

Humans flexibly integrate social information despite interindividual differences in reward

Alexandra Witt^{a,1}, Wataru Toyokawa^{b,c}, Kevin N. Lala^{d,+}, Wolfgang Gaissmaier^b, and Charley M. Wu^a

This manuscript was compiled on September 12, 2024

There has been much progress in understanding human social learning, including recent studies integrating social information into the reinforcement learning framework. Yet previous studies often assume identical payoffs between observer and demonstrator, overlooking the diversity of social information in real-world interactions. We address this gap by introducing a socially correlated bandit task that accommodates payoff differences among participants, allowing for the study of social learning under more realistic conditions. Our novel Social Generalization (SG) model, tested through evolutionary simulations and two online experiments, outperforms existing models by incorporating social information into the generalization process, but treating it as noisier than individual observations. Our findings suggest that human social learning is more flexible than previously believed, with the SG model indicating a potential resource-rational trade-off where social learning partially replaces individual exploration. This research highlights the flexibility of humans social learning, allowing us to integrate social information from others with different preferences, skills, or goals.

social learning | generalization | exploration | computational modelling | evolutionary simulations |

Imagine you are in a foreign city, trying to decide on a restaurant to visit for dinner. You check reviews within a certain radius. Do you go for the best-rated restaurant no matter what, trusting the majority judgement? Or do you assume your taste may differ from everyone else's in this city, and discount ratings based on your personal preferences, integrating what is popular with what you know about your own tastes? It may seem obvious that you would not generally assume that everyone you could possibly rely on for their opinion will share your exact tastes. However, much of the literature in social learning has focused on this idea of how we use information from others who are just like us (1–6).

Research on the use of social information has identified various *social learning strategies* (SLS) commonly deployed by humans (7–11). These SLS serve to selectively limit imitation to cases in which it would be beneficial to imitate others, and can be categorized into *when-*, *what-*, and *who-*strategies. When-strategies determine when social learning should be used, e.g. when an agent is uncertain (12–14), or when individual learning is costly (14, 15). What-strategies specify what is preferentially learnt from others, e.g. emotionally evocative content (16, 17), and information relevant for survival (18–20), or about social relationships (18, 21). Who-strategies determine who should be learnt from, e.g. prestigious (22, 23) or successful (24–26) individuals, or the majority (1, 13, 14, 27). Prior research has also sought to understand *how* people imitate, that is, what mechanism underlies their use of social information, e.g. stimulus enhancement (28), decision biasing (1, 3), or value shaping (2).

However, even when selectively limiting when and how imitation should be used for social learning, individuals may need to share the same goals or preferences as whoever they are imitating for imitation to yield favourable outcomes. In previous research, it has commonly been the case that demonstrators had the same payoff function as the participant (1, 2, 4–6, 29). Only few studies have considered social information use in matters of taste (30–32), and they have largely focused on the normative question of how best to craft social recommendations. Therefore, our understanding of the human ability to learn from others has been limited to settings in which imitation is optimal.

In real life, however, people can rarely assume that any stranger they may choose to imitate will share their exact goals. For instance, if the goal is to get home after work, following the first car in view is unlikely to lead to the desired outcome. Conversely, if an individual notices a usually bustling street deserted during rush hour, they may correctly choose to dodge the roadwork that has

Significance Statement

Social learning is one of humankind's most remarkable cognitive abilities, and underlies much of our success as a species. While prior research has uncovered much about the way people learn from others, it has usually investigated contexts in which demonstrator and learner shared the same goal. In our study, we investigate how humans can still use social information when learning from others who have similar, but not identical, goals. We find that they do so by treating social information as less reliable individual information. They also use social learning to partially replace costly individual exploration. Our research highlights the flexibility of human social learning, allowing us to integrate social information from others with different preferences, skills, or goals.

Author affiliations: ^aHuman and Machine Cognition Lab, University of Tübingen, Tübingen, 72074 Germany; ^bSocial Psychology and Decision Sciences, University of Konstanz, Konstanz, 78464 Germany; ^cRIKEN Center for Brain Science, RIKEN, Wako, 351-0198 Japan; ^dSchool of Biology, University of St Andrews, St Andrews, KY16 9AJ United Kingdom; ⁺formerly Laland

CMW, WT, and AW conceived the experiments. WT and AW conducted the experiments. AW and CMW analysed the results and wrote the manuscript. All authors reviewed the manuscript.

The authors declare no competing interests.

¹To whom correspondence should be addressed. E-mail: alexandra.witt@gmx.net

caused everyone's paths to change while still getting to the right house after a quick detour. This difference in exact goals (e.g. the destination of a trip) despite some shared preferences (e.g. avoiding traffic or closed roads) is commonly seen in many choice domains, like food selection, fashion, career choices, holiday planning, or scheduling, in which we can commonly learn from others. Thus, there must be more to social learning than just imitation: some consideration must be made of whether the interests of the imitating and imitated individuals are aligned.

The question of how humans learn socially from demonstrators with differing preferences is sometimes answered with Theory of Mind inference (33–35), i.e. the ability to infer others' mental states, like their goals or preferences, from their behaviour. Research in this domain has uncovered much about people's ability to infer mental-state information from others' behaviours (36–38), and specifically people's ability to infer others' preferences (39, 40). After such inference, people might indeed be able to determine whether they share another person's reward function, and thus whether exact imitation is a promising option. However, even if reward functions are not perfectly aligned, people may still be able to glean valuable insights that enhance their individual decision-making. Moreover, people in the modern world often make choices using social information that is merely an aggregate rating of others' opinions, with no way of inferring how similar each individual may be to them. Thus, there is an open question of how we can use inferred or otherwise-gained value information from others who do not share our exact preferences, which is what the current study aims to address.

Goals and scope. To this end, we introduce the *socially correlated bandit task*, which lets us investigate learning and exploration dynamics in social settings where exact imitation is not optimal. The task is based on the spatially correlated bandit (41), which uses spatially correlated rewards to allow for individual generalization, and is typically used to investigate asocial learning. We add social correlations to this setup, enabling the generalization of not only individual, but also social information (Fig. 1a-b).

In our socially correlated bandit, participants search individualized environments, which are correlated with one another. Thus, the highest rewards are generally in the same region for all participants, but directly copying another participant's best choice will not lead to the maximum payoff for oneself. This emulates the relationship of social information and diverse individual preferences and circumstances in the real world: while there are some standards that apply to everyone, not everyone would agree on the same option being optimal. While the spatially correlated multi-armed bandit has previously been used to investigate social learning, it was either in individual settings (42) or with both participants in the same environment (43), not with correlated rewards across participants.

Participants explored these socially correlated environments in groups of four. In group rounds, they had full information about other participant's choices and outcomes (Fig. 1c), thus sidestepping the actual social inference. We ran evolutionary simulations with multiple candidate models to find the normatively best strategy, which was our novel "Social Generalization" (henceforth "SG") model. We then fit these models to the behavioural data collected in three

online experiments. In Exp. 1, which consisted only of group rounds, we studied whether humans would be able to utilize social information in this novel setting, and if so, how they integrated it into their decision-making. We found that participants were able to use social information to their benefit, with search behaviour being significantly influenced by other participants finding high rewards. Their behaviour was most accurately predicted by SG. In Exp. 1R, we lowered the social correlations to a minimum, but kept all other aspects the same to ensure that these results are not artifacts induced by the instructions biasing experimental subjects to use social information inflexibly. We found that participants do not blindly use social information, being best fit by AS. Finally, we ran Exp. 2 as a preregistered replication (44) interleaving solo rounds and group rounds (Fig. 1d). This allowed us to disentangle behavioural signatures stemming from the correlated task structure from actual social learning. It also let us delve deeper into differences between individual and social learning in the task by comparing baseline learning model's parameters between conditions. Again, we find adaptive use of social information, with SG being the best fit model. Differences in exploration behaviour indicate that social learning may function as an exploration mechanism when available (9). Taken together, we find that humans can integrate social information with more nuance than what previous task designs implied, potentially using it to partially replace individual exploration.

Results

We use the socially correlated bandit (Fig. 1a) for this study. Each agent explores a multi-armed bandit arranged as a grid with spatial correlations (41), and can observe the other agents of their group doing the same. We generated sets of four positively correlated bandits (for details, see Methods), so that social information can be valuable, but is less so when used verbatim (Fig. 1b). In the experiments, this was framed as collecting salt samples in alien oceans as a team of scientists, with each scientist being interested in a different salt.

In the following, we first introduce four candidate models that differ in how they integrate social information into the reinforcement learning process (Fig. 2a, top panel). We then use these models in evolutionary simulations to find the best normative strategy. Finally, we report results from three online experiments. In Exp. 1, we investigated whether and how humans would be able to use social information in this new setup to enhance their decision-making, with social correlations set to $r = .6 \pm .05$. We repeated this experiment with social correlations of $r = .1 \pm .05$ in Exp. 1R. Finally, we expand on these results in Exp. 2 as a preregistered replication, where we interleave solo and group rounds to investigate how social learning influences individual exploration patterns.

Models. We first introduce an asocial baseline model (Asocial Learner; AS), followed by our candidate social models. We consider three social models, all of which build on the asocial baseline model. Each social model integrates social information into a different stage of the individual decision-making process: the policy (Decision Biasing; DB), value function (Value Shaping; VS), or reward generalization (Social Generalization; SG). All models are illustrated in Fig. 2a.

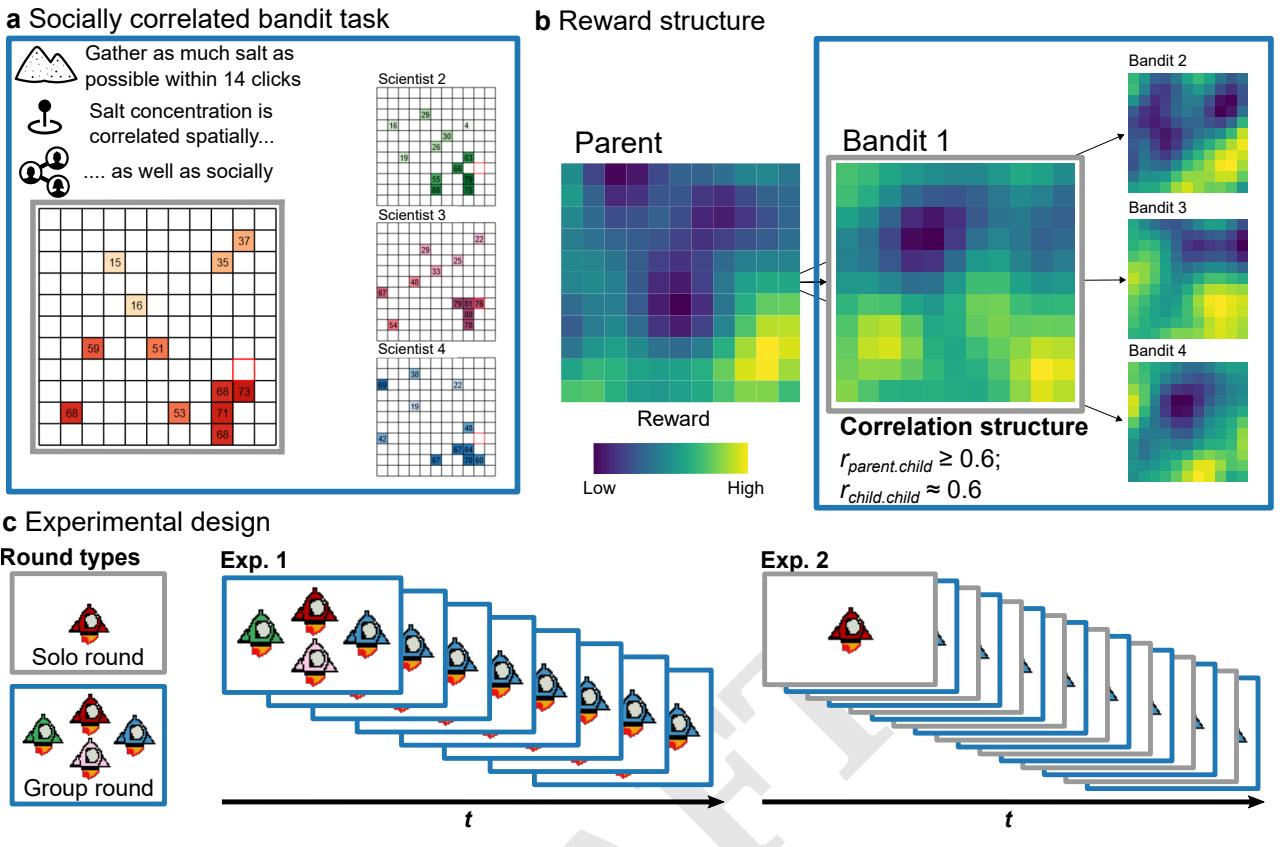


Fig. 1. Experiment overview. **a)** Screenshot of the socially correlated bandit task. Participants completed the task either individually (solo rounds, gray border) or in groups of four (group rounds, blue border). In the group condition, they had access to choice and outcome information of other group members. Participants were instructed they would collect salt samples on alien oceans with other scientists to explain the spatial and social correlation structure. For details, see Methods. **b)** Reward structure of the socially correlated bandit. Individual payoffs are generated from a common parent grid and are positively correlated. This leads to high and low payoffs being in the same general area across participants, while global optima are still distinct, limiting the effectiveness of exact imitation. **c)** Experimental design. Exp. 1 only included group rounds, while Exp. 2 had alternating group and solo rounds, in counterbalanced order.

Asocial Learner (AS). We use a Gaussian Process Upper Confidence Bound (GP-UCB) model (41) as a commonly used (42, 43, 45, 46) asocial baseline for the spatially correlated bandit problem. Gaussian Process regression is used to model expectations about the reward r associated with each action by generalizing from reward observations. For some novel option \mathbf{x}_* (i.e. a tile on the grid), and given past observations $\mathcal{D}_t = \{\mathbf{X}_t, \mathbf{y}_t\}$ of choices $\mathbf{x}_1, \dots, \mathbf{x}_t$ and rewards y_1, \dots, y_t , the posterior reward distribution is a multivariate Gaussian:

$$p(r(\mathbf{x}_* | \mathcal{D}_t) \sim \mathcal{N}(m(\mathbf{x}_* | \mathcal{D}_t), v(\mathbf{x}_* | \mathcal{D}_t)) \quad [1]$$

The posterior is thus defined by its mean $m(\mathbf{x}_* | \mathcal{D}_t)$ and variance $v(\mathbf{x}_* | \mathcal{D}_t)$:

$$\begin{aligned} m(\mathbf{x}_* | \mathcal{D}_t) &= \mathbf{K}_{*,t}^\top (\mathbf{K}_{t,t} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y}_t \\ v(\mathbf{x}_* | \mathcal{D}_t) &= \mathbf{K}_{*,*} - \mathbf{K}_{*,t}^\top (\mathbf{K}_{t,t} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{K}_{*,t}, \end{aligned} \quad [2]$$

Here, \mathbf{K} is the covariance matrix between different subsets of observations (* for new inputs and t for prior observations), σ_ϵ^2 is the observation noise, and \mathbf{I} is the identity matrix.

The assumed covariance depends on the kernel function \mathbf{k} , which determines how the model generalizes. We use a Radial Basis Function (RBF) kernel $k_{RBF}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x}-\mathbf{x}'\|^2}{2\lambda^2}\right)$. The length-scale λ determines the decay rate of the covariance

between two points as a function of distance, with higher values of λ assuming stronger spatial correlations. Thus, λ controls the range of generalization, with higher values leading to broader generalization, as a single data point affects more of the surrounding data. This follows the same principle as the generating function, presenting a reasonable solution to the individual generalization process. The GP also models the environment's observation noise σ_ϵ^2 , which allows for the model to not overfit noise.

After inferring reward, upper confidence bound (UCB) sampling is used to balance exploration and exploitation tendencies. This combines posterior mean and variance resulting in a *UCB value*.

$$UCB(\mathbf{x}) = m(\mathbf{x} | \mathcal{D}_t) + \beta \sqrt{v(\mathbf{x} | \mathcal{D}_t)} \quad [3]$$

The uncertainty-directed exploration parameter β trades off the value of an option against the uncertainty of that estimate: as it approaches 0, an agent will preferentially exploit the best known option, whereas higher β values induce more exploratory behaviour, by optimistically inflating the value of more uncertain options.

