

1 Exploration and generalization in vast spaces

2 Charley M. Wu^{1,2,*,+}, Eric Schulz^{3,+}, Maarten Speekenbrink³, Jonathan D. Nelson^{1,4}, and
3 Björn Meder^{1,5}

4 ¹Center for Adaptive Behavior & Cognition, Max Planck Institute for Human Development; Lentzeallee 94, 14195
5 Berlin, Germany

6 ²Center for Adaptive Rationality, Max Planck Institute for Human Development; Lentzeallee 94, 14195 Berlin,
7 Germany

8 ³Department of Experimental Psychology, University College London; 26 Bedford Way, London WC1H 0AP,
9 United Kingdom

10 ⁴School of Psychology, University of Surrey, 388 Stag Hill, Guildford GU2 7XH, UK

11 ⁵MPRG iSearch, Max Planck Institute for Human Development; Lentzeallee 94, 14195 Berlin, Germany

12 *cwu@mpib-berlin.mpg.de

13 +these authors contributed equally to this work

14 ABSTRACT

Foraging for food, developing new medicines, and learning complex games are search problems with vast numbers of possible actions. Under time or resource constraints, optimal solutions are generally unobtainable. How do humans generalize and learn which actions to take when not all outcomes can be explored? We present two behavioural experiments and competitively test 27 models for predicting individual search decisions. We find that a Bayesian function learning model, combined with an optimistic sampling strategy, robustly captures how humans use generalization to guide search behaviour. Taken together, these two form a model of exploration and generalization that leads to reproducible and psychologically meaningful parameter estimates, providing novel insights into the nature of human search in vast spaces. Importantly, our modelling results and parameter estimates are recoverable, and can be used to simulate human-like performance, bridging a critical gap between human and machine learning.

16 Introduction

17 From engineering proteins for medical treatment¹ to mastering a game like Go², many complex tasks
18 can be described as search problems³. Frequently, these tasks come with a vast space of possible actions,
19 each corresponding to some reward that can only be observed through experience. In such problems,
20 one must learn to balance the dual goals of exploring unknown options, while also exploiting existing
21 knowledge for immediate returns. This frames the *exploration-exploitation dilemma*, typically studied
22 using the multi-armed bandit framework^{*4,5}, with the assumption that each option has its own reward
23 distribution to be learned independently. Yet under real-world constraints of limited time or resources,
24 it is not enough to know *when* to explore, but also *where*. How could an intelligent agent, biological or
25 machine, learn which actions to take when not all outcomes can be explored?

26 There is an intriguing gap between human and machine learning, since humans are able to quickly
27 learn and adapt to unfamiliar environments, where the same situation is rarely encountered twice^{6,7}. This
28 contrasts with traditional approaches to reinforcement learning, which learn about the distribution of
29 rewards for each state independently⁸. Such an approach falls short in more realistic scenarios where it is

*The multi-armed bandit is a metaphor for a row of slot machines in a casino, where each slot machine has an independent payoff distribution. Solutions to the problem propose different policies for how to learn about which arms are better to play (exploration), while also playing known high-value arms to maximize reward (exploitation).

30 impossible to observe the outcomes of all possible states and actions^{9,10}. How could an agent efficiently
31 learn and make intelligent decisions about where to explore in problems with a vast space of possible
32 actions?

33 In computer science, one method for dealing with vast state spaces is to use *function learning* as a
34 mechanism to generalize prior experience to unobserved states¹¹. The function learning approach relates
35 different state-action contexts to each other by approximating a global value function over all contexts,
36 including unobserved outcomes⁷. This allows for generalization to vast and potentially infinite state spaces,
37 based on a small number of observations. Additionally, function learning scales to problems with complex
38 sequential dynamics and has been used in tandem with *restricted search* methods such as Monte Carlo
39 sampling for navigating intractably large search trees^{2,12}. While restricted search methods have been
40 proposed as models of human reinforcement learning in planning tasks^{13,14}, here we focus on situations in
41 which a rich model of environmental structure supports learning and generalization¹⁵.

42 Function learning has been successfully utilized for adaptive generalization in various machine learning
43 applications^{16,17}. However, relatively little is known about how humans generalize *in vivo* (e.g., in a
44 search task). Building on previous work exploring inductive biases in pure function learning contexts^{18,19},
45 and human behaviour in univariate function optimization²⁰, we present the first definitive research on how
46 people utilize generalization to effectively learn and search for rewards in large state spaces. Across two
47 studies using uni- and bivariate versions of a multi-armed bandit, we compare 27 different models in their
48 ability to predict individual human behaviour.

49 In both experiments, the vast majority of individual subjects are best captured by a model combining
50 function learning using Gaussian Process (*GP*) regression, with an optimistic Upper Confidence Bound
51 (UCB) sampling strategy that directly balances rewards and uncertainty. Importantly, we recover meaning-
52 ful and robust estimates of the nature of human generalization, showing the limits of traditional models
53 of associative learning²¹ in tasks where the environmental structure supports learning and inference.
54 Interestingly, the most predictive model of the behavioural data is also the only Bayesian optimization
55 algorithm with competitive guarantees²². This result has rich theoretical implications for reinforcement
56 learning.

57 The main contributions of this paper are threefold:

- 58 1. We introduce a novel paradigm, *the spatially correlated multi-armed bandit*, which allows us to
59 study the extent to which people use generalization to guide search through vast problem spaces.
- 60 2. We find that a Bayesian model of function learning robustly captures how humans generalize and
61 learn about the structure of the environment.
- 62 3. Participants solve the exploration-exploitation dilemma by optimistically inflating expectations of
63 reward by the underlying uncertainty, with recoverable evidence for the separate phenomena of
64 directed and undirected exploration.

65 Results

66 A useful inductive bias in many real world search tasks is to assume a spatial correlation between rewards
67 (i.e., clumpiness of resource distribution)²³. We present human data and modelling results from two
68 experiments using spatially correlated multi-armed bandits on univariate (Experiment 1) and bivariate
69 (Experiment 2) environments (Fig. 1). The spatial correlation of rewards provides a context to each arm of
70 the bandit, which can be learned and used to generalize to yet unobserved contexts, thereby guiding search

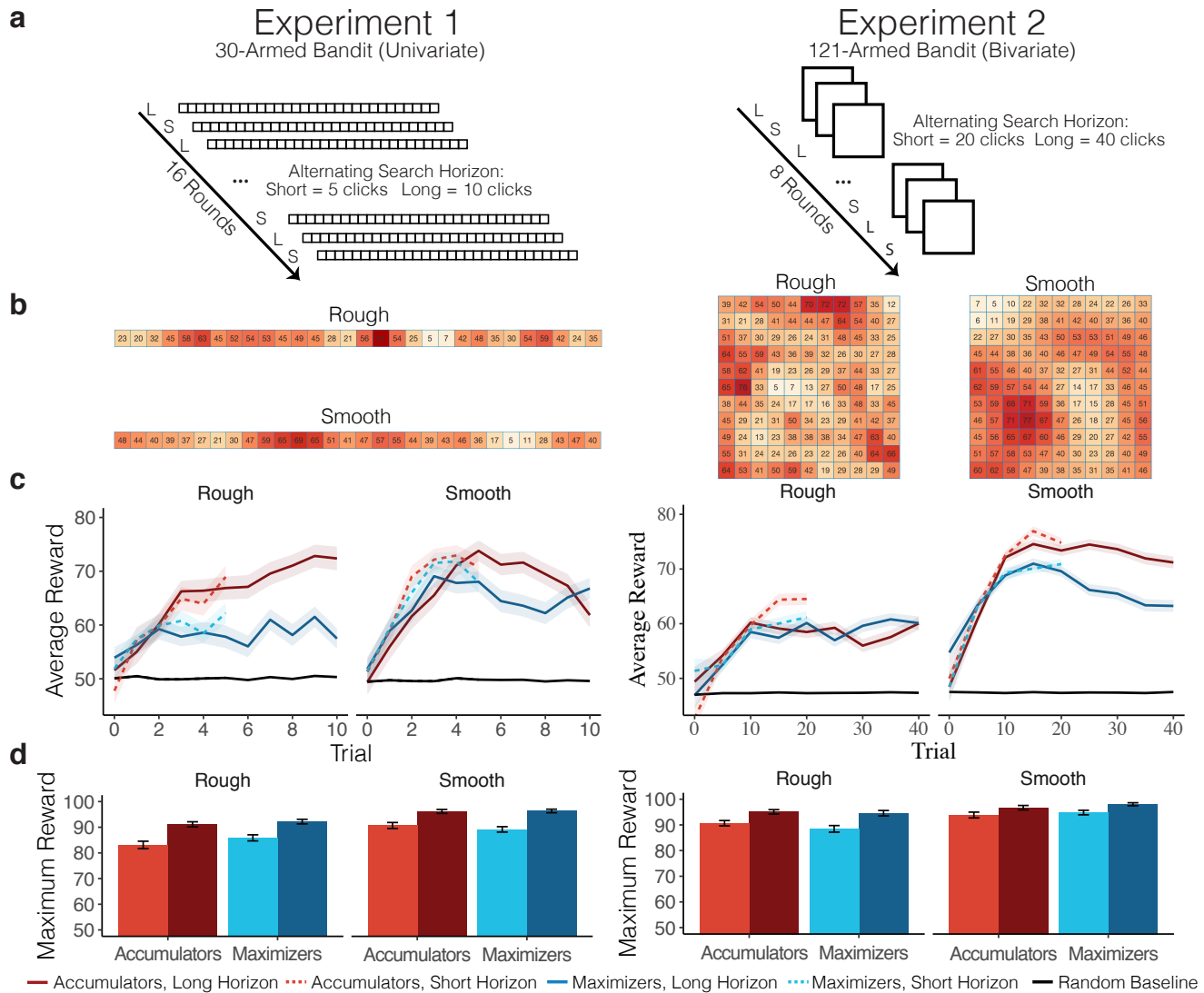


Figure 1. Procedure and behavioural results. Both experiments used a 2×2 between-subject design, manipulating the type of environment (Rough or Smooth) and the payoff condition (Accumulators or Maximizers). **a**, Experiment 1 used a 1D array of 30 possible options, while Experiment 2 used a 2D array (11×11) with 121 options. Experiments took place over 16 (Experiment 1) or 8 (Experiment 2) rounds, with a new environment sampled without replacement for each round. Search horizons alternated between rounds, with horizon order counter-balanced between subjects. **b**, Examples of fully revealed search environments, whereas tiles were initially blank at the beginning of each round, except for a single randomly revealed tile. Subjects were assigned to one of two different classes of environments, differing in the extent of spatial correlations (smoothness of the environment, see Methods). We describe behavioural results according to the average reward earned (**c**, Accumulator goal) and in terms of the maximum reward revealed up until the end of each round (**d**, Maximizer goal), where colours indicate the assigned payoff condition. Short horizon trials are indicated by lighter colours and/or dashed lines. Shaded regions (**c**) and error bars (**d**) show standard error of the mean. Black lines represent a random baseline simulated over 10,000 rounds.

71 decisions. Additionally, since recent work has connected both spatial and conceptual representations to

72 a common neural substrate²⁴, our results in a spatial domain provide potential pathways to other search
73 domains, such as contextual^{25,26} or semantic search^{27,28}.

74 Experiment 1

75 Participants searched for rewards on a 1×30 grid world, where each tile was a reward-generating arm of
76 the bandit (Fig. 1a). The mean rewards of each arm were spatially correlated, with stronger correlations in
77 smooth than in rough environments (between subjects; Fig. 1b). Participants were either assigned the goal
78 of accumulating the largest average reward (*Accumulators*), thereby balancing exploration-exploitation,
79 or of finding the best overall tile (*Maximizers*), an exploration goal directed towards finding the global
80 maximum. We hypothesized that, if search behaviour is guided by function learning, participants would
81 perform better and learn faster in smooth environments, since stronger spatial correlations reveal more
82 information about nearby tiles²⁹.

83 Participants in smooth environments obtained higher average rewards than participants in rough
84 environments ($t(79) = 3.58, p < .001, d = 0.8$), consistent with the hypothesis that spatial patterns in the
85 environment can be learned and used to guide search. The learning curves in Figure 1c show that longer
86 search horizons (solid lines) do not always lead to higher average reward ($t(80) = 0.60, p > .5, d = 0.07$).
87 We analysed both average reward and the maximum reward obtained for each subject, irrespective of their
88 payoff condition (Maximizer or Accumulator). Interestingly, while Accumulators performed better than
89 Maximizers on the average reward criterion ($t(79) = 2.89, p < .01, d = 0.65$), they performed equally
90 well when trying to find the highest overall reward ($t(79) = -0.43, p > .6, d = 0.09$; Fig. 1d). Thus, a
91 strategy balancing exploration and exploitation, at least for human learners, may find the global maximum
92 *en passant*.

