

A way around the exploration-exploitation dilemma.

Erik J Peterson^{a,b,1} and Timothy D Verstynen^{a,b,c,d}

^aDepartment of Psychology; ^bCenter for the Neural Basis of Cognition; ^cCarnegie Mellon Neuroscience Institute; ^dBiomedical Engineering, Carnegie Mellon University, Pittsburgh PA

For all animals the decision to explore comes with a risk of getting less. For example, a foraging bee might find less nectar, or hunting hawk less prey. This loss is often formalized as regret. It's been mathematically proven that exploring an uncertain world with a specific goal always has some regret. This is why exploration-exploitation can be a dilemma. Given this proof we wondered if the common advice to “focus on learning and not the goal” might have mathematical merit. So we re-imagined exploration in the dilemma as an open ended search for any new information. We then developed a new minimal description of information value, which generalizes existing ideas like curiosity, novelty and information gain. We use this description to model the dilemma as a competition between strategies that maximize reward and information independently. Here we prove this competition has a no regret solution. When we study this solution in simulation – using classic bandit tasks – it outperforms standard approaches, especially when rewards are sparse.

Introduction

Decision making in the natural world often leads to a dilemma. As an example let's imagine a bee foraging in a meadow (Figure 1A). The bee could go to the location of a flower it's been to before to gather nectar. Or the bee go somewhere new, and explore. Exploration comes though with the risk of getting less nectar. Perfectly optimizing away this risk is a mathematically intractable problem; there is no way to explore without enduring some regret (1–4), and so the decision can become a dilemma.

Resource gathering is not the only reason animals explore. Many animals, like our bee, explore out of curiosity (Figure 1B). This exploration lets them learn about their environment, developing an often simplified model that helps them in planning actions and making future decisions (5, 6). Borrowing from the field of artificial intelligence we refer to these models as *world models* (7–9). World models offer a principled explanation for why animals are intrinsically curious (10–15), and prone to explore even when no rewards are present or expected (16).

Curiosity raises the question of whether animals need to explore looking for specific goals or rewards are all. Perhaps we've misinterpreted their actions, and so misconceived of a fundamental problem in the learning and decision sciences. Here we explore a bold conjecture:

Exploration for reward is never needed. The only exploratory behavior an animal needs is that which builds its world model.

Our contribution is threefold. We define a new minimal (axiomatic) description for information value, which generalizes existing ideas like curiosity, novelty and information gain. In fact, the axioms let us formally disconnect information

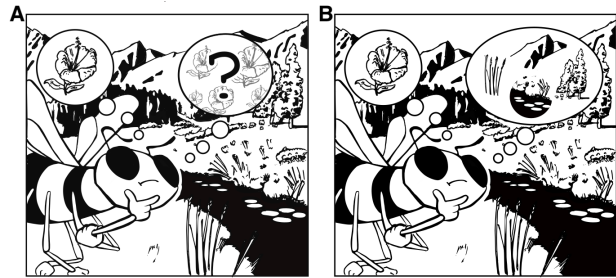


Fig. 1. Two views of exploration and exploitation. **A**. The classic dilemma: either exploit an action with a known reward (e.g., return to the previous plant) or explore other actions on the chance they will return a better outcome (e.g., find a plant with more flowers). **B**. Here we offer an alternative view of the dilemma, with two different competitive goals: maximize rewards (e.g., keep returning to known flower locations) or build a world model by learning new information (e.g., layout of the environment). Exploration here focused on learning in general, not on reward learning specifically. Artist credit: Richard Grant.

theory (17) from information value, suggesting we may have uncovered a new universal theory. Next we prove that the computer science method of dynamic programming (8, 18) provides an optimal way to maximize this kind of information value. Finally, we describe a simple winner-take-all scheduling algorithm that can optimally solve a competition between strategies which independently maximize information value and reward.

Results

Tangible rewards are a conserved resource, but learned information isn't. For example, if a rat shares potato chip with a cage-mate, she must necessarily split up the chip leaving less food for herself. Whereas if student shares the latest result from a scientific paper with a lab-mate, they do not necessarily forget a portion of that result. These differences make reward and information different concepts, and so considering information as a kind of reward isn't inconsistent.

If information value isn't a reward, we need another way to study and value it. To do this we first looked to the field of information theory (17), but the problem of information value is not based in the statistical problem of transmitting symbols, as was Shannon's goal. It is based on the problem of learning and remembering them.

A minimal model of memory. World models are memories with some amount of simplification (9, 19). They can range from simple novelty signals (20), to location or state counts (21, 22),

The authors have no conflicts of interest to declare.

¹To whom correspondence should be addressed. E-mail: Erik.Exists@gmail.com

state and action prediction (9, 15, 23), flow (24), learning progress (25), classic working or episodic memories (26, 27), Bayesian and hierarchical Bayesian models (23, 28–30), latent spaces (31) and recurrent neural networks (19, 32–34).

We have no mathematical reason to prefer any one kind of world model over any other. So we designed a new minimal definition, designed to overlap with all of them.

We must introduce some initial notation. We assume that time is a continuous value and denote increases in time using the differential quantity dt . We can then express changes in M (our world model, defined below) as a gradient, ∇M . We also assume that observations about the environment s are real numbers sampled from a finite state space $s \in S$, whose size is N (denoted S^N). Actions are also real numbers a , drawn from a finite space A^K . Rewards R_t – when they appear – are binary (0, 1) and are provided *only* by the external environment.

Definition 1. We can now formally define a world model M as a finite set of real numbers, whose maximum size is L (M^L). We say that every world model has a pair of functions f and g . Learning of s at time t (i.e. s_t) by M is done by the invertible encoder function f , $M_{t+dt} = f(M_t, s_t)$ and $M_t = f^{-1}(M_{t+dt}, s_t)$. Memories \hat{s}_t about s_t are recalled by the decoder function g , $\hat{s}_t = g(M_t, s_t)$.

The invertibility of f , denoted as f^{-1} , is a mathematical way to ensure that any observations encoded in the world model can also be forgotten. This is both an important aspect of real memory, and a critical point for our mathematical analysis.

The details of f and g define what kind of world model or memory M is. Let's consider some examples. If f adds states s_t to the memory, and g tests whether s_t is in M , then M is a model of novelty (20). If f counts states and g returns those counts, then M is a count-based heuristic (21, 22). If f follows Bayes rule and g decodes the probability of s_t , then M is a Bayesian memory (9, 15, 23, 29, 30). If the size of M is much smaller than the size of the state space S^N , then f can be seen as learning a latent or compressed representation in M (19, 28, 31, 33–37), and g decodes a reconstruction of s (\hat{s}_t) or future states (\hat{s}_{t+dt}).