We then use a softmax to convert the value function into the *policy*:

$$\pi(\mathbf{x}) \propto \exp(UCB(\mathbf{x}) / \tau) \quad [4]$$

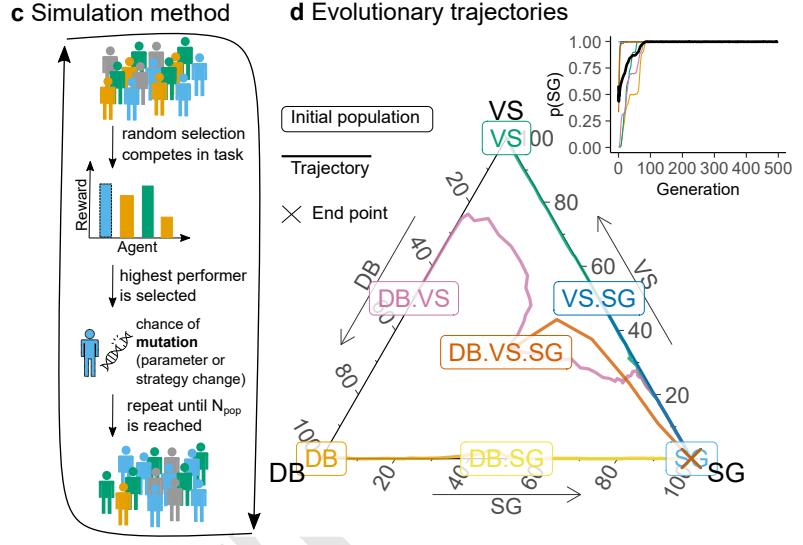
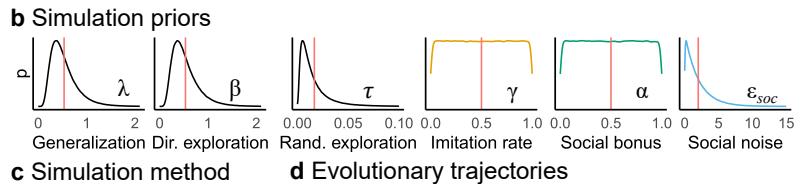
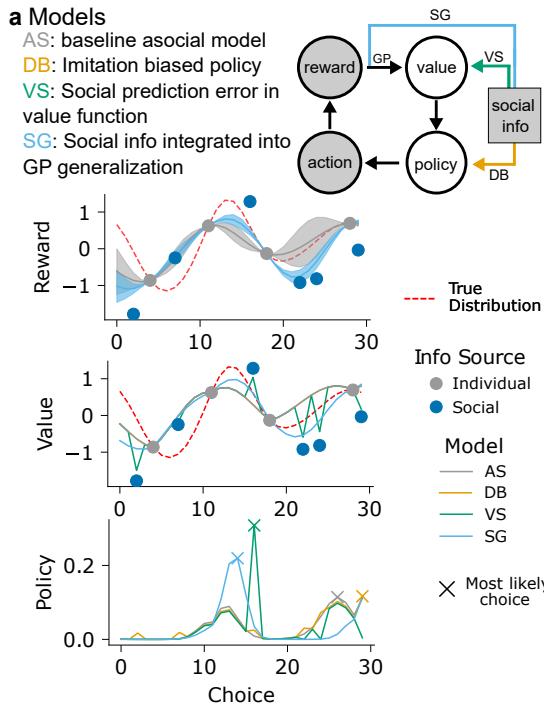


Fig. 2. Models and evolutionary simulations. **a)** Model overview. Top panel: Illustration of the individual decision-making circuit and the stages at which social information is integrated. Bottom panels: An illustrative 1D example of how models incorporate social information within the steps of the reinforcement learning circuit, where the x-axis is the discrete choice space. Reward: Only SG integrates social information into the GP posterior, whereas the other models (only AS shown for ease of reading) generalize only individual information. Value: VS integrates social information into the value function proportional to its deviation from expected value. Policy: DB integrates social information into the policy based on choice frequency. Crosses mark the most likely choice for each model. **b)** Simulation priors. Prior distribution densities used for model (including evolutionary) simulations. Red line shows the mean. **c)** Evolutionary simulation method. Agents were randomly selected to compete in one round of the task. The highest scoring agent was selected for the next generation, with a chance of parameter or type mutations. **d)** Results of the evolutionary simulations. Labels show the starting point of the various initial populations. Lines show evolutionary trajectory of model proportions, with crosses at the end point after 500 generations (all at bottom right vertex). For ease of reading, the inset plot shows only the development of SG over generations.

The temperature parameter τ controls how deterministically the model follows the value function: the higher it is, the more random the choices become. An agent's next action is chosen based on this policy.

Decision Biasing (DB) is the simplest social learning model, incorporating social information into the policy in a frequency-based manner (1, 3). This means that the choice probability for a given option is increased proportionally to how many agents have chosen that option. The policy becomes:

$$\pi = (1 - \gamma)\pi_{ind} + \gamma\pi_{soc}, \quad [5]$$

with the social policy π_{soc} tracking the other agents' choices in the previous trial such that $\pi_{soc}(x) \propto n_{x,soc,t-1}$. Here, n is the number of times an option was chosen. Individual and social policies are then combined, with the weight of social learning dependent on the mixing parameter γ .

Value Shaping (VS) incorporates social information into the value function. In previous studies, this was done by treating a social choice as a "pseudo-reward" (2). It can be seen as an implementation of either stimulus enhancement or local enhancement (28), in this case increasing the likelihood of choosing the same option one has seen chosen by the demonstrator by increasing its value. Previous implementations of this model had no reward information, as the action outcomes were not shown in their tasks. As outcomes are shown in our task, we augment the model to be value-sensitive by using a

simple prediction error approach:

$$V(x) = V_{x,ind} + \alpha(V_{x,soc} - V_{x,ind}) \quad [6]$$

with $V_{x,ind}$ being the individual UCB-value, and $V_{x,soc}$ the social value of a given option x . Thus, an observed action's value will be increased when it is better than individual expectation and decreased when it is worse. Social bonus parameter α governs the strength of this social influence. While including value information in VS improved it over a value-agnostic version (Fig. S2a), the same was not true for DB (Fig. S2b). Thus, we elected to keep using the simpler, equally good model for DB, but modified VS for better performance.

Social Generalization (SG) is a novel model that incorporates social information at the stage of the Gaussian Process regression. This means that, unlike in the other models, social information is generalized to surrounding options as well, which corresponds to a non-specific form of local enhancement (28) in a spatial context. However, social information is assumed to be noisier, and thus less reliable, than individual information. The formerly scalar noise term σ_e^2 (Eq. 2) becomes a vector, with its value depending on whether an observation was individual ($\delta_{soc}(x) = 0$) or social ($\delta_{soc}(x) = 1$).

$$\sigma_{e|x}^2 = \sigma_{e,ind}^2 + \delta_{soc}(x) \cdot \sigma_{e,soc}^2, \quad [7]$$

In addition to the term for the environment's observation noise σ_{eind}^2 , social noise $\sigma_{\epsilon_{\text{soc}}}^2$ is added to any social observations. This social noise (henceforth referred to as ϵ_{soc}) determines the reliance on social information. Higher social noise causes posterior means to deviate less from the prior mean and posterior variances to remain higher in social compared to individual observations. As social noise term ϵ_{soc} approaches 0, social information is relied on more and more, with the extreme case of $\epsilon_{\text{soc}} = 0$ treating social information as equally reliable as individual information.

Evolutionary simulations. We first used evolutionary simulations (see Methods) to determine which model achieves the best normative performance in this setting. Since social learning strategies have *frequency dependent fitness* (47), how well they perform depends on the frequency of other strategies in the population. Thus, evolutionary simulations are well-suited to evaluate frequency-based fitness without having to exhaustively evaluate every possible population composition (Fig. 2b-d). We first created different initial populations with all possible combinations of models in equal proportion, where each agent was parameterized by drawing from their model's respective prior distribution for parameters (Fig. 2b; for details see SI). We then used tournament selection (Fig. 2c) to iteratively sample groups of four agents (with replacement) to perform the task. The best performing agents in each group were selected to seed the next generation, with some chance of parameter and type mutation.

To ensure that our model implementations are in line with and extend results from the published literature, we first evaluate model performance replicating a setting where two agents are in the same environment ($r = 1$), with one making optimal choices as an expert (2). Our results replicate VS being the best model compared only against AS and DB (Fig. S4a). However, when we include SG, the results support a tie between SG and VS with no clear winner (average $p(\text{SG}) = .48$ and $p(\text{VS}) = .50$ in the final generation; Fig. S4b-d), due to the two models making the same predictions in identical reward environments with no spatial correlations.

Figure 2d shows the results of evolutionary simulations of the competing social learning models in our current task environment ($r \approx .6$). Here, due to individual differences in reward, all initial populations (even those that did not originally contain SG agents) evolve to be 100% SG agents (see Fig. S3a for starting populations including AS). This clearly suggests that SG is the normatively best model in our task.

As the parameters evolve throughout the simulations, we can also glean insight into what combinations of parameters were normatively optimal. Investigating the evolved parameters for SG (Fig. S3b), we find that λ nearly reaches the true underlying value of the environments, 2 ($\lambda = 1.96$). The random exploration parameter τ is fairly low at roughly 0.006, showing mostly deterministic choices based on the value function. The social noise parameter ϵ_{soc} shows considerable variation, but evolves to 3.2 on average. As this is higher than 0, we can see that indiscriminate social information use (like imitation) is not optimal in our task. The directed exploration parameter β evolves to lower values than what has previously been found in humans (41) with an average

of .19. This may indicate that directed exploration can be replaced by social information use in social learning settings.

Experiment 1. People flexibly use diverse social information.

Having determined SG to be the normatively best strategy for the socially correlated bandit, we now move on to online experiments using the task to see how human participants actually use social information with social correlations at $r = .6$. Exp. 1 consisted exclusively of group rounds, meaning that participants always had access to the choices and outcomes of the members of their group.

Behavioral results. Firstly, participants improved across trials (Fig. 3a), with the average performance being significantly higher than the chance level of 0.5 ($t(127) = 59.8, p < .001, d = 5.3, BF > 100$). There was a small but negligible learning effect over rounds (Fig. S1a). Average social search distance (the Euclidean distance between an option chosen at trial t and one chosen by another participant at t-1) decreased over time, indicating increased clustering of participant choices as the task progressed (Fig. 3b). In lieu of an asocial estimate, we find that social search distance was also significantly lower than what would be predicted by random choice (an average of 5.75; $t(127) = -29.4, p < .001, d = 2.6, BF > 100$), which provides further evidence for clustering.

This social clustering may have stemmed from a tendency to approach other participants who have earned a high reward in the previous trial. A regression of search distance over previous reward by information source (individual or social) shows that participants' individual search distance (the Euclidian distance from their previous choice) significantly decreased as previous individual reward increased (-6.95; Highest Density Interval [-7.24, -6.67]; Fig. 3c, black line). This is rational given the spatial correlations in the environment, and is consistent with predictions of all candidate models. However, they did not only show this tendency for individual, but also for social information: when another participant earned a high reward in the previous trial, they searched closer to this participant's position (-4.22 [-4.61, -3.88]; Fig. 3c blue line). It is worth noting that this effect of social reward on search distance is significantly lower than the effect of individual information (2.73 [2.63, 2.79]), reflecting the lower reliability of social information compared to individual information.

This value-sensitive social information use is in line with predictions made by VS and SG, which integrate social information based on their value. On the other hand, it does not match predictions made by AS or DB: AS predicts no reliance on social information at all, with social search distances at roughly chance level (5.75), while DB would predict lower social search distance regardless of previous social reward (Fig. S14c).

Further teasing apart the model predictions about social search distance, we consider the distinction of imitation (search distance = 0) vs. "innovation" (building on someone else's choice, search distance = 1; Fig. 3d). VS predicts value-sensitive imitation, that is, an increase of imitation rate as previous social reward increases, which we find in a linear model of social search distance frequency (0.04, 95%-CI: [0.03, 0.05], $p < .001$). However, this effect was even stronger for innovation (0.08, 95%-CI: [0.06, 0.10], $p < .001$), which

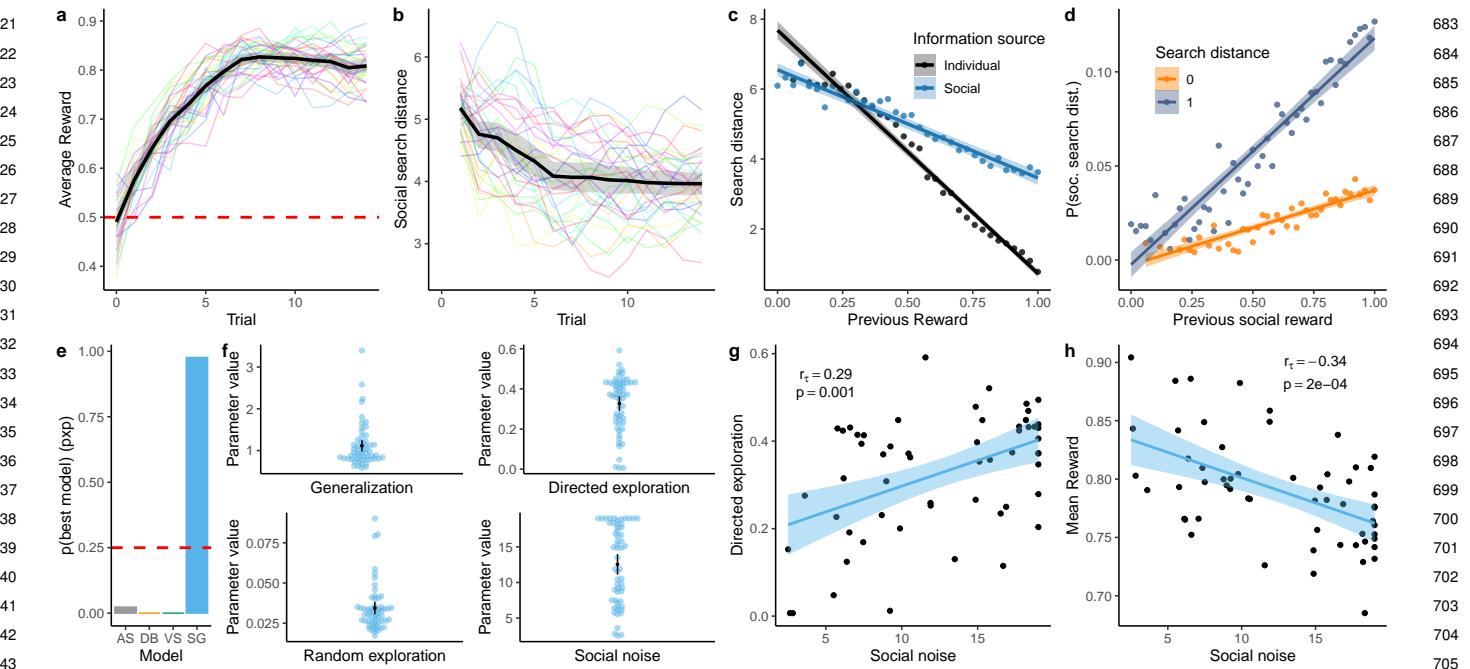


Fig. 3. Experiment 1 results. **a)** Learning curves. The average reward across participants is shown in black, with group averages as coloured lines. The red dashed line shows chance-level performance. **b)** Social search distance (average Euclidian distance from other participants) over trials. The black line is the population average, while group averages as shown as coloured lines. **c)** Search distance as a function of previous reward, split by information source. Lines are the posterior prediction of a Bayesian hierarchical regression, while points are data averaged across 20 bins. **d)** Value-dependent social search distance. A search distance of 0 is imitation, while a search distance of 1 means the participant explored an adjacent tile to a social observation. **e)** Model comparison, showing the protected exceedance probability (*p_{pxp}*), which describes the probability of a model best fitting the population, accounting for chance. Red dashed line shows chance level. **f)** Social Generalization (SG) parameters (limited to participants best fit by SG). **g)** β (directed exploration) over ϵ_{soc} . Higher values of ϵ_{soc} mean lower reliance on social information. Only participants best fit by SG are shown. **h)** Mean reward over ϵ_{soc} (social noise). Higher values of ϵ_{soc} mean lower reliance on social information. Only participants best fit by SG are shown.

only SG, the only model generalizing social information, could explain (Fig. S14d).

Modeling results. Turning to modelling, we find that SG did indeed fit human behaviour best, with hierarchical Bayesian model selection (48) showing it had the highest posterior probability of being the best model (protected exceedance probability: $p_{pxpSG} \approx .98$; Fig. 3e). In participants best described by SG (Fig. 3f), the generalization parameter was significantly lower than the ground truth of $\lambda = 2$ ($\lambda \approx 1.11$; $t(56) = -13.0$, $p < .001$, $d = 1.7$, $BF > 100$). This means that participants did not generalize their observations as broadly as would be optimal given the environment. The directed exploration parameter was significantly lower than values found for individually learning GP-UCB agents in the same task structure (41) ($\beta \approx 0.33$; $t(56) = -9.4$, $p < .001$, $d = 1.3$, $BF > 100$). The random exploration parameter $\tau \approx 0.03$. Social noise was significantly higher than the value of 3.29 found to be optimal in evolutionary simulations ($\epsilon_{soc} \approx 12.55$; $t(56) = 12.8$, $p < .001$, $d = 1.7$, $BF > 100$), meaning participants relied less on social information than optimal.

We also find a relationship between β and ϵ_{soc} : the more a participant relied on social information (lower ϵ_{soc}), the less they relied on directed exploration (lower β ; $r_\tau = .29$, $p = .001$, $BF = 28$; Fig. 3g). This might explain why β -values were lower than in previous, individual learning, settings (41). Additionally, participants best described by SG performed better when they showed higher reliance on social information ($r_\tau = -.34$, $p < .001$, $BF > 100$; Fig. 3h), where the negative

correlation reflects the fact that higher values of ϵ_{soc} mean lower reliance on social information.

In summary, Exp. 1 shows that participants could use social information to guide their decision-making even when it was not directly applicable to their own situation. Their behaviour followed the predictions of the SG model, implying that they used social information similarly to individual information, but treated it as more noisy, and thus less reliable. This method of integrating social information is optimal in our task environment (Fig. 2d), and lead to better results for participants the more they relied on social information. It is important to note that the linear relationship between social noise and reward only exists because participants relied on social information less than optimal. Indiscriminately using social information ($\epsilon_{soc} = 0$) is not beneficial in our task (Fig. S2c), so the expected relationship when the whole range of ϵ_{soc} is covered would be U-shaped with lower rewards for both higher and lower reliance on social information than optimal. In a replication of Exp. 1 with identical task structure and instructions, but social correlations set at $r = .1$, we found AS to be the dominant model ($p_{pxpAS} = 0.999$), and even higher social noise in the subset of participants best fit by SG ($\epsilon_{soc} \approx 14.29$; $t(141) = 2.8$, $p = .005$, $d = 0.6$, $BF = 7.1$; Fig. S6).

