.** Results for Task 4, whose best choice is initially “deceptive” in that it returns suboptimal reward value over the first 20 trials. **d.** Results for Task 5, which blends the information foraging task 1 with a larger version of Task 2. The yellow/blue stimuli are a max entropy distraction which do not predict the reward payout of each arm. **e.** Results for Task 6, which has 121 choices and a quite heterogeneous set of payouts but still with one best choice. **f.** Results for Task 7, which is identical to Task 6 except the best choice was changed to be the worst. The learners from Task 6 were trained on this Task beginning with the learned values from their prior experience – a test of robustness to sudden change in the environment. *Note:* To accommodate the fact that different tasks were run different numbers of trials, we normalized total reward by trial number. This lets us plot reward collection results for all tasks on the same scale.

Figure 7-Figure supplement 1. Summary of regret (*Tasks 2-7*)

Figure 7-Figure supplement 2. Summary of information value (*Tasks 2-7*)

Figure 7-Figure supplement 3. Examples of exploration (Task 1). Here we simulated different sets of hyper-parameters on the same task and random seed. We view each parameter set as a unique “animal” (*Prescott et al., 2006*). So this figure estimates the degree and kind of variability we might expect in a natural experiment.

Task 1. is an information gathering task, which we discussed already above. *Task 2* was designed to examine reward collection, in probabilistic reward setting very common in the decision making literature (Schönberg et al., 2007; Frank et al., 2004; Cavanagh et al., 2014; Jahfari et al., 2019; Collins and Frank, 2014; Collins et al., 2017; Gläscher et al., 2010). Rewards were 0 or 1. The best choice had a payout of $p(R = 1) = 0.8$. This is a much higher average payout than the others ($p(R = 1) = 0.2$). At no point does the task generate symbolic information. See, Fig. ??b.

The results for Task 1 were discussed above. As for Task 2, we expected all the exploration strategies to succeed. While this was indeed the case, our deterministic curiosity was the top-performer in terms of median rewards collected, though by a small margin (Fig ??a).

Task 3. was designed with very sparse rewards (Silver et al., 2016, 2018) and there were 10 choices, making this a more difficult task (Fig. ??c). Sparse rewards are a common problem in the natural world, where an animal may have to travel and search between each meal. This is a difficult but not impossible task for vertebrates (Anderson, 1984) and invertebrates (Westphal et al., 2006). That being said, most reinforcement learning algorithms will struggle in this setting because the thing they need to learn, that is rewards, are often absent. In this task we saw quite a bit more variation in performance, with the novelty-bonus strategy taking the top slot (this is the only time it does so).

Task 4. was designed with deceptive rewards. By deceptive we mean that the best long-term option presents itself initially with a decreasing reward value (Fig. ??d). Such small deceptions abound in many natural contexts where one must often make short-term sacrifices (Internicola and Harder, 2012). It is well known that classic reinforcement learning will often struggle in this kind of task (Lehman and Stanley, 2011b; Sutton and Barto, 2018). Here our deterministic curiosity is the only strategy that reaches above chance performance. Median performance of all other strategies are similar to the random control (Fig ??c).

Task 5 was designed to fool curiosity, our algorithm, by presenting information that was utterly irrelevant to reward collection, but had very high entropy and so “interesting” to our algorithm. We fully anticipated that this context would fool just our algorithm, with all other strategies performing well since they are largely agnostic to entropy. However, despite being designed to fail, deterministic curiosity still produced a competitive performance to the other strategies (Fig ??d).

Tasks 6-7 were designed as a pair, with both having 121 choices, and a complex payout structure. Tasks of this size are at the limit of human performance (Wu et al., 2018). We first trained all learners on Task 6, then tested them in Task 7 which identical to 6, except the best payout arm is reset to be worst (Fig. ??e-f). In other words Tasks 6 and 7 were joined to measure learning in a high dimensional, but shifting, environments.

In Task 6 deterministic curiosity performed well, securing a second place finish. We note the Bayesian strategy outperformed our approach. However, under the sudden non-stationarity in reward when switching tasks, the top Bayesian model became the worst on Task 7, and deterministic curiosity took the top spot. Compare Fig. ??e to f. This is the robustness that we’d expect for any curiosity algorithm, whose main goal is to learn everything unbiased by other objectives. The environment and the objectives do change in real life, and so we must be prepared for this. Note how the other less-biased-toward-reward-value exploration models (count-based, novelty, and entropy models) also saw gains, to a lesser degree.

Robustness

A good choice of boredom *eta* is critical for our definition of curiosity, and was optimized carefully. We show an example of this importance in Fig. 8. On the left hand side we did a random search over 10000 possible parameters to find a value for η that produced excellent results. (This was the same kind of search we did to achieve the results shown in Fig. ??). On the right hand side we choose a value for boredom we knew would lead to a poor relative result. We then ran the 100 different randomly generated simulations which gives us the results seen in Fig. 8).

Optimizing boredom for reward collection

100 simulated experiments

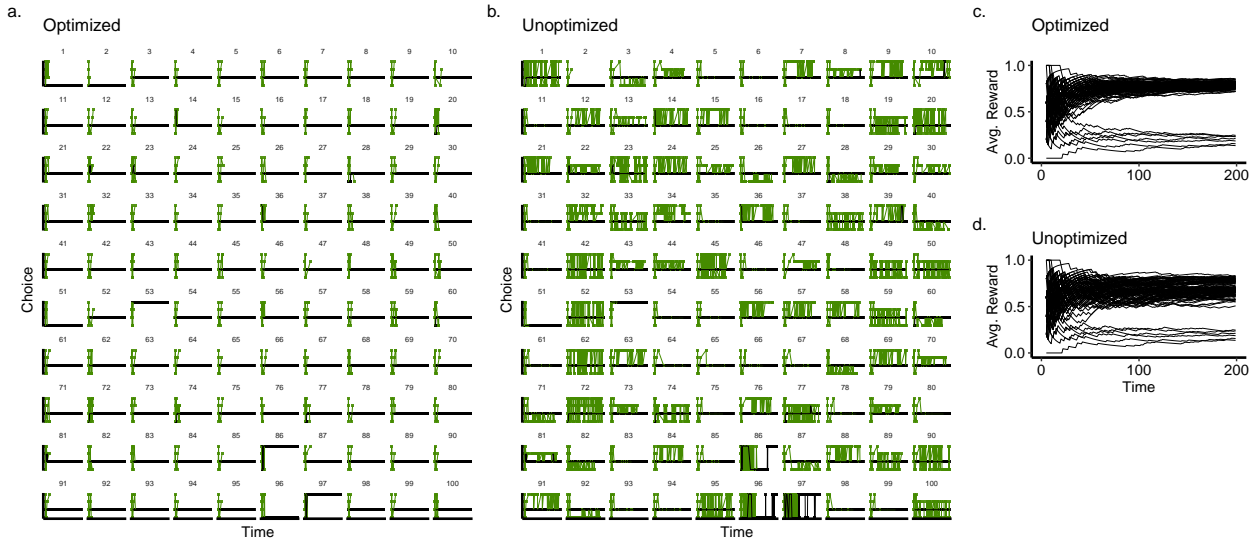


Figure 8. The importance of boredom during reward collection (Task 2). One hundred example experiments, each began with a different random seed. **a.** Behavioral choices with optimized boredom. **b.** Behavioral choices with unoptimized boredom. **c,d.** Average reward as a function of time for optimized (c) and unoptimized (d) boredom.

Carefully optimized boredom converged far more quickly (Fig. 8a-b), generating more reward over time (Fig. 8c-d), and showed far more in consistency during exploration (Fig. 8e-f).

These experiments in setting η are also useful in demonstrating the wide variance in behavior that is possible with a deterministic search. Exploration in these figures is deterministic, as prescribed by Eq. 7. The environment was stochastic however and so the variations come from the environment changing what is learned, which then changes what the best learner should choose to learn next.

Unlike idealised simulations, animals cannot pre-optimize their search parameters for every task. We therefore explored reward collection as a function of 1000 randomly chosen model parameters, and reexamined performance on Tasks 1 and 7=6. These were chosen to represent an “easy” task, and a “hard one”. These results are shown in Fig. 2.

As we observed in our other model comparisons, exploitation by deterministic curiosity produced top-3 performance on both tasks (Fig. 2a-b). Most interesting is the performance for the worst parameters. Here deterministic curiosity was markedly better, with substantially less variability across different parameter schemes. We can explain this in the following way. All the other exploration strategies use parameters to tune the degree of exploration. With deterministic curiosity, however, we tune only when exploration should stop. The degree of exploration, the measure by how close it is to maximum entropy, is fixed by the model itself. It seems that here at least it is far more robust to tune the stopping point, rather than the degree of exploration.

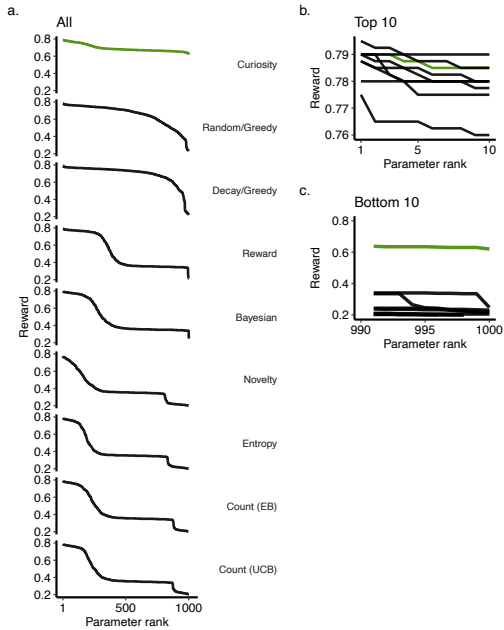
Discussion

We have offered a way around the exploration-exploitation dilemma. We proved the classic dilemma, found when optimizing exploration for reward value, can be resolved by exploring instead for curiosity. And that in this form the rational trade-off we are left with is solved with a classic strategy taken from game theory. The win-stay lose-shift strategy, that is.

