

1 Dynamic integration of forward planning and heuristic
2 preferences during multiple goal pursuit

3

4 Florian Ott^{1*}, Dimitrije Marković¹, Alexander Strobel¹, Stefan J. Kiebel¹

5 ¹Department of Psychology, Technische Universität Dresden, Germany

6

7 * Corresponding author

8 E-Mail: florian.ott@tu-dresden.de (FO)

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27 **Abstract**

28 Selecting goals and successfully pursuing them in an uncertain and dynamic environment is an
29 important aspect of human behaviour. In order to decide which goal to pursue at what point in time,
30 one has to evaluate the consequences of one's actions over future time steps by forward planning.
31 However, when the goal is still temporally distant, detailed forward planning can be prohibitively
32 costly. One way to select actions at minimal computational costs is to use heuristics. It is an open
33 question how humans mix heuristics with forward planning to balance computational costs with goal
34 reaching performance. To test a hypothesis about dynamic mixing of heuristics with forward
35 planning, we used a novel stochastic sequential two-goal task. Comparing participants' decisions
36 with an optimal full planning agent, we found that at the early stages of goal-reaching sequences, in
37 which both goals are temporally distant and planning complexity is high, on average 42% (SD =
38 19%) of participants' choices deviated from the agent's optimal choices. Only towards the end of the
39 sequence, participant's behaviour converged to near optimal performance. Subsequent model-based
40 analyses showed that participants used heuristic preferences when the goal was temporally distant
41 and switched to forward planning when the goal was close.

42 **Author summary**

43 When we pursue our goals, there is often a moment when we recognize that we did not make the
44 progress that we hoped for. What should we do now? Persevere to achieve the original goal, or
45 switch to another goal? Two features of real-world goal pursuit make these decisions particularly
46 complex. First, goals can lie far into an unpredictable future and second, there are many potential
47 goals to pursue. When potential goals are temporally distant, human decision makers cannot use an
48 exhaustive planning strategy, rendering simpler rules of thumb more appropriate. An important
49 question is how humans adjust the rule of thumb approach once they get closer to the goal. We
50 addressed this question using a novel sequential two-goal task and analysed the choice data using a
51 computational model which arbitrates between a rule of thumb and accurate planning. We found that
52 participants' decision making progressively improved as the goal came closer and that this
53 improvement was most likely caused by participants starting to plan ahead.

54 **Introduction**

55 Decisions of which goal to pursue at what point in time are central to everyday life [1-3]. Typically,
56 in our dynamic environment, the outcomes of our decisions are stochastic and one cannot predict
57 with certainty whether a preferred goal can be reached. Often, our environment also presents
58 alternative goals that may be less preferred but can be reached with a higher probability than the
59 preferred goal. For example, when working towards a specific dream position in a career, it may turn
60 out after some time that the position is unlikely to be obtained, while another less preferred position
61 can be secured. The decision to make is whether one should continue working towards the preferred
62 position, or switch goals and secure the less preferred position. The risk when pursuing the preferred
63 position is to lose out on both positions. This decision dilemma 'should I risk it and go after a big
64 reward or play it safe and gain less?' is typical for many decisions we have to make in real life.
65 Critically, for many such decisions, these binary choices do not emerge suddenly and unexpectedly,
66 but the decision maker is typically confronted with such decisions after some prolonged period of
67 time working towards enabling different options.

68 How would one choose one's actions during such a prolonged goal-reaching decision making
69 sequence? One way, if the rules of the dynamic environment and its uncertainties are known, is to use

70 forward planning to always choose the actions which maximize the gain (see [4, 5] reviewing
71 cognitive processes of forward planning). This would be the way one would program an optimal
72 agent in a game or experimental task environment. This approach is often used in cognitive
73 neuroscience to model the mechanism of how humans make decisions in temporally extended goal-
74 reaching scenarios, (e.g. [6-9]).

75 However, the implicit assumption made in these decision-making models, namely that humans use
76 detailed forward planning and compute the probabilities of reaching the goals, is difficult to justify,
77 because of the involved computational complexity. In a stochastic environment, forward planning in
78 artificial agents is typically achieved via sampling many possible policies (sequences of actions)
79 which requires substantial computing power that scales exponentially with the number of future
80 actions. In particular, when one is still temporally far from the goal, the computational burden of
81 simulating trajectories into the future is the largest, while the usefulness of the resulting action
82 selection is minimal: intuitively, in stochastic and sufficiently complex environments, anything may
83 yet happen on the long way to the goal so the gain of planning ahead at high cost may be small. The
84 importance of the balance between the benefits and its costs to better understand human decision
85 making became a recent research focus, e.g., [10-14]. The question is how one can select actions over
86 long stretches of time, without being exposed to the computational burden of forward planning or
87 similar dynamic programming schemes.

88 One obvious way to select actions at minimal computational costs is to use heuristics that do not
89 require forward planning towards a goal [15, 16], e.g. to always select the action towards a hard to
90 achieve and highly rewarded goal. Clearly, this and other heuristics come with the drawback that they
91 can be substantially suboptimal when close to the goal. For example, blindly working toward a hard
92 to achieve goal would ignore the risk of not reaching any goal. Another solution is to use habit-like
93 strategies to avoid computational costs [17]. However, habits are typically useful only when one
94 encounters exactly the same situation or context repeatedly, while goal reaching in uncertain
95 environments as presented here, often requires flexible behavioural control.

96 It is an open question how humans select their actions when the potentially reachable goals are still
97 far away and forward planning is complex. We hypothesized that people use a mixture of two
98 approaches to achieve an acceptable balance between outcome and computational costs. This mixture
99 changes with temporal distance to the goal: when far from the goal, people use a prior goal
100 preference to make their decision about which action to take. With this approach, one assumes that
101 one will eventually reach the preferred goal and selects the action that, if one looked backward in
102 time from the reached goal, is the most instrumental. When coming closer to the goal, one expects
103 that the influence of the goal preference should be progressively superseded by computationally more
104 expensive action selection using forward planning to optimally reach the preferred goal or, failing
105 that one, to pursue policies to reach an alternative goal.

106 To test whether participants used such an approach, we employed a novel behavioural task where
107 participants were placed in a dynamic and stochastic sequential decision task environment that
108 emulated reaching goals over an extended time period. In miniblocks of 15 trials, participants had to
109 make decisions to reach one or two goals, where reaching both goals was rewarded more than
110 reaching only one. In each miniblock, it was also possible, if blindly trying to obtain the higher
111 reward, to not reach any goal and not obtain any reward. While participants pass through the
112 miniblock, both the remaining trials to the end of the miniblock and the complexity of forward
113 planning decrease. This enables us to test and model whether participants switch from using
114 heuristics to forward planning during goal-reaching. To analyse the behavioural data of 89

115 participants and test hypotheses, we used stochastic variational inference, which provided posterior
116 beliefs about the goal strategy preference of each participant, among other free model parameters.
117 We show that the heuristic goal strategy preference parameter is key to explain participants' choices
118 when temporally distant from the goal, and how, when progressing towards a goal, this goal strategy
119 preference interacts with optimal forward planning to achieve near-optimal performance.

120 **Methods**

121 **Participants**

122 Eighty-nine participants took part in the experiment (58 women, mean age = 24.8, SD = 7.1).
123 Reimbursement was a fixed amount of 8€ or class credit plus a performance-dependent bonus (mean
124 bonus = 3.88€, SD = 13.6). The study was approved by the Institutional Review Board of the
125 Technische Universität Dresden and conducted in accordance to ethical standards of the Declaration
126 of Helsinki. All participants were informed about the purpose and the procedure of the study and
127 gave written informed consent prior to the experiment. All participants had normal or corrected-to-
128 normal vision.

129 **Table 1. Glossary of abbreviations**

Abbreviation	Explanation
A, B	Basic offers
Ab, aB	Mixed offers
Pts_t^A	A-points in trial t
Pts_t^B	B-points in trial t
g1	One-goal-choice = Sequential strategy choice = Choice that maximizes point difference
g2	Two-goal-choice = Parallel strategy choice = Choice that minimizes point difference
G1	One-goal-success = One point scale above threshold after 15 trials
G2	Two-goal-success = Both scales above threshold after 15 trials
Q(s,a)	Action value = Expected future reward of a choice
$Q_G(s,a)$	Goal choice value = Expected future reward of a goal strategy choice
DEV	Differential expected value = $Q_G(s, g2) - Q_G(s, g1)$

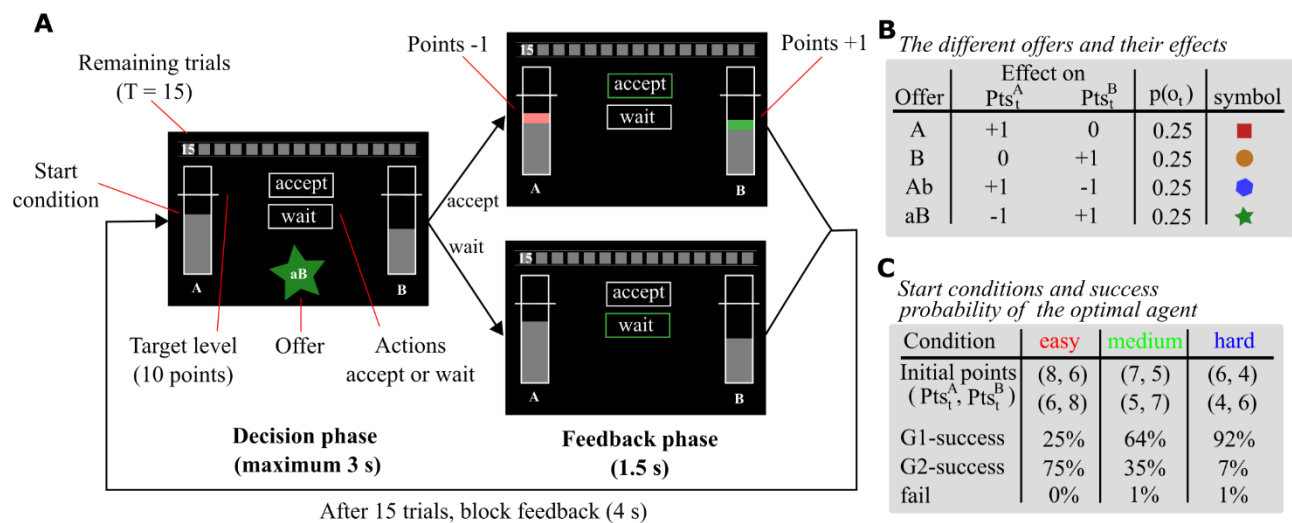
130 **Experimental Task**

131 The experiment included a training phase of 10 miniblocks, followed by the main experiment
132 comprising 60 miniblocks. The 60 miniblocks in the main experiment were subdivided into three
133 sessions of 20 miniblocks between which participants could make a self-determined pause. A
134 miniblock consisted of $T = 15$ trials in which participants had to accept or reject presented offers to
135 collect A-points (Pts_t^A) and B-points (Pts_t^B , see Table 1 for a glossary of abbreviations). If
136 participants reached the threshold of 10 points for either A- or B-point scale after 15 trials, they
137 received a reward of 5 cents. If participants reached the threshold for both point scales, they received
138 a reward of 10 cents. If none of the two thresholds was reached, no additional reward was provided.
139 In total, each participant completed 150 training trials and 900 trials in the main experiment.