93 Experiment 2

94 Experiment 2 had the same design as Experiment 1, but used a 11×11 grid representing an underlying
95 bivariate reward function (Fig. 1 right). Participants obtained higher rewards in smooth environments than
96 in rough environments ($t(78) = 6.55, p < .001, d = 1.47$), as in Experiment 1, but with a larger effect
97 size. As in Experiment 1, Accumulators were as good as Maximizers at discovering the highest rewards
98 ($t(78) = 0.01, p > .5, d = 0.002$). In Experiment 2, however, Accumulators did not perform substantially
99 better than Maximizers in terms of average reward ($t(78) = -1.31, p > .1, d = 0.29$). Again, short search
100 horizons led to the same level of performance as longer horizons, ($t(79) = -0.96, p > .34, d = 0.11$),
101 suggesting that for people, frugal search can be quite efficient. We also present full results for learning
102 over rounds and trials in the SI.

103 Modelling Generalization and Search

104 We competitively tested a diverse set of 27 different models in their ability to predict each subject's trial-
105 by-trial choices (Fig. S1, Table S1). These models include different combinations of models of learning
106 and sampling strategies, along with simple heuristics, which make predictions without maintaining a
107 model of the world. By far the most successful models used *Gaussian Process* (\mathcal{GP}) regression^{30,31} as a
108 method for learning an underlying value function relating all state-action contexts to each other, and *Upper*
109 *Confidence Bound* (UCB) sampling³² for predicting where to sample next. This maps onto the distinction
110 between belief and sampling models, central to theories in statistics³³, psychology³⁴, and philosophy of
111 science³⁵.

112 \mathcal{GP} s provide an expressive model for human function learning, and in contrast to neural network
113 function approximators³⁶, yield psychologically interpretable parameter estimates. \mathcal{GP} function learning

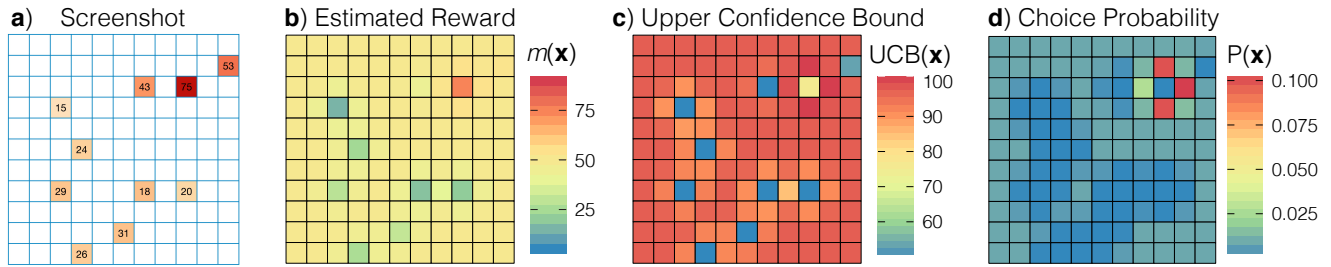


Figure 2. Overview of \mathcal{GP} -UCB specified using median participant parameter estimates (see Table S1). **a**, Screenshot of Experiment 2. Participants were allowed to select any tile until the search horizon was exhausted. **b**, Estimated reward as predicted by the \mathcal{GP} function learning engine, based on the sampled points in Panel **a**. (Not shown, the estimated uncertainty). **c**, Upper confidence bound of predicted rewards. **d**, Choice probabilities after a softmax choice rule, $P(\mathbf{x}) = \exp(\text{UCB}(\mathbf{x})/\tau) / \sum_{j=1}^N \exp(\text{UCB}(\mathbf{x}_j)/\tau)$, where τ is the temperature parameter (i.e., lower temperature values lead to more precise predictions).

114 can guide search by making normally distributed predictions about the expected mean $m(\mathbf{x})$ and the
 115 underlying uncertainty $s(\mathbf{x})$ (estimated as a standard deviation) for each option \mathbf{x} in the global state space
 116 (see Fig. 2b), conditioned on a finite number of previous observations of rewards $\mathbf{y}_T = [y_1, y_2, \dots, y_T]^\top$
 117 at inputs $\mathbf{X}_T = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$. Similarities between points are modeled by a *Radial Basis Function* (RBF)
 118 kernel:

$$119 \quad k_{RBF}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\lambda}\right), \quad (1)$$

120
 121 where λ governs how quickly correlations between points \mathbf{x} and \mathbf{x}' (e.g., two tiles on the grid) decay towards
 122 zero as their distance increases. We use λ as a free parameter, which can be interpreted psychologically as
 123 the extent to which people generalize spatially. This is similar to Shepard's gradient of generalization³⁷,
 124 which also models generalization as an exponentially decreasing function of distance.

125 Given estimates about expected rewards $m(\mathbf{x})$ and the underlying uncertainty $s(\mathbf{x})$ (expressed as a
 126 standard deviation), UCB sampling makes predictions about where participants search next (Fig. 2c) using
 127 a sum,

$$128 \quad \text{UCB}(\mathbf{x}) = m(\mathbf{x}) + \beta s(\mathbf{x}), \quad (2)$$

130 where β is a free parameter governing how much the reduction of uncertainty is weighted relative to
 131 expectations of reward. This trade-off between exploiting known high-value rewards and exploring to
 132 reduce uncertainty³⁸ can be interpreted as optimistically inflating expected rewards by their attached
 133 uncertainty, and can be decomposed into two separate components that only sample points based on high
 134 expected reward (Pure Exploitation) or high uncertainty (Pure Exploration).

$$135 \quad \text{PureExploit}(\mathbf{x}) = m(\mathbf{x}) \quad (3)$$

$$136 \quad \text{PureExplore}(\mathbf{x}) = s(\mathbf{x}) \quad (4)$$

138 Figure 2 shows how the \mathcal{GP} -UCB Function Learning Model makes inferences about the search space,
 139 and uses UCB sampling (with a softmax choice rule) to make probabilistic predictions about where the
 140 participant will sample next. We refer to this model as the *Function Learning Model* and contrast it with

141 a *Mean Tracker*. The Mean Tracker is a type of Kalman Filter without temporal dynamics⁵, and is a
 142 more traditional type of associative learning model that learns the distribution of rewards for each state
 143 independently. Like the \mathcal{GP} -UCB Function Learning Model, the Mean Tracker also generates normally
 144 distributed predictions $m(\mathbf{x})$ and $s(\mathbf{x})$, which we combine with the same set of sampling strategies to make
 145 probabilistic predictions about search.

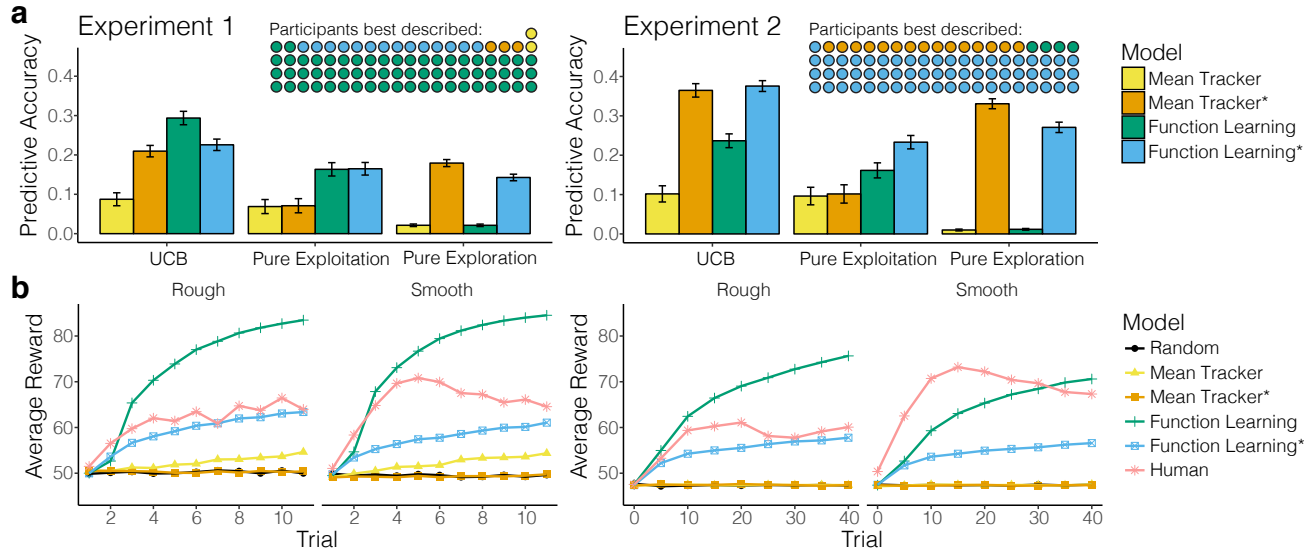


Figure 3. Modelling results. a, Cross-validated predictive accuracy of each model (bars), with the number of participants best described shown as an icon array (inset; aggregated by sampling strategies). Asterisks (*) indicate a localized variant of the Mean Tracker or Function Learning models, where predictions are weighted by the inverse distance from the previous choice (see Methods). **b**, Averaged learning curves of participants and models (UCB only) simulated over 10,000 replications using sampled participant parameter estimates. Learning curves (and parameter estimates) are separated by environment, but aggregated over payoff conditions and search horizons.

146 Modelling results

147 Experiment 1

148 Instead of learning rewards for each state independently, as assumed in the Mean Tracker and traditional
 149 associative learning models, participants were better described by the Function Learning Model ($t(80) =$
 150 14.01 , $p < .001$ $d = 1.56$; comparing cross-validated predictive accuracies, both using UCB sampling),
 151 providing evidence against the assumption of state independence. Furthermore, by decomposing the
 152 UCB sampling algorithm into Pure Exploitation or Pure Exploration components, we show that both
 153 high expectations of reward and the reduction of uncertainty are necessary components for the Function
 154 Learning Model to predict human search behaviour, with Pure Exploitation ($t(80) = 8.85$, $p < .001$,
 155 $d = 0.98$) and Pure Exploration ($t(80) = 16.63$, $p < .001$, $d = 1.85$) performing worse at predicting
 156 human behaviour than the combined UCB algorithm.

157 The distance between sequential samples was more localized than chance ($t(160) = 31.2$, $p < .0001$,
 158 $d = 1.92$; see Fig. S5), as has also been observed in semantic search²⁷ and causal learning³⁹ domains.
 159 Thus, we created a localized variant of both Mean Tracker and Function Learning models (indicated
 160 by an asterisk *; Fig. 3a), giving larger weight to options closer to the previous choice (see Methods).
 161 While localization improved predictive accuracy for Mean Tracker predictions ($t(80) = 16.13$, $p < .001$,

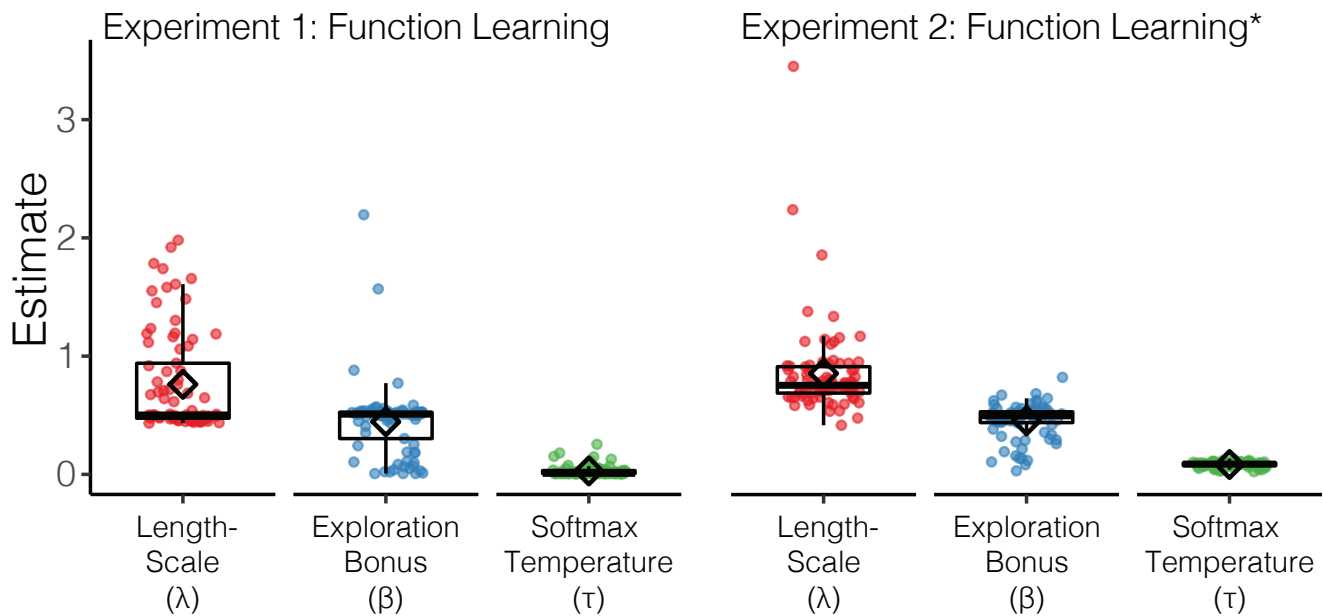


Figure 4. Parameter estimates of the best predicting model for each experiment. Each coloured point is the median estimate of a participant, with boxplots showing the median (line), mean (diamond), and interquartile range. λ is the length-scale of the RBF kernel reflecting the extent to which people generalize, β is the exploration bonus of the UCB sampling strategy, and τ is the temperature of the softmax choice rule.