A minimal description of information value. To formalize information value we use two axioms that define a real valued function, $E(s)$, that measures the value of any observation s_t given a world model M and a distance metric d .

Axiom 1 (Axiom of Change). *The value of information $E(s_t)$ depends only on the total distance M moves by making observation s_t .*

This axiom does three important things. It ensures information value depends only on the world model, that value is a distance in memory, and that value learning has the Markov property (8). Now, let's unpack it.

By distance we mean a function $\delta = d(m, m')$, where $m \in M$ and $m' \in M'$ are discrete memories drawn from two memories M and M' . We define d so $d \geq 0$ for all $s \in S$, and let $=0$ only if $M = M'$. Our definition of d does not require the distance in memories from M to M' be the same as from M' to M . Nor for the triangle inequality to hold. For the technically inclined, this definition makes d and so E a pre-metric.

By total distance we mean the norm $\|\Delta\|$, where $\Delta = \{\delta_1, \delta_2, \dots, \delta_L\}$.

In summary, Let $E \equiv \|\Delta\|$.

Different f and g pairs will naturally need different ways to measure distances in M . For example, in a novelty world model (20) either the hamming or Manhattan distance are applicable and would produce binary distance values, as would a count model (21, 22). A latent memory (9, 15) might instead use the euclidean norm of its own error gradient (38). While a probabilistic or Bayesian memory would likely use the Kullback–Leibler (KL) divergence (23, 28).

Axiom 2 (Axiom of Equilibrium). *To be valuable an observation s_t must be learnable by M*

By learnable we mean two things. First, with every (re)observation of s , M should change. Second, the change in M must eventually reach a learned equilibrium. To formalize these we constrain the average gradient of M , so $\mathbb{E}[\nabla^2 M] \leq 0$.

Most attempts to value information rest their definition on information theory. Value might rest on the intrinsic complexity of an observation (i.e., its entropy) (39) or on its similarity to the environment (i.e., mutual information) (40), or on some other salience signal (41). In our analysis, learning alone drives value. This is because learning might happen on a true world model or with a faulty world model, or be about a fictional narrative. The observation might be simple, or complex. From a subjective point of view, which is the right point of view for value, all of these are the same; value depends only on the total knowledge gained.

Exploration as a dynamic programming problem. Dynamic programming is a popular optimization method because it guarantees value is maximized using a simple algorithm that always chooses the largest option. In Theorem 1 (see *Mathematical Appendix*) we prove that our definition of memory has one critical property, optimal substructure, that is needed for an optimal dynamic programming solution (18, 42). The other two required properties, $E \geq 0$ and the Markov property (18, 42), are fulfilled by the *Axiom 1*. To write down our dynamic programming solution we introduce a little more notation. We let π denote an action policy, a function that takes a state s and returns an action a . We let δ denote the transition function, which takes a state-action pair (s_t, a_t) and returns a new state, s_{t+dt} . This function acts as an abstraction for the actual world. For notational consistency with the standard Bellman approach we also redefine $E(s)$ as a *payoff function*, $F(M_t, a_t)$ (18).

$$\begin{aligned} F(M_t, a_t) &= E(s) \\ \text{subject to the constraints} \\ a_t &= \pi(s_t) \\ s_{t+dt} &= \delta(s_t, a_t), \\ M_{t+dt} &= f(M_t, s_t) \end{aligned} \quad [1]$$

The value function for F is,

$$V_{\pi_E}(M_0) = \left[\max_{a \in A} \sum_{t=0}^{\infty} F(M_t, a_t) \mid M, d, S \right]. \quad [2]$$

And the recursive Bellman solution to learn this value function is,

$$V_{\pi_E}^*(M_t) = F(M_t, a_t) + \max_{a \in A} [F(M_{t+dt}, a_t)]. \quad [3]$$

For the full derivation of Eq 3 see the *Mathematical Appendix*, where we also prove that Eq 3 leads to exhaustive exploration of any finite space S (Theorems 2 and 3).

Scheduling a way around the dilemma. Remember that the goal of reinforcement learning is to maximize reward, an objective approximated by the value function $V_R(s)$ and an action policy π_R .

$$V_R^{\pi_R}(s) = \mathbb{E} \left[\sum_{k=0}^{\infty} R_{t+k+1} | s = s_t \right] \quad [4]$$

Remember too that our overall goal is to find an algorithm that maximizes both information and reward value. To do that we imagine the policies for exploration and exploitation are possible “jobs” competing to control behavior. We know that, by definition, each of these jobs produces non-negative values: E for information or R for reinforcement learning. So our goal is to find an optimal scheduler for these two jobs.

To do this we further simplify our assumptions. We assume each action takes a constant amount of time, and has no energetic cost. We assume the policy can only take one action at a time, and that those actions are exclusive. Most scheduling solutions also assume that the value of a job is fixed, while in our problem information value changes as the world model improves. In a general setting however, where one has no prior information about the environment, the best predictor of the next value is the last or most recent value (42, 43). We assume this precept holds in all of our analysis.

With these assumptions in place, the optimal solution to this kind of scheduling problem is known to be a purely local, winner-take-all, algorithm (18, 42). We state this winner-take-all solution here as a set of inequalities where R_t and E_t represent the value of reward and information at the last time-point.

$$\pi_{\pi}(s_t) = \begin{cases} \pi_E^*(s_t) & : E_t - \eta > R_t \\ \pi_R(s_t) & : E_t - \eta \leq R_t \end{cases} \quad [5]$$

subject to the constraints

$$\begin{aligned} p(\mathbb{E}[R]) &< 1 \\ E - \eta &\geq 0 \end{aligned}$$

To ensure that the default policy is reward maximization, Eq. 5 breaks ties between R_t and E_t in favor of π_R . In stochastic environments, M can show small continual fluctuations. To allow Eq. 5 to achieve a stable solution we introduce η , a boredom threshold for exploration. Larger values of η devalue information exploration and favor exploitation of reward.

The worst case algorithmic run time for Eq 5 is linear and additive in its policies. So if in isolation it takes T_E steps to earn $E_T = \sum_{T_E} E$, and T_R steps to earn $r_T = \sum_{T_R} R$, then the worst case training time for π_{π} is $T_E + T_R$. It is worth noting that this is only true if neither policy can learn from the other’s actions. There is, however, no reason that each policy cannot observe the transitions (s_t, a_t, R, s_{t+dt}) caused by the other. If this is allowed, worst case training time improves to $\max(T_E, T_R)$.