140 Each trial started with a response phase lasting until a response was made, but not more than 3 s (Fig
141 1, A). The current amount of A-points and B-points was visualized by two vertical bars flanking the
142 stimulus display. Horizontal white lines marked the threshold of 10 points. At the top of the screen, a
143 grey timeline informed the participants about the remaining trials in the miniblock. The current offer
144 was displayed at the bottom centre, and the two choice options were presented in the centre of the
145 screen by the framed words ‘accept’ and ‘wait’. Participants could accept an offer by an upwards
146 keypress and reject the offer by a downwards keypress. If participants did not respond within 3 s the
147 trial was aborted, and a message was displayed reminding the participant to pay attention. If
148 participants missed the response deadline more than 5 times in the whole main experiment, 50 cents
149 were subtracted from their final payoff (mean number of timeouts = 1.34, SD = 1.7). After the
150 response phase, feedback was displayed for 1.5 s. Response feedback included a change in colour of
151 the frame around the selected response from white to green. Additionally, the gain or loss of points
152 was visualized by colouring the respective area on the bar either green or red. After 15 trials,
153 feedback for the miniblock was displayed for 4 s informing the participants whether they won 5, 10
154 or 0 cents. Code for experimental control and stimulus presentation was custom written in Matlab
155 (MathWorks) with extensions from the Psychophysics toolbox [18].

156 Participants were presented with four different offers (A, B, Ab, and aB) that occurred with equal
157 probability on each trial of the miniblock (see Fig 1, B). We call A or B basic offers and Ab or aB
158 mixed offers. Accepting basic offers increased the corresponding point count, whereas accepting
159 mixed offers transferred a single point from one scale to the other. The basic offers introduce a
160 stochastic base rate of points, which allows participants to accumulate enough points on one or both
161 point scales. In contrast, mixed offers allow us to identify participants’ intention to reach a state in
162 which either both point scales are above threshold ($Pts_T^A \geq 10$ and $Pts_T^B \geq 10$) or only one point
163 scale is above threshold (e.g. $Pts_T^A < 10$ and $Pts_T^B \geq 10$; see below for more details). Rejecting an
164 offer did not have any effect on the current point count. All participants received the same sequence
165 of offers. We generated pseudorandomized lists for the training phase and for the three main
166 experimental phases such that the frequency of offers reflected an equal offer occurrence probability
167 in every list. We associated each offer with a coloured symbol to facilitate fast recognition.

168 Three different conditions modulated the difficulty to reach both thresholds by varying the number of
169 initial points (Fig 1, C). We chose the number of initial points such that an optimal agent’s
170 probability of reaching both thresholds was 75% in easy, 35% in medium and 7% in hard. The
171 agent’s goal reaching performance for each initial point configuration was based on 10,000 simulated
172 miniblocks with uniform offer probability (see below how we define the optimal agent). The same
173 sequence of start conditions was presented to all participants. Pseudorandomized lists with a balanced
174 frequency of initial point configurations were generated for the training phase and for the three main
175 experimental phases. Note that the observed agent behaviour in the results section deviates from what
176 we expected based on the experimental parametrization process. These discrepancies arise because
177 we used random offer sequences (offers with equal probability) for experimental parametrization, but
178 one specific offer sequence for the actual experiment. For example, in some miniblocks there were
179 only few basic offers (see S1-4 Fig for details about the used offer sequence).



180

181 **Fig 1. Experimental task.** (A) Depiction of trial timeline and stimulus features. Participants
 182 performed miniblocks of 15 trials in which they collected points to reach either one or two goals,
 183 rewarding them with additional 5 or 10 Cents. Each trial started with a decision phase (maximum 3s)
 184 in which participants had to accept or reject a presented offer. Depending on the offer, accepting
 185 increased or decreased A- and B-points. The current amount of points was displayed by two grey
 186 bars flanking the stimulus screen. In the feedback phase (1.5s), gained points were displayed as a
 187 green area and lost points as a red area on the bar. The horizontal lines crossing the bars indicated the
 188 threshold for reaching goal A and goal B. After 15 trials, feedback for the miniblock was displayed
 189 (4s) informing the participant about the reward gained. (B) Summary of offer types and their effect
 190 on point count. Offers occurred with equal probability in each trial of the miniblock. Basic offers (A
 191 and B) increased either A or B points. Mixed offers (Ab and aB) added one point on one side but
 192 subtracted one point on the other side. Only accepting an offer had an effect on points. (C) Three
 193 different conditions modulated the difficulty to reach both thresholds by varying the number of initial
 194 points. Using an optimal agent, we chose the number of initial points, such that the agent's
 195 probability of reaching both thresholds (G2-success) was 75% in easy, 35% in medium and 7% in
 196 hard.

197 Choice classification

198 In order to maximize reward, it was key for the participants to decide whether they should pursue the
 199 A- and B-goal in a sequential or in a parallel manner. A parallel strategy, i.e. balancing the two point
 200 scales, increases the likelihood that both goals (G2, see Table 1) will be reached at the end of the
 201 miniblock, but at the risk of failing. A sequential strategy, i.e. first secure one goal, then focus on the
 202 second one, might increase the likelihood to reach at least one goal (G1) within 15 trials, but
 203 decreases the likelihood to achieve G2.

204 To obtain a trial-wise measure of the pursued goal strategy, choices were classified based on the
 205 current point difference and the offer. Choices that minimized the difference between points were
 206 classified as two-goal-choice ($a_t = g2$), reflecting the intention to fill both bars using a parallel
 207 strategy. Choices that maximized the difference between points were classified as one-goal-choice
 208 ($a_t = g1$), reflecting the intention to pursue G1, or the intention to maintain one bar above threshold
 209 if G1-success has already been attained (see S1 Table). For example, if a participant has 8 A-points
 210 and 6 B-points and the current offer is Ab, accepting would be a g1-choice, whereas waiting would

211 be a g2-choice. Conversely, for an aB offer, accepting would be a g2-choice and waiting a g1-choice.
 212 If the difference between points ($Pts_t^A - Pts_t^B$) is 1 and the offer is aB, g-choice is not defined
 213 because the absolute point difference would not be changed. This also applies to the mirrored case,
 214 where the difference between points ($Pts_t^A - Pts_t^B$) is -1 and the offer is Ab. Note that, due to the
 215 experimental design, response (accept/wait) and g-choice (g2/g1) were weakly correlated ($r = 0.21$).
 216 Furthermore, g-choice classification is only defined for the mixed offers (Ab and aB). The basic
 217 offers (A and B) are not informative with respect to the participants' pursued goal strategy.
 218 Importantly, all trial-level analysis will be restricted to trials which can be related to g-choices.

219 Task model

220 Here we will formulate the task in an explicit mathematical form, which will help us clarify what
 221 implicit assumptions we make in the behavioural model [19]. We define a miniblock of the two-goal
 222 task as a tuple

$$(T, S, O, R, A, p(s_{t+1}|s_t, o_t, a_t), p(o_t), p(r_t|s_t)) \quad (1)$$

223 where

- 224 • $T = 15$ denotes the number of trials in a miniblock, hence $t = 1, \dots, 15$.
- 225 • $S = \{0, \dots, 20\}^2$ denotes the set of task states, corresponding to the point scale of the two
 226 point types (A, and B). Hence, a state s_t in trial t is defined as a tuple consisting of point
 227 counts along the two scales, $s_t = (Pts_t^A, Pts_t^B)$.
- 228 • $O = \{A, B, Ab, Ba\}$ denotes the set of four offer types, where the upper case letters denote an
 229 increase in points of a specific type and the lower case letters subtraction of points.
- 230 • $R = \{R_0, R_L, R_H\} = (0, 5, 10)$ denotes the set of rewards.
- 231 • $A = \{0, 1\}$ denotes the set of choices, where 0 corresponds to rejecting an offer and 1 to
 232 accepting an offer.
- 233 • $p(s_{t+1}|s_t, o_t, a_t)$ denotes state transitions which are implemented in a deterministic manner
 234 as $s_{t+1} = s_t + a_t * m(o_t)$, where $m(o_t)$ maps offer types into the point changes on the two
 235 point scales.
- 236 • $p(o_t = i) = \frac{1}{4}$ (for $\forall i \in O$) denotes a uniform distribution from which the offers are
 237 sampled.
- 238 • $p(r_t|s_t)$ denotes the state and trial dependent reward distribution defined as

$$\begin{aligned} p(r_t = R_0|s_t) &= 1, \text{ for } \forall t < T \\ p(r_T = R_L|Pts_T^A \geq 10 \oplus Pts_T^B \geq 10) &= 1 \\ p(r_T = R_H|Pts_T^A \geq 10 \wedge Pts_T^B \geq 10) &= 1 \end{aligned}$$

239 Note that in the experiment the participants are exposed to a pseudo-random sequence of offers,
 240 meaning that within one experimental block all participants observed the same sequence of offers
 241 pre-sampled from this uniform distribution (see S1-4 Fig. for additional information about the used
 242 offer sequence). For simulations and parameter estimates we use the same pseudo-random sequence
 243 of observations, hence in each trial t of a specific block b offers are selected from a predefined
 244 sequence $o_{1:T}^{1:B} = (o_1^1, \dots, o_T^1, \dots, o_1^B, \dots, o_T^B)$, initially generated from a uniform distribution.

245 Behavioural model

246 To build a behavioural model, we assume that participants have learned the task representation
 247 through the training session and initial instruction. Hence, the behavioural model is represented by
 248 the following tuple

$$(T, S, O, R_\kappa, A, p(s_{t+1}|s_t, o_t, a_t), p(o_t), p(r_t|s_t)) \quad (2)$$

249 where

- 250 • $T, S, O, A, p(s_{t+1}|s_t, o_t, a_t), p(o_t), p(r_t|s_t)$ are defined the same way as in the task model.
- 251 • $R_\kappa = \{0, 5, 10 \cdot \kappa\}$ denotes an agent-specific valuation of the rewarding states. Although the
 252 instructions for the experimental task clearly explained that participants receive a specific
 253 monetary reward depending on the final state reached during a miniblock, we considered a
 254 potential biased estimate of the ratio between G2 and G1 monetary rewards, quantified with
 255 the free model parameter $\kappa \in [0, 2]$. In other words, we assumed that the participants might
 256 overestimate or underestimate the value of a G2-success, relative to a G1-success.

257 Importantly, the process of action selection corresponds to following a behavioural policy that
 258 maximises expected value during a single miniblock. We classified as G2-success miniblocks in
 259 which both point scales were above threshold after the final trial ($Pts_T^A \geq 10$ and $Pts_T^B \geq 10$). We
 260 classified as G1-success miniblocks in which only one point scale was above threshold (e.g. $Pts_T^A <$
 261 10 or $Pts_T^B \geq 10$).

262 In what follows we derive the process of estimating choice values and subsequent choices based on
 263 dynamic programming applied to a finite horizon Markov decision process ([20]; for experimental
 264 studies see also [9, 21]).