162 $d = 1.79$; UCB), it led to worse performance for the \mathcal{GP} -UCB Function Learning Model ($t(80) = -5.05$,
 163 $p < .001$, $d = 0.56$).

164 Overall, 56 out of 81 participants were best described by the Function Learning Model combining \mathcal{GP}
 165 learning with UCB sampling (66 of 81 participants including the localized variant), again showing that
 166 generalization is important for predicting search behaviour. Figure 3b shows simulated learning curves of
 167 each model in comparison to human performance, where models were specified using parameters sampled
 168 from participant estimates. Whereas the Mean Tracker models achieve performance close to random, the
 169 Function Learning Model (and its localized variant) behaves sensibly, by utilizing the structure of the
 170 environment, and achieves similar or better performance than humans.

171 Looking more closely at the parameter estimates of the \mathcal{GP} -UCB Function Learning Model (Fig.
 172 4), we find that people tend to underestimate the extent of spatial correlations, with estimated λ values
 173 significantly lower than the ground truth ($\lambda_{Smooth} = 2$ and $\lambda_{Rough} = 1$) for both Smooth (mean estimate:
 174 $\hat{\lambda} = 0.82$, $t(41) = -17.60$, $p < .001$, $d = 2.71$) and Rough environments ($\hat{\lambda} = 0.78$, $t(38) = -3.89$,
 175 $p < .001$, $d = 0.62$), which could be interpreted as a tendency to avoid overgeneralization^{7,40}. Remarkably,
 176 simulations suggest that a tendency towards undergeneralization can benefit search performance (Fig. S4).
 177 The exploration bonus of UCB sampling (β) was robustly estimated above 0 ($\hat{\beta} = 0.47$, $t(80) = 12.78$,
 178 $p < .001$, $d = 1.42$, compared to the lower estimation bound), indicating participants valued the exploration
 179 of uncertain options, along with exploiting high expectations of reward. Additionally, we found very low
 180 estimates of the softmax temperature (τ), corresponding to more precise model predictions ($\hat{\tau} = 0.01$).
 181 The model comparison and parameter estimates were highly robust and recoverable (Figs. S2-S3).

182 Experiment 2

183 In a more complex bivariate environment, the Function Learning Model again predicted participants'
184 choices more accurately than the Mean Tracker Model, for both its standard ($t(79) = 9.99$, $p < .001$,
185 $d = 1.12$; UCB) and localized variants ($t(79) = 2.05$, $p < .05$, $d = 0.23$; UCB, Fig. 3a). In Experiment
186 2's two-dimensional search environment, adding localization improved predictions for both Mean Tracker
187 ($t(79) = 19.92$, $p < .001$, $d = 2.23$; UCB) and the Function Learning Model ($t(79) = 10.47$, $p < .001$,
188 $d = 1.17$; UCB), in line with the stronger tendency towards localized sampling behaviour in Experiment
189 2 (see Fig. S5). 61 out of 80 participants were best predicted by the Function Learning* Model with
190 UCB sampling, whereas only 12 participants were best described by the Mean Tracker* Model with UCB.
191 Again, both components of the UCB strategy—the expected reward ($t(79) = -6.44$, $p < .001$, $d = 0.72$)
192 and the attached uncertainty ($t(79) = -14.32$, $p < .001$, $d = 1.60$)—were necessary to predict choices.

193 As in Experiment 1, simulated learning curves of the Mean Tracker models performed close to random,
194 whereas both variants of the Function Learning Model achieved performance similar to participants (Fig.
195 3b). Median parameter estimates per participant from the Function Learning* Model (Fig. 4) showed
196 that participants again underestimated the strength of the underlying spatial correlation in both Smooth
197 ($\hat{\lambda} = 0.92$, $t(42) = -14.62$, $p < .001$, $d = 2.22$; comparison to $\lambda_{Smooth} = 2$) and Rough environments
198 ($\hat{\lambda} = 0.78$, $t(36) = -5.31$, $p < .001$, $d = 0.87$; comparison to $\lambda_{Rough} = 1$), suggesting a robust tendency
199 to undergeneralize. The estimated exploration bonus β was again greater than 0 ($\hat{\beta} = 0.45$, $t(79) = 27.02$,
200 $p < .001$, $d = 3.02$, compared to the lower estimation bound), while the estimated softmax temperature
201 parameter τ was slightly larger than in Experiment 1 ($\hat{\tau} = 0.09$; see Table S1), likely due to localization.
202 As in Experiment 1, the model comparison and parameter estimates were highly robust and recoverable
203 (Figs. S2-S3).

204 Experiment 2 therefore replicated the main findings of Experiment 1. Taken together, these results
205 provide strong evidence that human search behaviour is best explained by a combination of function
206 learning paired with an optimistic trade-off between exploration and exploitation.

207 Discussion

208 How should one behave in situations where the number of possible actions is vast and not all possibilities
209 can be explored? Function learning provides a possible mechanism for generalization, by relating different
210 options using spatial context. Gaussian Processes (\mathcal{GP}) combined with Upper Confidence Bound (UCB)
211 sampling have been successfully applied to problems in ecology⁴¹, robotics⁴², and biology⁴³, but there
212 has been little psychological research on human behaviour in such tasks.

213 We have presented the first study to apply cognitive modelling to predict individual decisions in such
214 a complex search task. Our rigorous comparison of 27 models yielded robust and recoverable model
215 comparisons (Fig. S2) and parameter estimates (Fig. S3). The spatial correlation of rewards made it
216 possible to generalize to unseen rewards by learning an approximate underlying value function based
217 on spatial context. Results show that participants capitalized on spatial context in all task variants, and
218 performed best in environments with the strongest spatial correlations.

219 Through multiple analyses, including trial-by-trial predictive cross-validation and simulated behaviour
220 using participant parameter estimates, we competitively studied which models best predicted human
221 behaviour. The vast majority of participants were best described by the \mathcal{GP} -UCB Function Learning
222 Model or its localized variant. Parameter estimates from the best-fitting \mathcal{GP} -UCB models suggest there
223 was a systematic tendency to undergeneralize the extent of spatial correlations, which can be a beneficial
224 bias for search (Fig. S4).

225 Whereas previous research on exploration bonuses has had mixed results^{5,10,44}, we find robustly

recoverable parameter estimates for the separate phenomena of directed exploration encoded in β and the random, undirected exploration encoded in the softmax temperature parameter τ , in the \mathcal{GP} -UCB Function Learning Model. Even though UCB sampling is both *optimistic* (always treating uncertainty as positive) and *myopic* (only planning the next timestep), it is nonetheless the only algorithm with known performance guarantees in a bandit setting (i.e., sublinear regret, or in other words, monotonically increasing average reward)²². This suggests a remarkable concurrence between intuitive human strategies and the state of the art in machine learning.

The \mathcal{GP} -UCB Function Learning Model also offers many opportunities for theory integration. The Mean Tracker models as specified here can be reformulated as special case of a \mathcal{GP} regression model⁴⁵ and hence implemented in our Function Learning Model. In addition, when the length-scale of the RBF kernel approaches zero ($\lambda \rightarrow 0$), the Function Learning Model effectively assumes state-independence, as in the Mean Tracker Model. Thus, there may be a continuum of reinforcement learning models, ranging from the traditional assumption of state independence to the opposite extreme, of complete state interdependence. Moreover, a \mathcal{GP} is also equivalent to a Bayesian Neural Network with infinite nodes⁴⁶, suggesting a further link to distributed function learning models⁴⁷. Indeed, one explanation for the impressive performance of Deep Reinforcement Learning¹² is that neural networks are specifically a powerful type of function approximator⁴⁸.

Lastly, recent findings have connected both spatial and conceptual representations to a common neural substrate²⁴, suggesting a potential avenue for applying the same \mathcal{GP} -UCB Function Learning Model for modelling human behaviour in domains such as contextual^{25,26} or semantic search^{27,28}. Marrying powerful yet interpretable function learning techniques with methods commonly applied in the reinforcement learning literature can further advance understanding of adaptive behaviour in complex and uncertain environments.

Methods

Participants

81 participants were recruited from Amazon Mechanical Turk for Experiment 1 (25 Female; mean \pm SD age 33 ± 11), and 80 for Experiment 2 (25 Female; mean \pm SD age 32 ± 9). In both experiments, participants were paid a participation fee of \$0.50 and a performance contingent bonus of up to \$1.50. Participants earned on average $\$1.14 \pm 0.13$ and spent 8 ± 4 minutes on the task in Experiment 1, while participants earned on average $\$1.64 \pm 0.20$ and spent 8 ± 4 minutes on the task in Experiment 2. Participants were only allowed to participate in one of the experiments, and were required to have a 95% HIT approval rate and 1000 previously completed HITs. The Ethics Committee of the Max Planck Institute for Human Development approved the methodology and all participants consented to participation through an online consent form at the beginning of the survey.

Design

Both experiments used a 2×2 between-subjects design, where participants were randomly assigned to one of two different payoff structures (*Accumulators* vs. *Maximizers*) and one of two different classes of environments (*Smooth* vs. *Rough*). Each grid world represented a (either uni- or bivariate) function, with each observation including normally distributed noise, $\epsilon \sim \mathcal{N}(0, 1)$. The task was presented over either 16 rounds (Exp. 1) or 8 rounds (Exp. 2) on different grid worlds drawn from the same class of environments. Participants had either a short or long search horizon (Exp. 1: [5,10]; Exp. 2: [20,40]) to sample tiles on the grid, including repeat clicks. The search horizon alternated between rounds (within subject), with initial horizon length counterbalanced between subjects.

269 **Materials and procedure**

270 Participants observed four fully revealed example environments and correctly completed three compre-
271 hension questions, prior to starting the task. Example environments were drawn from the same class
272 of environments assigned to the participant (Smooth or Rough). At the beginning of each round, one
273 random tile was revealed and participants could click any of the tiles in the grid until the search horizon
274 was exhausted, including re-clicking previously revealed tiles. Clicking an unrevealed tile displayed the
275 numerical value of the reward along with a corresponding colour aid, where darker colours indicated
276 higher point values. Per round, observations were scaled to a randomly drawn maximum value in the
277 range of 65 to 85, so that the value of the global optima could not be easily guessed (e.g., a value of 100).
278 Re-clicked tiles could show some variations in the observed value due to noise. For repeat clicks, the most
279 recent observation was displayed numerically, while hovering over the tile would display the entire history
280 of observation. The colour of the tile corresponded to the mean of all previous observations.

281 **Payoff conditions**

282 We compared performance under two different payoff conditions, requiring either a balance between
283 exploration and exploitation (*Accumulators*) or corresponding to consistently making exploration decisions
284 (*Maximizers*). In each payoff condition, participants received a performance contingent bonus of up to
285 \$1.50. *Accumulators* were given a bonus based on the average value of all clicks as a fraction of the global
286 optima, $\frac{1}{T} \sum (\frac{y_t}{y^*})$, where y^* is the global optimum, whereas *Maximizers* were rewarded using the ratio
287 of the highest observed reward to the global optimum, $(\frac{\max y_t}{y^*})^4$, taken to the power of 4 to exaggerate
288 differences in the upper range of performance and for between-group parity in expected earnings across
289 payoff conditions. Both conditions were equally weighted across all rounds and used noisy but unscaled
290 observations to assign a bonus of up to \$1.50. Subjects were informed in dollars about the bonus earned at
291 the end of each round.

292 **Smoothness of the environment**

293 We used two classes of environments, corresponding to different levels of smoothness. All environments
294 were sampled from a \mathcal{GP} prior with a RBF kernel, where the length-scale parameter (λ) determines the
295 rate at which the correlations of rewards decay over distance. *Rough* environments used $\lambda_{Rough} = 1$ and
296 *Smooth* environments used $\lambda_{Smooth} = 2$, with 40 environments (Exp. 1) and 20 environments (Exp. 2)
297 generated for each class (Smooth and Rough). Both example environments and task environments were
298 drawn without replacement from the assigned class of environments, where smoothness can be understood
299 as the extent of spatial correlations.

300 **Search horizons**

301 We chose two horizon lengths (Short=5 or 20 and Long=10 or 40) that were fewer than the total number
302 of tiles on the grid (30 or 121), and varied them within subject (alternating between rounds and counterbal-
303 anced). Horizon length was approximately equivalent between Experiments 1 and 2 as a fraction of the
304 total number of options (short $\approx \frac{1}{6}$; long $\approx \frac{1}{3}$).