This is a controversial claim. So we will address some of its criticisms as a series of questions.

Reward collection and parameter choice

Task 1 – Classic



Task 6 – High dimensional

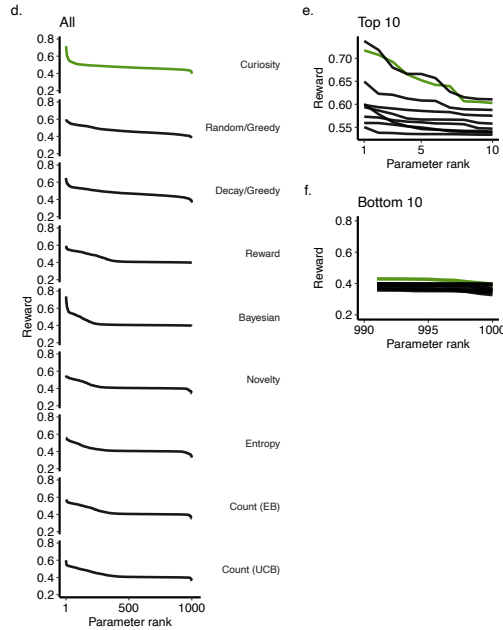


Figure 9. Parameter selection and search performance (*Task 2* and *6*). We wanted to evaluate how robust performance was to poor and random hyperparameters. A very robustness exploration algorithm would produce strong performance with both the best parameter choices, and the worst parameter choices. **a,d** Total reward for 1000 search hyperparameters, for each of strategies we considered. **b,e** Performance with the top 10 hyperparameters, curiosity (green) compared to all the others. **c,f** Performance with the bottom 10 hyperparameters.

Figure 9–Figure supplement 1. Adding noise to deterministic curiosity does degrade its performance. In this example control we added two levels of random noise. This only served to increase variability of convergence (a-c) and reduce the total rewards collected (d). Note: decision noise was modelled by sampling a Boltzmann distribution, as described in the *Methods*.

Figure 9–Figure supplement 2. Noise across strategies. Adding noise to degrades performance to nearly the same degree (*Task 2*).

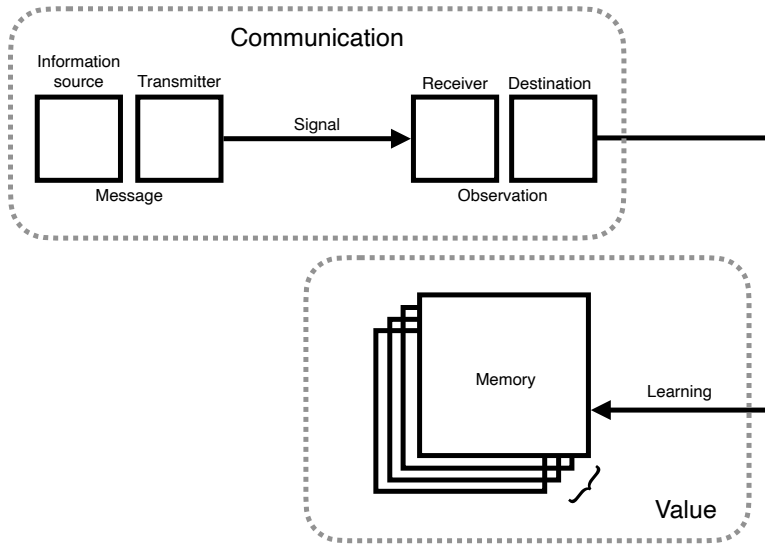


Figure 10. The relationship between the technical problem of communication with the technical problem of value. Note how value is not derived from the channel directly. Value is derived from learning about observations from the channel, which is in turn dependent on memory and its past history.

What about information theory?

In short, the problems of communication and value are different problems that require different theories.

Weaver (*Shannon, 1948*) in his classic introduction to the topic, describes information as a communication problem with three levels. 1. The technical problem of transmitting accurately. 2. The semantic problem of meaning. 3. The effectiveness problem of changing behavior. He then describes how Shannon's work addresses only problem A. It has turned out that Shannon's separation of the problem allowed the information theory to have an incredibly broad application.

Unfortunately valuing information is not, at its core, a communications problem. Consider the very simple example Bob and Alice, having a phone conversation and then Tom and Bob having the exact same conversation. The personal history of Bob and Alice will determine what Bob learns in their phone call, and so determines what he values in that phone call. These might be completely different from what he learns when talking to Tom, even if the conversations are identical. What we mean by this example is that the personal history defines what is valuable on a channel, and this is independent of the channel and the messages sent. We have summarized this diagrammatically, in Fig. 10.

There is, we argue, an analogous set of levels for information value as those Weaver describes. a. There is the technical problem of judging how much was learned. b. There is the semantic problem of both what this learning "means", and also what its consequences are. c. There is an effectiveness problem of using what was learned to some other effect. In many ways b and c are similar to those in the information theory. It is a. which is most different.

For the same reasons Shannon avoided b and c in his theory, we have avoided them also. Both are very hard to measure, which makes them very hard to study. We feel it necessary to first build a theory of information value that can act as a companion for technical problem of information theory. That is to say, Shannon's account. This is why we took an axiomatic approach, and why we have tried to ape Shannon's axioms, as it made sense to do so.

431 **Is this too complex?**

432 Perhaps turning a single objective into two, as we have, is too complex an answer. If this is true,
433 then it would mean our strategy is not a parsimonious solution to the dilemma, and so we could
434 or should reject it on that alone.

435 Questions about parsimony are resolved by considering the benefits versus the costs. The ben-
436 efits to curiosity-based search is that it leads to no regret solutions to exploration, to exploration-
437 exploitation, and that it so far seems to do very well in practice, at reward collection. At the same
438 time curiosity-as-exploration can also be building a model of the environment, useful for later
439 planning (*Ahilan et al., 2019; Poucet, 1993*), creativity, imagination (*Schmidhuber, 2010*), while also
440 building, in some cases, diverse action strategies (*Lehman and Stanley, 2011b; Lehman et al., 2013;*
441 *Mouret and Clune, 2015; Colas et al., 2020a*). In other words, we argue it is a necessary complexity.

442 **Is this too simple?**

443 Solutions to the regular dilemma are complex and require sophisticated mathematical efforts and
444 complex computations, when compared to our solution. Put another way, our solution may seem
445 too simple to some. So have we “cheated” by changing the problem in the way that we have?

446 The truth is we might have cheated, in this sense: the dilemma might have to be as hard as it
447 has seemed. But the general case for curiosity as useful in exploration is clear, and backed up by
448 the brute fact of its widespread presence. The question is: is curiosity so useful and so robust that
449 it can be sufficient for all exploration with learning.

450 The answer to this question is empirical. If our account does well in describing and predicting
451 animal behavior, that would be some evidence for it (*Sumner et al., 2019; Wang and Hayden, 2019;*
452 *Jaegle et al., 2019; Gottlieb and Oudeyer, 2018; Kidd and Hayden, 2015; Berlyne, 1950; Colas et al.,*
453 *2020b; Rahnev and Denison, 2018; Wilson et al., 2020; Cogliati Dezza et al., 2017; Berger-Tal et al.,*
454 *2014*). If it predicts neural structures (*Cisek, 2019; Kobayashi and Hsu, 2019*), that would be some
455 evidence for it (we discuss this more below). If this theory proves useful in machine learning and
456 artificial intelligence research, that would be some evidence for it (*Burda et al., 2018; Schmidhu-*
457 *ber, 1991; de Abril and Kanai, 2018; Fister et al., 2019; Lehman and Stanley, 2011b; Stanley and*
458 *Miikkulainen, 2004a; Colas et al., 2020a; Cully et al., 2015; Wilson et al., 2020; Pathak et al., 2019*).
459 These are empirical predictions that require further investigation.

460 **Is this a slight of hand, theoretically?**

461 Yes. It is often useful in mathematical studies to take one problem that cannot be solved, and
462 replace it with another related problem that can be. This is what we have done for exploration and
463 the dilemma. The regular view of the problem has no tractable solution, without regret. Ours has
464 a solution, with no regret.

465 **Does value as a technical problem even make sense?**

466 Information value has been taken up as a problem in meaning (*Kolchinsky and Wolpert, 2018*), or
467 usefulness (*Tishby et al., 2000*). These are based on what we will call *b-type* questions (a notion we
468 defined in the section above). It is not that *b* questions are not important or are not valuable. It is
469 that they are second order questions. The first order questions are the ‘technical’ or *a-type* ques-
470 tions (also defined above). Having a technical definition of value, free from any meaning, might
471 seem counterintuitive for a value measure. We argue that it is not more or less counterintuitive
472 than stripping information of meaning in the first place. Indeed, this is how Shannon got his infor-
473 mation theory. The central question for our account of value is, is it useful? It is enough to ensure a
474 complete but perfectly efficient search for information/learning in any finite space. It is enough to
475 play a critical part in solving/replacing the dilemma, which has for many years looked intractable.