265 **Forward Planning**

266 We start with a typical assumption used in reinforcement learning, namely that participants choose
 267 actions with the goal to maximize future reward. Starting from some state s_t at trial t , offer o_t , and
 268 following a behavioural policy π we define an expected future reward as

$$V[s_t, o_t|\pi] = \sum_{k=t+1}^T \gamma^{k-t-1} E[r_k|s_t, o_t, \pi] \quad (3)$$

269 where γ denotes a discount rate and $E[r_k|s_t, o_t, \pi]$ denotes expected reward at some future time step
 270 k . The behavioural policy sets the state-action probability $\pi(a_t, \dots, a_T|s_t, \dots, s_{T-1})$ over the current
 271 and future trials. Hence, we can obtain the expected reward as

$$E[r_k|s_t, \pi] = \sum_{r_k} r_k p(r_k|s_t, \pi) \quad (4)$$

272 where

$$p(r_k|s_t, \pi) = \sum_{s_{t+1:k}} \sum_{a_{t:k-1}} p(r_k|s_k) \prod_{\tau=t+1}^k p(s_\tau|s_{\tau-1}, o_{\tau-1}, a_{\tau-1}) p(o_{\tau-1}) \pi(a_{\tau-1}|s_{\tau-1}) \quad (5)$$

273 Note that we use $s_{t+1:k}$, and $a_{t:k-1}$ to denote a tuple of sequential variables, hence $x_{m:n} =$
 274 (x_m, \dots, x_n) . The key step in deriving the behavioural model was to find the policy which maximises
 275 the expected future reward, that is, the expected state-offer value. In practice, one obtains the optimal
 276 policy as

$$\pi^* = \underset{\pi}{\operatorname{argmax}} V[s_t, o_t | \pi] \quad (6)$$

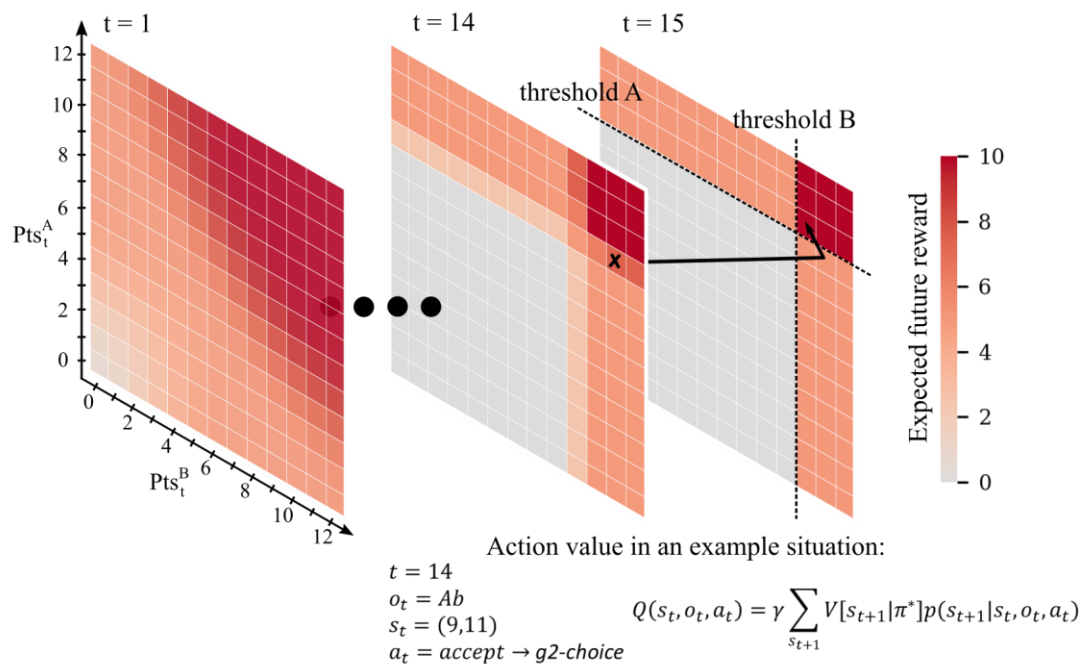
277 We solve the above optimization problem using the backward induction method of dynamic
 278 programming. The backward induction algorithm is defined in the following iterative steps:

- 279 (i) set the value of final state s_T as the reward obtained in that state $V[s_T | \pi^*] =$
 280 $\sum_{r_T \in R_\kappa} r_T p(r_T | s_T)$
- 281 (ii) compute state-offer-action value as $Q(s_k, o_k, a_k) = \gamma \sum_{s_{k+1}} V[s_{k+1} | \pi^*] p(s_{k+1} | s_k, o_k, a_k)$
- 282 (iii) set optimal choice for given state-offer pair as $a_k^* = \operatorname{argmax}_a Q(s_k, o_k, a)$
- 283 (iv) define the expected value of state s_k under optimal policy π^* as
 284 $V[s_k | \pi^*] = \sum_{o_k} Q(s_k, o_k, a_k^*) p(o_k)$
- 285 (v) repeat steps (ii) – (iv) until $k = t$

286 Hence, for a fixed value of the reward ratio (κ) an optimal choice at trial t corresponds to

$$a_t^* = \underset{a}{\operatorname{argmax}} Q(s_t, o_t, a) \quad (7)$$

287 We will define the optimal agent as an agent who has a correct representation of the reward ratio
 288 ($\kappa = 1$) and does not discount future reward ($\gamma = 1$). We illustrate in Fig 2 the Q-value to accept,
 289 estimated for the case of the optimal agent in an example trial ($Pts_t^A = 8$, $Pts_t^B = 11$, $o_t = Ab$).



290

291 **Fig 2. Illustration of the state space and associated expected future reward for the optimal**
 292 **agent ($\gamma = 1$, $\kappa = 1$).** The black arrow shows a hypothetical transition in the state space. In trial 14
 293 the participant has 9 A-points and 11 B-points (marked by the black cross) and accepts an offer Ab,
 294 gaining one A-point and losing one B-point (g2-choice). In the resulting state, both thresholds are
 295 reached; thus, the value of that state is 10 Cents. Similarly, the action that leads to that state has an
 296 associated Q-value of 10 Cents. In this example the agent would just have to wait in the last trial (15)
 297 to gain a 10 cents reward.

298 Response likelihood

299 Participants might compute expected values by mentally simulating and comparing sequences of
 300 actions towards the end of the miniblock. To illustrate the benefits of planning we consider the
 301 following example: There are 3 trials left in the current miniblock, and the participant has 9 A-points
 302 and 9 B-points (10 is threshold), and she receives offer Ab. Planning would, for example, allow to
 303 compute the probabilities for G2 when choosing either wait or accept. By waiting the participant
 304 would enter the second last trial with 9 A-points and 9 B-points. Receiving offer A or B in the
 305 second last trial (0.5 probability) followed by the complementary offer A or B in the last trial (0.25
 306 probability) would grant G2. When choosing accept, the participant will have in the second last trial
 307 10 A-points and 8 B-points. Consequently, she would need two consecutive B-offers (0.25 * 0.25
 308 probability) to achieve G2. Hence, by planning ahead one would conclude that wait gives the highest
 309 probability for a G2-success.

310 Still, planning an arbitrary number of future steps is complex and unrealistic. Hence, we make an
 311 assumption that the process of optimal action selection described above is perturbed by noise
 312 (planning noise, and response noise) which we quantify in the form of a parameter β , denoting
 313 response precision. Hence, this precision parameter is critical to characterize the participants'
 314 reliance on forward planning. Since the difference in expected future rewards of a g1- or g2- choice
 315 is high when the goal is close (S5 Fig), β is able to selectively capture g-choice performance at the
 316 end of the miniblock. Furthermore, instead of an elaborate planning process participants might use a
 317 simpler heuristic when deciding which action to select. We capture this heuristic in form of an

318 additional offer-state-action function $h(o_t, s_t, a_t, \theta)$ which evaluates choices relative to possible
 319 goals. We describe this heuristic evaluation below. Overall, we can express the response likelihood
 320 (the probability that a participant makes choice a_t) as

$$p(a_t | \beta, \theta, \gamma, \kappa) = s(\beta Q(o_t, s_t, a_t, \gamma, \kappa) + h(o_t, s_t, a_t, \theta)) \quad (8)$$

321 where $s(x)$ denotes the softmax function.

322 **Choice heuristic**

323 The choice heuristic is defined relative to the current offer o_t , current state s_t , and possible choices
 324 a_t . Importantly, we will interpret the choice heuristic in terms of participants' biases towards
 325 approaching both goals in a sequential or parallel manner. Hence, it is more intuitive to define the
 326 choice heuristic as choice biases relative to the goals, and not accept-reject choices. The choice
 327 heuristic is defined as follows

$$h(o_t, s_t, a_t, \theta) = \begin{cases} \infty, & \text{for } o_t \in \{A, B\}, \text{ and } a_t = 1 \\ \theta, & \text{for } o_t \in \{Ab, Ba\}, \text{ and } a_t \equiv g2 \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

328 where $a_t \equiv g2$ denotes choices (accept or reject) which can be classified as g2-choices (see
 329 subsection Choice classification for details). In summary, a choice which reduces the point difference
 330 ($Pts_t^A - Pts_t^B$), for the given offer and the current state, is classified as g2-choice and choice which
 331 increases the point difference as g1-choice. Essentially, the strategy preference parameter θ reflects
 332 participants' preference for pursuing a sequential (negative values) or parallel (positive values)
 333 strategy. For example, some participants might have a general tendency to pursue goals in a parallel
 334 manner, independent of the actual Q -values. Conversely, participants may prefer a more cautious
 335 sequential approach. Note that we expected this parameter to make the most significant contribution
 336 to participants' deviation from optimal behaviour, reflecting their reliance on decision heuristics early
 337 in the miniblock.

338 Finally, for those choices which can be classified as g2- or g1-choices, we can express the response
 339 likelihood in a simplified form, in terms of free model parameters $\beta, \theta, \gamma, \kappa$ (Table 2). We refer to the
 340 difference between Q -values for g-choice as the differential expected value (DEV),

$$DEV = Q_G(a_t = g2) - Q_G(a_t = g1) \quad (10)$$

341 Using DEV , we defined the probability of making a g2-choice as

$$p(g2) = \sigma(\beta \cdot DEV(\gamma, \kappa) + \theta) \quad (11)$$

342 where $\sigma(x) = \frac{1}{1+e^{-x}}$ denotes the logistic function. Note that the probability of g1-choice becomes
 343 $p(g1) = 1 - p(g2)$.

344 Table 2. Summary of four free model parameters, the variables, the transformations used to map
 345 values to unconstrained space and their function in modelling participant behaviour.

Name	Variable	Transform	Function
Precision	β	$x_1 = \ln \beta$	Captures the impact of <i>DEV</i> , derived by forward planning, on action selection
Strategy preference	θ	$x_2 = \theta$	Heuristic preference of pursuing a parallel ($\theta > 0$) or sequential ($\theta < 0$) strategy, independent of the actual <i>DEV</i>
Discount rate	γ	$x_3 = \ln \frac{\gamma}{1 - \gamma}$	Temporal discounting of <i>DEV</i> by the factor γ^{T-t} , where $T - t$ is the number of remaining trials
Reward ratio	κ	$x_4 = \ln \frac{\kappa}{2 - \kappa}$	Accounts for the possibility that participants may overweight ($\kappa > 1$) or underweight ($\kappa < 1$) the actual reward for G2-success relative to G1-success.

346

347 **Optimal agent comparison and general data analysis**

348 We compared participant behaviour with simulated behaviour of an optimal agent. To summarize, we
 349 denote the optimal agent as the agent which has a correct representation of the reward function
 350 ($\kappa = 1$), does not discount future rewards ($\gamma = 1$), is not biased in favour of any choice ($\theta = 0$),
 351 and who generates deterministic g-choices based on *DEV*-values (corresponding to $\beta \rightarrow \infty$ in the
 352 response likelihood, that is, the argmax operator). The optimal agent deterministically accepts A and
 353 B offers.

354 When simulating agent behaviour to evaluate successful goal reaching, the agent received the same
 355 sequence of offers and initial conditions as the participants. Analysis on the level of g-choices was
 356 performed by registering instances in which the g-choice of a participant differed from the g-choice
 357 the optimal agent would have made in the same context (Pts_t^A, Pts_t^B, o_t, t). Trials with A or B offers
 358 and trials in which G2 had already been reached, were excluded from the g-choice analysis.

359 The goal of this comparison between summary measures of both optimal agent and participants was
 360 two-fold: First, we used this comparison to visualize deviations from optimality and motivate the
 361 model-based analysis which was used to test the hypothesis that a shift from heuristics to forward
 362 planning may explain these deviations. Second, plotting suboptimal g-choices instead of g-choices
 363 (Fig. 4) makes behaviour between participants more comparable. Plotting the proportion of g-choices
 364 averaged across participants would have been mostly uninformative because the significance of a g-
 365 choice depends on the current state, which is a consequence of the individual history of past choices
 366 within a miniblock. By registering deviations from an optimal reference point, we circumvent this
 367 state dependence of g-choices.