In addition, we find low uncertainty directed exploration (lower β) correlates with greater reliance on social learning (lower ϵ_{soc}). This pattern suggests social learning may partially replace uncertainty-directed exploration, which we expand on in Exp. 2.

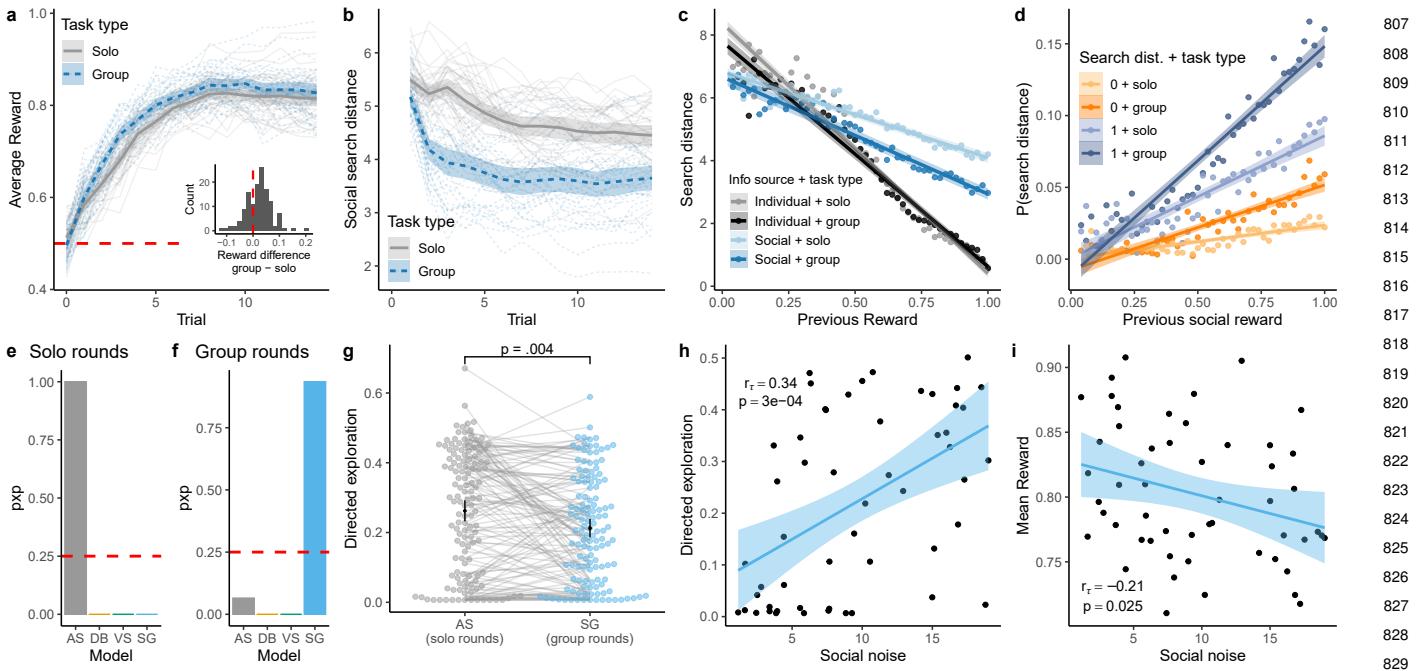


Fig. 4. Exp. 2 results. **a)** Learning curves by task type. Averages for solo (gray solid) and group (blue dashed) rounds shown in thick lines with shaded 95% CIs, while group averages are shown as thin lines. Red dashed line shows chance-level performance. Inset plot is an histogram of mean score in group - solo rounds, with the red dashed line showing no difference. **b)** Social search distance over trials by task type. Averages for solo (gray solid) and group (blue dashed) rounds thick and shaded, group level averages thin lines. **c)** Search distance over previous reward split by information source and task type. Lines are results of Bayesian hierarchical regression fixed effects, points are data averaged across 20 bins. Solo rounds in pale and group rounds in saturated colours, individual information in gray and social information in blue. **d)** Value dependent social search distance by task type. Solo rounds are solid lines, group rounds dashed. **e-f)** Protected exceedance probabilities (pxps) for solo (e) and group (f) rounds. Probability of a model best fitting the population, accounting for chance. Red dashed line shows chance level. **g)** β -parameter estimates for participants fit by AS in solo rounds (gray) and participants fit by SG in group rounds (blue). For all SG parameters, see Fig. S11. **h)** β (directed exploration) over ϵ_{soc} in group rounds. Higher values of ϵ_{soc} mean lower reliance on social information. Only participants best fit by SG shown. **i)** Mean reward over ϵ_{soc} in group rounds. Higher values of ϵ_{soc} mean lower reliance on social information. Only participants best fit by SG shown.

Experiment 2. Social learning partially replaces directed exploration. To understand the effects of social information on decision-making better, we conducted a preregistered replication of Exp. 1 with the addition of solo rounds (i.e. rounds where participants were still in correlated environments, but were not shown other participants' choices and outcomes) to provide an asocial baseline for each participant. This allows us to control for the generic effects of the correlated reward structure, and directly probe if directed exploration was actually lower in social learning settings than in individual. This was a pre-registered experiment (44), and any analyses that were not included in the preregistration are specified as exploratory.

Behavioral results. Participants improved throughout trials, with higher performance on average in group rounds compared to solo rounds ($t(131) = 6.0, p < .001, d = 0.5, BF > 100$; Fig. 4a). Again, there was a minimal learning effect over rounds, but no effect of condition order or their interaction on performance (Fig. S1b).

In following analyses, we compare social measures (like previous social reward, or social search distance) for both solo and group rounds despite no social information being provided in solo rounds. This serves as a baseline for effects that could be interpreted as social (e.g. lower search distances for high previous social rewards), which might also be explained by participants independently exploring correlated environments. Social search distance decreased over trials in both conditions

(Fig. 4b). However, it was significantly lower in group than in solo rounds ($t(131) = -14.8, p < .001, d = 1.7, BF > 100$), indicating that the clustering was not solely due to the social correlations between environments, but was influenced by social information.

Again, we investigate the effect of previous reward on search distance, splitting by information source (individual vs. social) and round type (solo vs. group). We replicate the results from Exp. 1 in group rounds: search distance was modulated by both individual (-7.75 [-8.00, -7.49]; Fig. 4c, black line) and social previous rewards (-5.44 [-5.75, -5.12]; Fig. 4c, dark blue line), with participants searching closer for higher values, and searching at greater distances for larger values. Again, social rewards influenced search distance to a lesser extent than individual rewards (2.31 [2.18, 2.44]). In solo rounds, we find the same effect for individual information (-7.97 [-8.22, -7.72]; Fig. 4c, gray line) with only a slight difference from group rounds (0.22 [0.02, 0.41]). This indicates that participants relied on previous individual reward slightly more in solo than in group rounds. Although previous social reward still significantly influenced social search distance in solo rounds, based on the correlated environmental structure only (-4.29 [-4.59, -4.00]; Fig. 4c, light blue line), it did so to a significantly lower degree than in group rounds (-1.15 [-1.33, -0.96]). This shows that, while the effect of social information on search distance can be partially explained by the socially correlated structure of the task, there is a significant component that can only be attributed to the use

of social information in how participants modulated their search. Again, this result is in line with the predictions of VS and SG, but not AS and DB, which predict either no or indiscriminate social information use, respectively.

Focusing on a finer delineation of social search distance in an exploratory analysis (Fig. 4d), we find a value-based increase in imitation frequency (0.023, 95% CI: [0.011, 0.035], $p < .001$) and even higher increase in innovation (0.062, 95% CI: [0.046, 0.078], $p < .001$) across round types. However, this increase in frequency was also significantly higher in group rounds compared to solo rounds (0.036, 95% CI: [0.019, 0.053], $p < .001$), and higher still for innovation in group rounds (0.038, 95% CI: [0.016, 0.061], $p = .001$). Thus, we replicate the value-sensitive increase in both imitation and innovation, which is only predicted by SG. The significant interaction of the effect with round type once again shows that, while the effects found in Exp. 1 can be partially explained by the correlation structure of the environments alone, they remain significant when controlling for this factor.

Modeling results. We again performed hierarchical Bayesian model comparisons, but separately for solo and group rounds. Here, we find that AS is the best fitting model for solo rounds, showing that our social models did not just exploit the correlated structure to improve fit (Fig. 4e, $p_{xp} \approx 1$). In group rounds, we again find that SG is the best fitting model (Fig. 4f; $p_{xp} \approx .94$).

In line with our finding of comparatively lower values of directed exploration parameter β in Exp. 1 than in previous individual learning literature, we find that β is indeed significantly lower in group than in solo rounds within participants (Wilcoxon signed-rank test; $Z = -2.7$, $p = .004$, $r = -.23$, $BF = 63$; Fig. 4g). In an exploratory analysis, we replicate the significant relationship between ϵ_{soc} and β from Exp. 1 ($r_\tau = .34$, $p < .001$, $BF > 100$; Fig. 4h), again suggesting a partial replacement of directed exploration with social learning when social information is available.

Regarding the relationship of social noise parameter ϵ_{soc} and average reward in group rounds in participants best fit by SG, we find a weakly significant correlation ($r_\tau = -.21$, $p = .025$, $BF = 2.1$, Fig. 4i). This might be explained by a ceiling effect of social learning that was not as strong in Exp. 1, when they had fewer rounds to familiarize themselves with the task.

In summary, Exp. 2 replicates the findings of experiment 1 in that participants use social information even when it is not directly applicable to their own situation, being best described by the SG model. This means that social information is used similarly to individual information, but treated as more noisy. The addition of an asocial baseline condition lets us compare individual and social strategies, where we confirm that the findings in Exp. 1 were actually indicative of social information use and not just a consequence of purely individual information use in a correlated environment. We replicate the finding that β -values in group rounds are lower than in previous literature. In comparison to solo rounds, we find that β -values are significantly lower, and significantly correlated with social noise, indicating that the use of social information replaces uncertainty-directed exploration to a degree. Use of social information was beneficial, shown both by the negative correlation of social noise and mean reward,

and the general higher scores in group compared to solo rounds.

Discussion

We introduce the socially correlated bandit as a task to study social learning with similar, but non-identical reward structures. Here, we find that the best normative and descriptive model, Social Generalization (SG), integrates social information into the generalization process, but in a noisy fashion. This expands our understanding of social learning mechanisms beyond settings where each individual has the exact same reward function (1, 2, 4, 43).

The socially correlated bandit has a spatially correlated reward structure within participant, which is also positively correlated across participants. We use this design to operationalize how humans may have similar but not perfectly identical preferences in matters of taste. The spatial correlations allow for social information to be integrated at a stage of the decision-making process unique to tasks in which reward information can be generalized. Thus, we introduce the SG model, which generalizes social information similarly to individual information, but treats it as more noisy (Eq. 7). We also show that the previously dominant Value Shaping (VS) model (2) can be seen as an edge case of SG when participants are in identical environments, there is no potential to generalize social information, and it is sensible to rely on social information as much as individual observations (Fig. S4).

In Exp. 1, we found that SG was the best descriptive model of human behaviour in this task, capturing behavioral patterns not predicted by other models (Fig. 3b-d) and providing the best predictions of trial-by-trial choices (Fig. 3e). We also found a relationship between the social noise parameter ϵ_{soc} and performance, suggesting participants were more successful the more they relied on social information (Fig. 3h). Additionally, we found that participants who relied more on social information displayed less uncertainty-directed exploration (Fig. 3g). This suggests a potential replacement of exploration with social learning, further corroborated by lower values of the directed exploration parameter β than found in previous works using an individual version of the task (41, 45), which motivated further investigation in Exp. 2.

In Exp. 2, we conducted a preregistered replication (44) of Exp. 1, which also added solo rounds to the task to assess how exploration behaviour changed within participants when social information is unavailable. The within-subject manipulation of solo vs. group rounds allowed us to ensure that none of the findings of Exp. 1 were solely the consequence of individual learning in correlated environments, and let us compare patterns and outcomes of individual and social learning. In this experiment, participants performed better in group than solo rounds, showing that they used available social information to their benefit (Fig. 4a). Again, we replicate the behavioral patterns corroborating SG as the best descriptive model even while accounting for the individual learning baseline (Fig. 4b-d), while model comparison shows that SG wins in group but not solo rounds (Fig. 4e-f). Consistent with our preregistered predictions, we found that β was significantly lower in group than in solo round, with β again being significantly correlated with social noise

993 ϵ_{soc} (Fig. 4h), further lending credibility to social learning
994 replacing directed individual exploration.

995 Across our experiments, we expand the scope for theories of
996 the integration of social information to cases where uncritical
997 imitation is not optimal, and find that humans can go beyond
998 imitation when the situation calls for it (10). Participants did
999 so by updating not only the social observations themselves,
1000 but also generalizing them when it was appropriate for their
1001 individual situation. Thus, the key difference between models
1002 is in *how* they use social information, with our results hinting
1003 at more sophisticated strategy use by humans than considered
1004 in previous studies.

1005 In sum, the findings of our study add to a rich literature
1006 showing that social learning is adaptive in stable environments
1007 (1, 43, 49–54), even with interindividual differences in reward.
1008 We show that this adaptive use of social information went
1009 hand in hand with a reduction of uncertainty-directed
1010 exploration, implying that social learning functioned as an
1011 exploration tool.

1012
1013 **Social learning and resource rationality.** While we find adaptive
1014 use of social information for our task setting, at the
1015 same time, we still find that participants underutilized social
1016 information compared to what would be optimal. Such
1017 underutilization of social information has also been found in
1018 other experimental settings (5, 24, 50, 55, 56). Participants'
1019 natural skill at social learning may have been limited by the
1020 artificial experimental setting. A part of this may be that
1021 some social learning strategies, like copying the expert, were
1022 unavailable due to lack of information, potentially reducing
1023 participants' inclination to rely on social information overall.

1024 However, besides social learning potentially being impeded
1025 by the artificial experimental setting, the discrepancy between
1026 adaptive social learning and underutilization of social information
1027 may also be explained by resource rationality. While
1028 it may be theoretically optimal to discount social information
1029 to a specific, low degree, this may also be significantly more
1030 complex than to rely on individual information more strongly,
1031 only referring back to presumably noisy social information
1032 when individual learning does not provide any promising
1033 options. In this regard, underutilizing social information
1034 may also be seen as resource-rational in that regard (57).
1035 Falling back on social information only when it is absolutely
1036 necessary also ties back to the social learning strategies (8, 10):
1037 Social Generalization agents generally copy (i.e. are strongly
1038 influenced by others' choices) when uncertain (i.e. their
1039 individual information does not outweigh social observations).

1040 The same resource-rationality based reasoning may be
1041 applied to our finding regarding directed exploration being
1042 partially replaced by social learning when possible. Directed
1043 exploration has been shown to be reduced by cognitive load
1044 (58, 59), indicating that individual exploration may be costly.
1045 Hence, our finding that social learning may have served as
1046 an exploration tool in our task hints at social learning as a
1047 method to let us offload these costs of directed exploration.
1048 It is also in line with prior research that suggests or shows
1049 exploration differences between asocial and social settings
1050 (1, 27, 42, 43, 51, 60, 61). It also provides empirical support
1051 for the outcome of the social learning strategies tournament,
1052 wherein winning models tended to almost exclusively use
1053 social learning for exploration (9).

1055 In our task, we show that this lower exploration is
1056 optimal for social learning using evolutionary simulations.
1057 However, simulations based on our participants' parameter
1058 estimates show that lower exploration would also lead to
1059 higher performance in asocial learning (Fig. S12a-b). This
1060 implies that social learning not only takes the function of
1061 uncertainty-directed exploration, but also helps participants
1062 avoid overexploration. This might be due to the need for
1063 exploration being diminished by the ability to gain more
1064 environmental information from others. It could also be a
1065 dynamic process wherein observing one teammate move from
1066 exploring to exploiting inspires participants to do the same,
1067 lowering overall exploration rates compared to individual
1068 settings. Such adjustments of strategy between individual and
1069 group settings, especially for individuals with low confidence,
1070 have been found before (49). However, the exact mechanism
1071 of this lowered exploration in a round-based task remains
1072 a subject for further research. In an exploratory analysis,
1073 we find a similar effect of social learning on generalization
1074 parameter λ being higher, and thus closer to the ground
1075 truth, in group rounds (Fig. S12c-d), which is in line with
1076 previous research (43).

1077 **Limitations and future directions.** Here, we used a task with
1078 spatially and socially correlated rewards, which allowed us to
1079 most easily communicate the individual and social learning
1080 structure to participants. However, this limits the scope of
1081 our current experiments to a spatial domain, whereas real-
1082 world social information may be conveyed on the basis of
1083 non-spatial features. While an investigation of SG in the
1084 non-spatial domain is beyond the scope of this work, our AS
1085 model, which forms the basis for all our social models, has
1086 been investigated in more abstract domains. In those works, it
1087 was successfully applied to non-spatial domains (e.g., abstract
1088 Gabor patch features; 62) and graph-structured environments
1089 (63). Thus, we would expect to find similar patterns for
1090 SG, which is computationally identical save for the added
1091 integration of social information. Testing this hypothesis, as
1092 well as investigating how exactly it would be parametrized,
1093 and if this would be different to generalization in spatial
1094 domains, remains a fruitful avenue for future research.

1095 Previous research often investigates the effects of demon-
1096 strator skill, contrasting one skilled and one unskilled
1097 demonstrator (2, 64). Following this reasoning, one might
1098 consider that the integration mechanism used depends on
1099 the skill level of teammates. However, value-sensitive models
1100 VS and SG benefit from any information about the structure
1101 of the environment, regardless of if it is positive or negative.
1102 Therefore, we did not investigate effects of participant skill
1103 directly. However, participants appeared more sensitive to
1104 choices of their peers that lead to high rewards (Fig. 3 and
1105 Fig. 4c), mirroring the human tendency to "copy" (here
1106 rather "learn from") the successful (24, 25). Nevertheless, it
1107 remains an open question how sensitive participants can be
1108 to others' perceived skill in this task, and in what way this
1109 would influence their decision-making.