305 **Models of Learning**

306 We use different *Models of Learning* (i.e., Function Learning and Mean Tracker), which combined with
307 a *Sampling Strategy* can make predictions about where a participant will search, given the history of
308 previous observations.

309 Function Learning

310 The Function Learning Model adaptively learns an underlying function mapping spatial locations onto
 311 rewards. We use Gaussian Process (\mathcal{GP}) regression as a Bayesian method of function learning³¹. A \mathcal{GP}
 312 is defined as a collection of points, any subset of which is multivariate Gaussian. Let $f : \mathcal{X} \rightarrow \mathbb{R}^n$ denote a
 313 function over input space \mathcal{X} that maps to real-valued scalar outputs. This function can be modelled as a
 314 random draw from a \mathcal{GP} :

$$315 \quad f \sim \mathcal{GP}(m, k), \quad (5)$$

317 where m is a mean function specifying the expected output of the function given input \mathbf{x} , and k is a kernel
 318 (or covariance) function specifying the covariance between outputs.

$$319 \quad m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] \quad (6)$$

$$320 \quad k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))] \quad (7)$$

322 Here, we fix the prior mean to the median value of payoffs, $m(\mathbf{x}) = 50$ and use the kernel function
 323 to encode an inductive bias about the expected spatial correlations between rewards (see Radial Basis
 324 Function kernel). Conditional on observed data $\mathcal{D}_t = \{\mathbf{x}_j, y_j\}_{j=1}^t$, where $y_j \sim \mathcal{N}(f(\mathbf{x}_j), \sigma^2)$ is drawn from
 325 the underlying function with added noise $\sigma^2 = 1$, we can calculate the posterior predictive distribution for
 326 a new input \mathbf{x}_* as a Gaussian with mean and variance given by:

$$327 \quad \mathbb{E}[f(\mathbf{x}_*)|\mathcal{D}_t] = m_t(\mathbf{x}_*) = \mathbf{k}_*^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}_t \quad (8)$$

$$328 \quad \mathbb{V}[f(\mathbf{x}_*)|\mathcal{D}_t] = v_t(\mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_*, \quad (9)$$

330 where $\mathbf{y} = [y_1, \dots, y_t]^\top$, \mathbf{K} is the $t \times t$ covariance matrix evaluated at each pair of observed inputs, and
 331 $\mathbf{k}_* = [k(\mathbf{x}_1, \mathbf{x}_*), \dots, k(\mathbf{x}_t, \mathbf{x}_*)]$ is the covariance between each observed input and the new input \mathbf{x}_* .

332 We use the Radial Basis Function (RBF) kernel as a component of the \mathcal{GP} function learning algorithm,
 333 which specifies the correlation between inputs.

$$334 \quad k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\lambda}\right) \quad (10)$$

336 This kernel defines a universal function learning engine based on the principles of Bayesian regression and
 337 can model any stationary function[†]. Intuitively, the RBF kernel models the correlation between points as
 338 an exponentially decreasing function of their distance. Here, λ modifies the rate of correlation decay, with
 339 larger λ -values corresponding to slower decays, stronger spatial correlations, and smoother functions. As
 340 $\lambda \rightarrow \infty$, the RBF kernel assumes functions approaching linearity, whereas as $\lambda \rightarrow 0$, there ceases to be any
 341 spatial correlation, with the implication that learning happens independently for each discrete input without
 342 generalization (similar to traditional models of associative learning). We treat λ as a hyper-parameter, and
 343 use cross-validated estimates to make inferences about the extent to which participants generalize.

344 Mean Tracker

345 The Mean Tracker is a type of traditional associative learning model, which assumes the average reward
 346 associated with each option is constant over time (i.e., no temporal dynamics, as opposed to the assumptions

[†]Note, sometimes the RBF kernel is specified as $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}\right)$ whereas we use $\lambda = 2l^2$ as a more psychologically interpretable formulation.

347 of a Kalman Filter or Temporal Difference Learning)⁵, as is the case in our experimental search tasks. In
 348 contrast to the Function Learning Model, the Mean Tracker learns the rewards of each option independently,
 349 by computing an independent posterior distribution for the mean μ_j for each option j . We implement a
 350 version that assumes rewards are normally distributed (as in the \mathcal{GP} Function Learning Model), with a
 351 known variance but unknown mean, where the prior distribution of the mean is again a normal distribution.
 352 This implies that the posterior distribution for each mean is also a normal distribution:

$$353 \quad p(\mu_{j,t} | \mathcal{D}_{t-1}) = \mathcal{N}(m_{j,t}, v_{j,t}) \quad (11)$$

355 For a given option j , the posterior mean $m_{j,t}$ and variance $v_{j,t}$ are only updated when it has been selected
 356 at trial t :

$$357 \quad m_{j,t} = m_{j,t-1} + \delta_{j,t} G_{j,t} [y_t - m_{j,t-1}] \quad (12)$$

$$358 \quad v_{j,t} = [1 - \delta_{j,t} G_{j,t}] v_{j,t-1} \quad (13)$$

360 where $\delta_{j,t} = 1$ if option j was chosen on trial t , and 0 otherwise. Additionally, y_t is the observed reward at
 361 trial t , and $G_{j,t}$ is defined as:

$$362 \quad G_{j,t} = \frac{v_{j,t-1}}{v_{j,t-1} + \theta_{\epsilon}^2} \quad (14)$$

364 where θ_{ϵ}^2 is the error variance, which is estimated as a free parameter. Intuitively, the estimated mean
 365 of the chosen option $m_{j,t}$ is updated based on the difference between the observed value y_t and the prior
 366 expected mean $m_{j,t-1}$, multiplied by $G_{j,t}$. At the same time, the estimated variance $v_{j,t}$ is reduced by a
 367 factor of $1 - G_{j,t}$, which is in the range $[0, 1]$. The error variance (θ_{ϵ}^2) can be interpreted as an inverse
 368 sensitivity, where smaller values result in more substantial updates to the mean $m_{j,t}$, and larger reductions
 369 of uncertainty $v_{j,t}$. We set the prior mean to the median value of payoffs $m_{j,0} = 50$ and the prior variance
 370 $v_{j,0} = 500$

371 Sampling Strategies

372
 373 Given the normally distributed posteriors of the expected rewards, which have mean $m_t(\mathbf{x})$ and variance
 374 $v_t(\mathbf{x})$, for each search option \mathbf{x} (for the Mean Tracker Model, we let $m_t(\mathbf{x}) = m_{j,t}$ and $v_t(\mathbf{x}) = v_{j,t}$, where j
 375 is the index of the option characterized by \mathbf{x}), we assess different sampling strategies that (with a softmax
 376 choice rule) make probabilistic predictions about where participants search next at time $t + 1$.

377 Upper Confidence Bound Sampling

378 Given the posterior predictive mean $m_t(\mathbf{x})$ and its attached standard deviation $s_t(\mathbf{x}) = \sqrt{v_t(\mathbf{x})}$, we calculate
 379 the upper confidence bound using a simple sum

$$380 \quad \text{UCB}(\mathbf{x}) = m_t(\mathbf{x}) + \beta s_t(\mathbf{x}), \quad (15)$$

382 where the exploration factor β determines how much reduction of uncertainty is valued (relative to
 383 exploiting known high-value options) and is estimated as a free parameter.

384 **Pure Exploitation and Pure Exploration**

385 Upper Confidence Bound sampling can be decomposed into a Pure Exploitation component, which only
386 samples options with high expected rewards, and a Pure Exploration component, which only samples
387 options with high uncertainty.

$$388 \text{PureExploit}(\mathbf{x}) = m_t(\mathbf{x}) \quad (16)$$

$$389 \text{PureExplore}(\mathbf{x}) = s_t(\mathbf{x}) \quad (17)$$

391 **Localization of Models**

392 To penalize search options by the distance from the previous choice, we weighted each option by the
393 inverse Manhattan distance (IMD) to the last revealed tile $\text{IMD}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^n |x_i - x'_i|$, prior to the softmax
394 transformation. For the special case where $\mathbf{x} = \mathbf{x}'$, we set $\text{IMD}(\mathbf{x}, \mathbf{x}') = 1$. Localized models are indicated
395 by an asterix (*).

396 **Model Comparison**

397 We use maximum likelihood estimation (MLE) for parameter estimation, and cross-validation to measure
398 out-of-sample predictive accuracy. A softmax choice rule transforms each model's prediction into a
399 probability distribution over options:

$$400 p(\mathbf{x}) = \frac{\exp(q(\mathbf{x})/\tau)}{\sum_{j=1}^N \exp(q(\mathbf{x}_j)/\tau)}, \quad (18)$$

402 where $q(\mathbf{x})$ is the predicted value of each option \mathbf{x} for a given model (e.g., $q(\mathbf{x}) = \text{UCB}(\mathbf{x})$ for the UCB
403 model), and τ is the temperature parameter. Lower values of τ indicate more concentrated probability
404 distributions, corresponding to more precise predictions. All models include τ as a free parameter.
405 Additionally, Function Learning models estimate λ (length-scale), Mean Tracker models estimate θ_ε^2
406 (error variance), and Upper Confidence Bound sampling models estimate β (exploration bonus).

407 **Cross Validation**

408 We fit all models—per participant—using cross-validated MLE, with either a Differential Evolution
409 algorithm⁴⁹ or a grid search if the model contained only a single parameter. Parameter estimates are
410 constrained to positive values in the range $[\exp(-5), \exp(5)]$.

411 Cross-validation is performed by first separating participant data according to horizon length, which
412 alternated between rounds within subject. For each participant, half of the rounds corresponded to a short
413 horizon and the other half corresponded to a long horizon. Within all rounds of each horizon length,
414 we use leave-one-out cross-validation to iteratively form a training set by leaving out a single round,
415 computing a MLE on the training set, and then generating out of sample predictions on the remaining
416 round. This is repeated for all combinations of training set and test set, and for both short and long horizon
417 sets. The cross-validation procedure yielded one set of parameter estimates per round, per participant,
418 and out-of-sample predictions for 120 choices in Experiment 1 and 240 choices in Experiment 2 (per
419 participant). In total, cross-validated model comparisons for both experiments required approximately
420 50,000 hours of computation, or about 3 days distributed across a 716 CPU cluster.

421 **Predictive Accuracy**

422 Prediction error (computed as log loss) is summed up over all rounds, and is reported as *predictive*
423 *accuracy*, using a pseudo- R^2 measure that compares the total log loss prediction error for each model to
424 that of a random model:

$$R^2 = 1 - \frac{\log \mathcal{L}(\mathcal{M}_k)}{\log \mathcal{L}(\mathcal{M}_{\text{rand}})}, \quad (19)$$

425

426

427 where $\log \mathcal{L}(\mathcal{M}_{\text{rand}})$ is the log loss of a random model (i.e., picking options with equal probability) and
428 $\log \mathcal{L}(\mathcal{M}_k)$ is the log loss of model k 's out-of-sample prediction error. Intuitively, $R^2 = 0$ corresponds
429 to prediction accuracy equivalent to chance, while $R^2 = 1$ corresponds to theoretical perfect prediction
430 accuracy, since $\log \mathcal{L}(\mathcal{M}_k)/\log \mathcal{L}(\mathcal{M}_{\text{rand}}) \rightarrow 0$ when $\log \mathcal{L}(\mathcal{M}_k) \ll \log \mathcal{L}(\mathcal{M}_{\text{rand}})$.