368 We used a sign test as implemented in the “sign_test” function of python’s “Statsmodels”[22]
 369 package to test whether participants total reward and success rates differed significantly from the
 370 optimal agent’s deterministic performance. We reported the p-value and the m-value $m = (N(+) -$
 371 $N(-))/2$, where $N(+)$ is the number of values above 0 and $N(-)$ is the number of values below
 372 and. To test for learning effects (in the main experimental phase), we used mixed effects models as

373 implemented in R [23] with the “lm4” package [24]. Intercepts and slopes were allowed to vary
 374 between participants. p-values were obtained using the “lmerTest” package [25].

375 Hierarchical Bayesian data analysis

376 To estimate the free model parameters (Table 2) that best match the behaviour of each participant, we
 377 applied an approximate probabilistic inference scheme over a hierarchical parametric model, so-
 378 called stochastic variational inference (SVI) [26].

379 As a first step, we define a generic (weakly informative) hierarchical prior over unconstrained space
 380 of model parameters. In Table 2 we summarize the roles of free model parameters of our behavioural
 381 model and the corresponding transforms that we used to map parameters into an unconstrained space.
 382 We use \mathbf{x}^n to denote a vector of free and unconstrained model parameters corresponding to the n th
 383 participant. Similarly, $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ will denote hyperpriors over group mean and variance for each free
 384 model parameter. We can express the hierarchical prior in the following form

$$\mu_i \sim N(m_i, s_i) \quad (12)$$

$$\sigma_i \sim C^+(0, 1) \quad (13)$$

$$x_i^n \sim N(\mu_i, \lambda \sigma_i) \quad (14)$$

$$\text{for } i \in [1, \dots, d], \text{ and } n \in [1, \dots, N] \quad (15)$$

385 where $C^+(0,1)$ denotes a Half-Cauchy prior with scale $s = 1$, d number of parameters, and N
 386 number of participants. Note that by using this form of a hierarchical prior we make an explicit
 387 assumption that parameters defining the behaviour of each participant are centred on the same mean
 388 and share the same prior uncertainty. Hence, both the prior mean and uncertainty for each parameter
 389 are defined at the group level. Furthermore, the hyper-parameters of the prior
 390 $\boldsymbol{\eta} = (m_1, \dots, m_d, s_1, \dots, s_d, \lambda)$ are also estimated from the data (Empirical Bayes procedure) in parallel
 391 to the posterior estimates of latent variables $\boldsymbol{\theta} = (\mu_1, \dots, \mu_d, \sigma_1, \dots, \sigma_d, \mathbf{x}^1, \dots, \mathbf{x}^N)$. For more details,
 392 see supporting information (S1 Notebook).

393 The behavioural model introduced above defines the response likelihood, that is, the probability of
 394 observing measured responses when sampling responses from the model, condition on the set of
 395 model parameters $(\mathbf{x}^1, \dots, \mathbf{x}^N)$. The response likelihood can be simply expressed as a product of
 396 response probabilities over all measured responses $A = (\mathbf{a}^1, \dots, \mathbf{a}^N)$, presented offers $O =$
 397 $(\mathbf{o}^1, \dots, \mathbf{o}^N)$, and states (point configurations) visited by each participant $S = (\mathbf{s}^1, \dots, \mathbf{s}^N)$ over the
 398 whole experiment

$$p(A|O, S, \mathbf{x}^1, \dots, \mathbf{x}^N) = \prod_{n=1}^N \prod_{b=1}^M \prod_{t=1}^T p(a_{b,t}^n | s_{b,t}^n, o_{b,t}^n, \mathbf{x}^n) \quad (16)$$

399 where b denotes experimental block, and t a specific trial within the block.

400 To estimate the posterior distribution (per participant) over free model parameters, we applied the
 401 following approximation to the true posterior

$$p(\mathbf{x}^1, \dots, \mathbf{x}^N, \boldsymbol{\mu}, \boldsymbol{\sigma} | A, S, O) \approx Q(\boldsymbol{\mu}, \boldsymbol{\sigma}) \prod_n^N Q(\mathbf{x}^n) \quad (17)$$

$$Q(\boldsymbol{\mu}, \boldsymbol{\sigma}) = \frac{1}{\sigma_1 \dots \sigma_d} \mathcal{N}_{2d}(\mathbf{z}; \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g) \text{ for } \mathbf{z} = (\mu_1, \dots, \mu_d, \ln \sigma_1, \dots, \ln \sigma_d) \quad (18)$$

$$Q(\mathbf{x}^n) = \mathcal{N}_d(\mathbf{x}^n; \boldsymbol{\mu}_x^n, \boldsymbol{\Sigma}_x^n) \quad (19)$$

402 Note that the approximate posterior captures posterior dependencies between free model parameters
 403 (in the true posterior) on both levels of the hierarchy using the multivariate normal and multivariate
 404 log-normal distributions. However, for practical reasons, we assume statistical independence between
 405 different levels of the hierarchy, and between participants. Independence between participants is
 406 justified by the structure of both response likelihood (responses are modelled as independent and
 407 identically distributed samples from conditional likelihood) and hierarchical prior (a priori statistical
 408 independence between model parameters for each participant).

409 Finally, to find the best approximation of the true posterior given the functional constraints of our
 410 approximate posterior, we minimized the variational free energy $F[Q]$ with respect to the parameters
 411 of the approximate posterior.

$$-\ln p(A|S, O) = F[Q] - D_{KL}(Q||p) \leq F[Q] = f(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g, \boldsymbol{\mu}_x^1, \boldsymbol{\Sigma}_x^1, \dots, \boldsymbol{\mu}_x^N, \boldsymbol{\Sigma}_x^N) \quad (20)$$

$$F[Q] = \int d\mathbf{x}^1 \dots d\mathbf{x}^N d\boldsymbol{\mu} d\boldsymbol{\sigma} Q(\boldsymbol{\mu}, \boldsymbol{\sigma}) \prod_n^N Q(\mathbf{x}^n) \ln \frac{Q(\boldsymbol{\mu}, \boldsymbol{\sigma}) \prod_n^N Q(\mathbf{x}^n)}{p(A|O, S, \mathbf{x}^1, \dots, \mathbf{x}^N) p(\mathbf{x}^1, \dots, \mathbf{x}^N, \boldsymbol{\mu}, \boldsymbol{\sigma})} \quad (21)$$

412 The optimization of the variational free energy $F[Q]$ is based on the SVI implemented in the
 413 probabilistic programming language Pyro [27] and the automatic differentiation module of PyTorch
 414 [28], an open source deep learning platform.

415 As a final remark, we would like to point out that it is possible to use a different hierarchical prior
 416 [29], different parametrization of the hierarchical model [30] or different factorization of the
 417 approximate posterior (e.g., mean-field approximation). However, through extensive comparison of
 418 posterior estimates on simulated data, we have determined that the presented hierarchical model and
 419 the corresponding approximate posterior provide the best posterior estimate of free model parameters
 420 among the set of parametric models we tested (S1 Notebook).

421 Results

422 To investigate how the balance between computationally costly forward planning and heuristic
 423 preferences changes as a function of temporal distance from the goals, participants performed
 424 sequences of actions in a novel sequential decision-making task. The task employed a two-goal
 425 setting, where participants had to decide between approaching the two goals in a sequential or in a
 426 parallel manner. We first performed a standard behavioural analysis, followed by a model-based

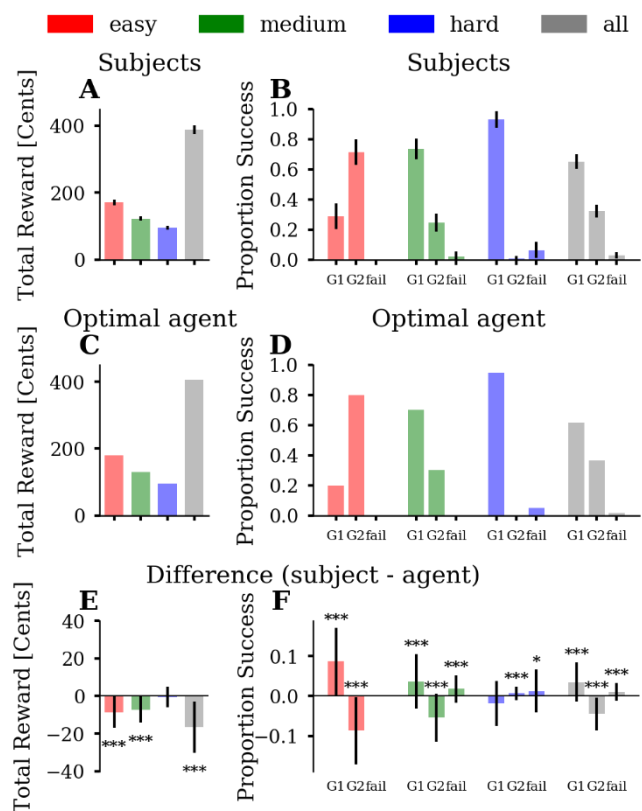
427 approach showing that participants use a mixture of strategy preference and forward planning to
428 select their action.

429 **Standard behavioural analysis**

430 We first analysed the general performance of all participants and – for each miniblock and trial –
431 compared it to the behaviour of an optimal agent possessing perfect knowledge of the task and
432 performing full forward planning to derive an optimal policy that maximizes total reward. The
433 motivation of this comparison was to detect differences between how the optimal agent and
434 participants perform the task. These differences will motivate our model-based analysis below. To
435 compute and compare optimal vs individual policies, all participants and the agent received exactly
436 the same sequence of offers and start conditions. The difference in total reward between participants
437 and agent was significant ($m = -35.5$, $p < 0.001$), where participants earned 388.5 Cents ($SD = 13.6$)
438 and the agent earned 405 Cents. As expected, both participants and agent earned more money in the
439 easy condition than in the medium condition and least in the hard condition (Fig 3, A, C). In the easy
440 and medium condition, the agent earned significantly more than the participants (easy: $M = 8.7$
441 Cents, $SD = 8.4$, $m = -33$, $p < 0.001$; medium: $M = 7.2$ Cents, $SD = 7.0$, $m = -30$, $p < 0.001$). In the
442 hard condition, the total reward did not differ significantly between the participants and agent, $m =$
443 0.5 , $p > 0.99$ (Fig 3, E). These results show that participant performance was generally close to the
444 optimal agent but differed significantly in the easy and medium condition.

445 Next, we analysed participants' goal reaching success and compared it to the optimal agent. There
446 were three possible outcomes in a miniblock: Achieving G1 (goal A or B), achieving G2 (A & B) or
447 fail (neither A nor B). The main experiment comprised 20 miniblocks of each difficulty level
448 modulating difficulty to reach G2. As expected, participants reached on average G2 more often in the
449 easy ($M = 71\%$, $SD = 8\%$) than in the medium condition ($M = 25\%$, $SD = 6\%$), $m = 44.5$, $p < 0.001$.
450 In the hard condition, participants reached G2 in only 1% ($SD = 2\%$) of the miniblocks. Participants
451 failed to reach any goal in 2% ($SD = 3\%$) of the miniblocks in the medium and in 6% ($SD = 5\%$) of
452 the miniblocks in the hard condition. They never failed in the easy condition (Fig 3, B). The agent
453 reached G2 in 80% in the easy, in 30% in the medium and in 0% in the hard condition (Fig 3, D).
454 Note that G2 cannot be reached in all miniblocks. We simulated all possible choice sequences ($n =$
455 2^{15}) for a given miniblock and evaluated whether G2 was theoretically possible. According to
456 these simulations, 90% G2 performance can be reached in the easy, 35% in the medium and 5% in
457 the hard condition.