1110 With the focus of this study being on the mechanisms of
1111 integrating social information as an individual, we left group
1112 dynamics of exploration throughout the experiment largely
1113 unexplored. While we find social clustering on a group level
1114 (Fig. 3b and Fig. 4b), we can explain this using individual-level
1115 mechanisms. In our task design, participants are incentivized
1116

1117 to maximize individual gain, which limits the benefit of active
1118 coordination. Based on participant's reported strategies,
1119 it is unlikely that they coordinated their behaviour to
1120 maximize information gain as well. However, it may be
1121 interesting to investigate coordination strategies in similar
1122 task settings, for example by changing the incentive (65) from
1123 maximizing individual reward to maximizing knowledge of
1124 the environment.

1125 Our main experiments limited the environmental correlation
1126 to $r = 0.6$. This is due to the fact that it was both harder
1127 to generate many environments with higher correlations, and
1128 the results would be less insightful, likely converging on
1129 imitation as they approach 1. However, our task using only
1130 one specific correlation of environments leads to a number of
1131 new questions: For which range of correlations humans are
1132 still sensitive to the optimal strategy? When (if ever) do they
1133 stop integrating social information altogether? Are humans
1134 able to make use of negatively correlated environments as well
1135 as positively correlated ones? While a large-scale battery of
1136 experiments across a range of social correlations is beyond the
1137 scope of this paper, we use evolutionary simulations to show
1138 that SG has the best normative performance for correlations
1139 as low as $r = .2 \pm 0.05$ (Fig. S5). Indeed, when the task
1140 environments had lower social correlations of $r = .1 \pm .05$
1141 (Exp. 1R), AS became dominant model ($p_{xpAS} = .999$),
1142 implying that humans are sensitive to the relevance of social
1143 information (Fig. S6). However, the threshold at which they
1144 stop relying on social information, whether they would over-
1145 or underweight social information before then, and how this
1146 could be influenced by the framing of the task, remains a
1147 question for further research.

1148 Additionally, in real life, we would expect to find some
1149 people with more similar tastes to us and others with more
1150 different tastes. How such varying correlations between
1151 participants would affect how social information is used
1152 remains an open question. Given prior research showing that
1153 humans are quite capable of adjusting their social learning
1154 based on the skill of the observed individual (2, 25, 26), it
1155 seems reasonable to assume they could adjust to higher or
1156 lower levels of correlation as well. It would be interesting to
1157 see if this would lead to only learning from the most closely
1158 correlated individual, or from all sources but with higher
1159 assumed noise for lower correlations.

1160 Given the novel task setting of social learning in pos-
1161itively correlated environments, we chose to investigate the
1162 naturalistic interactions of groups of four real participants.
1163 We would not have been able to make an informed choice of
1164 model for a more controlled setting where humans are placed
1165 in groups of artificial agents a priori. Thus, having humans
1166 do the task in groups ensured that we were not affecting
1167 their behaviour through unnatural model choices. In the
1168 future, more granular insights into the exact usage of social
1169 information could be gained by placing participants in groups
1170 with Social Generalization agents, which can be used to more
1171 precisely manipulate the usefulness of social information and
1172 control group dynamics.

1173 Finally, our task structure only maps to real-world
1174 scenarios in which it is not beneficial to learn about others'
1175 preferences. When an individual is making a choice for
1176 themselves in a matter of preference, their own preference
1177 has the most weight. However, when the individual lacks
1178

experience, relying on information from others can be more
1179 helpful than having to explore without any further guidance.
1180 However, in the real world, social information can be used
1181 for a multitude of inferences, many of which may go beyond
1182 using its use as an exploration tool. Knowing how much a
1183 friend likes a restaurant may not be helpful for individual
1184 decision-making once the individual has tried the restaurant
1185 themselves. But their friend's experience may provide more
1186 information about them, or the restaurant, which can be
1187 useful in other ways. Here, we model our task as a matter
1188 of preference, where the most relevant piece of information
1189 for reward is knowing one's own salt distribution, and leave
1190 collective information gathering about the world as a subject
1191 for future research.

1192 **Conclusions.** Across two experiments, we found that people
1193 used social information more flexibly than previously ac-
1194 counted for, successfully integrating information from others
1195 with diverse reward functions, but taking it "with a grain
1196 of salt". Our model captures this, by integrating social
1197 information as inherently noisier, since it is not directly
1198 applicable to one's own circumstances. Social learning
1199 also functioned as an exploration tool, partially replacing
1200 uncertainty-directed exploration, and potentially helping
1201 participants behave more optimally.

Materials and Methods

1202 **Experiment design.** Across all experiments, participants explored
1203 spatially correlated multi-armed bandits (41) with social cor-
1204 relations across participants. The bandits were displayed as
1205 grids consisting of 121 tiles. Environments were structured
1206 identically across studies and conditions. Each tile yielded normally
1207 distributed rewards: $r(\mathbf{x}) \sim \mathcal{N}(f(\mathbf{x}), \sigma_e^2)$ where the expected
1208 reward across all tiles was sampled from a Gaussian Process (GP)
1209 prior to induce spatial correlations $f \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$ and the
1210 variance was fixed to $\sigma_e^2 = .0001$. To generate the environments,
1211 we sampled a set of 11x11 parent grids from a Gaussian process
1212 prior with an RBF-kernel with a length scale of $\lambda = 2$. We then
1213 used these parent environment's means as the prior means to
1214 sample a set of child environments. To facilitate correlations across
1215 environments, the child environments were filtered to only include
1216 those which correlated with the parent environment by at least
1217 $r = .6$. This subset of child environments was then filtered to only
1218 include sets of 4 environments which had correlation coefficients of
1219 $r = .6 \pm 0.05$ with each other to use in the task. We generated 40
1220 sets of correlated environments for the experiments this way. At the
1221 start of the experiment, we sampled a number of environment sets
1222 corresponding to the number of rounds without replacement for
1223 each group. Exp. 1R, the replication of E1 with lower correlations,
1224 used the same strategy to sample environments with correlation
1225 coefficients of $r = .1 \pm 0.05$.

1226 For each round of the experiment, each participant was assigned
1227 an environment from such a correlated set of four. Exp.s 1 and 1R
1228 consisted of 8 rounds, and Exp. 2 consisted of 8 solo and 8 group
1229 rounds, totalling 16 rounds. The search horizon was 14 for all
1230 experiments, with one tile being revealed at random at the start of
1231 each round. To prevent participants from getting used to the same
1232 reward structure, including its global maximum, environments
1233 were rescaled to a randomly selected maximum value between
1234 60 and 80 for each round. This rescaling was consistent across
1235 participants. To prevent a single participant from holding up a
1236 group, a random tile would be selected if they did not make a
1237 choice within 10 seconds. Such random choice trials were excluded
1238 from analysis. After selecting a tile, they would wait for all other
1239 participants to make their selection as well. Once all participants
1240 made a choice, the task would move on to the next trial. In
solo rounds, participants would only see their own bandit. In

group rounds, participants were also permanently shown all other participants' bandits, including choices and outcomes. This was the only difference between the two conditions. Choice and outcome information in group rounds were updated for all participants at once after all group members had made a choice.

Participants and design. Participants for all experiments were recruited via Prolific and assigned to groups of four based on access time to the experiment. They were paid a base rate for expected experiment duration, and could earn a bonus of maximum the same amount based on performance. All experiments were approved by the Ethics Committee of the University of Konstanz ("Collective learning and decision-making study"), and participants provided informed consent prior to participation.

Exps. 1 and 1R were observational studies with only the group condition. For Exp. 1, we recruited $N=188$ participants. After eliminating all groups with drop-out, the final sample size was $N=128$ (mean age: 38.5 ± 12.7 SD; 44 females). On average, participants spent 20.8 ± 0.5 minutes on the task and earned £ 7.19 ± 0.04 . For Exp. 1R, we recruited $N=200$ participants. After eliminating all groups with drop-out, the final sample size was $N=156$ (mean age: 36.9 ± 10.7 SD; 87 females). On average, participants spent 22.4 ± 0.6 minutes on the task and earned £ 7.58 ± 0.07 .

For Exp. 2, which varied solo vs. group conditions within-subject in interleaved order, we recruited 220 participants. Condition order (solo round first vs. group round first) was counterbalanced across groups. After eliminating all groups with drop-out, the final sample size was $N=132$ (mean age: 35.8 ± 11.2 SD, 46 females). On average, participants spent 31.0 ± 0.6 minutes on the task and earned £ 10.4 ± 0.08 .

Materials and procedure. In all experiments, participants took part in groups of four, which they were assigned to based on access time. After giving informed consent, participants were instructed that they were embarking on a scientific mission to collect salt samples from alien oceans on other planets, and that their goal was to collect as many salt samples as possible. They were informed that they could revisit the same area to get a similar reward, with salt not depleting from repeated sampling. They were also told that other scientists on their team would collect different salts, so there would be no competition for resources, but that the salts were generated by the same process, and locations with high salt concentrations were thus correlated across the salts. In Exp. 2, participants were additionally told that they would be sent on both solo and group missions, with no information from their teammates being available in solo missions. However, participants were never instructed about how to use social information.

After the instructions, participants in all experiments were shown fully revealed example environments to ensure they understand the structure and how social information usage may benefit them. After passing a comprehension check, participants moved on to a waiting room, which would launch the task once four people had joined. If there was no group of four after 3 minutes of a room being open, all participants in that room were redirected to the post experiment questionnaire.

Once in the task, participants would be presented with their bandit grid with one tile revealed. In the group condition, they would additionally see the bandits of all other participants in

the group with the rewards revealed as well. While participants' bandits were still correlated in solo rounds, they could not see other group members' bandits or choices in this condition.

Evolutionary simulations. For any possible combination of models, we generated an initial population consisting of an equal proportion of all the models. For the three-way mixes, which lack one agent when all models have exactly equal numbers, the final agent was randomly selected to be any of the three models. Initial populations were generated based on a common set of priors (Fig. 2b, see also SI). We used tournament selection to select agents for the next generation: groups of four agents were randomly drawn with replacement to compete in one round of the task. The selection probability of agents thus selected was lowered to prevent the same agents from being sampled too often. The agent with the highest score in a group was selected to seed the next generation. This procedure was repeated until the full population size of $N=100$ was reached. Each agent was thus sampled about 4 times. Before repeating the process for the next generation, mutations were applied to a part of the population. There was a 2% chance of *parameter mutations*, in which a parameter would have Gaussian noise $\sim \mathcal{N}(0, 0.2)$ added. If this caused the parameter to go out of bounds, it was resampled from prior. There was a 0.2% chance of a *type mutation*, in which the agent's model would be randomly resampled. The new model could be one that was not initially present in the population. To allow for invasions, we kept the baseline (GP-UCB) parameters of the mutating agent stable, and only modified the social parameter, which determines the model. Simulations were run this way for 500 generations. Simulations of all initial populations were repeated 10 times to ensure stability of the results.

Model comparisons. We fit models based on cross-validated maximum-likelihood estimation. We iteratively formed the training sets by leaving one round out, computing parameter estimates on this set, and evaluating model predictions on the out-of-sample round. Overall goodness of fit was evaluated based on the sum of the prediction error on each of the out-of-sample predictions. For Exp. 2, participant data was split into solo and group rounds before fitting. We used the summed out-of-sample log likelihood as an approximation of the model evidence to perform hierarchical Bayesian model comparison (48).

Data and Code Availability. All data and code are publicly available at <https://github.com/AlexandraWitt/socialGeneralization>.

ACKNOWLEDGMENTS. AW and CMW are supported by the German Federal Ministry of Education and Research (BMBF): Tübingen AI Center, FKZ: 01IS18039A and funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy—EXC2064/1–390727645. The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting AW. WT and WG are supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy – EXC 2117 – 42203798. An early version of Exp. 1 was presented at the 45th Annual Conference of the Cognitive Science Society (66).

1. W Toyokawa, A Whalen, KN Laland, Social learning strategies regulate the wisdom and madness of interactive crowds. *Nat. Hum. Behav.* **3**, 183–193 (2019).
2. A Najar, E Bonnet, B Bahrami, S Palminteri, The actions of others act as a pseudo-reward to drive imitation in the context of social reinforcement learning. *PLoS biology* **18**, e3001028 (2020).
3. R McElreath, et al., Beyond existence and aiming outside the laboratory: estimating frequency-dependent and pay-off-biased social learning strategies. *Philos. Transactions Royal Soc. B: Biol. Sci.* **363**, 3515–3528 (2008).
4. SA Park, S Goiame, DA O'Connor, JC Dreher, Integration of individual and social information for decision-making in groups of different sizes. *PLoS biology* **15**, e2001958 (2017).
5. L Molleman, et al., Strategies for integrating disparate social information. *Proc. Royal Soc. B* **287**, 20202413 (2020).
6. CJ Charpentier, K ligaya, JP O'Doherty, A neuro-computational account of arbitration between choice imitation and goal emulation during human observational learning. *Neuron* **106**, 687–699 (2020).
7. L Rendell, et al., Cognitive culture: theoretical and empirical insights into social learning strategies. *Trends cognitive sciences* **15**, 68–76 (2011).
8. KN Laland, Social learning strategies. *Animal Learn. & Behav.* **32**, 4–14 (2004).
9. L Rendell, et al., Why copy others? Insights from the social learning strategies tournament. *Science* **328**, 208–213 (2010).
10. RL Kendal, et al., Social learning strategies: Bridge-building between fields. *Trends Cogn. Sci.* **22**, 651–665 (2018).
11. W Hoppitt, KN Laland, *Social learning: an introduction to mechanisms, methods, and models.* (Princeton University Press), (2013).
12. W Toyokawa, Y Saito, T Kameda, Individual differences in learning behaviours in humans: Asocial exploration tendency does not predict reliance on social learning. *Evol. Hum. Behav.* **38**, 325–333 (2017).
13. D Deffner, V Kleinow, R McElreath, Dynamic social learning in temporally and spatially variable environments. *Royal Soc. open science* **7**, 200734 (2020).