Exploration without regret. Suboptimal exploration strategies will lead to a loss of potential rewards by wasting time on actions that have a lower expected value. Regret G measures the value loss caused by such exploration. $G = \hat{V} - V_a$, where \hat{V} represents the maximum value and V_a represents the value found by taking an exploratory action rather than an exploitative one (8).

Optimal strategies for a solution to the exploration-exploitation dilemma should maximize total value with zero total regret.

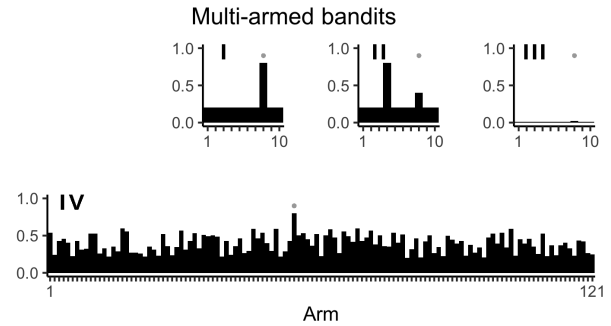


Fig. 2. Bandits. Reward probabilities for each arm in bandit tasks I-IV. Grey dots highlight the optimal (i.e., highest reward probability) arm. See main text for a complete description.

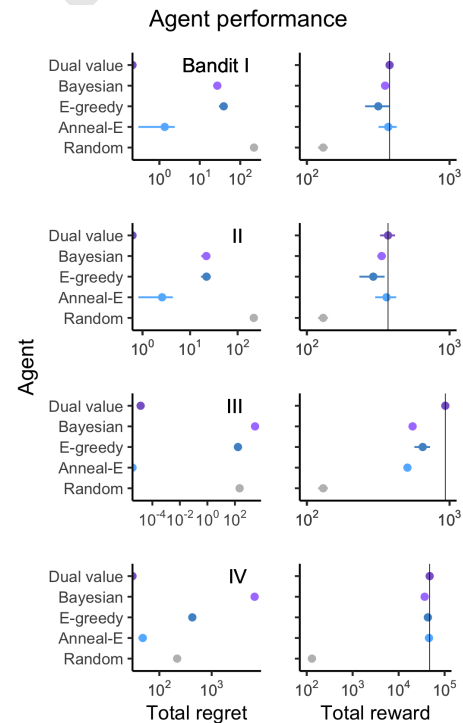


Fig. 3. Regret and total accumulated reward across models and bandit task. Median total regret (left column) and median total reward (right column) for simulations of each model type ($N = 100$ experiments per model). See main text and Table 1 for description of each model. Error bars in all plots represent median absolute deviation.

To evaluate dual value learning (Eq. 5) we compared total reward and regret across a range of both simple, and challenging multi-armed bandit tasks. Despite its apparent simplicity, the essential aspects of the exploration-exploitation dilemma

Table 1. Artificial agents.

Agent	Exploration mechanism
Dual value	Our algorithm (Eq 5).
E-greedy	With probability $1 - \epsilon$ follow a greedy policy. With probability ϵ follow a random policy.
Annealed e-greedy	Identical to E-greedy, but ϵ is decayed at fixed rate.
Bayesian reward	Use the KL divergence as a weighted intrinsic reward, sampling actions by a soft-max policy. $\sum_{\mathcal{T}} R_t + \beta E_t$
Random	Action are selected with a random policy (no learning)

exist in the multi-armed bandit task (8). Here the problem to be learned is the distribution of reward probabilities across arms (Figure 2). To estimate the value of any observation s_t , we compare sequential changes in this probabilistic memory, M_{t+dt} and M_t using the KL divergence (i.e. relative entropy; Figure 4A-B). The KL divergence is a standard way to measure the distance between two distributions (44) and is, by design, consistent with our axioms (see the *Supplementary Materials* for a more thorough discussion).

We start with a simple experiment involving a single high value arm. The rest of the arms have a uniform reward probability (Bandit I). This represents a trivial problem. Next we tried a basic exploration test (Bandit II), with one winning arm and one distractor arm whose value is close to but less than the optimal choice. We then move on to a more difficult sparse exploration problem (Bandit III), where the world has a single winning arm, but the overall probability of receiving any reward is very low ($p(R) = 0.02$ for the winning arm, $p(R) = 0.01$ for all others). Sparse reward problems are notoriously difficult to solve, and are a common feature of both the real world and artificial environments like Go, chess, and class Atari video games (45–47). Finally, we tested a complex, large world exploration problem (Bandit IV) with 121 arms, and a complex, randomly generated reward structure. Bandits of this type and size are near the limit of human performance (48).

We compared the reward and regret performance of 6 artificial agents. All agents used the same temporal difference learning algorithm (TD(0), (8)); see *Supplementary materials*. The only difference between the agents was their exploration mechanism (Table 1). The e-greedy algorithm is a classic exploration mechanism (8). Its annealed variant is common in state-of-the-art reinforcement learning papers, like Mnih *et al* ((45)). Other state-of-the-art exploration methods are models that treat Bayesian information gain as an intrinsic reward and the goal of all exploration is to maximize total reward (extrinsic plus intrinsic) (9, 49). To provide a lower bound benchmark of performance we included an agent with a purely random exploration policy.

All of the classic and state-of-the-art algorithms performed well at the different tasks in terms of accumulation of rewards (right column, Figure 3). The one exception to this being the sparse low reward probability condition (Bandit III), where the dual value algorithm consistently returned more rewards than the other models. In contrast, most of the traditional models still had substantial amounts of regret in most of the

tasks, with the exception of the annealed variant of the e-greedy algorithm during the sparse, low reward probability task (left column, Figure 3). In contrast, the dual value learning algorithm consistently was able to maximize total reward with zero or near zero (Bandit III) regret, as would be expected by an optimal exploration policy.

Discussion

Past work. We are certainly not the first to quantify information value (40, 50), or use that value to optimize reward learning (2, 9, 29, 51, 52). Information value though is typically framed as a means to maximize the amount of tangible rewards (e.g., food, water, money) accrued over time (8). This means that information is treated as an analog of these tangible or external rewards (i.e., an *intrinsic reward*) (9, 12, 23, 29). This approximation does drive exploration in a practical and useful way, but doesn't change the intractability of the dilemma (1–4).