476 **What about mutual information, or truth?**

477 In other prototypical examples, information value is based on how well what is learned relates to
 478 the environment (*Behrens et al., 2007; Kolchinsky and Wolpert, 2018; Tishby et al., 2000*), that is the
 479 mutual information the internal memory and the external environment. Colloquially, one might
 480 call this “truth”. The larger the mutual information between the environment and the memory, the
 481 more value it is said that we should then expect. This view seems to us at odds with actual learning
 482 behavior, especially in humans.

483 As a people we pursue information which is fictional (*Sternisko et al., 2020*), or based on analogy
 484 (*Gentner and Holyoak, 1997*), or outright wrong (*Loftus and Hoffman, 1989*). Conspiracy theories,
 485 disinformation, and our many more mundane errors, are far too commonplace for value to be
 486 based on fidelity alone. This does not mean that holding false beliefs cannot harm survival, but
 487 this is a second order question as far as learning goes. The fact that humans consistently seek to
 488 learn false things calls out us to accept that learning alone is a motivation for behavior.

489 **So you suppose there is always positive value for learning of all fictions, disinfor-** 490 **mation, and deceptions?**

491 We must. This is our most unexpected prediction. Yet, logically, it is internally consistent with our
 492 work we believe qualitatively consistent with the behavior of humans and animals alike.

493 **What about information foraging?**

494 Our work is heavily inspired by information foraging (*Inglis et al., 2001; Reddy et al., 2016*). Our
 495 motivation in developing our novel strategy was that information value should not depend on only
 496 a probabilistic reasoning, as it does in information foraging. This is for the simple reason that
 497 many of the problems animals need to learn are not probabilistic problems *from their point of*
 498 *view*. Of course, most any learning problem can be described as probabilistic from an outsider’s
 499 perspective; the one scientists’ rely on. But that means probabilistic approaches are descriptive
 500 approximations. We also feel curiosity is so important it needed a theoretical account that is both
 501 mechanistic and descriptive at the same time. This is why we developed an axiomatic approach
 502 that ended up addressing the technical problem of information value.

503 **Was it necessary to build a general theory for information value just to describe** 504 **curiosity?**

505 No. It was an idealistic choice that worked out. The field of curiosity studies has shown that there
 506 are many kinds of curiosity. At the extreme limit of this diversity is a notion of curiosity defined
 507 for any kind of observation, and any kind of learning. If we were able to develop a theory of value
 508 driven search at this limit, we can be sure it will apply in all possible examples of curiosity that we
 509 seem sure to discover in future. At this limit we can also be sure that the ideas will span fields,
 510 addressing the problem as it appears in computer science, psychology, neuroscience, biology, and
 511 perhaps economics.

512 **Doesn’t your definition of regret ignore lost rewards when the curiosity policy is in** 513 **control?**

514 Yes. Regret is calculated for the policy which is in control of behavior, not the policy which *could*
 515 have been in control of behavior.

516 **Is information a reward?**

517 It depends. If reward is defined as more-or-less any quantity that motivates behavior, then our
 518 definition of information value is a reward. This does not mean, however, that information value
 519 and environmental rewards are interchangeable. In particular, if you add reward value and infor-
 520 mation value to motivate search you can drive exploration in practical and useful ways. But this

doesn't change the intractability of the dilemma proper (*Thrun and Möller, 1992; Dayan and Sejnowski, 1996; Findling et al., 2018; Gershman, 2018*). So exploration will still generate substantial regret, as we show in Fig ?? . But perhaps more important is that these kinds of intrinsic rewards (*Schmidhuber, 1991; Berger-Tal et al., 2014; Itti and Baldi, 2009; Friston et al., 2016; Kobayashi and Hsu, 2019*) introduce a bias that will drive behavior away from what, could otherwise be, ideal skill learning (*Ng et al., 1999; Şimşek and Barto, 2006*).

But isn't curiosity impractical?

It does seem curiosity is just as likely to lead away from a needed solution as towards it. Especially so if one watches children who are the prototypical curious explorers (*Sumner et al., 2019; Kidd and Hayden, 2015*). This is why we limit curiosity with boredom, and counter it with a competing drive for reward collecting (i.e., exploitation).

We believe there is a useful analogy between science and engineering. Science can be considered an open-ended inquiry, and engineering a purpose driven enterprise. They each have their own pursuits but they also learn from each other, often in alternating iterations (*Gupta et al., 2006*). It is their different objectives though which make them such good long-term collaborators.

But what about curiosity on big problems?

A significant drawback to curiosity, as an algorithm, is that it will struggle to deliver timely results when the problem is very large, or highly detailed. First, it is worth noting that most learning algorithms struggle with this setting (*MacKay, 2003; Sutton and Barto, 2018*), that is unless significant prior knowledge can be brought to bear (*Zhang et al., 2020; Sutton and Barto, 2018*) or the problem can be itself compressed (*Ha and Schmidhuber, 2018; Fister et al., 2019*). Because "size" is a widespread problem in learning, there are a number of possible solutions and because our definition of learning, memory and curiosity is so general, we can take advantage of most. This does not guarantee curiosity will work out for all problems; all learning algorithms have counterexamples (*Wolpert and Macready, 1997*).

But what about other forms of curiosity?

We did not intend to dismiss the niceties of curiosity that others have revealed. It is that these prior divisions parcel out the idea too soon. Philosophy and psychology consider whether curiosity is internal or external, perceptual versus epistemic, or specific versus diversive, or kinesthetic (*Kidd and Hayden, 2015; Berlyne, 1950; Zhou et al., 2020*). Computer science tends to define it as needed, for the task at hand (*Stanley and Miikkulainen, 2004b; Friston et al., 2016; Lehman and Stanley, 2011b; Lehman et al., 2013; Mouret and Clune, 2015; Colas et al., 2020a*). Before creating a taxonomy of curiosity it is important to first define it, mathematically. Something we have done here.

What is boredom, besides being a tunable parameter?

A more complete version of the theory would let us derive a useful or even optimal value for boredom, if given a learning problem or environment. This would also answer the question of what boredom is computationally. We cannot do this. It is a next problem we will work on, and it is very important. The central challenge is deriving boredom for any kind of learning algorithm.

Does this mean you are hypothesizing that boredom is actively tuned?

We are predicting exactly that.

But do people subjectively feel they can tune their boredom?

From our experience, no. Perhaps it happens unconsciously? For example we seem to alter physiological learning rates without subjectively experiencing it (*Behrens et al., 2007*).

How can an animal tune boredom in a new setting?

The short answer is that one would need to guess and check. Our work in Fig 2 suggests that this can be somewhat robust.

Do you have any evidence in animal behavior and neural circuits?

There is some evidence for our theory of curiosity in psychology and neuroscience, however, in these fields curiosity and reinforcement learning have largely developed as separate disciplines (Berlyne, 1950; Kidd and Hayden, 2015; Sutton and Barto, 2018). Indeed, we have highlighted how they are separate problems, with links to different basic needs: gathering resources to maintain physiological homeostasis (Keramati and Gutkin, 2014; Juechems and Summerfield, 2019) and gathering information to decide what to learn and to plan for the future (Valiant, 1984; Sutton and Barto, 2018). Here we suggest that though they are separate problems, they are problems that can, in large part, solve one another. This insight is the central idea to our view of the explore-exploit decisions.

Yet there are hints of this independent cooperation of curiosity and reinforcement learning out there. Cisek (2019) has traced the evolution of perception, cognition, and action circuits from the Metazoan to the modern age (Cisek, 2019). The circuits for reward exploitation and observation-driven exploration appear to have evolved separately, and act competitively, exactly the model we suggest. In particular he notes that exploration circuits in early animals were closely tied to the primary sense organs (i.e. information) and had no input from the homeostatic circuits (Keramati and Gutkin, 2014; Cisek, 2019; Juechems and Summerfield, 2019). This neural separation for independent circuits has been observed in “modern” high animals, including zebrafish (Marques et al., 2019) and monkeys (White et al., 2019; Wang and Hayden, 2019).

What about aversive values?

We do not believe this is complete theory. For example, we do not account for any consequences of learning, or any other aversive events. Accounting for this is another important next step.

The need to add other values fits well within our notion of competitive simple policies, in a win-stay loose-shift game. The need for taking other values into account is why we choose the notation we have. For example, if we added a fear value F our method could be expanded to read,

$$\pi^\pi = [\pi_E, \pi_R, \pi_F] \quad (9)$$

$$\operatorname{argmax}_{\pi^\pi} [E_{t-1} - \eta, R_{t-1}, F_{t-1}]_{(2, 3, 1)} \quad (10)$$

Where we assume the fear/aversive term should take precedence when tied. Another more mundane example would be to add a policy for a grooming drive, G .

$$\pi^\pi = [\pi_E - \eta, \pi_R, \pi_G] \quad (11)$$

$$\operatorname{argmax}_{\pi^\pi} [E_{t-1}, R_{t-1}, G_{t-1}]_{(1, 2, 3)} \quad (12)$$

What about fitting deterministic models?

Exploration in the lab is often described as probabilistic (Calhoun et al., 2014; Song et al., 2019; Gershman, 2018; Schulz et al., 2018). By definition probabilistic models do not make exact predictions of behavior, only statistical ones. Our approach is deterministic and so can make exact predictions. If, that is, we can fit it and it turns out to be correct. In species ranging from human to *Caenorhabditis elegans*, there are hundreds of exploration-exploitation experiments. A deterministic theory can, in principle, open up entirely new avenues for reanalysis. Testing and exploring these is a critical next step.

Does regret matter?