458 When comparing participants' goal reaching success with the agent, we found that, on average, there
459 was a consistent pattern of deviations in the easy and medium conditions (Fig 3, F). In the easy
460 condition, participants reached G2 on average 9% ($SD = 8\%$) less often than the agent ($m = -33$, $p <$
461 0.001), but reached G1 9% ($SD = 8\%$) more often ($m = 33$, $p < 0.001$). In the medium condition,
462 participants reached G2 on average 6% ($SD = 6\%$) less often than the agent ($m = -26$, $p < 0.001$) but
463 reached G1 4% ($SD = 7\%$) more often ($m = 16.5$, $p < 0.001$). While the agent never failed,
464 participants had a 2% ($SD = 3\%$) fail rate ($m = 11.5$, $p < 0.001$). In the hard condition, participants
465 reached G2 on average 0.6% ($SD = 1.6\%$) more often than the agent ($m = 5.5$, $p < 0.001$). G1 ($m = -$
466 7 , $p = 0.087$) and fail-rate ($m = 3.5$, $p = 0.42$) did not differ significantly between participants and
467 agent. In summary, these differences in successful goal reaching between participants and the agent
468 explains the difference in accumulated total reward: Participants obtained less reward than the agent
469 because on average they missed some of the opportunities to reach G2 in the easy and medium
470 condition and sometimes even failed to achieve any goal in the medium and hard condition.



471

472 **Fig 3. Standard analyses of total reward and comparison to the optimal agent.** (A) Average total
 473 reward across participants. The three conditions are colour-coded (easy = red, medium = green, blue
 474 = hard) and the average over conditions is shown in grey. Error bars depict the standard deviation
 475 (SD). (B) Proportion of successful goal-reaching averaged across participants, for each of the three
 476 conditions. We plot the proportion of reaching, at the end of a miniblock, a single goal (G1), both
 477 goals (G2), or no goal (fail). The fourth block of bars in grey represents the proportions averaged
 478 over all three conditions. Error bars depict SD. (C) Simulated total reward of the optimal agent. (D)
 479 The goal-reaching proportions of the optimal agent. (E) Average difference between participants and
 480 agent with error bars depicting SD. (F) Averaged difference of proportion success between
 481 participants and agent with error bars depicting SD. One can see that the average goal-reaching
 482 proportions of participants were close to the agent's proportions. However, participants, on average,
 483 reached G2 less often than the agent. Asterisks indicate differences significantly greater than zero
 484 (Sign-test, * $\hat{=}$ $p < 0.05$, ** $\hat{=}$ $p < 0.01$, *** $\hat{=}$ $p < 0.001$).

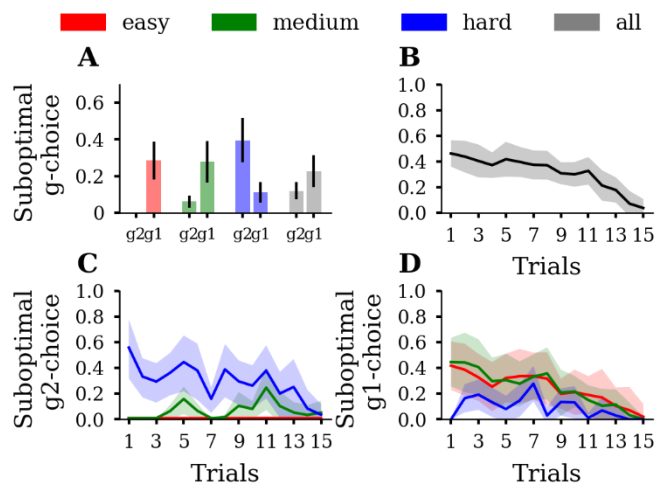
485 How can these differences in goal-reaching success be explained? To address this, we used the
 486 mixed-offer trials to identify which strategy a participant was pursuing in a given trial and compared
 487 the strategy choice to what the agent would have done in this trial. We classified strategy choices as
 488 evidence either of a parallel or a sequential strategy. With the parallel strategy (g2), participants make
 489 choices to pursue both goals in a parallel manner, while with a sequential strategy (g1), participants
 490 make choices to reach first a single goal and then the other. We inferred that participants used a g2-
 491 choice for a specific mixed-offer trial when the difference between the points of the two bars was
 492 minimized, while we inferred a g1-choice when the difference between points was maximized (see
 493 Methods). We categorized a participant's g2-choice as suboptimal when the optimal agent would
 494 have made a g1-choice in a specific trial and vice versa. Fig 4, A-D shows the proportions of
 495 suboptimal g-choices in mixed-offer trials. In the easy condition, participants made barely any

496 suboptimal g2-choice (mean = 0%, SD = 0.001%), but 29% (SD = 10%) suboptimal g1-choices (Fig
497 4, A). This means that participants, on average, preferred a sequential strategy more often than would
498 have been optimal. In the medium condition participants made on average 6% (SD = 3%) suboptimal
499 g2-choices and 28% (SD = 11%) suboptimal g1-choices. Similar to the easy condition, participants,
500 on average, preferred a sequential strategy where a parallel strategy would have been optimal. In the
501 hard condition, this pattern reversed. Participants made on average 40% (SD = 12%) suboptimal g2-
502 choices, relative to the agent, and 11% (SD = 6%) suboptimal g1-choices. Participants' suboptimal g-
503 choices were also reflected in goal reaching success. In the easy and medium condition, suboptimal
504 g1-choices, relative to the agent, resulted in a higher proportion of reaching G1, and a lower
505 proportion of reaching G2. In the hard condition, suboptimal g2-choices led to occasional fails and a
506 tiny margin of reaching G2. However, despite suboptimal g2-choices, participants still reached G1 in
507 93% (SD = 6%) of the miniblocks.

508 As the first test of our prediction that participants tend to use more forward planning when
509 temporally proximal to the goal, we analysed suboptimal decisions as a function of trial time. As
510 expected, suboptimal decisions, relative to the agent, decreased over trial time (Fig 4, B). While in
511 the first trial, 42% (SD = 19%) of participants' g-choices deviated from the agent's g-choices,
512 participant behaviour converged to almost optimal performance towards the end of the miniblock,
513 with only 4% deviating g-choices (SD = 7%). We also simulated a random agent that accepts all
514 basic A or B offers but guesses on mixed offers (S6-7 Fig). S7 Fig B shows that the random agent
515 makes approximately 50 % suboptimal g-choices across all trials in the miniblock. That means
516 participants used non-random response strategies, i.e. planning or heuristics, since their pattern of
517 suboptimality across trials deviated from the straight-line pattern of the random agent.

518 In the hard condition, the number of suboptimal g2-choices similarly decreased, but not in the easy
519 and medium condition (Fig 4, C). The number of suboptimal g1-choices decreased across trials in the
520 easy and medium, but not in hard condition (Fig 4, D). Note that in easy and the medium conditions,
521 opportunities to make suboptimal g2-choices are generally scarce, because the difference between
522 action values $DEV = Q_G(g2) - Q_G(g1)$ was mostly positive, which means that a g2-choice was
523 mostly optimal. Similarly, in the hard condition, as there was a low number of opportunities to make
524 suboptimal g1-choices, there was no clear decrease in the number of suboptimal g1-choices.

525 Although these findings of diminishing suboptimal choices over the course of miniblocks may be
526 explained by the participants' initial employment of a suboptimal heuristic, there is an alternative
527 explanation because we used an optimal agent, which uses a max operator to select its action: If this
528 agent computes, by using forward planning, a tiny advantage in expected reward of one action over
529 the other, the agent will always choose in a deterministic fashion the action with the slightly higher
530 expected reward. Therefore, at the beginning of the miniblock, where the distance to the final trial is
531 largest, the difference between goal choice values $DEV = Q_G(g2) - Q_G(g1)$ (S5 Fig) is close to 0.
532 The reason for this is that a single g2-choice at the beginning of the miniblock does not increase the
533 probability for G2-success by much. However, when only few trials are left, a single g2-choice might
534 make the difference between winning or losing G2. Since DEV s are close to 0 at the initial trials we
535 cannot exclude the possibility yet that participants actually may have used optimal forward planning
536 just like the agent but did not use a max operator. Instead, participants may have sampled an action
537 according to the computed probabilities of each action to reach the greater reward in the final trial.
538 Such a sampling procedure to select actions would also explain the observed pattern of diminishing
539 suboptimal g-choices over the miniblock (Fig. 4 B-C). To answer the question, whether there is
540 actually evidence that participants use heuristics, when far from the goal, even in the presence of
541 probabilistic action selection of participants, we now turn to a model-based analysis.



542

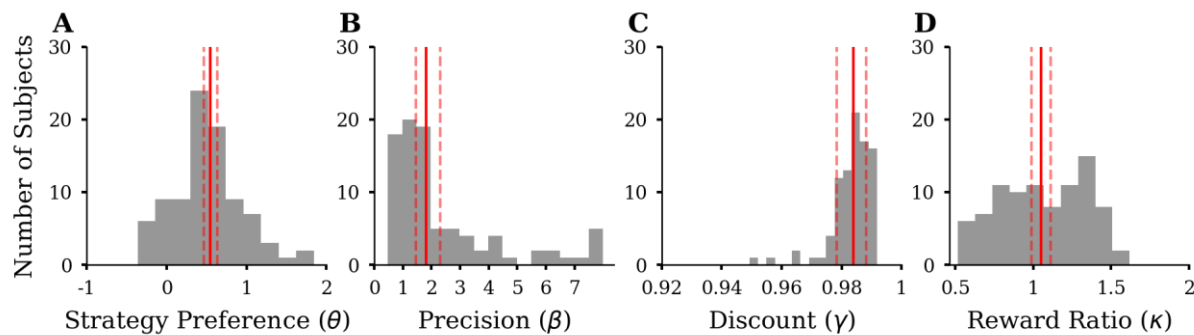
543 **Fig 4. (A)** Proportions of suboptimal g1-choices (g1) and suboptimal g2-choices (g2), averaged over
544 participants. Participants tend to make suboptimal g1-choices in the easy and medium condition
545 while this pattern reverses in the hard condition. Error bars depict SD. Conditions are colour coded.
546 **(B)** Suboptimal g-choices as a function of trial averaged over participants. Shaded areas depict SD.
547 **(C)** Suboptimal g2-choices as a function of trial averaged over participants. **(D)** Suboptimal g1-
548 choices as a function of trial averaged over participants. In both C and D, one can see that
549 participants made more suboptimal g-choices at the beginning of the miniblock than close to the final
550 trial. Shaded areas depict SD.

551 **Model-based behavioural analysis**

552 To infer the contributions of participants' forward planning and heuristic preferences, we conducted a
553 model-based analysis. If we find that participants' strategy preference θ is smaller or larger than zero,
554 we can conclude that participants indeed used a heuristic component to complement any forward
555 planning. This is especially relevant for choices early in the miniblock as *DEV* values are typically
556 close to zero. Indeed, when inferring the four parameters for all 89 participants using hierarchical
557 Bayesian inference, we found that participants' g-choices were influenced by a heuristic strategy
558 preference in addition to a forward planning component (Fig 5, A). For 74 out of 89 participants, we
559 found that the 90% credibility interval (CI) of the posterior over strategy preference did not include
560 zero. 68 of these participants had a positive strategy preference, meaning they preferred an overall
561 strategy of pursuing both goals in parallel. Six of these participants had a negative strategy
562 preference, meaning they preferred to pursue both goals sequentially. The median group
563 hyperparameter of strategy preference was 0.55 (90% CI = [0.47, 0.63]). For example, a participant
564 with this median strategy preference, in a mixed-offer trial where *DEV* = 0, would make a g2-choice
565 with 63% probability, whereas a participant without a strategy preference bias, i.e. $\theta = 0$, would
566 make a g2-choice with 50% probability. After the experiment, we had asked participants whether
567 they used any specific strategies to solve the task and to give a verbal description of the used
568 strategy. Reports reflected three main patterns: Pursuing one goal after the other (sequential strategy),
569 promoting both goals in a balanced way (parallel strategy), and switching between sequential and
570 parallel strategy, depending on context (mixed strategy). Reported strategies are in good qualitative
571 agreement with the estimated strategy preference parameter (S8 Fig), supporting our interpretation of
572 this parameter. Notably, the task instructions, given to the participants prior to the experiment, did
573 not point to any specific heuristic (S1 Text). Altogether, the non-zero strategy preference in 83% of

574 participants indicates that suboptimal decisions within a miniblock (see Fig 4) are not only caused by
575 probabilistic sampling for action selection, but also by the use of a heuristic strategy preference.