431 References

- 432 1. Romero, P. A., Krause, A. & Arnold, F. H. Navigating the protein fitness landscape with Gaussian
433 Processes. *Proc. Natl. Acad. Sci.* **110**, E193–E201 (2013).
- 434 2. Silver, D. *et al.* Mastering the game of Go with deep neural networks and tree search. *Nat.* **529**,
435 484–489 (2016).
- 436 3. Hills, T. T. *et al.* Exploration versus exploitation in space, mind, and society. *Trends Cogn. Sci.* **19**,
437 46–54 (2015).
- 438 4. Steyvers, M., Lee, M. D. & Wagenmakers, E.-J. A Bayesian analysis of human decision-making on
439 bandit problems. *J. Math. Psychol.* **53**, 168–179 (2009).
- 440 5. Speekenbrink, M. & Konstantinidis, E. Uncertainty and exploration in a restless bandit problem. *Top.*
441 *Cogn. Sci.* **7**, 351–367 (2015).
- 442 6. Lee, S. W., Shimojo, S. & O'Doherty, J. P. Neural computations underlying arbitration between
443 model-based and model-free learning. *Neuron* **81**, 687–699 (2014).
- 444 7. Gershman, S. J. & Daw, N. D. Reinforcement learning and episodic memory in humans and animals:
445 An integrative framework. *Annu. Rev. Psychol.* **68**, 101–128 (2017).
- 446 8. Sutton, R. S. & Barto, A. G. *Reinforcement learning: An introduction*, vol. 1 (MIT press Cambridge,
447 1998).
- 448 9. Lake, B. M., Ullman, T. D., Tenenbaum, J. B. & Gershman, S. J. Building machines that learn and
449 think like people. *Behav. Brain Sci.* 1–101 (2016).
- 450 10. Wilson, R. C., Geana, A., White, J. M., Ludvig, E. A. & Cohen, J. D. Humans use directed and
451 random exploration to solve the explore–exploit dilemma. *J. Exp. Psychol. Gen.* **143**, 2074 (2014).
- 452 11. Tesauro, G. Practical issues in temporal difference learning. *Mach. learning* **8**, 257–277 (1992).
- 453 12. Mnih, V. *et al.* Human-level control through deep reinforcement learning. *Nat.* **518**, 529–533 (2015).
- 454 13. Huys, Q. J. *et al.* Interplay of approximate planning strategies. *Proc. Natl. Acad. Sci.* **112**, 3098–3103
455 (2015).
- 456 14. Solway, A. & Botvinick, M. M. Evidence integration in model-based tree search. *Proc. Natl. Acad.*
457 *Sci.* **112**, 11708–11713 (2015).
- 458 15. Guez, A., Silver, D. & Dayan, P. Scalable and efficient Bayes-adaptive reinforcement learning based
459 on monte-carlo tree search. *J. Artif. Intell. Res.* **48**, 841–883 (2013).
- 460 16. Rasmussen, C. E. & Kuss, M. Gaussian processes in reinforcement learning. In *Advances in Neural*
461 *Information Processing Systems*, vol. 4, 1 (2003).

- 462 **17.** Sutton, R. S. Generalization in reinforcement learning: Successful examples using sparse coarse
463 coding. *Adv. Neural Inf. Process. Syst.* 1038–1044 (1996).
- 464 **18.** Lucas, C. G., Griffiths, T. L., Williams, J. J. & Kalish, M. L. A rational model of function learning.
465 *Psychon. Bull. & Rev.* **22**, 1193–1215 (2015).
- 466 **19.** Schulz, E., Tenenbaum, J. B., Duvenaud, D., Speekenbrink, M. & Gershman, S. J. Compositional
467 inductive biases in function learning. *bioRxiv* (2016). DOI 10.1101/091298.
- 468 **20.** Borji, A. & Itti, L. Bayesian optimization explains human active search. In *Advances in Neural*
469 *Information Processing Systems*, 55–63 (2013).
- 470 **21.** Dayan, P. & Niv, Y. Reinforcement learning: Ttfe good, the bad and the ugly. *Curr. opinion*
471 *neurobiology* **18**, 185–196 (2008).
- 472 **22.** Srinivas, N., Krause, A., Kakade, S. & Seeger, M. W. Gaussian process optimization in the bandit
473 setting: No regret and experimental design. In *Proceedings of the 27th International Conference on*
474 *Machine Learning*, 1015–1022 (2010).
- 475 **23.** Wilke, A. *et al.* A game of hide and seek: Expectations of clumpy resources influence hiding and
476 searching patterns. *PloS one* **10**, e0130976 (2015).
- 477 **24.** Constantinescu, A. O., O’Reilly, J. X. & Behrens, T. E. Organizing conceptual knowledge in humans
478 with a gridlike code. *Sci.* **352**, 1464–1468 (2016).
- 479 **25.** Stojic, H., Analytis, P. P. & Speekenbrink, M. Human behavior in contextual multi-armed bandit
480 problems. In *Proceedings of the 37th Annual Meeting of the Cognitive Science Society*, 2290–2295
481 (Cognitive Science Society, 2015).
- 482 **26.** Schulz, E., Konstantinidis, E. & Speekenbrink, M. Putting bandits into context: How function learning
483 supports decision making. *J. Exp. Psychol. Learn. Mem. Cogn.* (in press).
- 484 **27.** Hills, T. T., Jones, M. N. & Todd, P. M. Optimal foraging in semantic memory. *Psychol. review* **119**,
485 431 (2012).
- 486 **28.** Abbott, J. T., Austerweil, J. L. & Griffiths, T. L. Random walks on semantic networks can resemble
487 optimal foraging. *Psychol. Rev.* **122**, 558–569 (2015).
- 488 **29.** Schulz, E., Tenenbaum, J. B., Reshef, D. N., Speekenbrink, M. & Gershman, S. Assessing the
489 perceived predictability of functions. In *Proceedings of the 37th Annual Meeting of the Cognitive*
490 *Science Society*, 2116–2121 (Cognitive Science Society, 2015).
- 491 **30.** Rasmussen, C. & Williams, C. *Gaussian Processes for Machine Learning*. Adaptive Computation
492 and Machine Learning (MIT Press, 2006).
- 493 **31.** Schulz, E., Speekenbrink, M. & Krause, A. A tutorial on Gaussian process regression with a focus on
494 exploration-exploitation scenarios. *bioRxiv* 095190 (2016).
- 495 **32.** Auer, P. Using confidence bounds for exploitation-exploration trade-offs. *J. Mach. Learn. Res.* **3**,
496 397–422 (2002).
- 497 **33.** Lindley, D. V. On a measure of the information provided by an experiment. *The Annals Math. Stat.*
498 986–1005 (1956).
- 499 **34.** Nelson, J. D. Finding useful questions: On bayesian diagnosticity, probability, impact, and information
500 gain. *Psychol. Rev.* **112** (2005).

- 501 **35.** Crupi, V. & Tentori, K. State of the field: Measuring information and confirmation. *Stud. Hist. Philos.*
502 *Sci. Part A* **47**, 81–90 (2014).
- 503 **36.** Le Roux, N. & Bengio, Y. Deep belief networks are compact universal approximators. *Neural*
504 *computation* **22**, 2192–2207 (2010).
- 505 **37.** Shepard, R. N. Toward a universal law of generalization for psychological science. *Sci.* **237**, 1317–
506 1323 (1987).
- 507 **38.** Kaufmann, E., Cappé, O. & Garivier, A. On Bayesian upper confidence bounds for bandit problems.
508 In *Artificial Intelligence and Statistics*, 592–600 (2012).
- 509 **39.** Bramley, N. R., Dayan, P., Griffiths, T. L. & Lagnado, D. A. Formalizing neurath’s ship: Approximate
510 algorithms for online causal learning. *Psychol. review* **124**, 301 (2017).
- 511 **40.** Schulz, E., Speekenbrink, M., Hernández-Lobato, J. M., Ghahramani, Z. & Gershman, S. J. Quantify-
512 ing mismatch in Bayesian optimization. In *NIPS Workshop on Bayesian Optimization: Black-box*
513 *Optimization and beyond* (2016).
- 514 **41.** Gotovos, A., Casati, N., Hitz, G. & Krause, A. Active learning for level set estimation. In *International*
515 *Joint Conference on Artificial Intelligence (IJCAI)*, 1344–1350 (2013).
- 516 **42.** Cully, A., Clune, J., Tarapore, D. & Mouret, J.-B. Robots that can adapt like animals. *Nat.* **521**,
517 503–507 (2015).
- 518 **43.** Sui, Y., Gotovos, A., Burdick, J. & Krause, A. Safe exploration for optimization with Gaussian
519 processes. In *International Conference on Machine Learning*, 997–1005 (2015).
- 520 **44.** Daw, N. D., O’doherly, J. P., Dayan, P., Seymour, B. & Dolan, R. J. Cortical substrates for exploratory
521 decisions in humans. *Nat.* **441**, 876–879 (2006).
- 522 **45.** Reece, S. & Roberts, S. An introduction to Gaussian processes for the Kalman filter expert. In *13th*
523 *Conference on Information Fusion (FUSION)*, 1–9 (IEEE, 2010).
- 524 **46.** Neal, R. M. *Bayesian learning for neural networks* (Springer Science & Business Media, 1996).
- 525 **47.** LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nat.* **521**, 436–444 (2015).
- 526 **48.** Schölkopf, B. Artificial intelligence: Learning to see and act. *Nat.* **518**, 486–487 (2015).
- 527 **49.** Mullen, K., Ardia, D., Gil, D., Windover, D. & Cline, J. DEoptim: An R package for global
528 optimization by differential evolution. *J. Stat. Softw.* **40**, 1–26 (2011). URL [http://www.jstatsoft.](http://www.jstatsoft.org/v40/i06/)
529 [org/v40/i06/](http://www.jstatsoft.org/v40/i06/).
- 530 **50.** Bonawitz, E., Denison, S., Gopnik, A. & Griffiths, T. L. Win-stay, lose-sample: A simple sequential
531 algorithm for approximating bayesian inference. *Cogn. psychology* **74**, 35–65 (2014).

532 Acknowledgements

533 We thank Peter Todd, Neil Bramley, Henrik Singmann, and Mehdi Moussaïd for feedback. This work was
534 supported by the International Max Planck Research School on Adapting Behavior in a Fundamentally
535 Uncertain World (CMW), and DFG grants ME 3717/2-2 to BM and NE 1713/1-2 to JDN as part of
536 the New Frameworks of Rationality (SPP 1516) priority program. Code and data available at [https:](https://github.com/charleywu/gridsearch)
537 [//github.com/charleywu/gridsearch](https://github.com/charleywu/gridsearch)

538 **Author contributions statement**

539 C.M.W. and E.S. designed the experiments, collected and analysed the data, and wrote the paper. M.S.,
540 J.D.N., and B.M. designed the experiments and wrote the paper.

Supplemental Materials: Exploration and generalization in vast spaces

Full Model Comparison

We report the full model comparison of 27 models, of which 12 (i.e., four learning models and three sampling strategies) are included in the main text. We use different *Models of Learning* (i.e., Function Learning and Mean Tracking), which combined with a *Sampling Strategy* can make predictions about where a participant will search, given the history of previous observations. We also include comparisons to *Simple Heuristic Strategies*, which make predictions about search decisions without maintaining a representation of the world (i.e., with no learning model). Table S1 shows the predictive accuracy, the number of participants best described, and the median parameter estimates of each model. Figure S1 shows a more detailed assessment of predictive accuracy, with participants separated by payoff condition and environment type.

Additional Sampling Strategies

Expected Improvement

At any point in time t , the best observed outcome can be described as $\mathbf{x}^+ = \arg \max_{\mathbf{x}_i \in \mathbf{x}_{1:t}} m_t(\mathbf{x}_i)$. Expected Improvement (EXI) evaluates each option by *how much* (in the expectation) it promises to be better than the best observed outcome \mathbf{x}^+ :

$$\text{EXI}(\mathbf{x}) = \begin{cases} \Phi(Z)(m_t(\mathbf{x}) - m_t(\mathbf{x}^+)) + s_t(\mathbf{x})\phi(Z), & \text{if } s_t(\mathbf{x}) > 0 \\ 0, & \text{if } s_t(\mathbf{x}) = 0 \end{cases} \quad (\text{S1})$$

where $\Phi(\cdot)$ is the normal CDF, $\phi(\cdot)$ is the normal PDF, and $Z = (m_t(\mathbf{x}) - m_t(\mathbf{x}^+))/s_t(\mathbf{x})$.

Probability of Improvement

The Probability of Improvement (POI) strategy evaluates an option based on *how likely* it will be better than the best outcome (\mathbf{x}^+) observed so far:

$$\begin{aligned} \text{POI}(\mathbf{x}) &= P(f(\mathbf{x}) \geq f(\mathbf{x}^+)) \\ &= \Phi\left(\frac{m_t(\mathbf{x}) - m_t(\mathbf{x}^+)}{s_t(\mathbf{x})}\right) \end{aligned} \quad (\text{S2})$$

Probability of Maximum Utility

The Probability of Maximum Utility (PMU) samples each option according to the probability that it results in the highest reward of all options in a particular context⁵. It is a form of probability matching and can be implemented by sampling from each option's predictive distribution once, and then choosing the option with the highest sampled pay-off.

$$\text{PMU}(\mathbf{x}) = P(f(\mathbf{x}_j) > f(\mathbf{x}_{i \neq j})) \quad (\text{S3})$$

We implement this acquisition function by Monte Carlo sampling from the posterior predictive distribution of a learning model for each option, and evaluating how often a given option turns out to be the maximum over 1,000 generated samples.