Another way around the dilemma is to ignore regret altogether. This is the domain of pure exploration methods, like the Gitten's Index (*Gittins, 1979*), or probably-approximately correct approaches (*Valiant, 1984*), or other bounded forms of reinforcement learning (*Brafman and Tennenholtz, 2002*). Such methods can guarantee that you find the best value, eventually. These methods do not consider regret. But regret is critical for making sure that an animal's actions are efficient when resources and time are scarce, as they are in almost every natural setting.

Is the algorithmic run time practical?

It is common in computer science to study the run time of an algorithm as a measure of its efficiency. So what is the run time of our win stay, loose switch solution? There is unfortunately no fixed answer for the average or best case. The worst case algorithmic run time is linear and additive in the independent policies. If it takes T_E steps for π_E to converge, and T_R steps for π_R , then the worst case run time for π_π is $T_E + T_R$. Of course this is the worst case run time and would still be within reasonable ranges for learning in most naturalistic contexts.

Summary

There will certainly be cases in an animal's life when their exploration must focus on tangible physical rewards. Even here, searching for information about physical rewards, rather than searching for their value, will make the search more robust to deception, and more reliable in the future. We think of this as curiosity about reward. There are many more cases though where curiosity can include observations of the environment and the self. We have argued one should put aside intuitions that suggest an open-ended curious search is too inefficient to be practical. We have shown it can be made very practical. We have argued one should put aside traditional intuitions for what exploratory behavior in animals represents. We have shown how there is a zero-regret way to solve the dilemma, a problem that has seemed unsolvable. Simply be curious.

Mathematical Appendix.

Information value as a dynamic programming problem

To find a dynamic programming solution based on the Bellman equation (Eq ??) we must prove our memory Θ has optimal substructure. This is because the normal route, which assumes the problem rests in a Markov Space, is closed to us. By optimal substructure we mean that the process of learning in Θ can be partitioned into a collection, or series of memories, each of which is itself a solution. That is, it lets us find a kind of recursive structure in memory learning, which is what we need to apply the Bellman. Proving optimal substructure also allows for simple proof by induction (*Roughgarden, 2019*), a property we will take advantage of.

Theorem 1 (Optimal substructure). *Let $X, A, \Theta, f_{\Theta},$ (Def. 1), π_E and δ be given. Assuming transition function δ is deterministic, if $V_{\pi_E}^*$ is the optimal information value given by policy π_E , a memory Θ_{t+1} has optimal substructure if the the last observation x_t can be removed from Θ_t , by $\Theta_t = f^{-1}(\Theta_{t+1}, x_t)$ such that the resulting value $V_{t-1}^* = V_t^* - E_t$ is also optimal.*

Proof. Given a known optimal value V^* given by π_E we assume for the sake of contradiction there also exists an alternative policy $\hat{\pi}_E \neq \pi_E$ that gives a memory $\hat{M}_{t-1} \neq M_{t-1}$ and for which $\hat{V}_{t-1}^* > V_{t-1}^*$. To recover the known optimal memory M_t we lift \hat{M}_{t-1} to $M_t = f(\hat{M}_{t-1}, x_t)$. This implies $\hat{V}^* > V^*$ which in turn contradicts the purported original optimality of V^* and therefore $\hat{\pi}_E$. \square

This proof required two things. First, it was required we extend the idea of memory. We need to introduce a mechanism for forgetting of a very particular kind. We must assume that any last learning step $f(x, \theta) \rightarrow \theta'$ can be undone by a new vector valued function f^{-1} , such that $f(x, \theta') \rightarrow \theta$. In other words we must assume what was last remembered, can always be forgotten.

648 Second we must also assume the environment δ is deterministic. Determinism is consistent
 649 with the natural world, which does evolve in a deterministic way, at the scales we concerned with.
 650 This assumption is however at odds with much of reinforcement learning theory (?) and past ex-
 651 perimental work (?). Both tend to study stochastic environments. We'll address this discrepancy
 652 later on using numerical simulations.

653 Bellman solution

654 Knowing the optimal substructure of Θ and given an arbitrary starting value E_0 , the Bellman solu-
 655 tion to curiosity optimization of E is given by,

$$V_E(x) = E_0 + \operatorname{argmax}_a \left[E_t \right] \quad (13)$$

656 To explain this derivation we'll begin simply with the definition of a value function and work from
 657 there. A value function $V(x)$ is defined as the best possible value of the objective for x . Finding a
 658 Bellman solution is discovering what actions a to take to arrive at $V(x)$. Bellman's insight was to
 659 break the problem down into a series of smaller problems, leading a recursive solution. Problems
 660 that can be broken down this way have optimal substructure. The value function for a finite time
 661 horizon T and some payout function F is given by Eq 14. It's Bellman form is given by Eq. 15

$$V(x) = \operatorname{argmax}_a \left[\sum_T F(x, a) \right] \quad (14)$$

$$V(x) = \operatorname{argmax}_a \left[F(x_0, a_0) + V(x_1) \right] \quad (15)$$

662 In reinforcement learning the value function is based on the sum of future rewards giving Eq. 16-
 663 17,

$$V_R(x) = \operatorname{argmax}_a \left[\sum_T R_t \right] \quad (16)$$

$$V_R(x) = \operatorname{argmax}_a \left[R_0 + V(x_1) \right] \quad (17)$$

664 In our objective the best possible value is not found by temporal summation as it is during
 665 reinforcement learning. E is necessarily a temporally unstable function. We expect it to vary sub-
 666 stantially before contracting to approach 0. It's best possible value is well approximated then by
 667 only looking at the maximum value for the last time step, making our optimization myopic, that is
 668 based on only last values encountered (*Hocker and Park, 2019*).

$$V_E(x) = \operatorname{argmax}_a \left[E_t \right] \quad (18)$$

$$\begin{aligned} V_E(x) &= \operatorname{argmax}_a \left[E_0 + V(x_1) \right] \\ &= E_0 + \operatorname{argmax}_a \left[E_t \right] \end{aligned} \quad (19)$$

669 Good exploration by curiosity optimization

670 Recall from the main text we consider that a good exploration should,

- 671 1. Visit all available states of the environment at least once.
- 672 2. Should cease only once learning about the environment has plateaued.
- 673 3. Should take as few steps as possible to achieve criterion 1 and 2.

Limiting π_E to a deterministic policy makes proving these three properties amounts to solving sorting problems on E . If a state is visited by our algorithm it must have the highest value. So if every state must be visited under a deterministic greedy policy every state must, at one time or another, generate maximum value. This certain if we know that all values will begin contracting towards zero. *Violations of this assumption are catastrophic*. We discuss this limit in the Discussion but note that things are not as dire as they may seem.

Definitions. Let Z be the set of all visited states, where Z_0 is the empty set $\{\}$ and Z is built iteratively over a path P , such that $Z \rightarrow \{x|x \in X \text{ and } x \notin Z\}$.

To formalize the idea of ranking we take an algebraic approach. Give any three real numbers (a, b, c) ,

$$a \leq b \Leftrightarrow \exists c; b = a + c \quad (20)$$

$$a > b \Leftrightarrow (a \neq b) \wedge (b \leq a) \quad (21)$$

Theorem 2 (State search: breadth). *Given some arbitrary value E_0 , an exploration policy governed by π_E^* will visit all states $x \in X$ in finite number of steps T .*

Proof. Let $\mathbf{E} = (E_1, E_2, \dots)$ be ranked series of E values for all states X , such that $(E_1 \geq E_2, \geq \dots)$. To swap any pair of values $(E_i \geq E_j)$ so $(E_i \leq E_j)$ by Eq. 20 $E_i - c = E_j$.

Therefore, again by Eq. 20, $\exists \int \delta E(s) \rightarrow -c$.

Recall: $\nabla^2 f_\Theta(x) < 0$ after a finite time T .

However if we wished to instead swap $(E_i \leq E_j)$ so $(E_i \geq E_j)$ by definition $\nexists c; E_i + c = E_j$, as $\nexists \int \delta \rightarrow c$.

To complete the proof, assume that some policy $\hat{\pi}_E \neq \pi_E^*$. By definition policy $\hat{\pi}_E$ can be any action but the maximum, leaving $k - 1$ options. Eventually as $t \rightarrow T$ the only possible swap is between the max option and the k th, but as we have already proven this is impossible as long as Axiom 2 holds. Therefore, the policy $\hat{\pi}_E$ will leave at least 1 option unexplored and $S \neq Z$. \square

Theorem 3 (State search: depth). *An exploration policy governed by π_E^* will only revisit states for where learning is possible.*

Proof. *Recall:* Axiom 2. Each time π_E^* visits a state s , so $\Theta \rightarrow \Theta', F(\Theta', a_{t+dt}) < F(\Theta, a_t)$

In Theorem 2 we proved only a deterministic greedy policy will visit each state in X in T trials.

By induction, if $\pi^* E$ will visit all $x \in X$ in T trials, it will revisit them at most $2T$, therefore as $T \rightarrow \infty$, $E \rightarrow \eta$. \square

We assume in the above the action policy π_E can visit all possible states in X . If for some reason π_E can only visit a subset of X then proofs above apply only to that subset.

These proofs come with some fine print. E_0 can be any positive and finite real number, $E_0 > 0$. Different choices for E_0 will not change the proofs, especially their convergence. So in that sense one can chosen it in an arbitrary way. Different choices for E_0 can however change individual choices, and their order. This can be quite important in practice, especially when trying to describe some real data. This choice will not change the algorithm's long-term behavior *but* different choices for E_0 will strongly change the algorithm's short-term behavior. This can be quite important in practice.