576 As expected, we found that the *DEV* (see Table 1) derived by forward planning influenced action
577 selection (median group hyperparameter of the inferred precision $\beta = 1.82$, 90% CI = [1.45, 2.3], Fig
578 5, B). For example, a hypothetical participant with parameters similar to the group hyperparameters
579 ($\theta = 0.55$ and $\beta = 1.82$), when encountering a *DEV* = 0.5, would make a g2-choice with 82%
580 probability. Increasing *DEV* by 1 would increase the g2-choice probability to 96%. In contrast, a
581 participant with low precision but the same median strategy preference ($\theta = 0.55$ and $\beta = 0.5$),
582 when encountering a *DEV* = 0.5, would make a g2-choice with 69% probability. Increasing *DEV* by
583 1 would increase g2-choice probability to 79%. We found evidence only for weak discounting of
584 future rewards, as for most participants the inferred discount was close to 1 (median of the inferred
585 discount parameter $\gamma = 0.984$, 90% CI = [0.978, 0.988], Fig 5, C). We found that some participants
586 used a reward ratio different from the objective value of 1 (CI not containing 1). Twelve participants
587 had a reward ratio greater than 1 and 17 participants had a reward ratio smaller than 1. However, the
588 median group hyperparameter of the inferred reward ratio was close to the objective value of 1 ($\kappa =$
589 1.05, 90% CI = [0.99, 1.11], Fig 5, D). A reward ratio of 1.2 means, that participants behaved as if
590 the value of achieving G2 would be 2.4 times the value of achieving G1 (when in reality the reward is
591 only double as high). While strategy preference has its greatest influence during the first few trials of
592 a miniblock, the reward ratio has an influence only when forward planning, i.e. changes the *DEV*,
593 and will therefore affect action selection most during the final trials of a miniblock. In addition, we
594 found only low posterior correlation between the strategy preference and reward ratio parameter,
595 indicating that these two parameters model distinct influences on goal reaching behaviour.



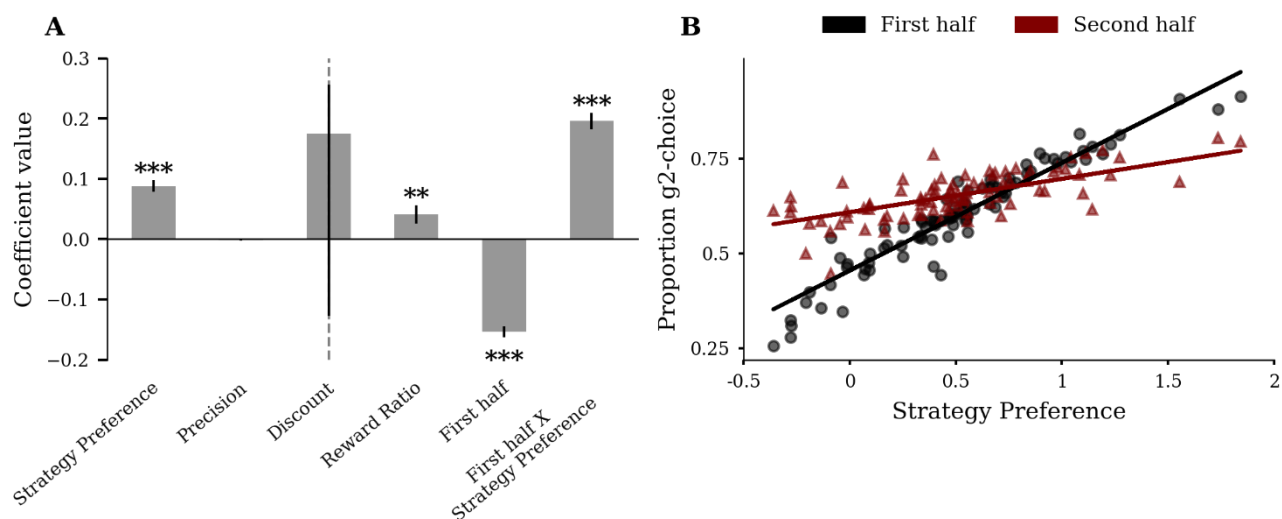
596

597 **Fig 5. Summary of inferred parameters of the four-parameter model for all 89 participants.** We
598 show histograms of the median of the posterior distribution, for each participant. Solid red lines
599 indicate the median of the group hyperparameter posterior estimate with dashed lines indicating 90%
600 credibility intervals (CI). (A) Histogram of strategy preference parameter θ . (B) Histogram of
601 precision parameter β (last bin containing values > 8). (C) Histogram of discount parameter γ . (D)
602 Histogram of reward ratio parameter κ .

603 To show that our model with constant parameters is able to capture a dynamic shift from heuristic
604 decision making to forward planning we conducted two sets of simulations where we systematically
605 varied the response precision β and the strategy preference parameter θ . First, we simulated
606 behaviour where we varied β between 0.25 and 3 with θ , γ , and κ sampled from their fitted
607 population mean (S1-2 Movie). S2 Movie, B shows that the higher β , the fewer suboptimal g-choices
608 are made towards the end of the miniblock. Second, we simulated behaviour where we varied θ
609 varied between -1 and 1 with β , γ , and κ sampled from their fitted population mean (S3-4 Movie). S4

610 Movie, B shows that a change in θ affects the number of suboptimal g-choices made at the beginning
611 but not at the end of the miniblock. These two results support the argument that the θ parameter is
612 able to capture heuristic decision making at the beginning of the miniblock while the β parameter is
613 able to capture planning behaviour at the end of the miniblock. The reason for this interaction
614 between parameter effect and trial number is that differential expected value (*DEV*) computed by
615 forward planning is close to zero at the beginning of the miniblock but increases towards the end of
616 the miniblock (S5 Fig). For small *DEVs*, the influence of β on choice probability is marginal;
617 therefore, the relative influence of the strategy preference parameter θ is high, and behaviour is
618 explained by using the heuristic. For higher trial numbers, i.e. closer to the end of the miniblock,
619 *DEVs* tend to be high so that the influence of the response precision β is high, and the relative
620 influence of θ is low; therefore, towards the end of the miniblock behaviour is explained by forward
621 planning with a shift in between, depending on the dynamics of the *DEV*. We also implemented a
622 model with changing parameters over trials and compared it to the constant model. Parameters were
623 fit separately for three partitions of the miniblock, i.e. early (trials 1- 5), middle (trials 6-10) and late
624 trials (11-15). Model comparisons showed that this model with changing parameters had lower model
625 evidence compared to the model with constant parameters (S9 Fig). We interpret these results as
626 further evidence that the described constant parameterization is sufficient to describe a hidden shift
627 from using a heuristics to forward planning.

628 Finally, as an additional test of the hypothesis that participants rely more on heuristic preferences
629 when the goal is temporally distant, we conducted a multiple regression analysis (Fig 6, A). To do
630 this, we divided the data into the first (first 7 trials) and the second half (last 8 trials) of miniblocks,
631 and computed, for each participant the proportion of g2-choices in the mixed-offer trials. We fitted,
632 across participants, these proportions of g2-choices against 6 regressors: strategy preference,
633 precision, discount rate, reward ratio, a dummy variable coding for the first and second miniblock
634 half and interaction between strategy preference and miniblock half. We found a significant
635 interaction between strategy preference and miniblock-half ($p < 0.001$), demonstrating that strategy
636 preference is more predictive for the proportion of g2-choices in the first half of the miniblock than in
637 the second half. Fig 6, B visualizes the interaction effect showing that the slope of the marginal
638 regression line for the first half of the miniblock is greater than the slope of the marginal regression
639 line for the second half of the miniblock. This finding provides additional evidence that participants
640 rely on heuristic preferences when the goal is temporally far away but use differential expected
641 values (*DEV*) derived by forward planning when the goal is closer.



642

643 **Fig 6. Strategy preference is more predictive for participant's proportion of g2-choices in the**
 644 **first than in the second half of the miniblock. (A)** Linear regression of proportion g2-choice
 645 against parameters from the four-parameter model, a dummy variable coding for miniblock-half and
 646 interaction between miniblock-half and strategy preference. The significant interaction term supports
 647 the hypothesis that the influence of strategy preference on g2-choice proportion is greater in the first
 648 than in the second half of the miniblock. Error bars represent SE. Asterisks indicate coefficients
 649 significantly different from 0 (t-test, * \triangleq $p < 0.05$, ** \triangleq $p < 0.01$, *** \triangleq $p < 0.001$). **(B)** Strategy
 650 preference plotted against the proportion of g2-choices in the first half of the miniblock (black) and
 651 in the second half of the miniblock (red). Solid lines represent marginal regression lines.

652 In addition, we conducted model comparisons, posterior predictive checks and parameter recovery
 653 simulations to test whether our model is an accurate and parsimonious fit to the data. First, we
 654 compared variants of our model, where we fixed individual parameters (S9 Fig). Adding θ and β
 655 increased model evidence, confirming their importance in explaining participant behaviour. The
 656 three-parameter model (θ , β , κ) had the highest model evidence among all 16 models. Adding γ did
 657 not increase model evidence. This result is consistent since we found only little evidence for
 658 discounting when fitting the parameters, see Fig. 5C. To test whether participants used condition-
 659 specific response strategies (e.g., use heuristics in the easy and hard but plan forward in the medium
 660 difficult condition) we estimated model parameters separately for conditions. However, the
 661 condition-wise model had lower model evidence compared to the conjoint model, indicating that
 662 participants use a condition-general approach to arbitrate between using a heuristic and planning
 663 ahead. Second, we simulated data using the group mean parameters as inferred from the participants'
 664 data and compared it to the observed data. Visual inspection shows that both the simulated
 665 performance pattern (S10 Fig) and the simulated frequency of suboptimal g-choices (S11 Fig) closely
 666 resemble the experimentally observed patterns (Fig. 3 and 4). Third, we simulated data using
 667 participants' posterior mean and tested whether we could reliably infer parameters (S1 Notebook).
 668 Results showed that the inferred β , θ and κ align with the true parameter value, but simulation-based
 669 calibration [31] suggests that estimates of γ are biased. Taken together, our model provides a good fit
 670 to the data, where the data are informative about the three parameters β , θ and κ .

671

672 We also tested whether participants showed learning effects in the main experimental phase. In a first
 673 linear model, the depended variable was the total reward and the predictor was the experimental
 674 block number (miniblock 1-20, miniblock 21-40, miniblock 41-60). The analysis revealed a

675 significant but small main effect of experiment block ($\beta = 5.4$, $SE = 0.5$, $p < 0.001$). In a second
676 logistic model the dependent variable was suboptimal goal choice (1 = suboptimal, 0 = optimal) and
677 the predictor was experiment block. The second analysis revealed a significant but small main effect
678 of experiment block on the probability to make a suboptimal g-choice ($\beta = -0.084$, $SE = 0.02$, $p <$
679 0.001). Furthermore, we fitted the three parameter model (θ, β, κ) separately for experiment blocks.
680 Model comparisons revealed that the experiment block-wise model had lower model evidence
681 compared to the conjoint model (S9 Fig.).

682
683 As a final control analysis, we used logistic regression to establish how the absolute difference
684 between A- and B-points affects goal choice as a function of the number of trials remaining in the
685 miniblock. If participants rely on a fixed strategy preference when far from the goal, there should be
686 no effect of absolute score difference on goal choice at the start of miniblocks. In this model the
687 depended variable was goal choice (1 = g2, 0 = g1) and the predictors were absolute score difference
688 ($|Pts_t^A - Pts_t^B| \in [0..15]$), miniblock-half (1 = trial 1-7, 0 = trial 8-15) and the interaction term
689 absolute score difference*miniblock-half. There was a significant main effect of absolute score
690 difference ($\beta = 0.14$, $SE = 0.008$, $p < 0.001$) and miniblock-half ($\beta = 0.29$, $SE = 0.039$, $p < 0.001$).
691 Importantly, the analysis revealed a significant interaction between miniblock-half and absolute score
692 difference ($\beta = -0.2$, $SE = 0.013$, $p < 0.001$). This means that goal choice was more affected by the
693 absolute score difference in the second half the miniblock compared to the first half. The analysis
694 supports our conclusion that participants relied on a heuristic strategy preference when far from the
695 goal.