At the other extreme from reinforcement learning are pure exploration methods, like curiosity (15, 49, 53) or PAC approaches (54). Curiosity learning is not generally known to converge on rewarding actions with certainty, but nevertheless can be an effective heuristic (15, 55, 56). Within some bounded error, PAC learning is certain to converge (54). For example, it will find the most rewarding arm in a bandit, and do so with a bounded number of samples (57). However, the number of samples is fixed and based on the size of the environment (but see (58, 59)). So while PAC will give the right answer, eventually, its exploration strategy also guarantees high regret.

Cost. It is not fair to talk about benefits without talking about costs. The worst-case run-time of a dual value algorithm is $\max(T_E, T_R)$, where T_E and T_R represent the time to learn to some criterion (see *Results*). In the unique setting where minimizing regret, maximizing data efficiency, exploration efficiency, and transfer do not matter, dual value learning can be a suboptimal choice.

Animal behavior. In psychology and neuroscience, curiosity and reinforcement learning have developed as separate disciplines (8, 53, 60). And they are separate problems, with links to different basic needs: gathering resources to maintain physiological homeostasis (61, 62) and gathering information to plan for the future (8, 54). Here we suggest that though they are separate problems, they are problems that can, in large part, solve one another.

The theoretical description of exploration in scientific settings is probabilistic (4, 63–65). By definition probabilistic models can't make exact predictions of behavior, only statistical ones. Our approach is deterministic, and so does make exact predictions. Our theory predicts that it should be possible to guide exploration in real-time using, for example, optogenetic methods in neuroscience, or well timed stimulus manipulations in economics or other behavioral sciences.

Artificial intelligence. Progress in reinforcement learning and artificial intelligence research is limited by three factors: data efficiency, exploration efficiency, and transfer learning (19). Our algorithm speaks directly to all three of these limits. By treating exploration as a problem in building a world model,

our algorithm always ensures high quality exploration. The focus on the world model also means it can be naturally integrated with data efficient model-based reinforcement learning (8, 66). Finally, as it builds a world model that is free of any task specific bias and so is ideal for later transfer or fine-tuning (67, 68).

We describe here a simple and optimal algorithm to combine nearly any world model with any reinforcement learning algorithm. This effectively joins the two approaches to reinforcement learning – model-free and model-based – into an advantageous whole where exploration is model-based, but exploitation and reward learning is algorithmically model-free.

Everyday life. The uncertainty of the unknown can always be recast as an opportunity to learn. But rather than being a trick of positive psychology, we prove this view is (in the narrow sense of our formalism, anyway) mathematically optimal.

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

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.