Optimality of π_π

Recall that in the main text we introduce the equation below as a candidate with zero regret solution to the exploration-exploitation dilemma.

$$\pi^\pi = \left[\pi_E, \pi_R \right] \quad (22)$$

$$\operatorname{argmax}_{\pi_\pi} \left[E_{t-1} - \eta, R_{t-1} \right]_{(2, 1)} \quad (23)$$

In the following section we prove two things about the optimality of π_π . First, if π_R had any optimal asymptotic property for value learning before their inclusion into our scheduler, they retain that optimal property under π_π when $\eta = 0$, or is otherwise sufficiently small. Second, show that if both π_R and π_E are greedy, and π_π is greedy in its definition, then Eq 22 is certain to maximize total value. The total value of R and E is the exact quantity to maximize if information seeking and reward seeking are equally important, overall. This is, as the reader may recall, one of our key assumptions. Proving this optimality is analogous to the classic activity selection problem from the job scheduling literature (Roughgarden, 2019).

Theorem 4 (π_π is unbiased). *Let any $S, A, \Theta, \pi_R, \pi_E$, and δ be given. Assuming an infinite time horizon, if π_E is optimal and π_R is optimal, then π_π is also optimal in the same sense as π_E and π_R .*

Proof. The optimality of π_π can be seen by direct inspection. If $p(R = 0) > 0$ we are given an infinite horizon, then π_E will have a unbounded number of trials meaning the optimality of P^* holds. Likewise, $\sum E < \eta$ as $T \rightarrow \infty$, ensuring π_R will dominate π_π therefore π_R will asymptotically converge to optimal behavior. \square

In proving this optimality of π_π we limit the probability of a positive reward to less than one, denoted by $p(R_t = 1) < 1$. Without this constraint the reward policy π_R would always dominate π_π when rewards are certain. While this might be useful in some circumstances, from the point of view π_E it is extremely suboptimal. The model would never explore. Limiting $p(R_t = 1) < 1$ is a reasonable constraint, as rewards in the real world are rarely certain. A more naturalistic way to handle this edge case is to introduce reward satiety, or a model physiological homeostasis (Keramati and Gutkin, 2014; Juechems and Summerfield, 2019).

Optimal scheduling for dual value learning problems

In classic scheduling problems the value of any job is known ahead of time (Bellmann, 1954; Roughgarden, 2019). In our setting, this is not true. Reward value is generated by the environment, *after* taking an action. In a similar vein, information value can only be calculated *after* observing a new state. Yet Eq. 22 must make decisions *before* taking an action. If we had a perfect model of the environment, then we could predict these future values accurately with model-based control. In the general case though we don't know what environment to expect, let alone having a perfect model of it. As a result, we make a worst-case assumption: the environment can arbitrarily change–bifurcate–at any time. This is a highly nonlinear dynamical system (Strogatz, 1994). In such systems, myopic control–using only the most recent value to predict the next value– is known to be an robust and efficient form of control (Hocker and Park, 2019). We therefore assume that last value is the best predictor of the next value, and use this assumption along with Theorem 5 to complete a trivial proof that Eq. 22 maximizes total value.

Optimal total value

If we prove π_π has optimal substructure, then using the same replacement argument (Roughgarden, 2019) as in Theorem 5, a greedy policy for π_π will maximize total value.

Theorem 5 (Total value maximization of π_π). *Let any S, A, Θ, π_R , and δ be given. If π_R is defined on a Markov Decisions, then π_π is Bellman optimal and will maximize total value.*

Proof. We assume Reinforcement learning algorithms are embedded in Markov Decisions space, which by definition has the same decomposition properties as that found in optimal substructure.

Recall: The memory Θ has optimal substructure (Theorem 1).

Recall: The asymptotic behavior of π_R and π_E are independent under π_π (Theorem 5)

Recall: The controller π_π is deterministic and greedy.

If both π_R and π_E have optimal substructure and are independent, then π_π must also have optimal substructure. If π_π has optimal substructure, and is greedy-deterministic then it is Bellman optimal. \square

Acknowledgments

We wish to thank Jack Burgess, Matt Clapp, Kyle “Donovank” Dunovan, Richard Gao, Roberta Klatzky, Jayanth Koushik, Alp Muyesser, Jonathan Rubin, and Rachel Storer for their comments on earlier drafts. We also wishes to thank Richard Grant for his illustration work in Figure 1, and Jack Burgess for his assitance with Figure ??.

The research was sponsored by the Air Force Research Laboratory (AFRL/AFOSR) award FA9550-18-1-0251. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. government. TV was supported by the Pennsylvania Department of Health Formula Award SAP4100062201, and National Science Foundation CAREER Award 1351748.

References

- Ahilan S**, Solomon RB, Breton YA, Conover K, Niyogi RK, Shizgal P, Dayan P. Learning to Use Past Evidence in a Sophisticated World Model. *PLoS Comput Biol*. 2019 Jun; 15(6):e1007093. doi: [10.1371/journal.pcbi.1007093](https://doi.org/10.1371/journal.pcbi.1007093).
- Anderson O**. Optimal foraging by largemouth bass in structured environments. *Ecology*. 1984; 65(3):851–861.
- Auersperg AMI**. Exploration Technique and Technical Innovations in Corvids and Parrots. In: *Animal Creativity and Innovation* Elsevier; 2015.p. 45–72. doi: [10.1016/B978-0-12-800648-1.00003-6](https://doi.org/10.1016/B978-0-12-800648-1.00003-6).
- Behrens TEJ**, Woolrich MW, Walton ME, Rushworth MFS. Learning the Value of Information in an Uncertain World. *Nat Neurosci*. 2007 Sep; 10(9):1214–1221. doi: [10.1038/nn1954](https://doi.org/10.1038/nn1954).
- Bellmann R**. The Theory of Dynamic Programming. *Bull Amer Math Soc*. 1954; 60(6):503–515.
- Bench S**, Lench H. On the Function of Boredom. *Behavioral Sciences*. 2013 Aug; 3(3):459–472. doi: [10.3390/bs3030459](https://doi.org/10.3390/bs3030459).
- Berger-Tal O**, Nathan J, Meron E, Saltz D. The Exploration-Exploitation Dilemma: A Multidisciplinary Framework. *PLoS ONE*. 2014 Apr; 9(4):e95693. doi: [10.1371/journal.pone.0095693](https://doi.org/10.1371/journal.pone.0095693).
- Berlyne D**. Novelty and Curiosity as Determinants of Exploratory Behaviour. *British Journal of Psychology*. 1950; 41(1):68–80.
- Brafman RI**, Tennenholtz M. R-Max – A General Polynomial Time Algorithm for Near-Optimal Reinforcement Learning. *Journal of Machine Learning Research*. 2002; 2:19.
- Burda Y**, Edwards H, Pathak D, Storkey A, Darrell T, Efros AA. Large-Scale Study of Curiosity-Driven Learning. arXiv:180804355 [cs, stat]. 2018 Aug; .
- Calhoun AJ**, Chalasani SH, Sharpee TO. Maximally Informative Foraging by *Caenorhabditis Elegans*. *eLife*. 2014 Dec; 3:e04220. doi: [10.7554/eLife.04220](https://doi.org/10.7554/eLife.04220).
- Cavanagh JF**, Masters SE, Bath K, Frank MJ. Conflict acts as an implicit cost in reinforcement learning. *Nature communications*. 2014; 5(1):1–10.
- Cisek P**. Resynthesizing Behavior through Phylogenetic Refinement. *Atten Percept Psychophys*. 2019 Jun; doi: [10.3758/s13414-019-01760-1](https://doi.org/10.3758/s13414-019-01760-1).
- Cogliati Dezza I**, Yu AJ, Cleeremans A, Alexander W. Learning the Value of Information and Reward over Time When Solving Exploration-Exploitation Problems. *Sci Rep*. 2017 Dec; 7(1):16919. doi: [10.1038/s41598-017-17237-w](https://doi.org/10.1038/s41598-017-17237-w).
- Colas C**, Huizinga J, Madhavan V, Clune J. Scaling MAP-Elites to Deep Neuroevolution. *Proceedings of the 2020 Genetic and Evolutionary Computation Conference*. 2020 Jun; p. 67–75. doi: [10.1145/3377930.3390217](https://doi.org/10.1145/3377930.3390217).