- 1365 14. TJ Morgan, LE Rendell, M Ehn, W Hoppitt, KN Laland, The evolutionary basis of human
1366 social learning. *Proc. Royal Soc. B: Biol. Sci.* **279**, 653–662 (2012). 1427
- 1367 15. T Kameda, D Nakanishi, Cost–benefit analysis of social/cultural learning in a nonstationary
1368 uncertain environment: An evolutionary simulation and an experiment with human subjects.
Evol. Hum. Behav. **23**, 373–393 (2002). 1428
- 1369 16. C Heath, C Bell, E Sternberg, Emotional selection in memes: the case of urban legends. *J.
1370 personality social psychology* **81**, 1028 (2001). 1429
- 1371 17. K Eriksson, JC Coulas, Corpses, maggots, poodles and rats: Emotional selection operating
1372 in three phases of cultural transmission of urban legends. *J. Cogn. Cult.* **14**, 1–26 (2014). 1430
- 1373 18. JM Stubbersfield, JJ Tehrani, EG Flynn, Serial killers, spiders and cybersex: Social and
1374 survival information bias in the transmission of urban legends. *Br. journal psychology* **106**,
288–307 (2015). 1431
- 1375 19. T Blaine, P Boyer, Origins of sinister rumors: A preference for threat-related material in the
1376 supply and demand of information. *Evol. Hum. Behav.* **39**, 67–75 (2018). 1432
- 1377 20. JM Stubbersfield, JJ Tehrani, EG Flynn, Chicken tumours and a fishy revenge: Evidence for
1378 emotional content bias in the cumulative recall of urban legends. *J. Cogn. Cult.* **17**, 12–26
(2017). 1433
- 1379 21. A Mesoudi, A Whiten, R Dunbar, A bias for social information in human cultural
1380 transmission. *Br. journal psychology* **97**, 405–423 (2006). 1434
- 1381 22. C Brand, S Heap, T Morgan, A Mesoudi, The emergence and adaptive use of prestige in an
1382 online social learning task. *Sci. reports* **10**, 1–11 (2020). 1435
- 1383 23. C Brand, A Mesoudi, TJ Morgan, Trusting the experts: The domain-specificity of
1384 prestige-biased social learning. *PLoS one* **16**, e0255346 (2021). 1436
- 1385 24. A Mesoudi, An experimental comparison of human social learning strategies: payoff-biased
1386 social learning is adaptive but underused. *Evol. Hum. Behav.* **32**, 334–342 (2011). 1437
- 1387 25. R Watson, TJ Morgan, RL Kendal, J Van de Vyver, J Kendal, Social learning strategies and
1388 cooperative behaviour: Evidence of payoff bias, but not prestige or conformity, in a social
1389 dilemma game. *Games* **12**, 89 (2021). 1438
- 1390 26. CM Wu et al., Visual-spatial dynamics drive adaptive social learning in immersive
1391 environments. *bioRxiv* (2023). 1439
- 1392 27. W Toyokawa, W Gaissmaier, Conformist social learning leads to self-organised prevention
1393 against adverse bias in risky decision making. *Elife* **11**, e75308 (2022). 1440
- 1394 28. BG Galef, Imitation in animals: history, definition, and interpretation of data from the
1395 psychological laboratory in *Social learning*. (Psychology Press), pp. 15–40 (2013). 1441
- 1396 29. G Biele, J Rieskamp, LK Krugel, HR Heekeken, The neural basis of following advice. *PLoS
1397 biology* **9**, e1001089 (2011). 1442
- 1398 30. PP Analytis, D Barkoczi, SM Herzog, Social learning strategies for matters of taste. *Nat.
1399 human behaviour* **2**, 415–424 (2018). 1443
- 1400 31. J Müller-Trede, S Choshen-Hillel, M Barneron, I Yaniv, The wisdom of crowds in matters of
1401 taste. *Manag. Sci.* **64**, 1779–1803 (2018). 1444
- 1402 32. I Yaniv, S Choshen-Hillel, M Milyavsky, Receiving advice on matters of taste: Similarity,
1403 majority influence, and taste discrimination. *Organ. Behav. Hum. Decis. Process.* **115**,
111–120 (2011). 1445
- 1404 33. J Jara-Ettinger, Theory of mind as inverse reinforcement learning. *Curr. Opin. Behav. Sci.*
1405 **29**, 105–110 (2019). 1446
- 1406 34. J Jara-Ettinger, LE Schulz, JB Tenenbaum, The naive utility calculus as a unified,
1407 quantitative framework for action understanding. *Cogn. Psychol.* **123**, 101334 (2020). 1447
- 1408 35. P Shafto, ND Goodman, MC Frank, Learning from others: The consequences of
1409 psychological reasoning for human learning. *Perspectives on Psychol. Sci.* **7**, 341–351
(2012). 1448
- 1410 36. S Croom, H Zhou, C Firestone, Seeing and understanding epistemic actions. *Proc. Natl.
1411 Acad. Sci.* **120**, e2303162120 (2023). 1449
- 1412 37. M Berke, J Jara-Ettinger, Integrating experience into bayesian theory of mind in
1413 *Proceedings of the Annual Meeting of the Cognitive Science Society*. Vol. 44, (2022). 1450
- 1414 38. RD Hawkins, et al., Flexible social inference facilitates targeted social learning when
1415 rewards are not observable. *Nat. Hum. Behav.* **7**, 1767–1776 (2023). 1451
- 1416 39. CL Baker, J Jara-Ettinger, R Saxe, JB Tenenbaum, Rational quantitative attribution of
1417 beliefs, desires and percepts in human mentalizing. *Nat. Hum. Behav.* **1**, 0064 (2017). 1452
- 1418 40. A Jern, CG Lucas, C Kemp, People learn other people's preferences through inverse
1419 decision-making. *Cognition* **168**, 46–64 (2017). 1453
- 1420 41. CM Wu, E Schulz, M Speekenbrink, JD Nelson, B Meder, Generalization guides human
1421 exploration in vast decision spaces. *Nat. human behaviour* **2**, 915–924 (2018). 1454
- 1422 42. RC Plate, H Ham, AC Jenkins, When uncertainty in social contexts increases exploration
1423 and decreases obtained rewards. *J. Exp. Psychol. Gen.* (2023). 1455
- 1424 43. A Naito, K Katahira, T Kameda, Insights about the common generative rule underlying an
1425 information foraging task can be facilitated via collective search. *Sci. Reports* **12**, 1–12
(2022). 1456
- 1426 44. A Witt, W Toyokawa, W Gaissmaier, P Lala, Kevin N, CM Wu, Social learning in a spatially
1427 correlated multi- armed bandit (2024). 1457
- 1428 45. AP Giron, et al., Developmental changes in exploration resemble stochastic optimization.
1429 *Nat. Hum. Behav.* (2023). 1458
- 1430 46. CM Wu, B Meder, E Schulz, Unifying principles of generalization: past, present, and future.
1431 *Annu. Rev. Psych* **76**, 1–33 (2024). 1459
- 1432 47. AR Rogers, Does biology constrain culture? *Am. Anthropol.* **90**, 819–831 (1988). 1460
- 1433 48. L Rigoux, KE Stephan, KJ Friston, J Daunizeau, Bayesian model selection for group
1434 studies—revisited. *Neuroimage* **84**, 971–985 (2014). 1461
- 1435 49. AN Tump, TJ Pleskac, RH Kurvers, Wise or mad crowds? The cognitive mechanisms
1436 underlying information cascades. *Science Advances* **6**, eabb0266 (2020). 1462
- 1437 50. AN Tump, M Wolf, J Krause, RH Kurvers, Individuals fail to reap the collective benefits of
1438 diversity because of over-reliance on personal information. *J. Royal Soc. Interface* **15**,
20180155 (2018). 1463
- 1439 51. W Toyokawa, Hr Kim, T Kameda, Human collective intelligence under dual
1440 exploration-exploitation dilemmas. *PLoS one* **9**, e95799 (2014). 1464
- 1441 52. M Wolf, J Krause, PA Carney, A Bogart, RH Kurvers, Collective intelligence meets medical
1442 decision-making: the collective outperforms the best radiologist. *PLoS one* **10**, e0134269
(2015). 1465
- 1443 53. D Bang, CD Frith, Making better decisions in groups. *Royal Soc. open science* **4**, 170193
(2017). 1466
- 1444 54. B Bahrami, et al., Optimally interacting minds. *Science* **329**, 1081–1085 (2010). 1467
- 1445 55. O Morin, PO Jacquet, K Vaesen, A Acerbi, Social information use and social information
1446 waste. *Philos. Transactions Royal Soc. B* **376**, 20200052 (2021). 1468
- 1447 56. A Acerbi, C Tennie, A Mesoudi, Social learning solves the problem of narrow-peaked search
1448 landscapes: experimental evidence in humans. *Royal Soc. open science* **3**, 160215 (2016). 1469
- 1449 57. R Bhui, L Lai, SJ Gershman, Resource-rational decision making. *Curr. Opin. Behav. Sci.*
1450 **41**, 15–21 (2021). 1470
- 1451 58. I Cogliati Dezza, A Cleeremans, W Alexander, Should we control? the interplay between
1452 cognitive control and information integration in the resolution of the exploration-exploitation
1453 dilemma. *J. Exp. Psychol. Gen.* **148**, 977 (2019). 1471
- 1454 59. CM Wu, E Schulz, TJ Pleskac, M Speekenbrink, Time pressure changes how people
1455 explore and respond to uncertainty. *Sci. reports* **12**, 4122 (2022). 1472
- 1456 60. JD Cohen, SM McClure, AJ Yu, Should i stay or should i go? how the human brain
1457 manages the trade-off between exploitation and exploration. *Philos. Transactions Royal
1458 Soc. B: Biol. Sci.* **362**, 933–942 (2007). 1473
- 1459 61. N Fleischhut, FM Artinger, S Olschewski, R Hertwig, Not all uncertainty is treated equally:
1460 Information search under social and nonsocial uncertainty. *J. Behav. Decis. Mak.* **35**, e2250
(2022). 1474
- 1460 62. CM Wu, E Schulz, MM Garvert, B Meder, NW Schuck, Similarities and differences in spatial
1461 and non-spatial cognitive maps. *PLOS Comput. Biol.* **16**, e1008149 (2020). 1475
- 1461 63. CM Wu, E Schulz, SJ Gershman, Inference and search on graph-structured spaces.
1462 *Comput. Brain & Behav.* **4**, 125–147 (2021). 1476
- 1462 64. I Selbing, B Lindström, A Olsson, Demonstrator skill modulates observational aversive
1463 learning. *Cognition* **133**, 128–139 (2014). 1477
- 1463 65. D Deffner, et al., Collective incentives reduce over-exploitation of social information in
1464 unconstrained human groups. *Nat. Commun.* (2024). 1478
- 1464 66. A Witt, W Toyokawa, K Lala, W Gaissmaier, CM Wu, Social learning with a grain of salt in
1465 *Proceedings of the 45th Annual Conference of the Cognitive Science Society*, eds. M
1466 Goldwater, F Anggoro, B Hayes, D Ong. (Cognitive Science Society, Sydney, Australia),
1467 (2023). 1479
- 1465 67. S Croom, H Zhou, C Firestone, Seeing and understanding epistemic actions. *Proc. Natl.
1466 Acad. Sci.* **120**, e2303162120 (2023). 1480
- 1466 68. M Berke, J Jara-Ettinger, Integrating experience into bayesian theory of mind in
1467 *Proceedings of the Annual Meeting of the Cognitive Science Society*. Vol. 44, (2022). 1481
- 1467 69. RD Hawkins, et al., Flexible social inference facilitates targeted social learning when
1468 rewards are not observable. *Nat. Hum. Behav.* **7**, 1767–1776 (2023). 1482
- 1468 70. CL Baker, J Jara-Ettinger, R Saxe, JB Tenenbaum, Rational quantitative attribution of
1469 beliefs, desires and percepts in human mentalizing. *Nat. Hum. Behav.* **1**, 0064 (2017). 1483
- 1469 71. A Jern, CG Lucas, C Kemp, People learn other people's preferences through inverse
1470 decision-making. *Cognition* **168**, 46–64 (2017). 1484
- 1470 72. CM Wu, E Schulz, M Speekenbrink, JD Nelson, B Meder, Generalization guides human
1471 exploration in vast decision spaces. *Nat. human behaviour* **2**, 915–924 (2018). 1485
- 1471 73. RC Plate, H Ham, AC Jenkins, When uncertainty in social contexts increases exploration
1472 and decreases obtained rewards. *J. Exp. Psychol. Gen.* (2023). 1486
- 1472 74. A Naito, K Katahira, T Kameda, Insights about the common generative rule underlying an
1473 information foraging task can be facilitated via collective search. *Sci. Reports* **12**, 1–12
(2022). 1487
- 1473 75. W Toyokawa, W Gaissmaier, P Lala, Kevin N, CM Wu, Social learning in a spatially
1474 correlated multi- armed bandit (2024). 1488
- 1474 76. AP Giron, et al., Developmental changes in exploration resemble stochastic optimization.
1475 *Nat. Hum. Behav.* (2023). 1489
- 1475 77. CM Wu, B Meder, E Schulz, Unifying principles of generalization: past, present, and future.
1476 *Annu. Rev. Psych* **76**, 1–33 (2024). 1490
- 1476 78. AR Rogers, Does biology constrain culture? *Am. Anthropol.* **90**, 819–831 (1988). 1491
- 1477 79. L Rigoux, KE Stephan, KJ Friston, J Daunizeau, Bayesian model selection for group
1478 studies—revisited. *Neuroimage* **84**, 971–985 (2014). 1492
- 1478 80. AN Tump, TJ Pleskac, RH Kurvers, Wise or mad crowds? The cognitive mechanisms
1479 underlying information cascades. *Science Advances* **6**, eabb0266 (2020). 1493
- 1479 81. AN Tump, M Wolf, J Krause, RH Kurvers, Individuals fail to reap the collective benefits of
1480 diversity because of over-reliance on personal information. *J. Royal Soc. Interface* **15**,
20180155 (2018). 1494
- 1480 82. W Toyokawa, Hr Kim, T Kameda, Human collective intelligence under dual
1481 exploration-exploitation dilemmas. *PLoS one* **9**, e95799 (2014). 1495

Supplementary Information for Humans flexibly integrate social information despite interindividual differences in reward.

Alexandra Witt, Wataru Toyokawa, Kevin N. Lala, Wolfgang Gaissmaier & Charley M. Wu

Learning and ordering effects

In Exp. 1, there was a small learning effect over rounds ($0.004, 95\% - CI : [0.001, 0.007], p = 0.0027$). The same was true for Exp. 2 ($0.002, 95\% - CI : [0.0003, 0.003], p = 0.012$). In Exp. 2, there was no effect of block order ($0.005, 95\% - CI : [-0.02, 0.015], p = 0.63$), or the interaction of round and block order ($-0.0001, 95\% - CI : [-0.002, 0.002], p = 0.94$) on performance (Fig. S1).

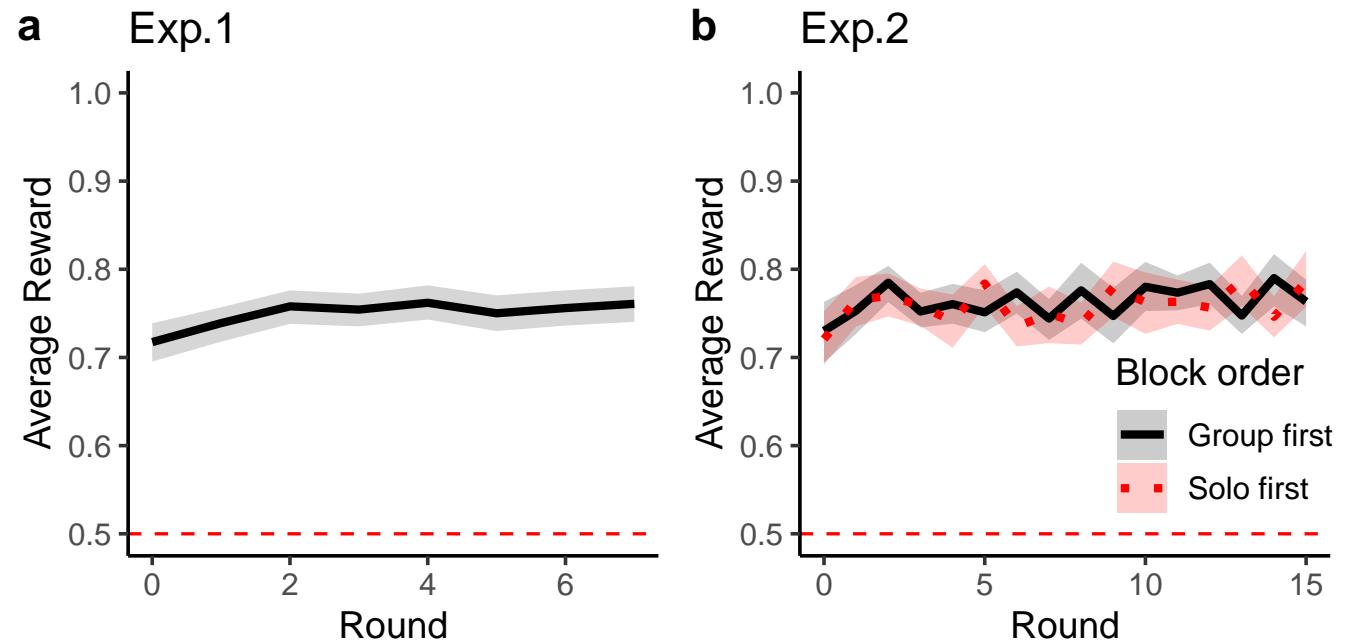


Fig. S1. Learning over rounds. a) Learning over rounds in Exp. 1. Red dashed line gives chance level performance. b) Learning over rounds and by block order in Exp. 2. Red dashed line gives chance level performance.

Model simulations and variants

Priors. For initial agent-based and evolutionary simulations, we drew parameters from a set of parameter priors. Priors for the baseline asocial learner were based on the values found in prior studies(41). Since none of the parameter values can be negative, but have no upper bound, we used log-normal distributions around the reported average participant estimates. This resulted in the following parameter priors:

$$\begin{aligned} \lambda, \beta &\sim \text{LogNormal}(-0.75, 0.5) \\ \tau &\sim \text{LogNormal}(-4.5, 0.75) \end{aligned} \quad [8]$$

For the social parameters, no prior empirical results existed, so we used priors that covered as much of the theoretical space as possible. While we are able to cover the entire possible range for Decision Biassing and Value Shaping, since $\alpha, \gamma \in [0, 1]$, the Social Generalization noise parameter ϵ_{soc} cannot be negative, but can grow infinitely large. Therefore, we chose to centre an exponential distribution around $\epsilon_{soc} = 2$, which we found to be good, but not optimal, in simulations, resulting in the following priors:

$$\begin{aligned} \alpha, \gamma &\sim \text{Uniform}(0, 1) \\ \epsilon_{soc} &\sim \text{Exponential}(0.5) \end{aligned} \quad [9]$$

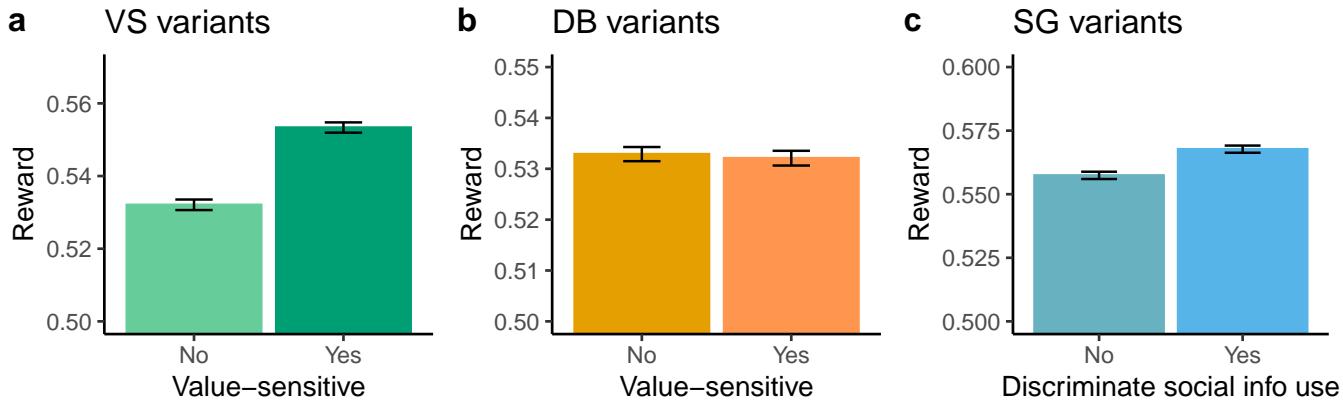
Unless otherwise stated, simulations were run using these priors.

Model variants. For model variants, we simulated groups of two asocial agents as well as one canonical and one modified agent. This was done to be able to directly compare the model's performance given identical information. Asocial agents were used to prevent Roger's paradox(47), i.e. the frequency-dependent fitness of social models, from affecting the results. We ran these simulations with task settings (search horizon of 14, 8 rounds) with 1000 different parameter sets.