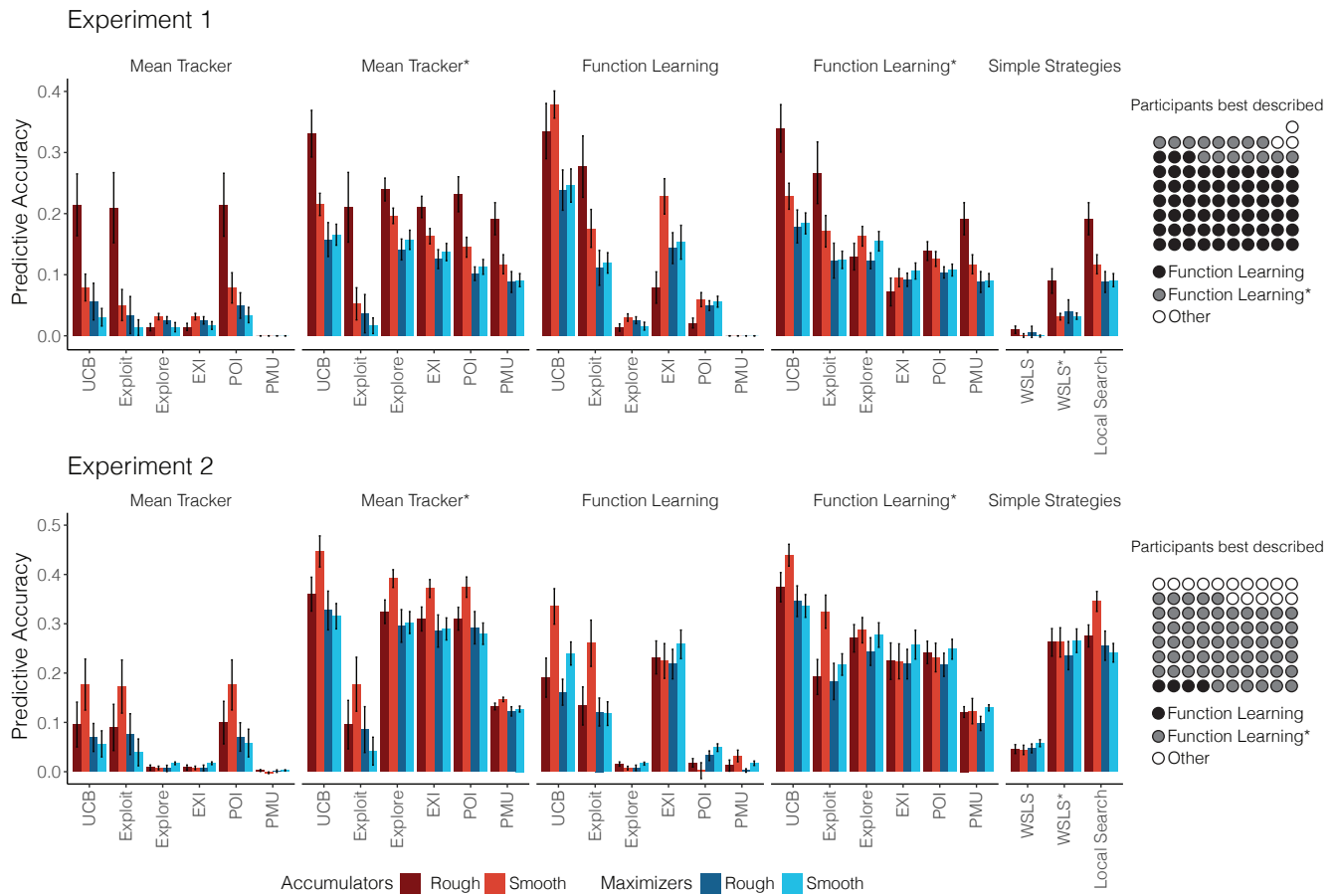


Figure S1. Full model comparison of all 27 models, with the learning model indicated above (or lack of in the case of simple heuristic strategies), and sampling strategy along the x-axis. Bars indicate predictive accuracy along with standard error, and are separate by payoff condition (colour) and environment type (darkness). Icon array (right) shows the number participants best described (out of the full 27 models) and are aggregated over payoff conditions, environment types, and sampling strategy. Table S1 provides more detail about the number of participants best described by each model.

576 Simple Heuristic Strategies

577

578 We also compare various simple heuristic strategies that make predictions about search behaviour without
579 learning about the distribution of rewards.

580 Local Search

581 Local search predicts that search decisions have a tendency to stay local to the previous choice. We use
582 inverse Manhattan distance (IMD) to quantify locality:

$$583 \text{IMD}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^n |x_i - x'_i|, \quad (S4)$$

584

585 where \mathbf{x} and \mathbf{x}' are vectors in \mathbb{R}^n . For the special case where $\mathbf{x} = \mathbf{x}'$, we set $\text{IMD}(\mathbf{x}, \mathbf{x}') = 1$.

586 Win-Stay Lose-Sample

587 We also consider a form of a win-stay lose-sample (WLSL) heuristic⁵⁰, where a *win* is defined as finding a
588 payoff with a higher or equal value than the previous best. When the decision-maker “wins”, we assume

589 that any tile with a Manhattan distance ≤ 1 is chosen (i.e., a repeat or any of the four cardinal neighbours)
 590 with equal probability. *Lossing* is defined as the failure to improve, and results in sampling any unrevealed
 591 tile with equal probability.

592 Localization of Models

593

594 With the exception of the *Local Search* model, all other models include a localized variant, which
 595 introduced a locality bias by weighting the predicted value of each option $q(\mathbf{x})$ by the inverse Manhattan
 596 distance (IMD) to the previously revealed tile. This is equivalent to a multiplicative combination with the
 597 Local Search model, without the introduction of any additional free parameters. Localized models are
 598 indicated with an asterisk (e.g., Function Learning*). See *Locality of sampling behaviour* for a behavioural
 599 analysis of locality.

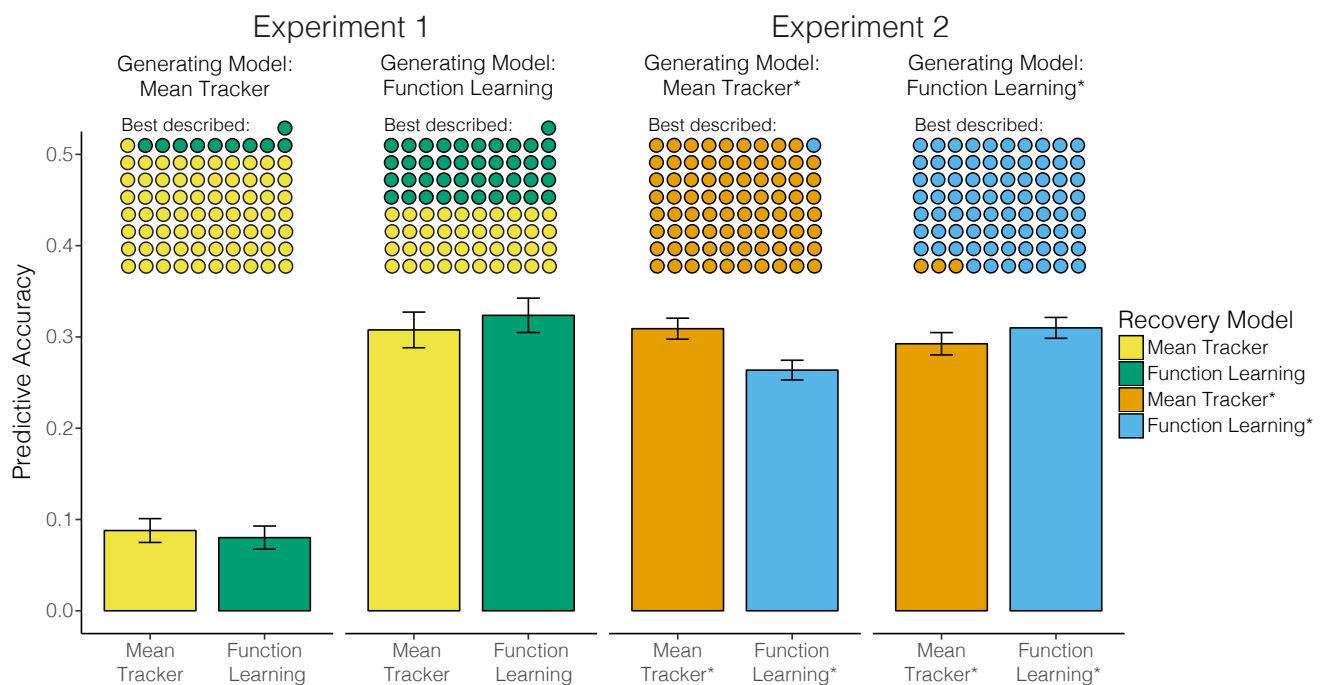


Figure S2. Model recovery results, where data was generated by the specified generating model using individual participant parameter estimates. The recovery process used the same cross-validation method used in the model comparison. We report the predictive accuracy of each candidate recovery model (bars including standard error) and the number of simulated participants best described (icon array). For both generating and recovery models, we used UCB sampling. Table S1 reports the median values of the cross-validated parameter estimates used to specify each generating model.

600 Model recovery

601 We present model recovery results that assess whether or not our predictive model comparison procedure
 602 allows us to correctly identify the true underlying model. To assess this, we generated data based on
 603 each individual participant's parameter estimates. More specifically, for each participant and round,
 604 we use the cross-validated parameter estimates to specify a given model, and then generate new data
 605 resembling participant data. We generate data using the Mean Tracker and the Function Learning model
 606 for Experiment 1 and the Mean Tracker* model and the Function Learning* model for Experiment 2. In all

607 cases, we use the UCB sampling strategy in conjunction with the specified learning model. We then utilize
608 the same cross-validation method as before in order to determine if we can successfully identify which
609 model has generated the underlying data. Figure S2 shows the cross-validated predictive performance
610 (bars) for the simulated data, along with the number of simulated participants best described (icon array).

611 **Experiment 1**

612 In the simulation for Experiment 1, our predictive model comparison procedure shows that the Mean
613 Tracker model is a better predictor for data generated from the same underlying model, whereas the
614 Function Learning model is only marginally better at predicting data generated from the same underlying
615 model. This suggests that our main model comparison results are robust to Type I errors, and provides
616 evidence that the better predictive accuracy of the Function Learning model on participant data is unlikely
617 due to overfitting.

618 When the Mean Tracker model generates data using participant parameter estimates, the same Mean
619 Tracker model achieves an average predictive accuracy of $R^2 = .1$ and describes 71 out of 81 simulated
620 participants best. On the same generated data, the Function Learning model achieves an average predictive
621 accuracy of $R^2 = .08$ and only describes 10 out of 81 simulated participants best.

622 When the Function Learning model has generated the underlying data, the same Function Learning
623 model achieves a predictive accuracy of $R^2 = .4$ and describes 41 out of 81 simulated participants best,
624 whereas the Mean Tracker model achieves a predictive accuracy of $R^2 = .39$ and describes 40 participants
625 best. This makes our finding of the Function Learning as the best predictive model even stronger as
626 –technically– the Mean Tracker model could mimic parts of the Function Learning behaviour.

627 **Experiment 2**

628 In the simulations for Experiment 2, we used the localized version of each type of learning model for
629 both generation and recovery, since in both cases, localization improved predictive accuracy of human
630 participants (Table S1). Here, we find very clear recoverability in all cases, with the recovering model best
631 predicting the vast majority of simulated participants when it is also the generating model (Fig. S2).

632 When the Mean Tracker* model generated the data, the Mean Tracker* model achieves a predictive
633 accuracy of $R^2 = .32$ and predicts 79 out of 80 simulated participants best, whereas the Function Learning*
634 model predicts only a lone simulated participant better, with an average predictive accuracy of $R^2 = .26$.

635 If the Function Learning* model generated the underlying data, the same Function Learning* model
636 achieves a predictive accuracy of $R^2 = .34$ and describes 77 out of 80 simulated participants best, whereas
637 the Mean Tracker* model only describes 3 out of 80 simulated participants better, with a average predictive
638 accuracy of $R^2 = .32$.

639 In all of the these simulations, the model that has generated the underlying data is also the best
640 performing model, as assessed by its predictive accuracy and the number of simulated participants
641 predicted best. Thus, we can confidently say that our cross-validation procedure distinguishes between the
642 two assessed model classes. Moreover, in the cases where the Function Learning or Function Learning*
643 model has generated the underlying data, the predictive accuracy of the same model is not perfect (i.e.,
644 $R^2 = 1$), but rather close to the predictive accuracies we found for participant data (Table S1).