1. Thrun S, Möller K (1992) Active Exploration in Dynamic Environments. *Advances in neural information processing systems* pp. 531–538.
2. Dayan P, Sejnowski TJ (1996) Exploration bonuses and dual control. *Machine Learning* 25(1):5–22.
3. Findling C, Skvortsova V, Dromnelle R, Palminteri S, Wyart V (2018) Computational noise in reward-guided learning drives behavioral variability in volatile environments, (Neuroscience), Preprint.
4. Gershman SJ (2018) Deconstructing the human algorithms for exploration. *Cognition* 173:34–42.
5. Ahilan S, et al. (2019) Learning to use past evidence in a sophisticated world model. *PLoS Computational Biology* 15(6):e1007093.
6. Poucet B (1993) Spatial cognitive maps in animals: New hypotheses on their structure and neural mechanisms. *Psychological Review* 100(1):162–183.
7. Schmidhuber J (2019) Unsupervised Minimax: Adversarial Curiosity, Generative Adversarial Networks, and Predictability Minimization. *arXiv:1906.04493 [cs]*.
8. Sutton RS, Barto AG (2018) *Reinforcement Learning: An Introduction*, Adaptive Computation and Machine Learning Series. (The MIT Press, Cambridge, Massachusetts), Second edition.
9. Schmidhuber (1991) 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* pp. 222–227.
10. Mehlhorn K, et al. (2015) Unpacking the exploration–exploitation tradeoff: A synthesis of human and animal literatures. *Decision* 2(3):191–215.
11. Gupta AA, Smith K, Shalley C (2006) The Interplay between Exploration and Exploitation. *The Academy of Management Journal* 49(4):693–706.
12. Berger-Tal O, Nathan J, Meron E, Saltz D (2014) The Exploration-Exploitation Dilemma: A Multidisciplinary Framework. *PLoS ONE* 9(4):e95693.
13. Gottlieb J, Oudeyer PY (2018) Towards a neuroscience of active sampling and curiosity. *Nature Reviews Neuroscience* 19(12):758–770.
14. Schwartenbeck P, et al. (2019) Computational mechanisms of curiosity and goal-directed exploration. *eLife* (e41703):45.
15. Pathak D, Agrawal P, Efros AA, Darrell T (2017) Curiosity-Driven Exploration by Self-Supervised Prediction in 2017 *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. (IEEE, Honolulu, HI, USA), pp. 488–489.
16. Hughes RN (1997) Intrinsic exploration in animals: Motives and measurement. *Behavioural Processes* 41(3):213–226.
17. Shannon C (1948) A Mathematical Theory of Communication. *The Bell System Technical Journal* 27:379–423, 623–656.
18. Bellman R (1954) The theory of dynamic programming. *Bull. Amer. Math. Soc* 60(6):503–515.
19. Ha D, Schmidhuber J (2018) World Models. *arXiv:1803.10122 [cs, stat]*.
20. Kakade S, Dayan P (2002) Dopamine: Generalization and bonuses. *Neural Networks* 15(4-6):549–559.
21. Bellemare MG, et al. (2016) Unifying Count-Based Exploration and Intrinsic Motivation. *arXiv:1606.01868 [cs, stat]*.
22. Dayan P (1993) Improving Generalisation for Temporal Difference Learning: The Successor Representation. *Neural Computation* 5(4):613–624.
23. Friston K, et al. (2016) Active inference and learning. *Neuroscience & Biobehavioral Reviews* 68:862–879.
24. Yang HK, Chiang PH, Hong MF, Lee CY (2019) Exploration via Flow-Based Intrinsic Rewards. *arXiv:1905.10071 [cs, stat]*.
25. Lopes M, Lang T, Toussaint M, Oudeyer PY (2012) Exploration in Model-based Reinforcement Learning by Empirically Estimating Learning Progress. *NIPS* 25:1–9.
26. Miller G (1956) The Magical Number Seven, Plus or Minus Two: Some Limits on Our Capacity for Processing Information. *The Psychological Review* 63:81–97.
27. Tulving E (2002) Episodic Memory: From Mind to Brain. *Annual Review of Psychology* 53(1):1–25.
28. Park IM, Pillow JW (2017) Bayesian Efficient Coding, (Neuroscience), Preprint.
29. Itti L, Baldi P (2009) Bayesian surprise attracts human attention. *Vision Research* 49(10):1295–1306.
30. Tenenbaum JB, Griffiths TL, Kemp C (2006) Theory-based Bayesian models of inductive learning and reasoning. *Trends in Cognitive Sciences* 10(7):309–318.
31. Kingma DP, Welling M (2013) Auto-Encoding Variational Bayes. *arXiv:1312.6114 [cs, stat]*.
32. Ganguli S, Huh D, Sompolinsky H (2008) Memory traces in dynamical systems. *Proceedings of the National Academy of Sciences* 105(48):18970–18975.
33. Schmidhuber J (2015) On Learning to Think: Algorithmic Information Theory for Novel Combinations of Reinforcement Learning Controllers and Recurrent Neural World Models. *arXiv:1511.09249 [cs]*.
34. Mante V, Sussillo D, Shenoy KV, Newsome WT (2013) Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature* 503(7474):78–84.
35. Schmidhuber J (2008) Driven by Compression Progress: A Simple Principle Explains Essential Aspects of Subjective Beauty, Novelty, Surprise, Interestingness, Attention, Curiosity, Creativity, Art, Science, Music, Jokes. *arXiv:0812.4360 [cs]*.
36. Levi-Aharoni H, Shriki O, Tishby N (2019) Surprise response as a probe for compressed memory states, (Neuroscience), Preprint.
37. Ganguli S, Sompolinsky H (2010) Short-term memory in neuronal networks through dynamical compressed sensing. p. 9.
38. Pascanu R, Bengio Y (2013) Revisiting Natural Gradient for Deep Networks. *arXiv:1301.3584 [cs]*.
39. Haarnoja T, et al. (2018) Soft Actor-Critic Algorithms and Applications. *arXiv:1812.05905 [cs, stat]*.
40. Kolchinsky A, Wolpert DH (2018) Semantic information, autonomous agency and non-equilibrium statistical physics. *Interface Focus* 8(6):20180041.
41. Tishby N, Pereira FC, Bialek W (2000) The Information Bottleneck Method. *Arxiv* 0004057:11.
42. Roughgarden T (2019) *Algorithms Illuminated (Part 3): Greedy Algorithms and Dynamic Programming*. Vol. 1.
43. Hocker D, Park IM (2019) Myopic control of neural dynamics. 15(3):24.
44. MacKay DJC (2003) *Information Theory, Inference, and Learning Algorithms*. 2 edition.
45. Mnih V, et al. (year?) Playing Atari with Deep Reinforcement Learning. p. 9.
46. Silver D, et al. (2016) Mastering the game of Go with deep neural networks and tree search. *Nature* 529(7587):484–489.
47. Silver D, et al. (2018) Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm. *Science* 362(6419):1140–1144.
48. Wu CM, Schulz E, Speekenbrink M, Nelson JD, Meder B (2018) Generalization guides human exploration in vast decision spaces. *Nature Human Behaviour* 2(12):915–924.
49. Jaegle A, Mehrpour V, Rust N (2019) Visual novelty, curiosity, and intrinsic reward in machine learning and the brain. *Arxiv* 1901.02478:13.
50. Cogliati Dezza I, Yu AJ, Cleeremans A, Alexander W (2017) Learning the value of information and reward over time when solving exploration-exploitation problems. *Scientific Reports* 7(1).
51. Kelly JL (1956) A New Interpretation of Information Rate. *the bell system technical journal* p. 10.
52. de Abrisil IM, Kanai R (2018) Curiosity-Driven Reinforcement Learning with Homeostatic Regulation in 2018 *International Joint Conference on Neural Networks (IJCNN)*. (IEEE, Rio de Janeiro), pp. 1–6.
53. Berylyne D (1950) Novelty and curiosity as determinants of exploratory behaviour. *British Journal of Psychology* 41(1):68–80.
54. Valiant L (1984) A theory of the learnable. *Communications of the ACM* 27(11):1134–1142.
55. Burda Y, et al. (2018) Large-Scale Study of Curiosity-Driven Learning. *arXiv:1808.04355 [cs, stat]*.
56. Colas C, Fournier P, Sigaud O, Chetouani M, Oudeyer PY (2019) CURIOUS: Intrinsically Motivated Multi-Task Multi-Goal Reinforcement Learning. *Arxiv* 1810.06284v3:1–15.
57. Even-Dar E, Mannor S, Mansour Y (2002) PAC Bounds for Multi-armed Bandit and Markov Decision Processes in *Computational Learning Theory*, eds. Goos G, Hartmanis J, van Leeuwen J, Kivinen J, Sloan RH. (Springer Berlin Heidelberg, Berlin, Heidelberg) Vol. 2375, pp. 255–270.
58. Even-Dar E, Mannor S, Mansour Y (2006) Action Elimination and Stopping Conditions for the Multi-Armed Bandit and Reinforcement Learning Problems. *Journal of Machine Learning Research* 7:1–27.
59. Strehl AL, Li L, Littman ML (2009) Reinforcement Learning in Finite MDPs: PAC Analysis. *Journal of Machine Learning Research* 10:1–32.
60. Kidd C, Hayden BY (2015) The Psychology and Neuroscience of Curiosity. *Neuron* 88(3):449–460.