Value-agnostic Value Shaping. As mentioned in the main text, Value Shaping benefitted from including social value information (Fig. S2a) compared to the unbiased version common in previous literature, where value is generically boosted for options selected by others. We implemented this as $V(\mathbf{x}) = V_{x,ind} + \alpha \cdot n_{x,soc,t-1}$. Our canonical prediction error implementation of VS significantly outperforms this alternative ($t(1998) = 8.4, p < .001, d = 0.4, BF > 100$) and was thus chosen as the main implementation.

1613 **Value-sensitive Decision Biasing.** We also tried to adapt Decision Biasing to our task by including value information. This might have been
 1614 a beneficial update, since outcome information was generally not available in previous studies(1, 2), but was in ours. We modified the
 1615 social policy so that it increased proportionally to the frequency of a social choice, weighted by how much higher the social reward was
 1616 than the average experienced individual reward $\pi_{soc}(x) \propto n_{x_{soc_{t-1}}} \cdot (m(x_{soc_{t-1}}) - m(\bar{x}_{ind}))$. In cases where $m(x_{soc_{t-1}}) < m(\bar{x}_{ind})$, the
 1617 social information was ignored to prevent negative probabilities in the policy. Despite this added information (and the large sample size),
 1618 there was no significant difference between the two models' scores ($t(1998) = 0.3$, $p = .730$, $d = 0.02$, $BF = .05$; Fig. S2a). Thus, we chose
 1619 to use the simpler model.
 1620

1621 **Indiscriminate Social Generalization.** There is an edge case of Social Generalization for $\epsilon_{soc} = 0$. It means that social and individual
 1622 information are treated identically. While technically not a separate model, we show that discriminate use of social information ($\epsilon_{soc} \neq 0$)
 1623 is significantly better than indiscriminate use in our task ($t(1998) = 3.4$, $p < .001$, $d = 0.2$, $BF = 18$; Fig. S2b).
 1624



1625 **Fig. S2. Model variants.** a) Value Shaping with (canonical) and without value sensitivity added. b) Decision Biasing with and without (canonical)
 1626 value sensitivity added. c) SG with ϵ_{soc} set to 0 and not (canonical).
 1627

1628 Detailed evolutionary simulations

1629 Visualization as a ternary plot as in the main text only allows for comparisons between 3 models at a time. As this paper put a focus on
 1630 social learning, we chose to compare the three candidate social models in this manner. Figure S3a shows the evolutionary trajectories for
 1631 all starting populations, including AS. As reported in the main text, SG takes over and dominates all populations.
 1632

1633 When it comes to parameter evolution, we can gain insight into the “optimal” SG agent based on evolutionary simulations as well
 1634 (Fig. S3b). We only report SG parameters as parameters of other models are unstable due to their low population size. The parameter
 1635 evolutions are discussed in the main text.
 1636

1637 We additionally ran evolutionary simulations in a setting analogous to previous literature(2) with one expert choosing the correct
 1638 option and the (social) learning agent being in an identical environment. This served to show that the spatially correlated bandit setup is
 1639 not inherently different from simpler bandits used in previous literature. As reported in the main text, we replicate VS being the dominant
 1640 model in such settings when comparing only the previously established models (Fig. S4a), and find an equilibrium between VS and SG
 1641 when considering all of our candidate models (Fig. S4b-c). This is because when fully socially reliant ($\alpha=1$ for VS, or $\epsilon_{soc} = 0$ for SG), as
 1642 is optimal when learning from an expert in the same environment, and in the same environment as said expert, VS and SG make identical
 1643 choice predictions, only differing at the stage at which social information is integrated (Fig. S4d). SG can be viewed as an extension of VS
 1644 to cases where one has to learn from others in non-identical environments, not a completely new model.
 1645

1646 Normative strategy across environmental correlations

1647 We only test environments with a social correlation of $r = .6 \pm 0.05$, which is somewhat arbitrary. To further investigate how Social
 1648 Generalization would perform in different social correlation settings, we ran evolutionary simulations as described in the main text methods
 1649 for environments with a range of social correlations (from $r = 0 \pm 0.05$ to $r = .9 \pm 0.05$ in .1 increments). While Asocial Learning performs
 1650 better in uncorrelated environments (as was to be expected), Social Generalization appears to be able to exploit even correlations in the
 1651 range of $r = .1 \pm 0.05$, and taking over as the prevalent model in the final generation starting at correlations of $r = .2 \pm 0.05$. (Fig. S5).
 1652

1653 Human behaviour in environments with $r = 1$

1654 Following the insight provided by the evolutionary simulations across different social correlations, we followed up with another replication
 1655 of Exp. 1, changing the social correlation to $r = .1 \pm .05$ to assess if humans would be sensitive to the lowered correlation and how it
 1656 would affect their behaviour (N=156). We kept instructions identical to the .6-version of the task, which also let us assess whether the
 1657 instruction influenced on use of social information. In environments with low correlations, participants were predominantly in line with
 1658 asocial behaviour ($p_{xp} = .999$). In the subset of participants which were still best fit by SG, the average value of ϵ_{soc} was 14.29, which
 1659 was significantly higher than the values observed in environments with higher correlations ($t(141) = 2.8$, $p = .005$, $d = 0.6$, $BF = 7.1$).
 1660

1661 Figure S6 compares p_{xpSG} (a) and social noise (b) across correlations. Humans appear to be able to adapt to varying correlations,
 1662 reducing the propensity of social learning (SG in our case) and their reliance on social information in relation to the relevance of such
 1663 information. This adaptation seems to be unrelated to instructions, which remained identical across correlations.
 1664

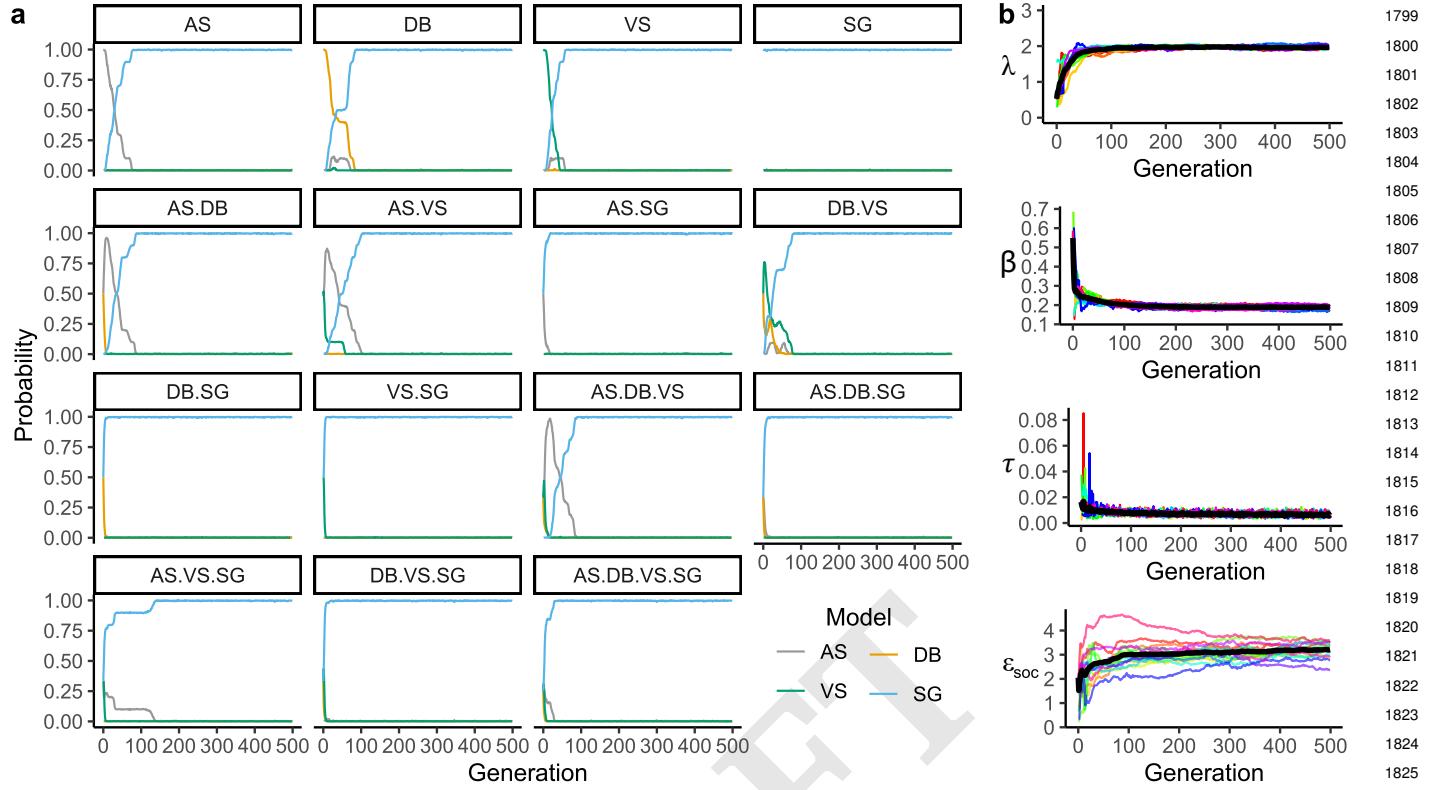


Fig. S3. Evolutionary simulations in correlated environments in detail. **a)** Evolutionary trajectories across all starting populations. Facet labels show initial population, and lines show the probability of a given model in the population. Social Generalization dominates across all initial populations. **b)** Evolved parameters for Social Generalization. Thick black line is the average.

Reward improvement

In the analysis of experiment 1, we investigated how participants used both individual and social information to guide their exploration. To this end, we analyzed the influence of improvement potential (the difference between previous individual and previous social reward in the case of social information, and the difference between maximum possible reward and previous individual reward for individual information) on reward improvement (the difference between current and previous individual reward). While the data corroborated no effect of negative social information (improvement potential < 0), there seemed to be a strong relationship between positive social improvement potential and reward improvement (Fig. S7a). When modelling this relationship, we not only found a general relationship between improvement potential and improvement (0.53 [0.47, 0.60]), but also both a significant positive effect of social information (0.06 [0.04, 0.07]) and its interaction with improvement potential (0.12 [0.09, 0.15]; Fig. S7b). This seemed to indicate that social information was even more effective than individual information in guiding participants' exploration.

However, while the relationship remained similar in experiment 2 (Fig. S7c), and the baseline effects replicated (improvement potential: 0.53 [0.51, 0.55]; social information: 0.05 [0.04, 0.06]; their interaction: 0.12 [0.09, 0.15]), we found none of their interactions with task type were significant (improvement potential*group round: 0.00 [-0.03, 0.03]; social info*group round: 0.01 [-0.01, 0.02]; improvement potential*social info*group round: -0.01 [-0.05, 0.03]; Fig. S7d). This shows that the effect we found in experiment 1 was solely based on the task structure, and not any actual benefit of social information usage.

Model bounding

As the baseline asocial learning model is nested in all social models, we determined bounds for the social models to minimize model mimicry, and thus improve recovery. This serves to make the modelling more stringent compared to previous work (66). The bounds were determined based on the social mechanisms of the respective models. Since DB effectively only changes imitation rate (mixing parameter γ effectively trades off between individual learning and imitation, which makes it interpretable as an average imitation rate per trial), we chose to determine the bound based on expected average imitation based on individual learning in correlated environments. We simulated AS with priors from previous literature, and set the lower bound at 95%-quantile of the resulting Poisson distribution based on imitation counts (Fig. S8a). This meant, that agents fit by DB were expected to imitate at least as much as the 5% tail-end of the asocial population.

Since VS affects the value function at a social observation, we determined the lower bound based on the minimum effect of a maximum reward social observation on a naive social learner (no individual observations) across a range of β -values. The criterion was a minimum of 5% change from the individual value (Fig. S8b).

For SG, ϵ_{soc} affects how strong of an effect social information has on the posterior of the GP. Hence, we set the bound at the social observation retaining at least 5% of its value given a naive social learner (Fig. S8c).

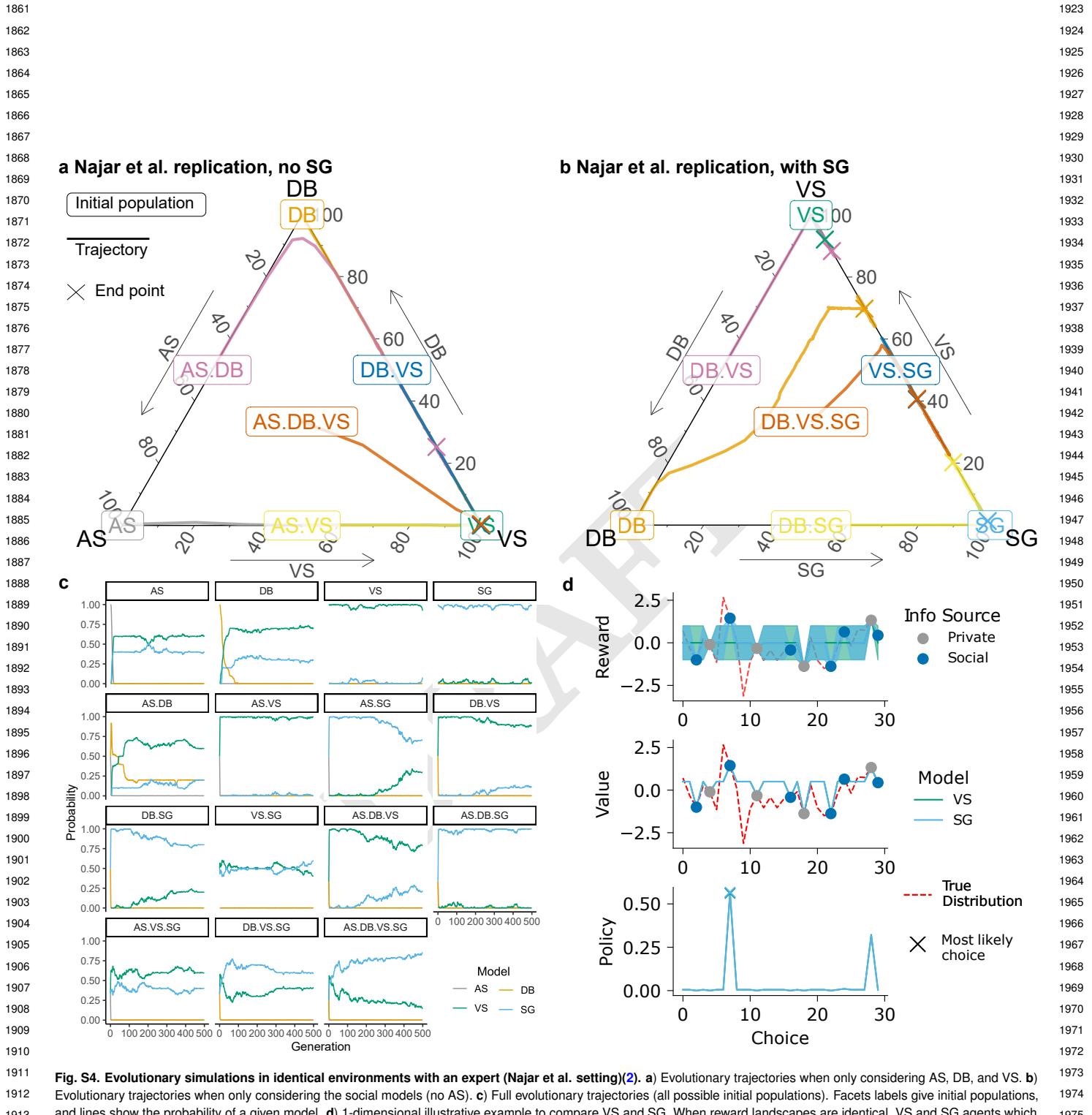


Fig. S4. Evolutionary simulations in identical environments with an expert (Najar et al. setting)(2). **a**) Evolutionary trajectories when only considering AS, DB, and VS. **b)** Evolutionary trajectories when only considering the social models (no AS). **c**) Full evolutionary trajectories (all possible initial populations). Facets labels give initial populations, and lines show the probability of a given model. **d**) 1-dimensional illustrative example to compare VS and SG. When reward landscapes are identical, VS and SG agents which fully rely on social information ($\alpha = 1$ or $\epsilon_{soc} = 0$) behave identically.