- Colas C**, Karch T, Lair N, Dussoux JM, Moulin-Frier C, Dominey PF, Oudeyer PY. Language as a Cognitive Tool to Imagine Goals in Curiosity-Driven Exploration. *arXiv:200209253 [cs]*. 2020 Jun; .
- Collins AG**, Albrecht MA, Waltz JA, Gold JM, Frank MJ. Interactions among working memory, reinforcement learning, and effort in value-based choice: A new paradigm and selective deficits in schizophrenia. *Biological Psychiatry*. 2017; 82(6):431–439.
- Collins AG**, Frank MJ. Opponent actor learning (OpAL): Modeling interactive effects of striatal dopamine on reinforcement learning and choice incentive. *Psychological review*. 2014; 121(3):337.
- Cully A**, Clune J, Tarapore D, Mouret JB. Robots That Can Adapt like Animals. *Nature*. 2015 May; 521(7553):503–507. doi: 10.1038/nature14422.
- Dayan P**, Sejnowski TJ. Exploration Bonuses and Dual Control. *Mach Learn*. 1996 Oct; 25(1):5–22. doi: 10.1007/BF00115298.
- de Abril IM**, Kanai R. Curiosity-Driven Reinforcement Learning with Homeostatic Regulation. In: *2018 International Joint Conference on Neural Networks (IJCNN)* Rio de Janeiro: IEEE; 2018. p. 1–6. doi: 10.1109/IJCNN.2018.8489075.
- Findling C**, Skvortsova V, Dromnelle R, Palminteri S, Wyart V. Computational Noise in Reward-Guided Learning Drives Behavioral Variability in Volatile Environments. *Neuroscience*; 2018.
- Fister I**, Iglesias A, Galvez A, Del Ser J, Osaba E, Fister I, Perc M, Slavinec M. Novelty Search for Global Optimization. *Applied Mathematics and Computation*. 2019 Apr; 347:865–881. doi: 10.1016/j.amc.2018.11.052.
- Frank MJ**, Seeberger LC, O'reilly RC. By carrot or by stick: cognitive reinforcement learning in parkinsonism. *Science*. 2004; 306(5703):1940–1943.
- Friston K**, FitzGerald T, Rigoli F, Schwartenbeck P, O'Doherty J, Pezzulo G. Active Inference and Learning. *Neuroscience & Biobehavioral Reviews*. 2016 Sep; 68:862–879. doi: 10.1016/j.neubiorev.2016.06.022.
- Gentner D**, Holyoak KJ. Reasoning and learning by analogy: Introduction. *American psychologist*. 1997; 52(1):32.
- Gershman SJ**. Deconstructing the Human Algorithms for Exploration. *Cognition*. 2018 Apr; 173:34–42. doi: 10.1016/j.cognition.2017.12.014.
- Gittins JC**. Bandit Processes and Dynamic Allocation Indices. *Journal of the Royal Statistical Society Series B (Methodological)*. 1979; 41(2):148–177.
- Gläscher J**, Daw N, Dayan P, O'Doherty JP. States versus rewards: dissociable neural prediction error signals underlying model-based and model-free reinforcement learning. *Neuron*. 2010; 66(4):585–595.
- Gottlieb J**, Oudeyer PY. Towards a Neuroscience of Active Sampling and Curiosity. *Nat Rev Neurosci*. 2018 Dec; 19(12):758–770. doi: 10.1038/s41583-018-0078-0.
- Gupta AA**, Smith K, Shalley C. The Interplay between Exploration and Exploitation. *The Academy of Management Journal*. 2006; 49(4):693–706.
- Ha D**, Schmidhuber J. Recurrent World Models Facilitate Policy Evolution. *NeurIPS*. 2018; p. 13.
- Hocker D**, Park IM. Myopic Control of Neural Dynamics. *PLOS Computational Biology*. 2019; 15(3):24.
- Inglis IR**, Langton S, Forkman B, Lazarus J. An Information Primacy Model of Exploratory and Foraging Behaviour. *Animal Behaviour*. 2001 Sep; 62(3):543–557. doi: 10.1006/anbe.2001.1780.
- Internicola AI**, Harder LD. Bumble-bee learning selects for both early and long flowering in food-deceptive plants. *Proceedings of the Royal Society B: Biological Sciences*. 2012; 279(1733):1538–1543.
- Ishii S**, Yoshida W, Yoshimoto J. Control of Exploitation–Exploration Meta-Parameter in Reinforcement Learning. *Neural Networks*. 2002 Jun; 15(4-6):665–687. doi: 10.1016/S0893-6080(02)00056-4.
- Itti L**, Baldi P. Bayesian Surprise Attracts Human Attention. *Vision Research*. 2009 Jun; 49(10):1295–1306. doi: 10.1016/j.visres.2008.09.007.
- Jaegle A**, Mehrpour V, Rust N. Visual Novelty, Curiosity, and Intrinsic Reward in Machine Learning and the Brain. *Arxiv*. 2019; 1901.02478:13.

- Jahfari S**, Ridderinkhof KR, Collins AG, Knapen T, Waldorp LJ, Frank MJ. Cross-task contributions of frontobasal ganglia circuitry in response inhibition and conflict-induced slowing. *Cerebral Cortex*. 2019; 29(5):1969–1983.
- Juechems K**, Summerfield C. Where Does Value Come From? *PsyArXiv*; 2019.
- Kelly JL**. A New Interpretation of Information Rate. *the bell system technical journal*. 1956; p. 10.
- Keramati M**, Gutkin B. Homeostatic Reinforcement Learning for Integrating Reward Collection and Physiological Stability. *eLife*. 2014 Dec; 3:e04811. doi: [10.7554/eLife.04811](https://doi.org/10.7554/eLife.04811).
- Kidd C**, Hayden BY. The Psychology and Neuroscience of Curiosity. *Neuron*. 2015 Nov; 88(3):449–460. doi: [10.1016/j.neuron.2015.09.010](https://doi.org/10.1016/j.neuron.2015.09.010).
- Kobayashi K**, Hsu M. Common Neural Code for Reward and Information Value. *Proc Natl Acad Sci USA*. 2019 Jun; 116(26):13061–13066. doi: [10.1073/pnas.1820145116](https://doi.org/10.1073/pnas.1820145116).
- Kolchinsky A**, Wolpert DH. Semantic Information, Autonomous Agency and Non-Equilibrium Statistical Physics. *Interface Focus*. 2018 Dec; 8(6):20180041. doi: [10.1098/rsfs.2018.0041](https://doi.org/10.1098/rsfs.2018.0041).
- Lee MD**, Zhang S, Munro M, Steyvers M. Psychological Models of Human and Optimal Performance in Bandit Problems. *Cognitive Systems Research*. 2011 Jun; 12(2):164–174. doi: [10.1016/j.cogsys.2010.07.007](https://doi.org/10.1016/j.cogsys.2010.07.007).
- Lehman J**, Stanley KO. Abandoning Objectives: Evolution Through the Search for Novelty Alone. *Evolutionary Computation*. 2011 Jun; 19(2):189–223. doi: [10.1162/EVCO_a_00025](https://doi.org/10.1162/EVCO_a_00025).
- Lehman J**, Stanley KO. Novelty Search and the Problem with Objectives. In: Riolo R, Vladislavleva E, Moore JH, editors. *Genetic Programming Theory and Practice IX* New York, NY: Springer New York; 2011.p. 37–56. doi: [10.1007/978-1-4614-1770-5_3](https://doi.org/10.1007/978-1-4614-1770-5_3).
- Lehman J**, Stanley KO, Miikkulainen R. Effective Diversity Maintenance in Deceptive Domains. In: *Proceeding of the Fifteenth Annual Conference on Genetic and Evolutionary Computation Conference - GECCO '13* Amsterdam, The Netherlands: ACM Press; 2013. p. 215. doi: [10.1145/2463372.2463393](https://doi.org/10.1145/2463372.2463393).
- Loftus EF**, Hoffman HG. Misinformation and memory: The creation of new memories. *Journal of experimental psychology: General*. 1989; 118(1):100.
- Mackay DJC**. *Information Theory, Inference, and Learning Algorithms*. Second ed.; 2003.
- Marques J**, Meng L, Schaak D, Robson D, Li J. Internal State Dynamics Shape Brainwide Activity and Foraging Behaviour. *Nature*. 2019; p. 27.
- Mehlhorn K**, Newell BR, Todd PM, Lee MD, Morgan K, Braithwaite VA, Hausmann D, Fiedler K, Gonzalez C. Unpacking the Exploration–Exploitation Tradeoff: A Synthesis of Human and Animal Literatures. *Decision*. 2015 Jul; 2(3):191–215. doi: [10.1037/dec0000033](https://doi.org/10.1037/dec0000033).
- Mouret JB**. Novelty-Based Multiobjectivization. In: Kacprzyk J, Doncieux S, Bredèche N, Mouret JB, editors. *New Horizons in Evolutionary Robotics*, vol. 341 Berlin, Heidelberg: Springer Berlin Heidelberg; 2011.p. 139–154. doi: [10.1007/978-3-642-18272-3_10](https://doi.org/10.1007/978-3-642-18272-3_10).
- Mouret JB**, Clune J. Illuminating Search Spaces by Mapping Elites. *arXiv:150404909 [cs, q-bio]*. 2015 Apr; .
- Münch D**, Ezra-Nevo G, Francisco AP, Tastekin I, Ribeiro C. Nutrient Homeostasis — Translating Internal States to Behavior. *Current Opinion in Neurobiology*. 2020 Feb; 60:67–75. doi: [10.1016/j.conb.2019.10.004](https://doi.org/10.1016/j.conb.2019.10.004).
- Ng A**, Harada D, Russell S. Policy Invariance under Reward Transformations: Theory and Application to Reward Shaping. In *Proceedings of the Sixteenth International Conference on Machine Learning*. 1999; p. 278–287.
- Nowak M**, Sigmund K. A Strategy of Win-Stay Lose-Shift That Outperforms Tit-for-Tat in the Prisoner's Dilemma Game. *Nature*. 1993; 364(1):56–58.
- Pathak D**, Agrawal P, Efros AA, Darrell T. Curiosity-Driven Exploration by Self-Supervised Prediction. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)* Honolulu, HI, USA: IEEE; 2017. p. 488–489. doi: [10.1109/CVPRW.2017.70](https://doi.org/10.1109/CVPRW.2017.70).
- Pathak D**, Gandhi D, Gupta A. Self-Supervised Exploration via Disagreement. *Proceedings of the 36th International Conference on Machine Learning*. 2019; p. 10.