61. Keramati M, Gutkin B (2014) Homeostatic reinforcement learning for integrating reward collection and physiological stability. *eLife* 3:e04811.
62. Juechems K, Summerfield C (2019) Where does value come from?. (PsyArXiv). Preprint.
63. Calhoun AJ, Chalasani SH, Sharpee TO (2014) Maximally informative foraging by *Caenorhabditis elegans*. *eLife* 3:e04220.
64. Song M, Bnaya Z, Ma WJ (2019) Sources of suboptimality in a minimalistic explore–exploit task. *Nature Human Behaviour* 3(4):361–368.
65. Schulz E, et al. (2018) Exploration in the wild, (Animal Behavior and Cognition), Preprint.
66. Shyam P, Jaśkowski W, Gomez F (2018) Model-Based Active Exploration. *arXiv:1810.12162 [cs, math, stat]*.
67. Yosinski J, Clune J, Bengio Y, Lipson H (2014) How transferable are features in deep neural networks? *arXiv:1411.1792 [cs]*.
68. Barreto A, Dabney W, Munos R, Hunt JJ, Schaul T (2018) Successor Features for Transfer in Reinforcement Learning. *Arxiv* 1606.05312v2:1–11.
69. López-Fidalgo J, Tommasi C, Trandafir PC (2007) An optimal experimental design criterion for discriminating between non-normal models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 69(2):231–242.
70. Still S, Sivak DA, Bell AJ, Crooks GE (2012) The thermodynamics of prediction. *Physical Review Letters* 109(12).
71. Ay N (2015) Information Geometry on Complexity and Stochastic Interaction. *Entropy* 17(4):2432–2458.
72. Li L, Jamieson K, DeSalvo G, Rostamizadeh A, Talwalkar A (2016) Hyperband: A Novel Bandit-Based Approach to Hyperparameter Optimization. *arXiv:1603.06560 [cs, stat]*.
73. Strogatz SH (1994) *Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry, and Engineering*, Studies in Nonlinearity. (Addison-Wesley Pub, Reading, Mass).

PREPRINT

Supplementary materials.

Dual value implementation.

Value initialization and tie breaking. The initial value E_0 for π_E^* can be arbitrary, with the limit $E_0 > 0$. In theory E_0 does not change π_E^* 's long term behavior, but different values will change the algorithm's short-term dynamics and so might be quite important in practice. By definition a pure greedy policy, like π_E^* , cannot handle ties. There is simply no mathematical way to rank equal values. Theorems 3 and 2 ensure that any tie breaking strategy is valid, however, like the choice of E_0 , tie breaking can strongly affect the transient dynamics. Viable tie breaking strategies taken from experimental work include, “take the closest option”, “repeat the last option”, or “take the option with the highest marginal likelihood”. We do suggest the tie breaking scheme is deterministic, which maintains the determinism of the whole theory. See *Information value learning* section below for concrete examples both these choices.

The rates of exploration and exploitation. In Theorem 4 we proved that π_π inherits the optimality of policies for both exploration π_E and exploitation π_R over infinite time. However this does proof does not say whether π_π will not alter the rate of convergence of each policy. By design, it does alter the rate of each, favoring π_R . As you can see in Eq. ??, whenever $r_t = 1$ then π_R dominates that turn. Therefore the more likely $p(r = 1)$, the more likely π_R will have control. This doesn't of course change the eventual convergence of π_E , just delays it in direct proportion to the average rate of reward. In total, these dynamics mean that in the common case where rewards are sparse but reliable, exploration is favored and can converge more quickly. As exploration converges, so does the optimal solution to maximizing rewards.

Re-exploration. The world often changes. Or in formal parlance, the world is non-stationary process. When the world does change, re-exploration becomes necessary. Tuning the size of ϵ in π_π (Eq ??) tunes the threshold for re-exploration. That is, once the π_E^* has converged and so π_R^* fully dominates π_π , if ϵ is small then small changes in the world will allow π_E to exert control. If instead ϵ is large, then large changes in the world are needed. That is, ϵ acts a hyper-parameter controlling how quickly rewarding behavior will dominate, and easy it is to let exploratory behavior resurface.

Bandits.

Design. Like the slot machines which inspired them, each bandit returns a reward according to a predetermined probability. As an agent can only chose one bandit (“arm”) at a time, so it must decide whether to explore and exploit with each trial.

We study four prototypical bandits. The first has a single winning arm ($p(R) = 0.8$, Figure 2A); denoted as bandit I. We expect any learning agent to be able to consistently solve this task. Bandit II has two winning arms. One of these (arm 7, $p(R) = 0.8$) though higher payout than the other (arm 3, $p(R) = 0.6$). The second arm can act as a “distractor” leading an to settle on this suboptimal choice. Bandit III also has a single winning arm, but the overall probability of receiving any reward is very low ($p(R) = 0.02$ for the winning arm, $p(R) = 0.01$ for all others). Sparse rewards problems like these are difficult to solve and are common feature of both the real

world, and artificial environments like Go, chess, and class Atari video games (45–47). The fourth bandit (IV) has 121 arms, and a complex randomly generated reward structure. Bandits of this type and size are probably at the limit of human performance (48).

World model and distance. All bandits share a simple basic common structure. They have a set of n -arms, each of which delivers rewards in a probabilistic fashion. This lends itself to simple discrete n -dimensional world model, with a memory slot for each arm/dimension. Each slot then represents the independent probability of receiving a reward (Supp. Fig 4A).

The Kullback–Leibler divergence (KL) is a widely used information theory metric, which measures the information gained by replacing one distribution with another. It is highly versatile and widely used in machine learning (?), Bayesian reasoning (23, 29), visual neuroscience (29), experimental design (69), compression (70?) and information geometry (71), to name a few examples. KL has seen extensive use in reinforcement learning.

The Kullback–Leibler (KL) divergence satisfies all five value axioms (Eq. 6).

Itti and Baladi (29) developed an approach similar to ours for visual attention, where our information value is identical to their *Bayesian surprise*. Itti and Baladi (2009) showed that compared to range of other theoretical alternative, information value most strongly correlates with eye movements made when humans look at natural images. Again in a Bayesian context, KL plays a key role in guiding *active inference*, a mode of theory where the dogmatic central aim of neural systems is make decisions which minimize free energy (14, 23).

Let E represent value of information, such that $E := KL(M_{t+dt}, M_t)$ (Eq. 6) after observing some state s .

$$KL(M_{t+dt}, M_t) = \sum_{s \in S} M_{t+dt}(s) \log \frac{M_{t+dt}(s)}{M_t(s)} \quad [6]$$

Axiom ?? is satisfied by limiting E calculations to successive memories. Axiom ??-?? are naturally satisfied by KL. That is, $E = 0$ if and only if $M_{t+dt} = M_t$ and $E \geq 0$ for all pairs (M_{t+dt}, M_t) .

To make Axiom 2 more concrete, in Figure 5 we show how KL changes between a hypothetical initial distribution (always shown in grey) and a “learned” distribution (colored). For simplicity's sake we use a simple discrete distribution representing a 10-armed bandit, though the illustrated patterns hold true for any pair of appropriate distributions. In Figure 5C we see KL increases substantially more for a local exchange of probability compared to an even global re-normalization (compare panels A. and B.).