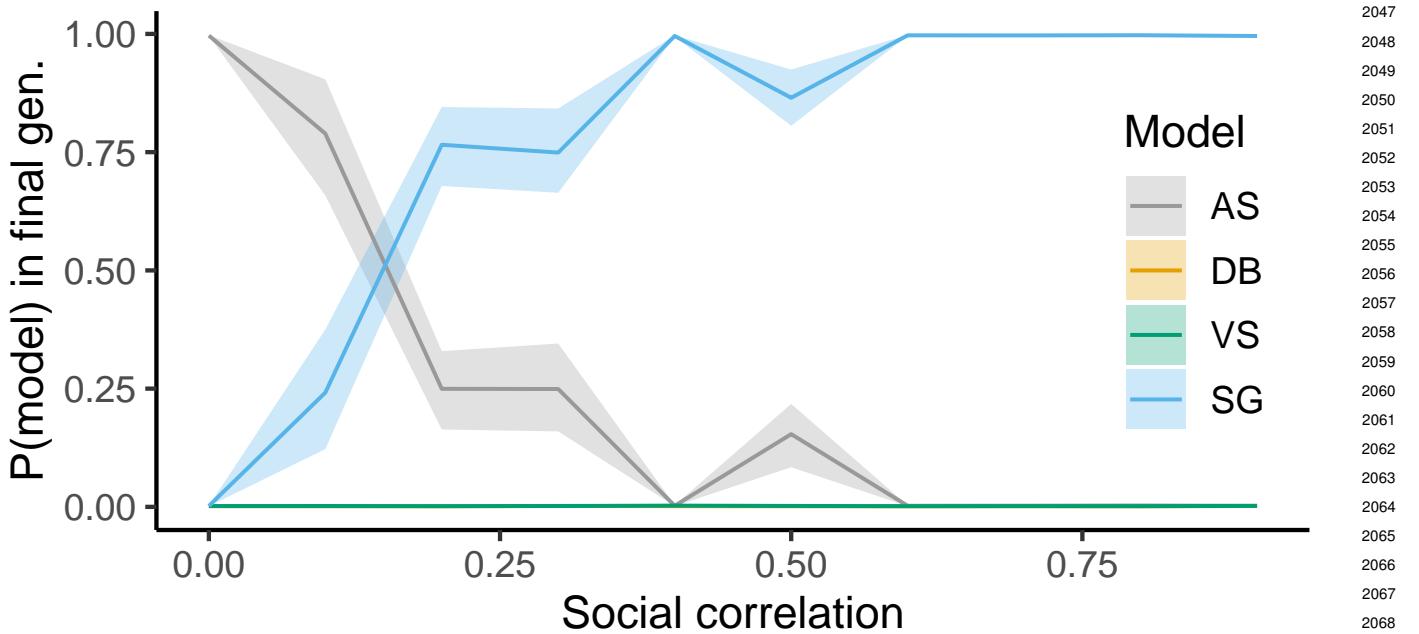


Fig. S5. Model performance in evolutionary simulations across social correlations. Depicted is the proportion of each model in the final generation of the evolutionary simulation over the environmental correlation across agents.

Model and parameter recovery

To assess model recovery, we simulated data using parameters fitted to participants for all models in experiment 1 and the group rounds of experiment 2. We then fit the simulated participants following the same procedure as used for actual participants, assigned each simulated participant a model based on best fit. We computed conditional probabilities for confusion and inversion matrices (Fig. S9). Despite the bounding, there is still some confusion potential between models, especially DB and VS, which hardly get fit with social parameters above the lower bound, and AS. There is also some confusion potential between AS and SG, but it is roughly balanced between the two, so overall fitting results should not be biased either way. In turn, the likelihood of a fit model being the generating one is highest for DB (0.85) and VS (0.9), but still high for AS (0.7) and SG (0.79).

When it comes to parameter recovery, we used the same procedure of simulating and fitting the data as for model recovery, but looked at the correlations between the parameters of the generating model and its fit instead of the best fitting model (Fig. S10). λ ($r_\tau = .87$, $p < .001$, $BF > 100$) and τ ($r_\tau = .86$, $p < .001$, $BF > 100$) correlate near perfectly with the generating parameters across models. When it comes to β and the social parameters, the issue of lower bound social parameters recurred, leading to worse fits for DB and VS. β correlations are still high overall ($r_\tau = .85$, $p < .001$, $BF > 100$), whereas the social parameter correlation is lower ($r_\tau = .27$, $p < .001$, $BF > 100$). However, given that neither DB nor VS fit participant data well, leading to them mimicking AS as much as possible, this lack of correlation is less concerning. Looking at only the correlation for ϵ_{soc} , it is noticeably higher ($r_\tau = .51$, $p < .001$, $BF > 100$).

Model performance

Models were fit using leave-one-round-out cross-validation. Negative log likelihoods were summed across test-rounds, with mean values being reported in Tables S1 and S2. Pseudo-R² was computed as $R^2 = 1 - (nLL_{model}/nLL_{random})$ where nLL_{random} is treating every choice as equally likely (1/121) for all non-random trials for that participant. Random trials were excluded from model fitting.

As the models were nested, and AS generally provides a good explanation even in social settings, especially once a high value option has been found, performance does not differ greatly between models. However, SG is consistently the best fit across both experiments with high social correlations.

Model	Mean nLL	Pseudo-R ²
AS	448.1	0.1576
DB	357.2	0.3268
VS	379.5	0.2845
SG	356.9	0.3273

Table S1. Model Performance Metrics for Exp. 1.

Exp. 2 parameters

To focus on the differences between β -values, we do not report all parameter values for the group rounds of Exp. 2 in the main text (Fig. S11). Generalization parameter $\lambda \approx 1.1$, which is significantly lower than the ground truth $\lambda = 2$ ($t(52) = -14.4$, $p < .001$, $d = 2.0$, $BF > 100$). Directed exploration parameter $\beta \approx 0.22$, and random exploration parameter $\tau \approx 0.06$. Social noise $\epsilon_{soc} \approx 9.5$, which is significantly higher than optimal the optimal value found in evolutionary simulations ($t(52) = 8.3$, $p < .001$, $d = 1.1$, $BF > 100$).

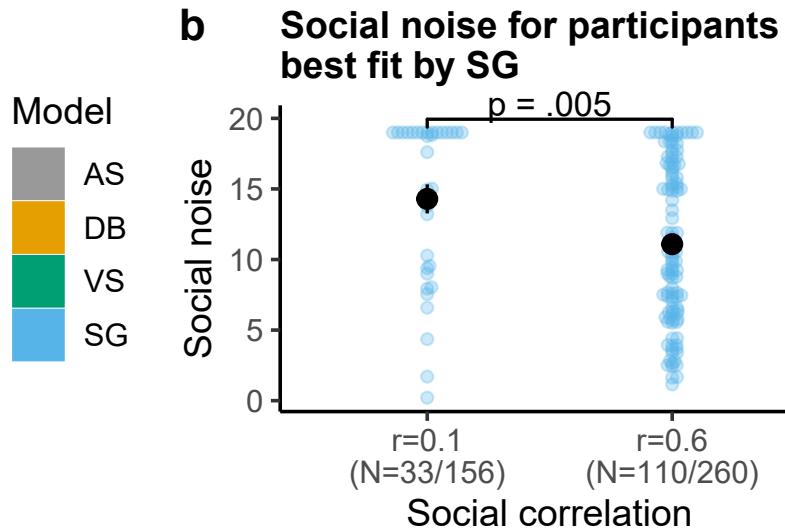
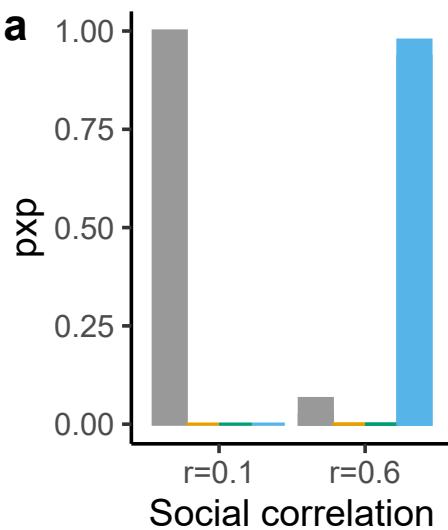


Fig. S6. Human sensitivity to varied social correlations. Values for correlation = 1 (hollow circles) are based on prior literature and parameter values making SG equivalent to winning models in these works. **a)** Protected exceedance probability of SG across social correlations. **b)** Social noise parameters across correlations.

Model	Mean nLL	Pseudo-R ²
AS	369.4	0.3106
DB	367.2	0.3147
VS	378.8	0.2932
SG	365.5	0.3178

Table S2. Model Performance Metrics for Exp. 2 group rounds.

Exploration optimality

Following up the comparatively lower β -parameter in experiment 1, we compare participant's β -parameters between the solo and group rounds. As reported in the main text, participants had significantly higher β -values in solo than in group rounds ($Z = -2.7$, $p = .004$, $r = -.23$, $BF = 63$, Fig. S12a). In simulations based on participant parameters while varying the parameter of interest, we find that such low values of β are actually optimal in the current task, both in solo and group rounds. We see that the group round value of β is closer to optimal than the solo round one, and both are lower and thus closer to optimal than the average found in previous literature(41).

As we also found higher values of λ in Exp. 1, we exploratively repeat these analyses. λ is significantly higher, and thus closer to ground truth, in group than in solo rounds. Again, group round λ is closest to the theoretical optimum, followed by solo round and previous study values. Taken together, this implies that social information may improve exploration behaviour closer to optimality in general.

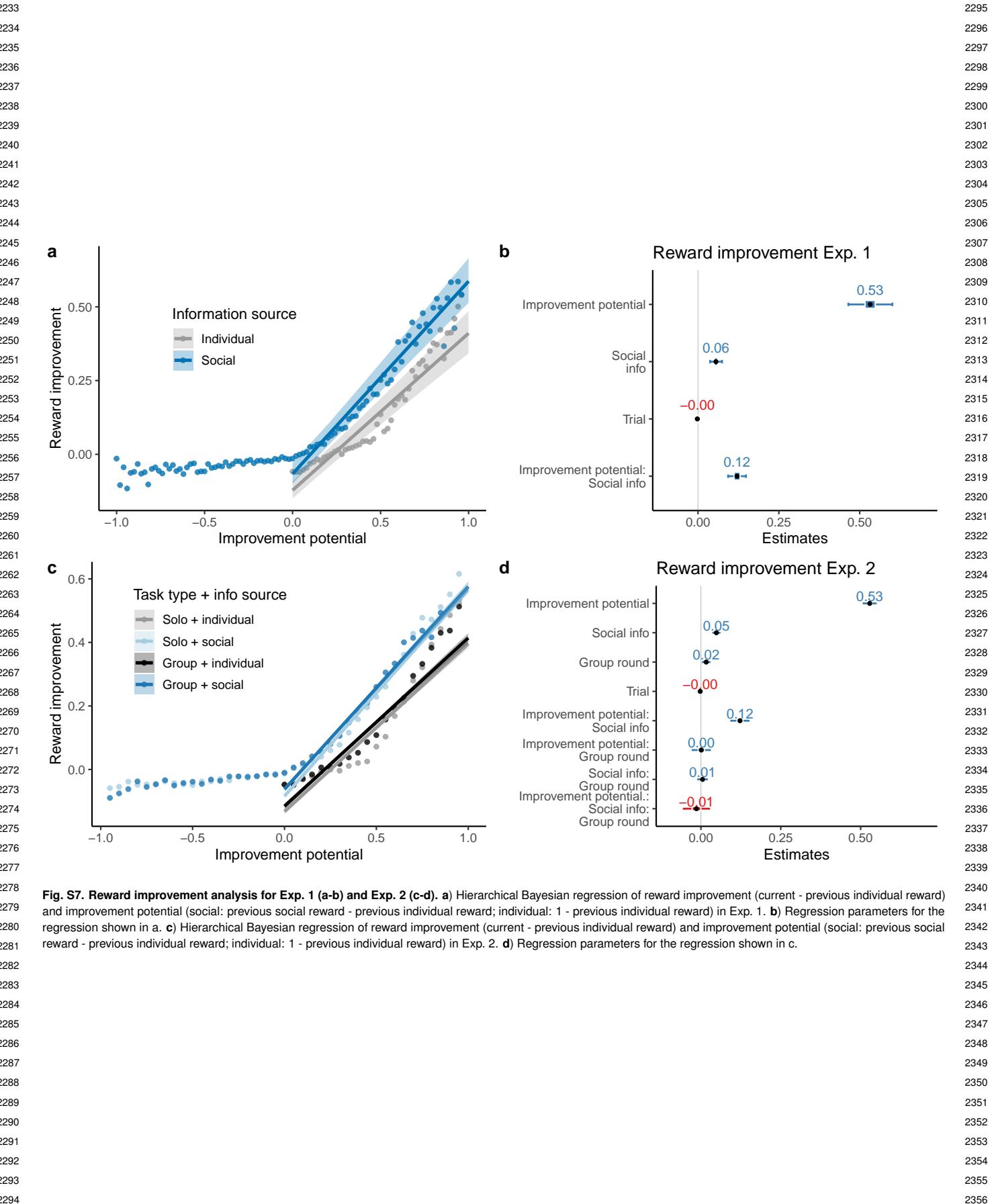
Exclusion analyses

To ensure that correlations were not spurious because of participants at the upper bound of ϵ_{soc} , we redid the correlation analyses for reward and β excluding any participants whose $\epsilon_{soc} > 18.9999$.

All relationships remained significant. In Exp. 1, there was a significant negative correlation between ϵ_{soc} and mean reward ($r_\tau = -.31$, $p = .001$, $BF = 24$), indicating higher reliance on social information leading to higher scores. There was also a significant positive correlation between ϵ_{soc} and β ($r_\tau = .34$, $p < .001$, $BF = 61$), indicating that participants using relatively more social learning used less directed exploration and vice versa. In Exp. 2, we replicate both the relationship of ϵ_{soc} with reward ($r_\tau = -.20$, $p = .034$, $BF = 1.6$), and β ($r_\tau = .35$, $p < .001$, $BF > 100$).

Model-predicted behaviour

For reference, we simulated model behaviour using the priors given in the priors section.



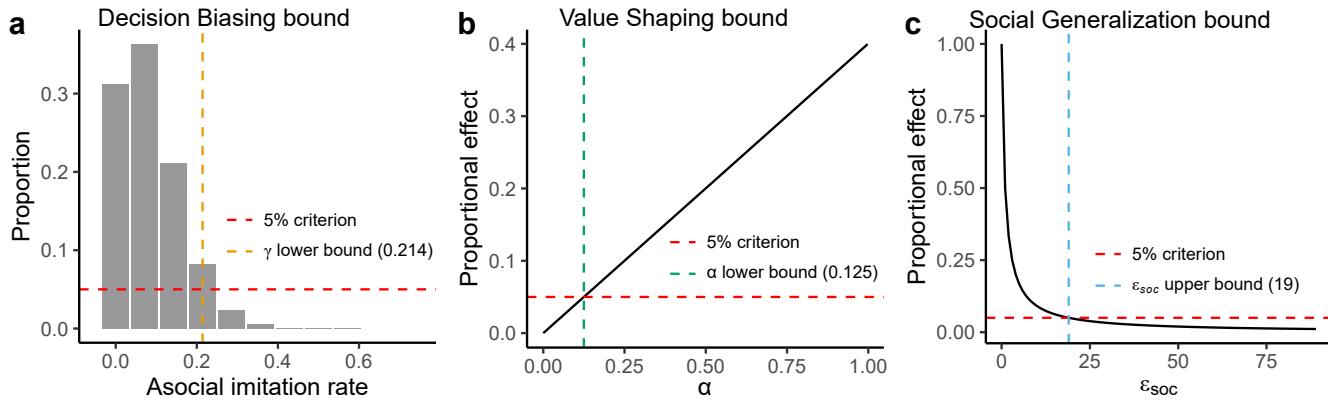


Fig. S8. Model bounding of the social models. Illustrations of the model bounding concepts, criteria, and parameter cut-offs for Decision Biasing (a), Value Shaping (b), and Social Generalization (c)

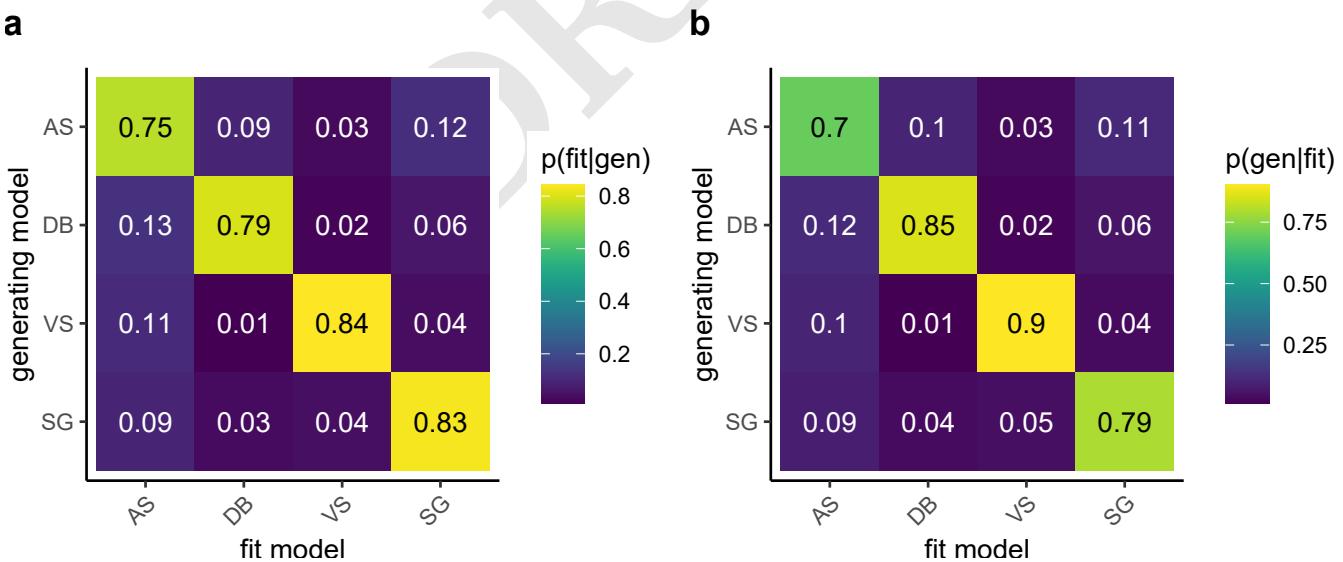


Fig. S9. Model recovery. a) Confusion matrix giving the conditional probability of a model being the best fit given the generating model. Values on the diagonal give the probability that the correct model is fit. b) Inversion matrix giving the conditional probability of a model being the generating model given it is fit. Values on the diagonal give the probability that the fit model is correct.