- Poucet B.** Spatial Cognitive Maps in Animals: New Hypotheses on Their Structure and Neural Mechanisms. *Psychological Review*. 1993; 100(1):162–183.
- Prescott TJ**, Montes González FM, Gurney K, Humphries MD, Redgrave P. A Robot Model of the Basal Ganglia: Behavior and Intrinsic Processing. *Neural Networks*. 2006 Jan; 19(1):31–61. doi: [10.1016/j.neunet.2005.06.049](https://doi.org/10.1016/j.neunet.2005.06.049).
- Rahnev D**, Denison RN. Suboptimality in Perceptual Decision Making. *Behavioral and Brain Sciences*. 2018; 41:e223. doi: [10.1017/S0140525X18000936](https://doi.org/10.1017/S0140525X18000936).
- Reddy G**, Celani A, Vergassola M. Infomax Strategies for an Optimal Balance between Exploration and Exploitation. *Journal of Statistical Physics*. 2016 Jun; 163(6):1454–1476. doi: [10.1007/s10955-016-1521-0](https://doi.org/10.1007/s10955-016-1521-0).
- Roughgarden T.** Algorithms Illuminated (Part 3): Greedy Algorithms and Dynamic Programming, vol. 1; 2019.
- Schmidhuber.** A Possibility for Implementing Curiosity and Boredom in Model-Building Neural Controllers. *Proc of the international conference on simulation of adaptive behavior: From animals to animats*. 1991; p. 222–227.
- Schmidhuber J.** Formal Theory of Creativity, Fun, and Intrinsic Motivation (1990–2010). *IEEE Transactions on Autonomous Mental Development*. 2010 Sep; 2(3):230–247. doi: [10.1109/TAMD.2010.2056368](https://doi.org/10.1109/TAMD.2010.2056368).
- Schönberg T**, Daw ND, Joel D, O'Doherty JP. Reinforcement learning signals in the human striatum distinguish learners from nonlearners during reward-based decision making. *Journal of Neuroscience*. 2007; 27(47):12860–12867.
- Schulz E**, Bhui R, Love BC, Brier B, Todd MT, Gershman SJ. Exploration in the Wild. *Animal Behavior and Cognition*; 2018.
- Shannon C.** A Mathematical Theory of Communication. *The Bell System Technical Journal*. 1948; 27:379–423, 623–656,.
- Silver D**, Huang A, Maddison CJ, Guez A, Sifre L, van den Driessche G, Schrittwieser J, Antonoglou I, Panneershelvam V, Lanctot M, Dieleman S, Grewe D, Nham J, Kalchbrenner N, Sutskever I, Lillicrap T, Leach M, Kavukcuoglu K, Graepel T, Hassabis D. Mastering the Game of Go with Deep Neural Networks and Tree Search. *Nature*. 2016 Jan; 529(7587):484–489. doi: [10.1038/nature16961](https://doi.org/10.1038/nature16961).
- Silver D**, Hubert T, Schrittwieser J, Antonoglou I, Lai M, Guez A, Lanctot M, Sifre L, Kumaran D, Graepel T, Lillicrap T, Simonyan K, Hassabis D. Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm. *Science*. 2018; 362(6419):1140–1144.
- Şimşek Ö**, Barto AG. An Intrinsic Reward Mechanism for Efficient Exploration. In: *Proceedings of the 23rd International Conference on Machine Learning - ICML '06* Pittsburgh, Pennsylvania: ACM Press; 2006. p. 833–840. doi: [10.1145/1143844.1143949](https://doi.org/10.1145/1143844.1143949).
- Song M**, Bnaya Z, Ma WJ. Sources of Suboptimality in a Minimalistic Explore–Exploit Task. *Nature Human Behaviour*. 2019 Apr; 3(4):361–368. doi: [10.1038/s41562-018-0526-x](https://doi.org/10.1038/s41562-018-0526-x).
- Stanley KO**, Miikkulainen R. Evolving a Roving Eye for Go. In: Kanade T, Kittler J, Kleinberg JM, Mattern F, Mitchell JC, Naor M, Nierstrasz O, Pandu Rangan C, Steffen B, Sudan M, Terzopoulos D, Tygar D, Vardi MY, Weikum G, Deb K, editors. *Genetic and Evolutionary Computation – GECCO 2004*, vol. 3103 Berlin, Heidelberg: Springer Berlin Heidelberg; 2004.p. 1226–1238. doi: [10.1007/978-3-540-24855-2_130](https://doi.org/10.1007/978-3-540-24855-2_130).
- Stanley KO**, Miikkulainen R. Evolving a Roving Eye for Go. In: Kanade T, Kittler J, Kleinberg JM, Mattern F, Mitchell JC, Naor M, Nierstrasz O, Pandu Rangan C, Steffen B, Sudan M, Terzopoulos D, Tygar D, Vardi MY, Weikum G, Deb K, editors. *Genetic and Evolutionary Computation – GECCO 2004*, vol. 3103 Berlin, Heidelberg: Springer Berlin Heidelberg; 2004.p. 1226–1238. doi: [10.1007/978-3-540-24855-2_130](https://doi.org/10.1007/978-3-540-24855-2_130).
- Stanton C**, Clune J. Deep Curiosity Search: Intra-Life Exploration Can Improve Performance on Challenging Deep Reinforcement Learning Problems. *arXiv:180600553 [cs]*. 2018 Nov; .
- Sternisko A**, Cichocka A, Van Bavel JJ. The dark side of social movements: Social identity, non-conformity, and the lure of conspiracy theories. *Current opinion in psychology*. 2020; 35:1–6.
- Strogatz SH.** Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry, and Engineering. *Studies in Nonlinearity*, Reading, Mass: Addison-Wesley Pub; 1994.

- 941 **Sumner ES**, Li AX, Perfors A, Hayes BK, Navarro DJ, Sarnecka BW. The Exploration Advantage: Children's Instinct
942 to Explore Allows Them to Find Information That Adults Miss. *PsyArxiv*. 2019; h437v:11.
- 943 **Sutton RS**, Barto AG. Reinforcement Learning: An Introduction. Second edition ed. Adaptive Computation and
944 Machine Learning Series, Cambridge, Massachusetts: The MIT Press; 2018.
- 945 **Thrun S**, Möller K. Active Exploration in Dynamic Environments. *Advances in neural information processing*
946 *systems*. 1992; p. 531–538.
- 947 **Thrun SB**. Efficient Exploration In Reinforcement Learning. *NIPS*. 1992; p. 44.
- 948 **Tishby N**, Pereira FC, Bialek W. The Information Bottleneck Method. *Arxiv*. 2000; 0004057:11.
- 949 **Valiant L**. A Theory of the Learnable. *Communications of the ACM*. 1984; 27(11):1134–1142.
- 950 **Wang MZ**, Hayden BY. Monkeys Are Curious about Counterfactual Outcomes. *Cognition*. 2019 Aug; 189:1–10.
951 doi: [10.1016/j.cognition.2019.03.009](https://doi.org/10.1016/j.cognition.2019.03.009).
- 952 **Westphal C**, STEFFAN-DEWENTER I, Tschardt T. Foraging trip duration of bumblebees in relation to
953 landscape-wide resource availability. *Ecological Entomology*. 2006; 31(4):389–394.
- 954 **White JK**, Bromberg-Martin ES, Heilbronner SR, Zhang K, Pai J, Haber SN, Monosov IE. A Neural Network for
955 Information Seeking. *Neuroscience*; 2019.
- 956 **Wilson RC**, Bonawitz E, Costa V, Ebitz B. Balancing Exploration and Exploitation with Information and Random-
957 ization. *PsyArXiv*; 2020.
- 958 **Wolpert DH**, Macready WG. No Free Lunch Theorems for Optimization. *IEEE Trans Evol Computat*. 1997 Apr;
959 1(1):67–82. doi: [10.1109/4235.585893](https://doi.org/10.1109/4235.585893).
- 960 **Woodgate JL**, Makinson JC, Lim KS, Reynolds AM, Chittka L. Continuous Radar Tracking Illustrates the Devel-
961 opment of Multi-Destination Routes of Bumblebees. *Scientific Reports*. 2017 Dec; 7(1). doi: [10.1038/s41598-](https://doi.org/10.1038/s41598-017-17553-1)
962 017-17553-1.
- 963 **Wu CM**, Schulz E, Speekenbrink M, Nelson JD, Meder B. Generalization Guides Human Exploration in Vast
964 Decision Spaces. *Nature Human Behaviour*. 2018 Dec; 2(12):915–924. doi: [10.1038/s41562-018-0467-4](https://doi.org/10.1038/s41562-018-0467-4).
- 965 **Zhang M**, Li H, Pan S, Liu T, Su S. One-Shot Neural Architecture Search via Novelty Driven Sampling. In: *Proceed-*
966 *ings of the Twenty-Ninth International Joint Conference on Artificial Intelligence* Yokohama, Japan: International
967 Joint Conferences on Artificial Intelligence Organization; 2020. p. 3188–3194. doi: [10.24963/ijcai.2020/441](https://doi.org/10.24963/ijcai.2020/441).
- 968 **Zhou D**, Lydon-Staley DM, Zurn P, Bassett DS. The Growth and Form of Knowledge Networks by Kinesthetic
969 Curiosity. *arXiv:200602949 [cs, q-bio]*. 2020 Jun; .