645 **Parameter Recovery**

646 Another important question is whether or not the reported parameter estimates of the two Function
647 Learning models are reliable and recoverable. We address this question by assessing the recoverability of
648 the three parameters of the Function Learning model, the length-scale λ , the exploration factor β , and
649 the temperature parameter τ of the softmax choice rule. We use the results from the model recovery
650 simulation described above, and correlate the empirically estimated parameters used to generate data (i.e.,

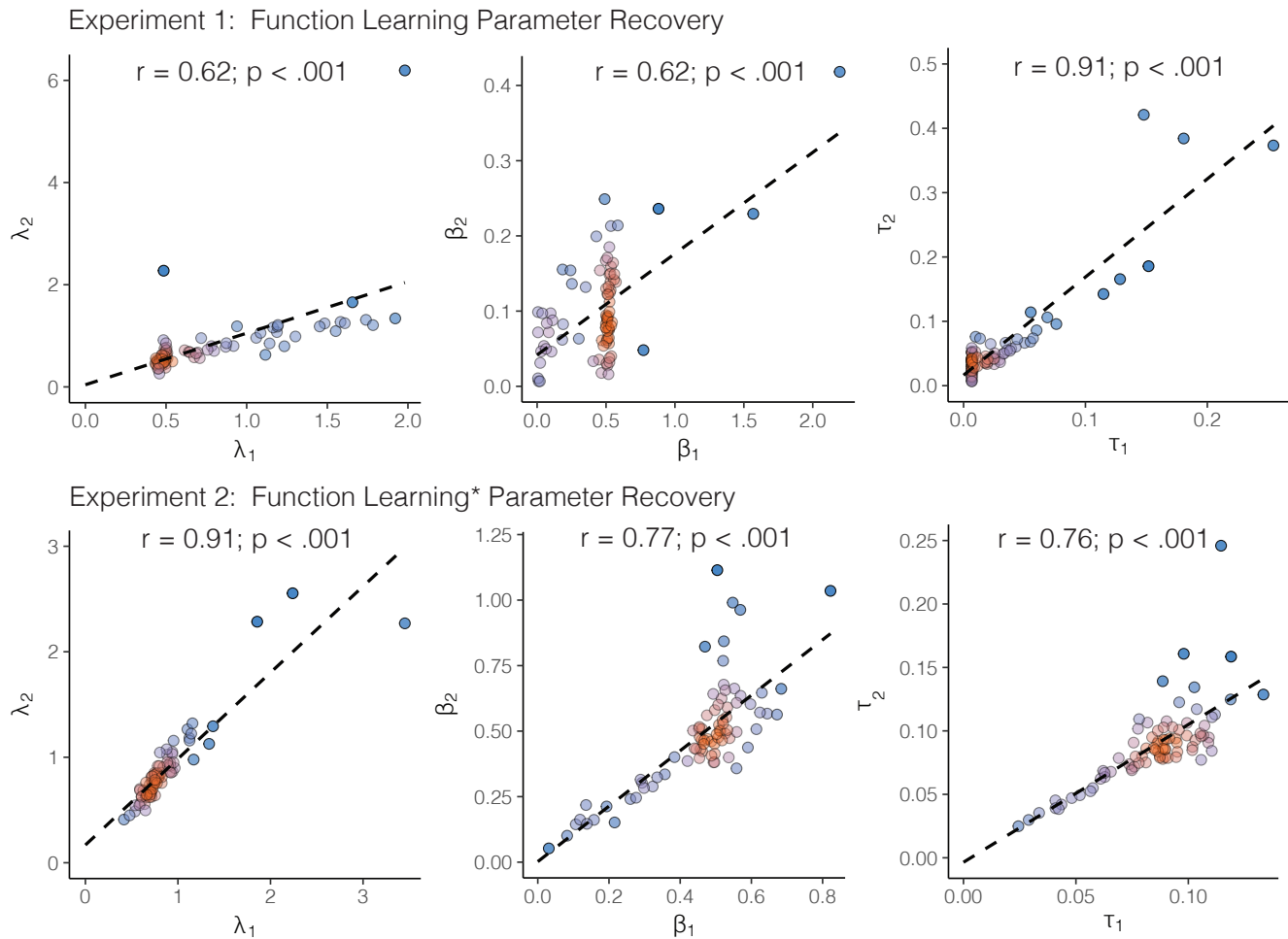


Figure S3. Parameter recovery results. The generating parameter estimate is on the x-axis and the recovered parameter estimate is on the y-axis. The generating parameter estimates are from the cross-validated participant parameter estimates, which were used to simulate data (see Model Recovery). Recovered parameter estimates are the result of the cross-validated model comparison (see Model Comparison) on the simulated data. While the cross-validation procedure yielded k -estimates per participant, one for each round ($k_{exp1} = 16; k_{exp2} = 8$), we show the median estimate per (simulated) participant. The dashed line shows a linear regression on the data, while the Pearson correlation and p-value is shown above the plot. For readability, colours represent the bivariate kernel density estimate, with red indicating higher density.

651 the estimates based on participants' data), with the parameter estimates of the recovering model (i.e.,
 652 the MLE from the cross-validation procedure on the simulated data). We assess whether the recovered
 653 parameter estimates are similar to the parameters that were used to generated the underlying data. We
 654 present parameter recovery results for the Function Learning model for Experiment 1 and the Function
 655 Learning* model for Experiment 2, both using the UCB sampling strategy. We report the results in Figure
 656 S3, with the generating parameter estimate on the x-axis and the recovered parameter estimate on the
 657 y-axis.

658 For Experiment 1, the correlation between the generating and the recovered length-scale λ is $r = .62$,
 659 $p < .001$, the correlation between the generating and the recovered exploration factor β is $r = 0.62$,
 660 $p < .001$, and the correlation between the generating and the recovered softmax temperature parameter

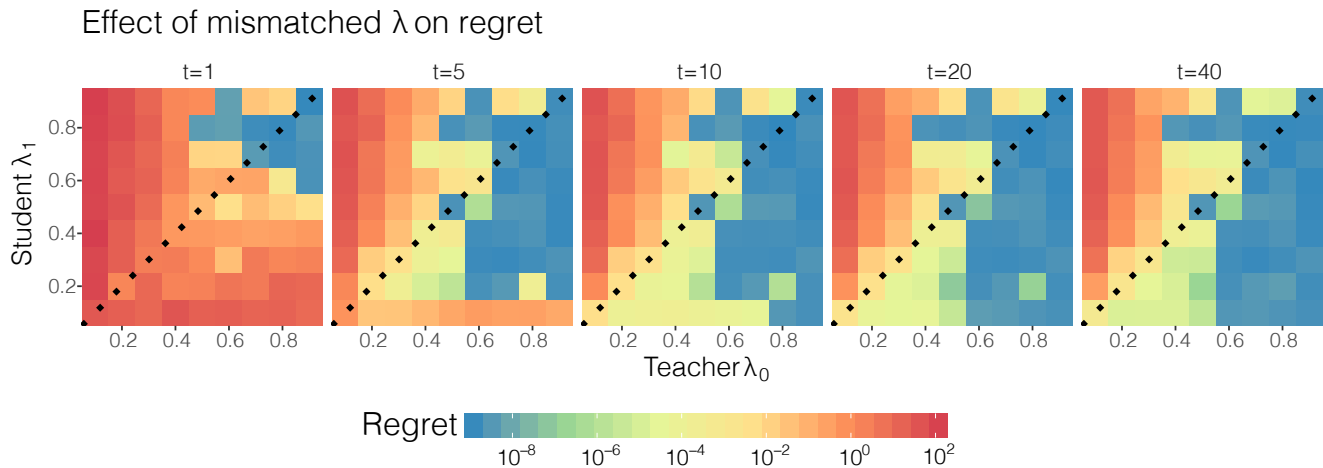


Figure S4. Mismatched length-scale simulation results. The generating teacher length-scale parameter is on the x-axis and the student length-scale parameter is on the y-axis. The teacher length-scale values were used to generating environments, while the student length-scale parameters were used to parameterize the \mathcal{GP} -UCB Function Learning Model to simulate human search performance. Each tile of the heat-map indicates the median regret of that particular λ_0 - λ_1 -combination, aggregated over 100 replications at trial numbers $t = \{1, 5, 10, 20, 40\}$. The dotted lines show where $\lambda_0 = \lambda_1$ and mark the difference between undergeneralization and overgeneralization, where points below the line show the regret produced by undergeneralizing student kernels. For readability, colours represent the log-regret of different λ_0 - λ_1 -combinations, with red indicating higher regret.

661 τ is $r = 0.91$, $p < .001$. For Experiment 2, the correlation between the generating and the recovered λ
 662 is $r = 0.91$, $p < .001$, for β the correlation is $r = 0.77$, $p < .001$, and for τ the correlation is $r = 0.76$,
 663 $p < .001$.

664 These results show that the correlation between the generating and the recovered parameters is high
 665 for both experiments and for all parameters. Thus, we have strong evidence to support the claim that
 666 the reported parameter estimates of the Function Learning model (Table S1) are recoverable, reliable,
 667 and therefore interpretable. Importantly, we find that estimates for β (exploration bonus) and τ (softmax
 668 temperature) are indeed recoverable, providing evidence for the existence of a *directed* exploration bonus¹⁰,
 669 as a separate phenomena from random, undirected exploration⁴⁴ in our behavioural data.

670 Mismatched generalization

671 We assess the effect of mismatched λ -estimates on the performance of the \mathcal{GP} -UCB Function Learning
 672 Model. A mismatch is defined as estimating a different level of spatial correlations (captured by the
 673 per participant λ -estimates) than the ground truth in the environment ($\lambda_{Smooth} = 2$, and $\lambda_{Rough} = 1$ for
 674 both experiments). In both experiments, we found that participant λ -estimates were systematically
 675 lower than the true value (Fig. 4), which can be interpreted as a tendency to undergeneralize about the
 676 spatial correlation of rewards in the world. In order to test how this tendency to undergeneralize (i.e.,
 677 underestimate λ) influences task performance, we present simulation results (Fig. S4) using different λ
 678 values in a teacher kernel (x-axis) and a student kernel (y-axis).

679 Both teacher and student kernels were RBF kernels, where the teacher kernel was parameterized with
 680 a length-scale λ_0 and the student kernel with a length-scale λ_1 . For situations in which $\lambda_0 \neq \lambda_1$, the
 681 smoothness assumptions can be seen as misaligned. The student *overgeneralizes* when $\lambda_1 > \lambda_0$ (Fig. S4

682 above the dotted line), and *undergeneralizes* when $\lambda_1 > \lambda_0$ (Fig. S4 below the dotted line), as was captured
683 by in our behavioural data.

684 We simulate every possible combination between $\lambda_0 = \{0.1, 0.2, \dots, 1\}$ and $\lambda_1 = \{0.1, 0.2, \dots, 1\}$,
685 leading to 100 different combinations of student-teacher scenarios. For each of these combinations, we
686 sample a bivariate target function from a \mathcal{GP} parameterized by λ_0 and then use the \mathcal{GP} -UCB Function
687 Learning Model parameterized by λ_1 to search for rewards. The exploration parameter β was set to 0.5 to
688 resemble participant behaviour (Table S1).

689 Figure S4 shows the median regret for 100 replications for all 100 λ_0 - λ_1 -combination at trial $t =$
690 $\{1, 5, 10, 20, 40\}$. Regret is defined as the difference between the reward obtained at trial t by sampling
691 x_t , and the best possible reward that could have been obtained by sampling the global optimum x_* , if the
692 reward distributions of all options was fully known (i.e., with perfect knowledge).

$$693 \quad R_t = f(x_*) - f(x_t) \quad (S5)$$

694

695 The simulations revealed several interesting results. First of all, regret is generally lower (blue values)
696 when the student undergeneralizes (below the dotted line) than when the student overgeneralizes (above
697 the dotted line). This effect is more pronounced over time, whereby a mismatch in the direction of
698 undergeneralization recovers over time (less regret for larger values of t). This is not the case for a
699 mismatch in the direction of overgeneralization, which continues to produce high regret, even at $t = 40$.
700 Thus, undergeneralization leads to better performance than overgeneralization.

701 Estimating the best possible alignment between λ_0 and λ_1 to produce the lowest regret revealed
702 that underestimating λ_0 by an average of about 0.21 produces the best regret over all scenarios. These
703 simulation results provide strong evidence that the systematically lower estimates of λ captured by our
704 model comparison procedure do not necessarily suggest a flaw or bias in human behaviour—but instead—
705 can sometimes lead to better performance. Undergeneralization, as it turns out, might not be a bug but
706 rather a feature of human behaviour.

707 Further Behavioural Analysis

708 Locality of sampling behaviour

709 Figure S5 shows the locality of participants' sampling behaviour compared to a random baseline. Locality
710 is assessed by the euclidean distance between two consecutively sampled points. This distance is compared
711 between participants' sampling behaviour and a fully random sampler. Whereas participants sample more
712 locally than a random sampler in both the univariate Experiment 1 ($t(160) = 31.2, p < .0001$) and the
713 bivariate Experiment 2 ($t(158) = 42.7, p < .0001$), this locality effect is much stronger for the latter
714 ($d = 4.47$ vs. $d = 1.92$), in line with our finding that localized models generate better predictions in
715 Experiment 2.