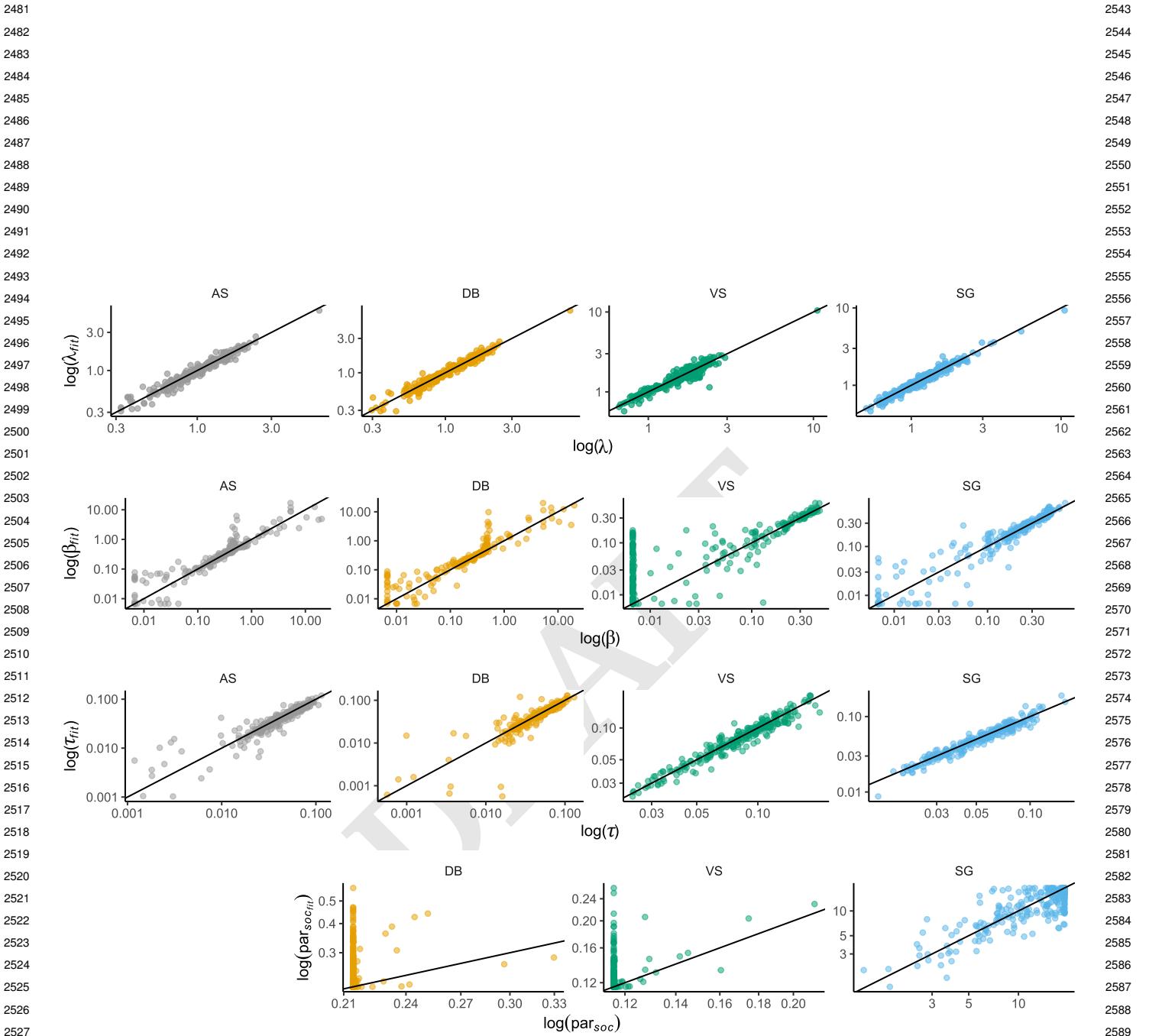


Fig. S10. Parameter recovery. Relationship between generating and recovered parameter by model (facet label).

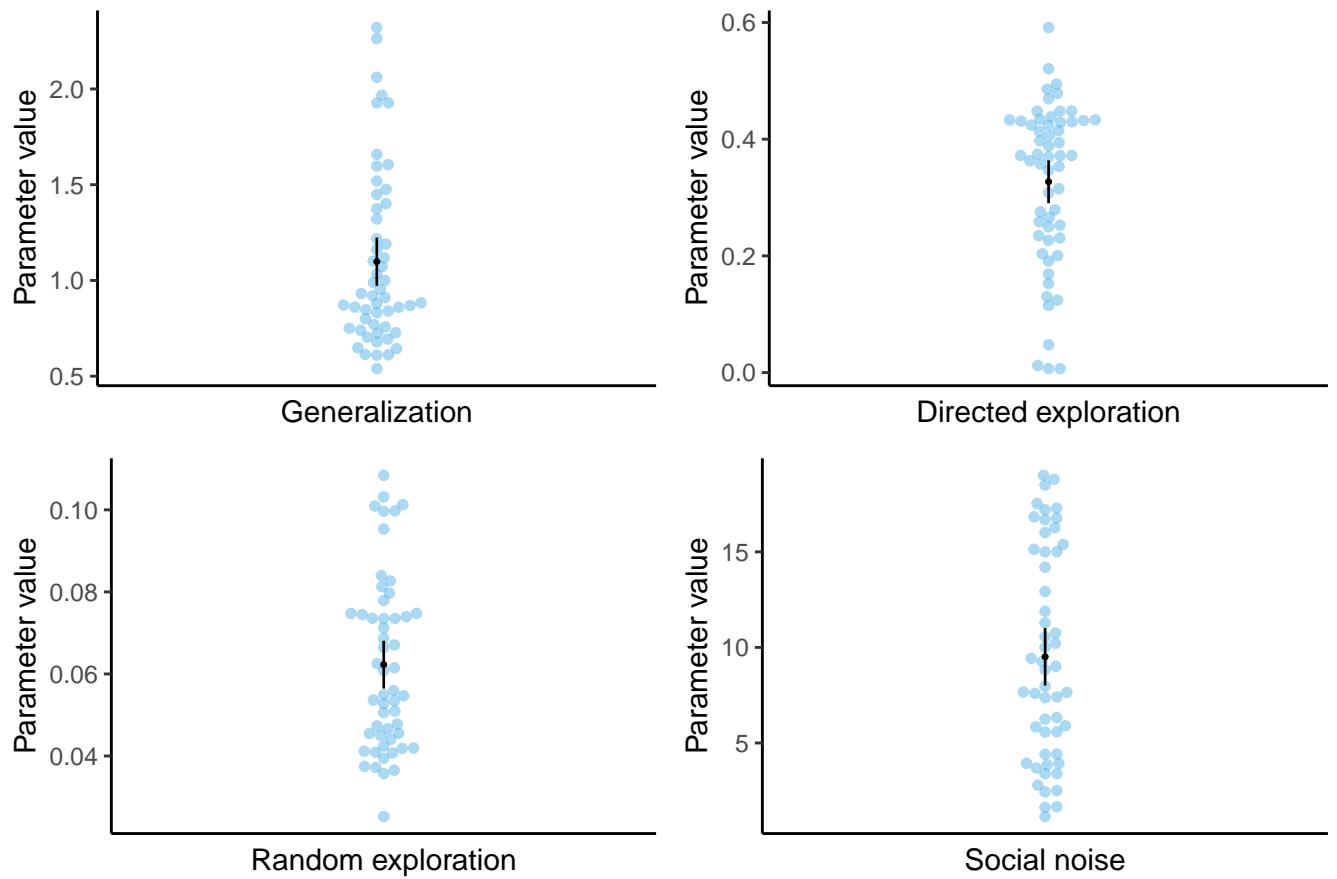
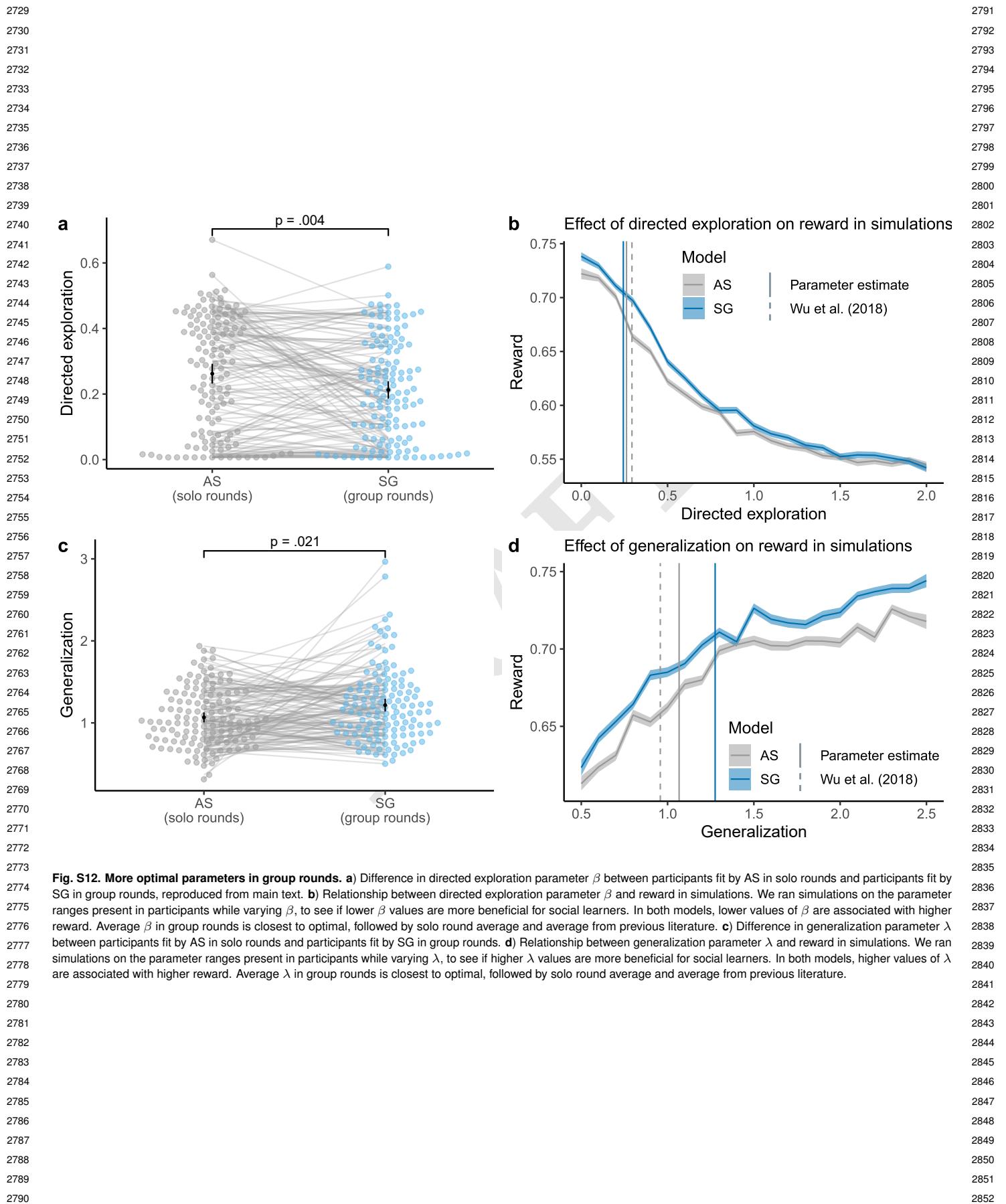


Fig. S11. Parameter fits for group rounds of Exp. 2. Only participants best fit by SG shown.

2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666

2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728



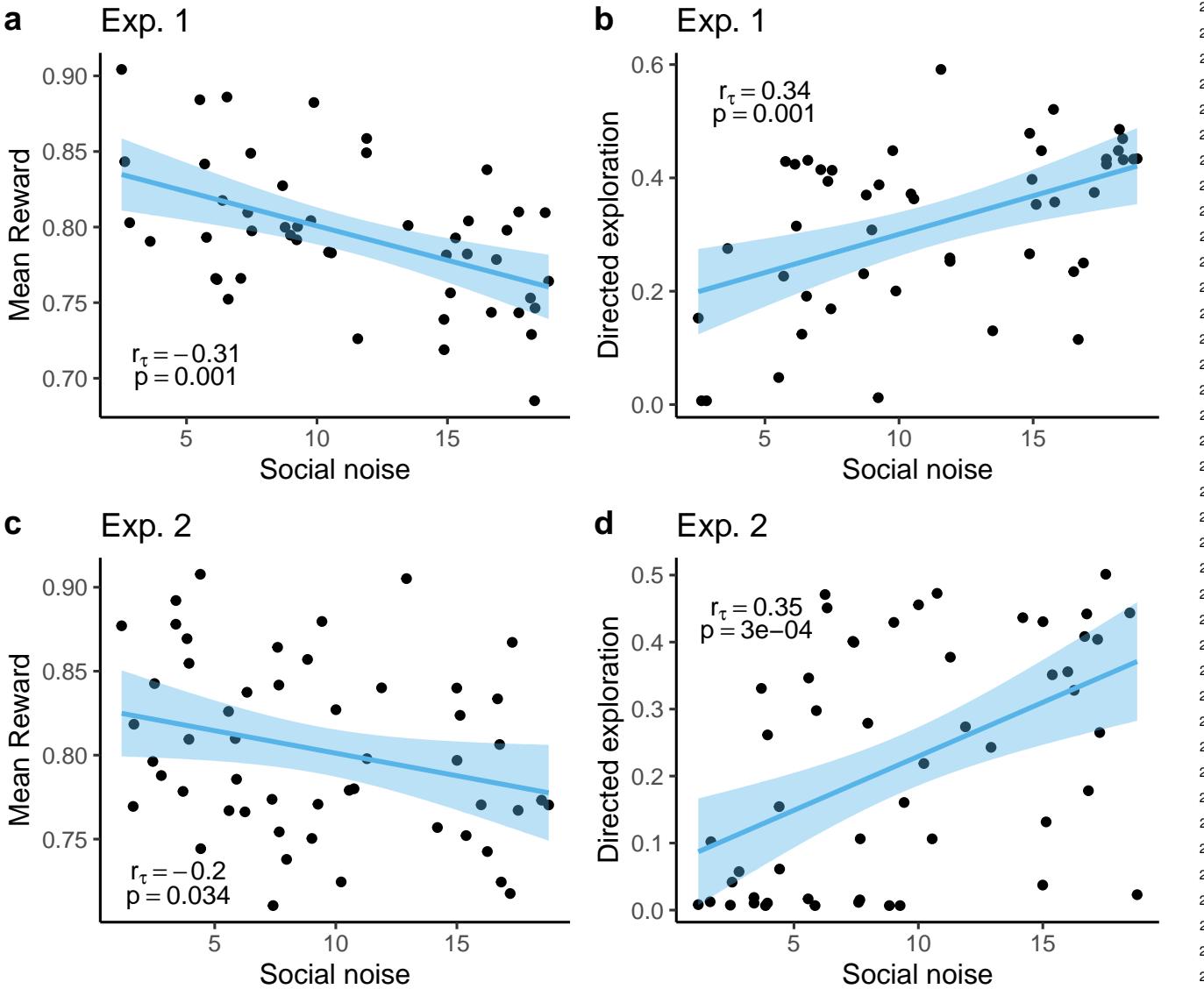


Fig. S13. Correlation analyses excluding ϵ_{soc} at the upper bounds. **a)** Relationship between ϵ_{soc} and mean reward in Exp. 1. **b)** Relationship between ϵ_{soc} and β in Exp. 1. **c)** Relationship between ϵ_{soc} and mean reward in Exp. 2. **d)** Relationship between ϵ_{soc} and β in Exp. 2.

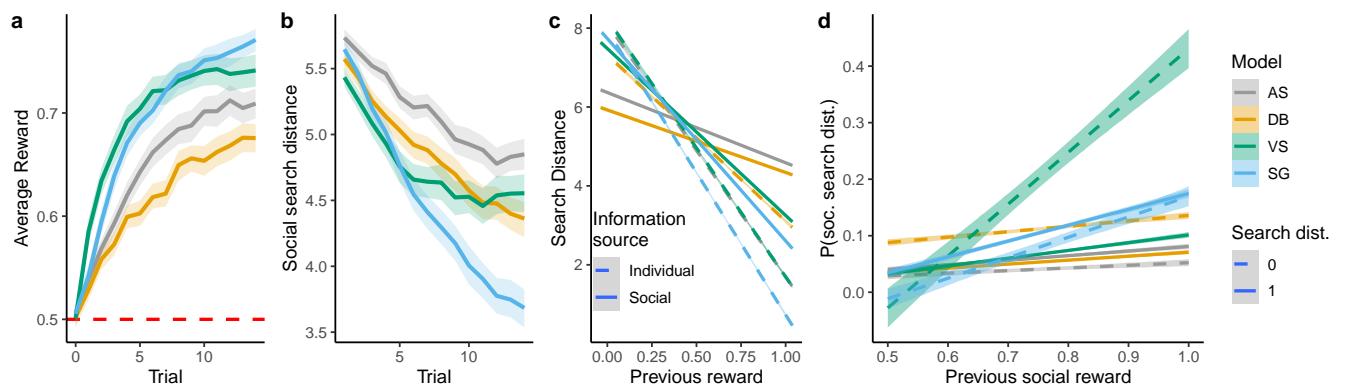


Fig. S14. Behaviour as predicted by the models. **a)** Learning curves. Red dashed line indicates chance-level performance. **b)** Social search distance over trials. **c)** Search distance over previous reward. Dashed lines indicate the relationship for individual previous reward and individual search distance; solid lines indicate the relationship for social previous reward and social search distance. **d)** Probability of imitation (social search distance = 0; dashed lines) and innovation (social search distance = 1; solid lines) over previous social reward.