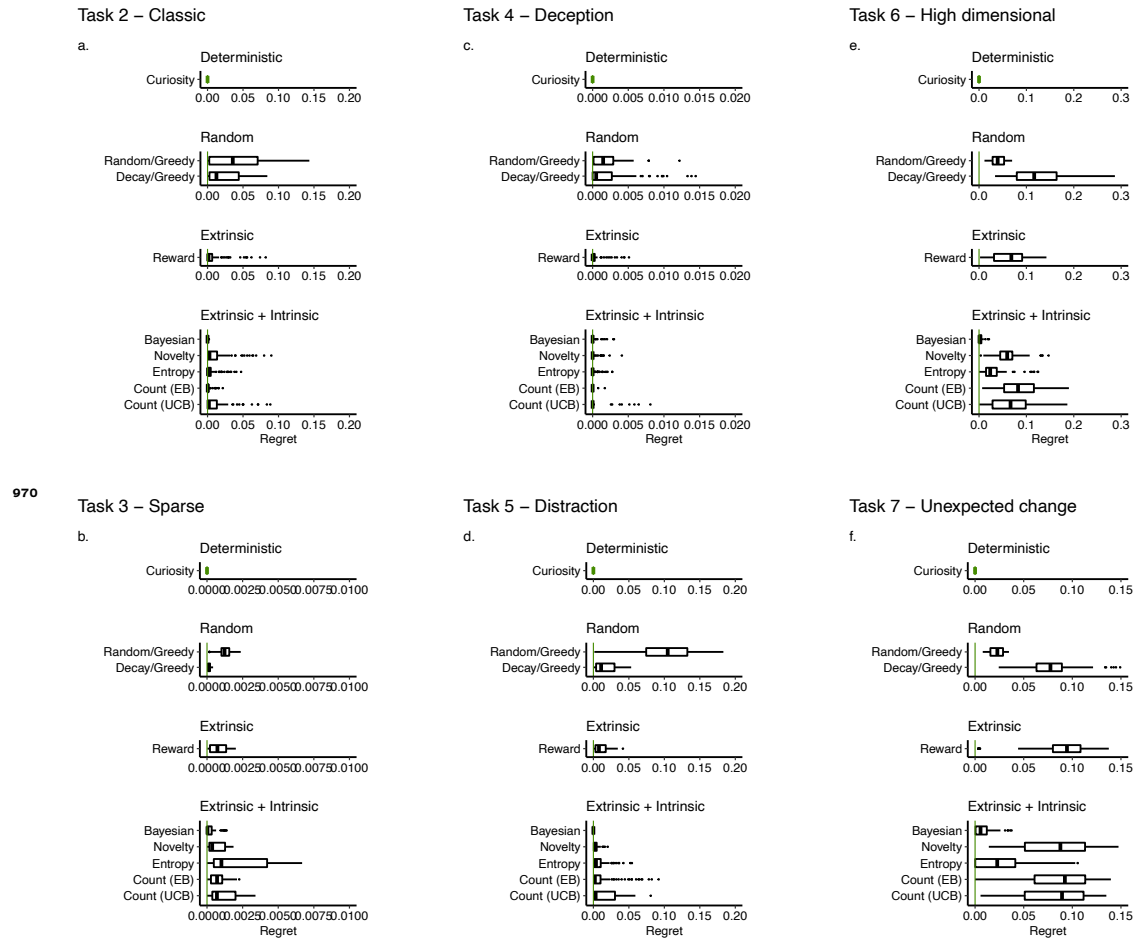
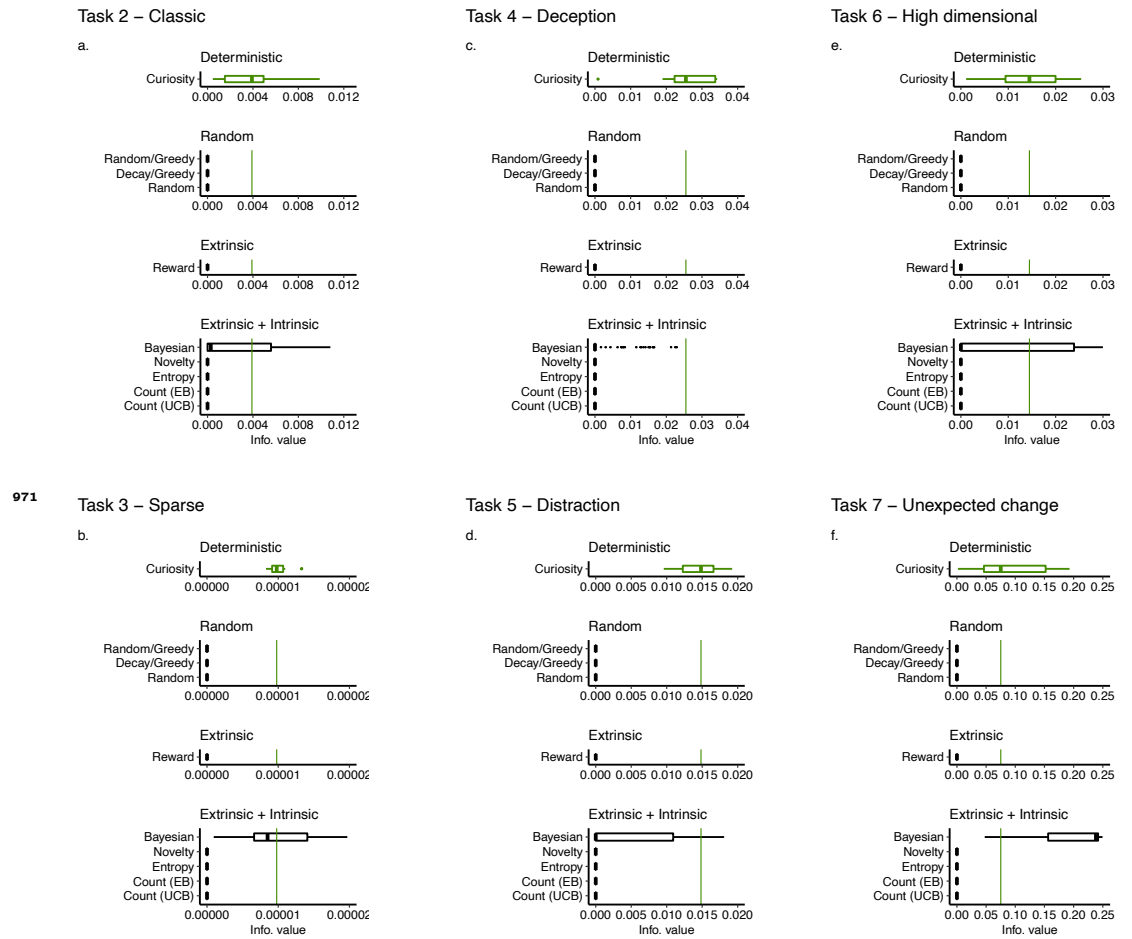


Figure 7–Figure supplement 1. Summary of regret (Tasks 2-7)



Exploration strategies Ten 'animals'



Figure 7-Figure supplement 3. Examples of exploration (Task 1). Here we simulated different sets of hyper-parameters on the same task and random seed. We view each parameter set as a unique “animal” (*Prescott et al., 2006*). So this figure estimates the degree and kind of variability we might expect in a natural experiment.

Reward collection with independent policies: deterministic versus stochastic

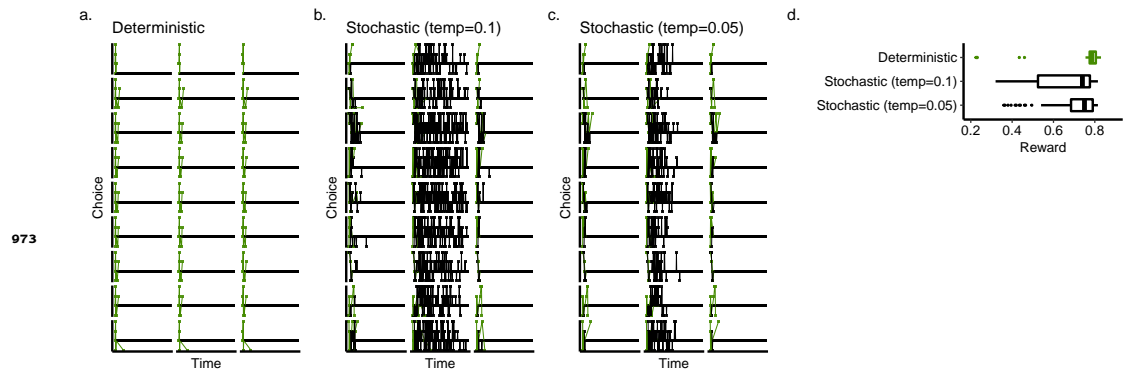


Figure 9-Figure supplement 1. Adding noise to deterministic curiosity does degrade its performance. In this example control we added two levels of random noise. This only served to increase variability of convergence (a-c) and reduce the total rewards collected (d). Note: decision noise was modelled by sampling a Boltzmann distribution, as described in the *Methods*.

The effect of noise on stochastic search

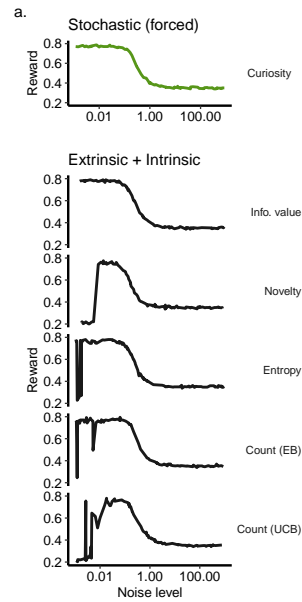


Figure 9-Figure supplement 2. Noise across strategies. Adding noise to degrades performance to nearly the same degree (Task 2).