696 Discussion

697 In the current study, we investigated how humans change the way they decide what goal to pursue
698 while approaching two potential goals. To emulate real life temporally extended decision making
699 scenarios of goal pursuit, we used a novel sequential decision making task. In this task environment,
700 decisions of participants had deterministic consequences, but the options given to participants on
701 each of the 15 trials were stochastic. This meant that especially during the first few trials, participants
702 could not predict with certainty what goal was achievable. Using model-based analysis of
703 behavioural data we find that most participants, during the initial trials, relied on computationally
704 inexpensive heuristics and switched to forward planning only when closer to the final trial.

705 We inferred the transition from a heuristic action selection to action selection based on forward
706 planning using a model parameter that captured participants' preference for pursuing both goals
707 either in a sequential or parallel manner. This strategy preference had its strongest impact for the first
708 few trials, when participants, due to the stochasticity of future offers, could not predict well which of
709 the two available actions in a mixed trial would enable them to maximize their gain. This can be seen
710 from Eq. 11 where two terms contribute to making a decision: the term containing the differential
711 expected value (DEV) and the strategy preference θ . In our computational model, the DEV is the
712 difference between the expected value of a sequential strategy choice and a parallel strategy choice.
713 The DEV enables the agent to choose actions which maximize the average reward gain in a
714 miniblock (see methods). Critically, this DEV is typically close to 0 in the first few trials, i.e. there is
715 high uncertainty on what action is the best one. In this situation, the strategy preference mostly
716 determines the action selection of the agent. In our model, we computed the DEV by using forward
717 planning, where the agent hypothetically runs simulations through all remaining future trials until the
718 end of a miniblock, i.e. to the 15th trial. The number of state space trajectories to be considered in
719 these simulations scales exponentially with the number of remaining trials – and so does in principle
720 the computational costs needed to simulate these trajectories. Therefore, full forward planning would

721 be both prohibitively costly and potentially useless when the deadline is far away, rendering simpler
722 heuristics [16] the more appropriate alternative.

723 It is an open question what heuristic participants actually used. In our model, the strategy preference
724 parameter simply quantifies a preference for a parallel or sequential strategy and biases a
725 participant's action selection accordingly. This may mean that participants had a prior expectation
726 whether they are going to reach G2 or just G1. Given this prior, participants could choose their action
727 without any forward planning. In other words, to select an action in a mixed trial, participants simply
728 assumed that they are going to reach, for example, G2. This simplifies action selection tremendously
729 because, under the assumption that G2 will be reached, the optimal action is to use the parallel
730 strategy at all times. To an outside observer, a participant with a strong preference for a parallel
731 strategy may be described as overly optimistic, as this participant would choose g2-choices even if
732 reaching G2 is not very likely, e.g. in the hard condition. Conversely, a participant with a strong
733 preference for a sequential strategy may be described as too cautious, e.g. because that participant
734 chooses one-goal actions in the easy condition (see S12 Fig for two example participants).
735 Importantly, the difference in total reward between the agent and the participants is only about 5%
736 (see Fig 3, E). This means that even though participants used a potentially suboptimal strategy
737 preference, the impact on total reward is not that large. This is because, as we have shown, later in
738 the miniblock, when *DEVs* become larger and are more predictive of what goal can be reached,
739 participants choose their actions accordingly. Although we do not quantify the relative costs of full
740 forward planning versus the observed mixture of heuristic and forward planning, we assume that an
741 average loss of 5% of the earnings is small as compared to the reduction of computational costs when
742 using heuristics.

743 There were two important features of our sequential decision making task: The first was that we used
744 a rather long series of 15 trials to model multiple goal pursuit, where typically sequential decision
745 making tasks would use fewer trials, e.g. 2 in the two-step task [32] with common values around 5
746 [21] to 8 trials [7, 8] per miniblock. The reason why we chose a rather large number of trials is that
747 this effectively precluded the possibility that participants can plan forward and ensure that
748 participants were exposed at least to some initial trials where they had to rely on other information
749 than forward planning. This initial period when participants have to select actions without an accurate
750 estimate of the future consequences of these actions is potentially most interesting for studying meta-
751 decisions about how we use heuristics when detailed information about goal reaching probabilities is
752 scarce. It is probably in this period of uncertainty during goal reaching, when internal beliefs and
753 preferences have their strongest influence.

754 The second important feature of our task was that participants had to prioritize between two goals.
755 This is a departure from most sequential decision making tasks, where there is typically a single goal,
756 e.g. to collect a minimum number of points, where the alternative is a fail [7]. In our task,
757 participants could reach one of two goals, which enables addressing questions about how participants
758 select and pursue a specific goal, see also [9]. Our findings complement work investigating
759 behavioural strategies for pursuing multiple goals, e.g. [33], showing that pursuit strategies depend
760 on environmental characteristics, subjective preferences and changes in context when getting closer
761 to the goal. In line with our findings, a recent study [34] showed that decisions whether to redress the
762 imbalance between two assets or to focus on a distinct asset during sequential goal pursuit were best
763 fit by a dynamic programming model with a limited time horizon of 7.5 trials (20 trials would be the
764 optimum). In future research, the pursuit of multiple goals in sequential decision making tasks may
765 also be a basis for addressing questions about cognitive control during goal-reaching, e.g. how

766 participants regulate the balance between stable maintenance and flexible updating of goal
767 representations [35].

768 Another important factor when modelling the use of forward planning is that complexity and time
769 can, in principle be dissociated. For example, a temporally distant goal might have only low planning
770 complexity because one must consider only a few decision sequences leading to the goal.
771 Conversely, a temporally proximate goal might have high planning complexity because of a large
772 number of potential actions sequences that may lead to the goal. In future research, by testing
773 sequential tasks with varying branching factor (number of potential actions in each trial) one could
774 selectively test how time to goal and planning complexity influence the arbitration of forward
775 planning and the use of heuristics.

776 It is unclear what mechanism made participants actually use a strategy preference different from zero
777 in our task. It is tempting to assume that participants might have used their usual approach, which
778 they might apply in similar real-life situations, to select their goal strategies when the computational
779 costs of forward planning are high and the prediction accuracy is low. In other words, participants
780 who had a preference for a parallel strategy might either show a tendency towards working on
781 multiple goals at the same time or entertain the belief that tasks should be approached with an
782 optimistic stance. Conversely, participants with a preference for a sequential strategy might have
783 made good experiences with using a more cautious approach and would tend to pursue one goal after
784 the other.

785 We would like to note that the proposed model does not explicitly model the arbitration between
786 forward planning and heuristic decision making. The computational model to fit participant
787 behaviour uses at its core full forward planning as the optimal agent does. The effect of strategy
788 preference just changes the action selection result, but the underlying computation to determine the
789 *DEV* is still based on forward planning. Clearly, if a real agent used our model, this agent would not
790 save any computations because forward planning is still used for all trials. The open question is how
791 an agent makes a meta-decision to not use goal-directed forward planning but to rely on heuristics
792 and other cost-efficient action selection procedures [11]. To make this meta-decision, an agent cannot
793 rely on the *DEV* because this value is computed by forward planning. An alternative way would be to
794 use an agent's prior experience to decide that the goal is still too temporally distant to make an
795 informed decision with an acceptable computational cost. Such a meta-decision would depend on
796 several factors, e.g. the relevance of reaching G2, intrinsic capability and motivation of planning
797 forward, or a temporal distance parameter which signals urgency to start planning forward. In the
798 future we plan to develop such meta-decision-making models and predict the moment at which
799 forward planning takes over the action selection process.

800 It is also possible that participants use, apart from simple heuristics, other approximate planning
801 strategies to reduce computational costs. For example, one could sample only a subset of sequences
802 to compute value estimates. Indeed, in another study it was found that participants prune a part of the
803 decision tree in response to potential losses, even if this pruning was suboptimal [36]. Another
804 important point is that the planning process itself might be error-prone and therefore value
805 calculations over longer temporal horizons may be noisier. This could presumably account for
806 temporal modulations of the precision parameter β . In future work one could test for evidence of
807 alternative planning algorithms that allow to sample subsets of (noisy) forward planning trajectories
808 to further delineate how humans deal with computational complexity in goal-directed decision
809 scenarios.

810 Taken together, the present research shows that over prolonged goal-reaching periods, individuals
811 tend to behave in a way that approaches the behaviour of an optimal agent, with noticeable
812 differences early in the goal-reaching period, but nearly optimal behaviour when the goal is close. It
813 also highlights the potential of computational modelling to infer the decision parameters individuals
814 use during different stages of sequential decision-making. Such models may be a promising means to
815 further elucidate the dynamics of decision-making in the pursuit of both laboratory and everyday life
816 goals.

817 **References**

- 818 1. Schmidt AM, DeShon RP. What to do? The effects of discrepancies, incentives, and time on dynamic
819 goal prioritization. *Journal of Applied Psychology*. 2007;92(4):928.
- 820 2. Schmidt AM, Dolis CM. Something's got to give: The effects of dual-goal difficulty, goal progress, and
821 expectancies on resource allocation. *Journal of Applied Psychology*. 2009;94(3):678.
- 822 3. Neal A, Ballard T, Vancouver JB. Dynamic Self-Regulation and Multiple-Goal Pursuit. *Annual Review of*
823 *Organizational Psychology and Organizational Behavior*. 2017;4:401-23.
- 824 4. Schacter DL, Addis DR, Hassabis D, Martin VC, Spreng RN, Szpunar KK. The future of memory:
825 remembering, imagining, and the brain. *Neuron*. 2012;76(4):677-94.
- 826 5. Hayes-Roth B, Hayes-Roth F. A cognitive model of planning. *Cognitive science*. 1979;3(4):275-310.
- 827 6. Economides M, Guitart-Masip M, Kurth-Nelson Z, Dolan RJ. Anterior cingulate cortex instigates
828 adaptive switches in choice by integrating immediate and delayed components of value in
829 ventromedial prefrontal cortex. *J Neurosci*. 2014;34(9):3340-9. Epub 2014/02/28. doi:
830 10.1523/JNEUROSCI.4313-13.2014. PubMed PMID: 24573291; PubMed Central PMCID:
831 PMCPMC3935091.
- 832 7. Kolling N, Wittmann M, Rushworth MFS. Multiple neural mechanisms of decision making and their
833 competition under changing risk pressure. *Neuron*. 2014;81(5):1190-202. Epub 2014/03/13. doi:
834 10.1016/j.neuron.2014.01.033. PubMed PMID: 24607236; PubMed Central PMCID:
835 PMCPMC3988955.
- 836 8. Schwartenbeck P, FitzGerald TH, Mathys C, Dolan R, Friston K. The Dopaminergic Midbrain Encodes
837 the Expected Certainty about Desired Outcomes. *Cereb Cortex*. 2015;25(10):3434-45. Epub
838 2014/07/25. doi: 10.1093/cercor/bhu159. PubMed PMID: 25056572; PubMed Central PMCID:
839 PMCPMC4585497.
- 840 9. Ballard T, Yeo G, Neal A, Farrell S. Departures from optimality when pursuing multiple approach or
841 avoidance goals. *Journal of Applied Psychology*. 2016;101(7):1056.
- 842 10. Gershman SJ, Horvitz EJ, Tenenbaum JB. Computational rationality: A converging paradigm for
843 intelligence in brains, minds, and machines. *Science*. 2015;349(6245):273-8.
- 844 11. Boureau Y-L, Sokol-Hessner P, Daw ND. Deciding how to decide: self-control and meta-decision
845 making. *Trends in cognitive sciences*. 2015;19(11):700-10.
- 846 12. Shenhav A, Botvinick MM, Cohen JD. The expected value of control: an integrative theory of anterior
847 cingulate cortex function. *Neuron*. 2013;79(2):217-40.
- 848 13. Shenhav A, Musslick S, Lieder F, Kool W, Griffiths TL, Cohen JD, et al. Toward a rational and
849 mechanistic account of mental effort. *Annual review of neuroscience*. 2017;40:99-124.
- 850 14. Lieder F, Griffiths TL. Strategy selection as rational metareasoning. *Psychological Review*.
851 2017;124(6):762.