716 Learning over trials and rounds

717 Next, we assess whether participants are more strongly improving over trials or over rounds (Fig. S6). If
718 they are improving over trials, this means that they are indeed finding better and better options, whereas if
719 they are improving over rounds, this would also suggest some kind of meta-learning as they would get
720 better at the task the more rounds they have performed previously. To test this, we fit a linear regression to
721 every participant's outcome individually, either only with trials or only with rounds as the independent
722 variable. Afterwards, we extract the mean standardized slopes for each participant including their standard

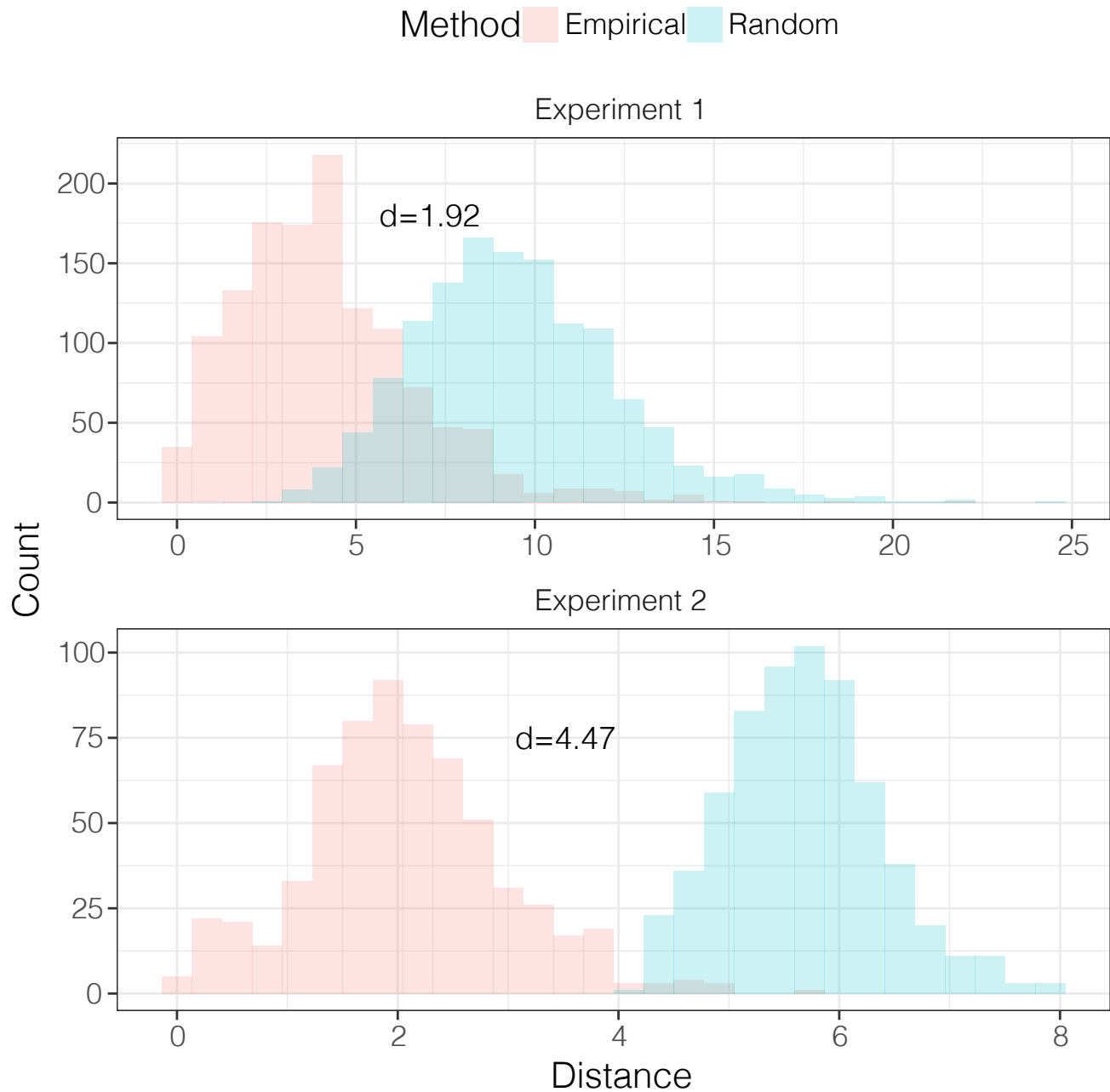


Figure S5. Locality of empirically observed and random sampling behaviour. Whereas participants consistently sample more locally than what would be expected by chance, this effect is much stronger for the bivariate Experiment 2 ($d=4.47$) than for the univariate Experiment 1 ($d=1.92$).

⁷²³ errors[‡]. Results (from one-sample t-tests with $\mu_0 = 0$) show that participants' scores improve significantly
⁷²⁴ over trials for both Experiment 1 ($t(80) = 5.57$, $p < .0001$, $d = 0.62$) and Experiment 2 ($t(79) = 2.78$,
⁷²⁵ $p < .001$, $d = 0.31$). There were no effects for rounds in either Experiment 1 ($t(80) = -2.7$, $p > .9$,
⁷²⁶ $d = -0.3$) or Experiment 2 ($t(79) = 0.21$, $p > 0.4$, $d = 0.02$).

[‡]Notice that these estimates are based on a linear regression, whereas learning curves are probably non-linear. Thus, this method might underestimate the true underlying effect of learning over time

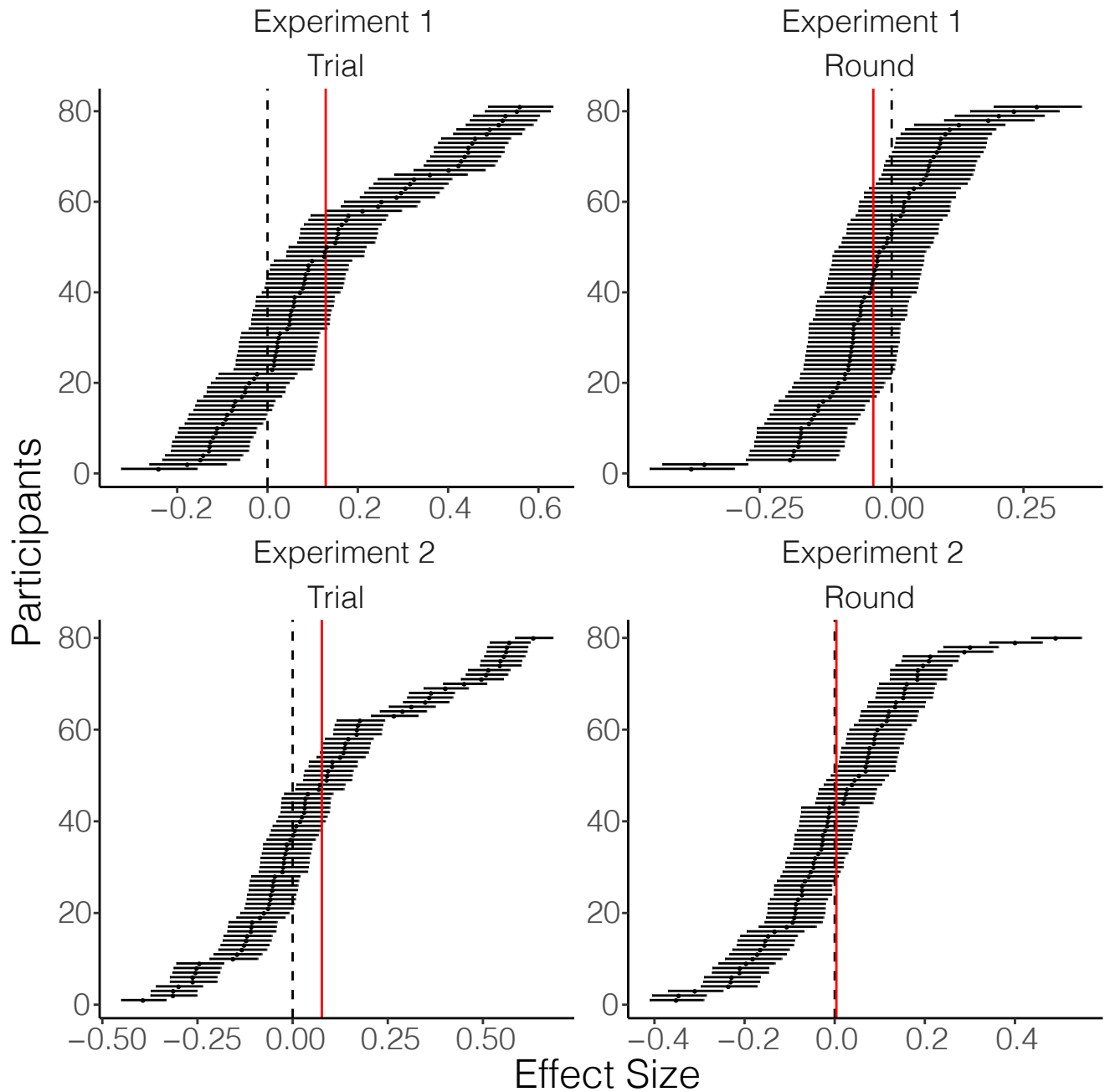


Figure S6. Average correlational effect size of trial and round on score per participant as assessed by a standardized linear regression. Participants are ordered by effect size in decreasing order. Dashed lines indicate no effect. Red lines indicate average effect size. Whereas participants consistently improve over trials, there is no effect over rounds.

Table S1. Modelling Results

Model	Experiment 1						Experiment 2					
	Model Comparison		Parameter Estimates				Model Comparison		Parameter Estimates			
	Predictive Accuracy	Participants Best Described	Length Scale (λ)	Exploration Bonus (β)	Error Variance ($\sqrt{\theta_2^2}$)	Softmax Temperature (τ)	Predictive Accuracy	Participants Best Described	Length Scale (λ)	Exploration Bonus (β)	Error Variance ($\sqrt{\theta_2^2}$)	Softmax Temperature (τ)
Mean Tracker												
Upper Confidence Bound	0.09	0	–	3.51	0.94	0.03	0.1	0	–	0.97	1.96	0.02
Pure Exploitation	0.07	1	–	–	54.6	54.6	0.1	0	–	–	148.41	148.41
Pure Exploration	0.02	0	–	–	0.32	0.02	0.01	0	–	–	15.9	0.03
Expected Improvement	0.02	0	–	–	0.37	0.01	0.01	0	–	–	1.56	0.02
Probability of Improvement	0.09	0	–	–	0.01	0.15	0.1	0	–	–	0.01	0.11
Probability of Maximum Utility	0.00	0	–	–	0.69	0.69	0	0	–	–	0.54	0.01
Mean Tracker*												
Upper Confidence Bound	0.21	1	–	44.7	0.01	28.07	0.36	12	–	44.08	0.07	15.79
Pure Exploitation	0.07	1	–	–	54.6	0.01	0.1	0	–	–	148.41	148.41
Pure Exploration	0.18	0	–	–	0.01	0.71	0.33	3	–	–	0.58	0.43
Expected Improvement	0.16	0	–	–	0.01	0.27	0.32	0	–	–	0.63	0.14
Probability of Improvement	0.14	0	–	–	0.01	0.19	0.32	0	–	–	0.01	0.09
Probability of Maximum Utility	0.12	0	–	–	0.67	0.46	0.13	0	–	–	0.36	0.01
Function Learning												
Upper Confidence Bound	0.29	48	0.5	0.51	–	0.01	0.24	4	0.54	0.47	–	0.02
Pure Exploitation	0.16	6	1.94	–	–	0.15	0.16	0	1.55	–	–	0.11
Pure Exploration	0.02	0	0.11	–	–	0.03	0.01	0	0.17	–	–	0.55
Expected Improvement	0.15	9	0.56	–	–	0.01	0.23	0	0.67	–	–	0.05
Probability of Improvement	0.05	0	3.43	–	–	0.18	0.02	0	0.87	–	–	0.09
Probability of Maximum Utility	0.00	0	0.69	–	–	7.17	0.02	0	0.49	–	–	0.01
Function Learning*												
Upper Confidence Bound	0.23	10	0.96	0.54	–	0.16	0.38	60	0.76	0.49	–	0.09
Pure Exploitation	0.16	1	7.13	–	–	0.12	0.23	0	14.4	–	–	0.06
Pure Exploration	0.14	3	0.08	–	–	0.32	0.27	0	0.17	–	–	.19
Expected Improvement	0.09	1	0.71	–	–	0.11	0.23	1	0.67	–	–	0.05
Probability of Improvement	0.12	0	7.14	–	–	0.2	0.24	0	0.84	–	–	0.09
Probability of Maximum Utility	0.12	0	0.67	–	–	0.46	0.12	0	0.46	–	–	0.01
Win-Stay Lose-Sample												
Win-Stay Lose-Sample	0.00	0	–	–	–	3.72	0.05	0	–	–	–	0.32
Win-Stay Lose-Sample*	0.05	0	–	–	–	0.73	0.26	0	–	–	–	0.22
Local Search	0.12	0	–	–	–	0.46	0.28	0	–	–	–	0.22

Note: Parameter estimates are the mean over all participants. There were 81 participants in Experiment 1 and 80 participants in Experiment 2. We have highlighted the best performing model for each experiment in boldface.