Initializing π_π . In these simulations we assume that at the start of learning an animal should have a uniform prior over the possible actions $A \in \mathbb{R}^K$. Thus $p(a_k) = 1/K$ for all $a_k \in A$. We transform this uniform prior into the appropriate units for our KL-based E using Shannon entropy, $E_0 = \sum_K p(a_k) \log p(a_k)$.

In our simulations we use a tie breaking “right next” heuristic which keeps track of past breaks, and in a round robin fashion iterates rightward over the action space.

Reinforcement learning. Reinforcement learning in all agent models was done with using the TD(0) learning rule (8) (Eq. 7). Where $V(s)$ is the value for each state (arm), \mathbf{R}_t is the *return*

Table 2. Hyperparameters for individual bandits (I-IV).

Agent	Parameter	I	II	III	IV
Dual value	η	0.053	0.017	0.003	5.8e-09
Dual value	α	0.34	0.17	0.15	0.0011
E-greedy	ϵ	0.14	0.039	0.12	0.41
E-greedy	α	0.087	0.086	0.14	0.00048
Annealed e-greedy	τ_E	0.061	0.084	0.0078	0.072
Annealed e-greedy	ϵ	0.45	0.98	0.85	0.51
Annealed e-greedy	α	0.14	0.19	0.173	0.00027
Bayesian	β	0.066	0.13	0.13	2.14
Bayesian	α	0.066	0.03	0.17	0.13
Bayesian	γ	0.13	0.98	0.081	5.045

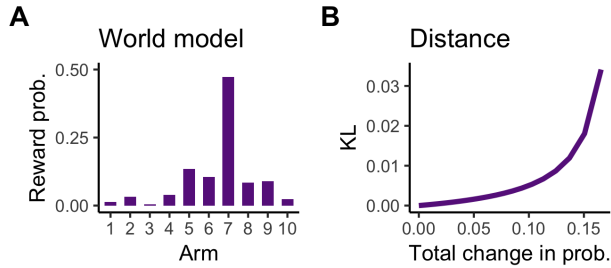


Fig. 4. A world model for bandits. **B.** Example of a single world model suitable for all bandit learning. **B** Changes in the KL divergence—our choice for the distance metric during bandit learning—compared to changes in world model, as by measured the total change in probability mass.

for the current trial, and α is the learning rate ($0 - 1$). See the *Hyperparameter optimization* section for information on how α chosen for each agent and bandit.

$$V(s) = V(s) + \alpha(\mathbf{R}_t - V(s)) \quad [7]$$

The return \mathbf{R}_t differed between agents. Our dual value agent, and both the variations of the e-greedy algorithm, used the reward from the environment R_t as the return. This value was binary. The Bayesian reward agent used a combination of information value and reward $\mathbf{R}_t = R_t + \beta E_t$, with the weight β tuned as described below.

Hyperparameter optimization. The hyperparameters for each agent were tuned independently for each bandit using a modified version of Hyperband (72). For a description of hyperparameters seen Table 1, and for the values themselves Table ??.

Exploration and value dynamics. While agents earned nearly equivalent total reward in Bandit I (Fig 3, top row), their exploration strategies were quite distinct. In Supp. Fig 6B-D) we compare three prototypical examples of exploration, for each major class of agent: ours, Bayesian, and E-greedy for Bandit I. In Supp. Fig 6A) we include an example of value learning in our agent.

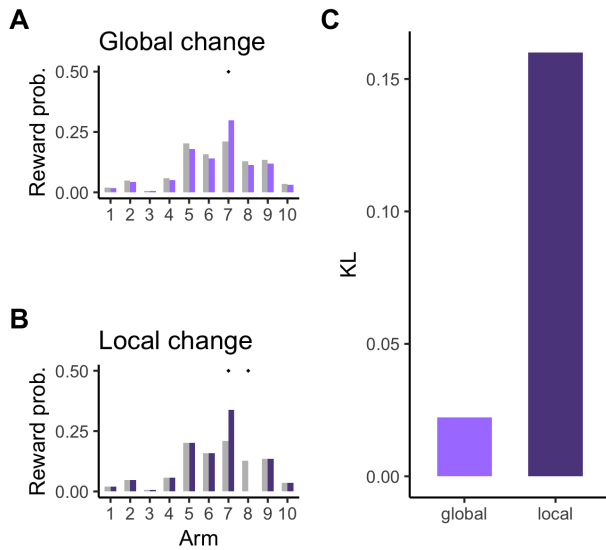


Fig. 5. An example of observation specificity during bandit learning. **A.** A initial (grey) and learned (distribution), where the hypothetical observation s increases the probability of arm 7 by about 0.1, and the expense of all the other probabilities. **B.** Same as A except that the decrease in probability comes only from arm 8. **C.** The KL divergence for local versus global learning.

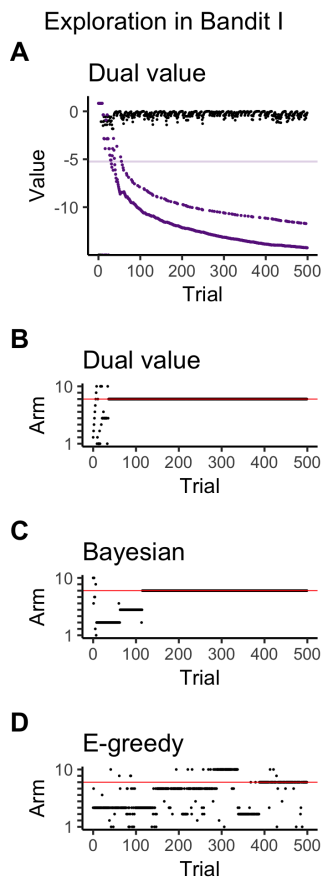


Fig. 6. Exploration and value dynamics. **A.** An example of our dual value learning algorithm during 500 trials on Bandit. The light purple line represents the boredom threshold η (Eq. 5). **B.** An example of exploration dynamics (i.e arm selection) on Bandit. Note how the search is structured, and initially sequential. **C-D.** Exploration dynamics for two other agents. **C.** The Bayesian agent, which like our algorithm uses active sampling, and values information. Note how this shows a mixture of structures and repeated choices, mixed with seemingly random behavior. **D.** The E-greedy agent, which uses purely random sampling. Note how here the agent is either greedy, repeating the same arm, or seemingly random.

Mathematical Appendix.

Information value as a dynamic programming problem. To find greedy dynamic programming (8, 42) answers we must prove our memory M has optimal substructure. By optimal substructure we mean that M can be partitioned into a small number, collection, or series of memories, each of which is itself a dynamic programming solution. In general by proving we can decompose some optimization problem into a small number of sub-problems whose optimal solution are known, or easy to prove, it becomes trivial to prove that we can also grow the series optimally. That is, proving optimal sub-structure nearly automatically allows for proof by induction (42).

Theorem 1 (Optimal substructure). *Assuming transition function δ is deterministic, if $V_{\pi_E}^*$ is the optimal information value given by π_E , a memory M_{t+dt} has optimal substructure if the the last observation s_t can be removed from M_t , by $M_{t+dt} = f^{-1}(M_{t+dt}, s_t)$ where the resulting value $V_{t-dt}^* = V_t^* - F(M_t, a_t)$ 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 $\hat{M}_{t-dt} \neq M_{t-dt}$ and for which $\hat{V}_{t-dt}^* > V_{t-dt}^*$.

To recover the known optimal memory M_t we lift \hat{M}_{t-dt} to $M_t = f(\hat{M}_{t-dt}, s_t)$. This implies $\hat{V}^* > V^*$ which in turn contradicts the purported original optimality of V^* and therefore $\hat{\pi}_E$. \square

Bellman solution. Armed with optimal substructure of M we want to do the next natural thing and find a recursive Bellman solution to maximize our value function for F (Eq. 1). (A Bellman solution of F is also a solution for E (Eq.2). We do this in the classic way by breaking up the series for F into an initial value F_0 , and the remaining series in the summation. We can then apply this same decomposition recursively (Eq 3) to arrive at a final “two-step” or recursive form which is shown Eq. 8).

$$\begin{aligned} V_{\pi_E}^*(M_0) &= \max_{a \in A} \left[\sum_{t=0}^{\infty} F(M_t, a_t) \right] \\ &= \max_{a \in A} \left[F(M_0, a_0) + \sum_{t=1}^{\infty} F(M_{t+dt}, a_{t+dt}) \right] \\ &= F(M_0, a_0) + \max_{a \in A} \left[\sum_{t=1}^{\infty} F(M_{t+dt}, a_{t+dt}) \right] \\ &= F(M_0, a_0) + V_{\pi_E}^*(M_{t+dt}) + V_{\pi_E}^*(M_{t+2}), \dots \end{aligned} \quad [8]$$

A greedy policy explores exhaustively. To prevent any sort of sampling bias, we need our exploration policy π_E (Eq.3) to visit each state s in the space S . As our policy for E is a greedy policy, proofs for exploration are really sorting problems. That is if a state is to be visited it must have highest value. So if every state must be visited (which is what we need to prove to avoid bias) then under a greedy policy every state’s value must, at one time or another, be the maximum value.

We assume implicitly here the action policy π_E can visit all possible states in S . If for some reason π_E can only visit a subset of S , then the following proofs apply only to exploration of that subset.

To begin our proof, some notation. 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 $Z_{t+} = \{s | s \in P \text{ and } s \notin Z_t\}$. As sorting requires ranking, we also need to formalize ranking. To do this we take an algebraic approach, are define inequality for any three real numbers (a, b, c) (Eq. 9).

$$\begin{aligned} a \leq b &\Leftrightarrow \exists c; b = a + c & [9] \\ a > b &\Leftrightarrow (a \neq b) \wedge (b \leq a) & [10] \end{aligned}$$

Theorem 2 (State search: breadth). *A greedy policy π is the only deterministic policy which ensures all states in S are visited, such that $Z = S$.*

Proof. Let $\mathbf{E} = (E_1, E_2, \dots)$ be ranked series of E values for all states S , 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. 9 $E_i - c = E_j$.

Therefore, again by Eq. 9, $\exists \int \delta E(s) \rightarrow -c$.

Recall: Axiom 5.

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

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 5 holds. Therefore, the policy $\hat{\pi}_E$ will leave at least 1 option unexplored and $S \neq Z$. \square

Theorem 3 (State search: depth). *Assuming a deterministic transition function Λ , a greedy policy π_E will resample S to convergence at $E_t \leq \eta$.*

Proof. Recall: Axiom 5.

Each time π_E^* visits a state s , so $M \rightarrow M'$, $F(M', a_{t+dt}) < F(M, a_t)$

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

By induction, if $\pi^* E$ will visit all $s \in S$ in T trials, it will revisit them in $2T$, therefore as $T \rightarrow \infty$, $E \rightarrow 0$. \square

Optimality of π_π . In the following section we prove two things about the optimality of π_π . First, if π_R and/or π_E had any optimal asymptotic property for value learning before their inclusion into our scheduler, they retain that optimal property under π_π . Second, we use this Theorem to show if both π_R and π_E are greedy, and π_π is greedy, then Eq 5 is certain to maximize total value. This is analogous to the classic activity selection problem (42).

Independent policy convergence.

Theorem 4 (Independence policy convergence under π_π). *Assuming an infinite time horizon, if π_E is optimal and π_R is optimal, then π_π is also optimal in the same senses as π_E and π_R .*

Proof. The optimality of π_π can be seen by direct inspection. If $p(R = 1) < 1$ and we have 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 $p_i R$ 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 reasonable constraint, as rewards in the real world are rarely certain. A more naturalistic to handle this edge case is to introduce reward satiety, or a model physiological homeostasis (61, 62).

Optimal scheduling for dual value learning problems. In classic scheduling problems the value of any job is known ahead of time (18, 42). 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. 5 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 what environment to expect, let alone having a perfect model of it. As result, we make a worst-case assumption: the environment can arbitrarily change—bifurcate—at any time. This is, it is a highly nonlinear dynamical system (73). 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 (43). We therefore assume that last value is the best predictor of the next value, and use this assumption along with Theorem 4 to complete a trivial proof that Eq. 5 maximizes total value.

Optimal total value. If we prove π_π has optimal substructure, then using the same replacement argument (42) as in Theorem 4, a greedy policy for π_π will maximize total value.

Theorem 5 (Total value maximization of π_π). π_π must have an optimal substructure.

Proof. Recall: Reinforcement learning algorithms are embedded in Markov Decisions space, which by definition have optimal substructure.

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

Recall: The asymptotic behavior of π_R and π_E are independent under π_π (Theorem 4)

If both π_R and π_E have optimal substructure, and are asymptotically independent, then π_π must also have optimal substructure. \square