- 852 15. Gigerenzer G, Gaissmaier W. Heuristic decision making. *Annual review of psychology*. 2011;62:451-
853 82.
- 854 16. Soltani A, Khorsand P, Guo C, Farashahi S, Liu J. Neural substrates of cognitive biases during
855 probabilistic inference. *Nature communications*. 2016;7:11393.
- 856 17. Keramati M, Smittenaar P, Dolan RJ, Dayan P. Adaptive integration of habits into depth-limited
857 planning defines a habitual-goal-directed spectrum. *Proceedings of the National Academy of
858 Sciences*. 2016;113(45):12868-73.
- 859 18. Kleiner M, Brainard D, Pelli D, Ingling A, Murray R, Broussard C. What's new in Psychtoolbox-3.
860 *Perception*. 2007;36(14):1.
- 861 19. Ostwald D, Bruckner R, Heekeren H. Computational mechanisms of human state-action-reward
862 contingency learning under perceptual uncertainty. *Conference on Cognitive Computational
863 Neuroscience*; 2018; Philadelphia, Pennsylvania, USA2018.
- 864 20. Puterman ML. *Markov decision processes: discrete stochastic dynamic programming*: John Wiley &
865 Sons; 2014.
- 866 21. Korn CW, Bach DR. Heuristic and optimal policy computations in the human brain during sequential
867 decision-making. *Nature communications*. 2018;9(1):325.
- 868 22. Seabold S, Perktold J, editors. *Statsmodels: Econometric and statistical modeling with python*.
869 *Proceedings of the 9th Python in Science Conference*; 2010: Scipy.
- 870 23. Team RC. *R: A language and environment for statistical computing*. 2013.
- 871 24. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4. *arXiv preprint
872 arXiv:14065823*. 2014.
- 873 25. Kuznetsova A, Brockhoff PB, Christensen RHB. lmerTest package: tests in linear mixed effects models.
874 *Journal of Statistical Software*. 2017;82(13).
- 875 26. Hoffman MD, Blei DM, Wang C, Paisley J. Stochastic variational inference. *The Journal of Machine
876 Learning Research*. 2013;14(1):1303-47.
- 877 27. Bingham E, Chen JP, Jankowiak M, Obermeyer F, Pradhan N, Karaletsos T, et al. Pyro: Deep universal
878 probabilistic programming. *arXiv preprint arXiv:181009538*. 2018.
- 879 28. Paszke A, Gross S, Chintala S, Chanan G, Yang E, DeVito Z, et al. Automatic differentiation in pytorch.
880 2017.
- 881 29. Polson NG, Scott JG. Shrink globally, act locally: Sparse Bayesian regularization and prediction.
882 *Bayesian statistics*. 2010;9:501-38.
- 883 30. Bernardo J, Bayarri M, Berger J, Dawid A, Heckerman D, Smith A, et al., editors. Non-centered
884 parameterisations for hierarchical models and data augmentation. *Bayesian Statistics 7: Proceedings
885 of the Seventh Valencia International Meeting*; 2003: Oxford University Press, USA.
- 886 31. Talts S, Betancourt M, Simpson D, Vehtari A, Gelman A. Validating Bayesian inference algorithms with
887 simulation-based calibration. *arXiv preprint arXiv:180406788*. 2018.
- 888 32. Daw ND, Gershman SJ, Seymour B, Dayan P, Dolan RJ. Model-based influences on humans' choices
889 and striatal prediction errors. *Neuron*. 2011;69(6):1204-15.
- 890 33. Orehek E, Vazeou-Nieuwenhuis A. Sequential and concurrent strategies of multiple goal pursuit.
891 *Review of General Psychology*. 2013;17(3):339.
- 892 34. Juechems K, Balaguer J, Castañón SH, Ruz M, O'Reilly JX, Summerfield C. A network for computing
893 value equilibrium in the human medial prefrontal cortex. *Neuron*. 2019;101(5):977-87. e3.

- 894 35. Goschke T. Dysfunctions of decision-making and cognitive control as transdiagnostic mechanisms of
 895 mental disorders: advances, gaps, and needs in current research. *International journal of methods in*
 896 *psychiatric research*. 2014;23(S1):41-57.
- 897 36. Huys QJ, Eshel N, O'Nions E, Sheridan L, Dayan P, Roiser JP. Bonsai trees in your head: how the
 898 Pavlovian system sculpts goal-directed choices by pruning decision trees. *PLoS computational*
 899 *biology*. 2012;8(3):e1002410.

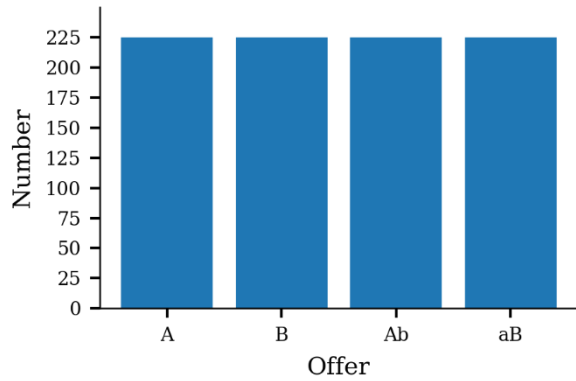
900

901 Supporting information

902 **S1 Table. Classification of accept-wait responses into either two-goal-choices (g2) or one-goal-**
 903 **choices (g1).**

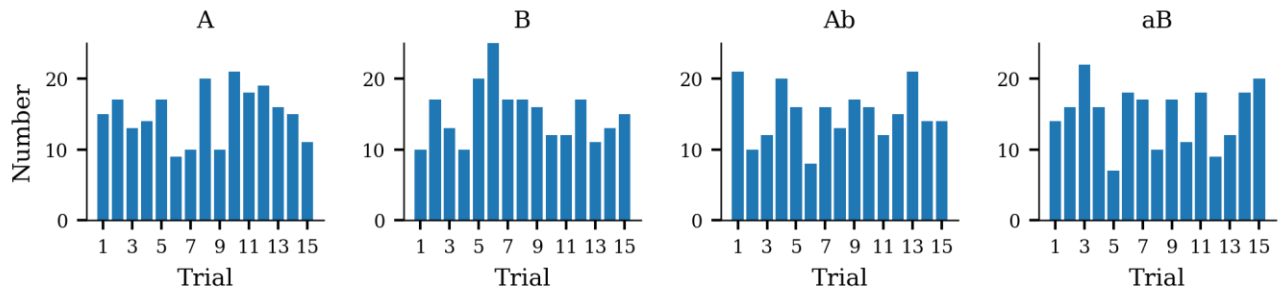
Offer	Points	Response	Classification
Ab	$Pts_t^A - Pts_t^B > 1$	accept	g1
Ab	$Pts_t^A - Pts_t^B > 1$	wait	g2
Ab	$Pts_t^A - Pts_t^B < -1$	accept	g2
Ab	$Pts_t^A - Pts_t^B < -1$	wait	g1
Ab	$Pts_t^A - Pts_t^B = 1$	accept	g1
Ab	$Pts_t^A - Pts_t^B = 1$	wait	g2
Ab	$Pts_t^A - Pts_t^B = -1$	accept	nan
Ab	$Pts_t^A - Pts_t^B = -1$	wait	nan
Ab	$Pts_t^A - Pts_t^B = 0$	accept	g1
Ab	$Pts_t^A - Pts_t^B = 0$	wait	g2
aB	$Pts_t^A - Pts_t^B > 1$	accept	g2
aB	$Pts_t^A - Pts_t^B > 1$	wait	g1
aB	$Pts_t^A - Pts_t^B < -1$	accept	g1
aB	$Pts_t^A - Pts_t^B < -1$	wait	g2
aB	$Pts_t^A - Pts_t^B = 1$	accept	nan
aB	$Pts_t^A - Pts_t^B = 1$	wait	nan
aB	$Pts_t^A - Pts_t^B = -1$	accept	g1
aB	$Pts_t^A - Pts_t^B = -1$	wait	g2
aB	$Pts_t^A - Pts_t^B = 0$	accept	g1
aB	$Pts_t^A - Pts_t^B = 0$	wait	g2

904



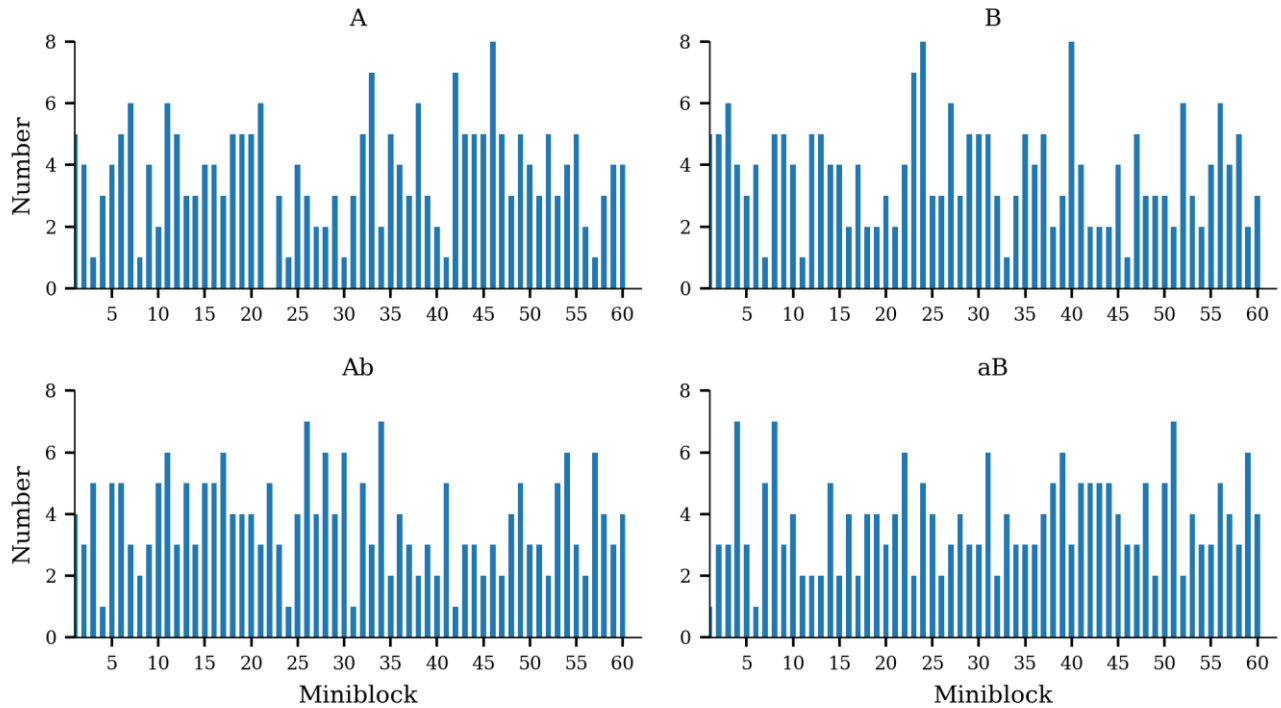
905

906 **S1 Fig. Occurrence of offer types across all 900 trials.**



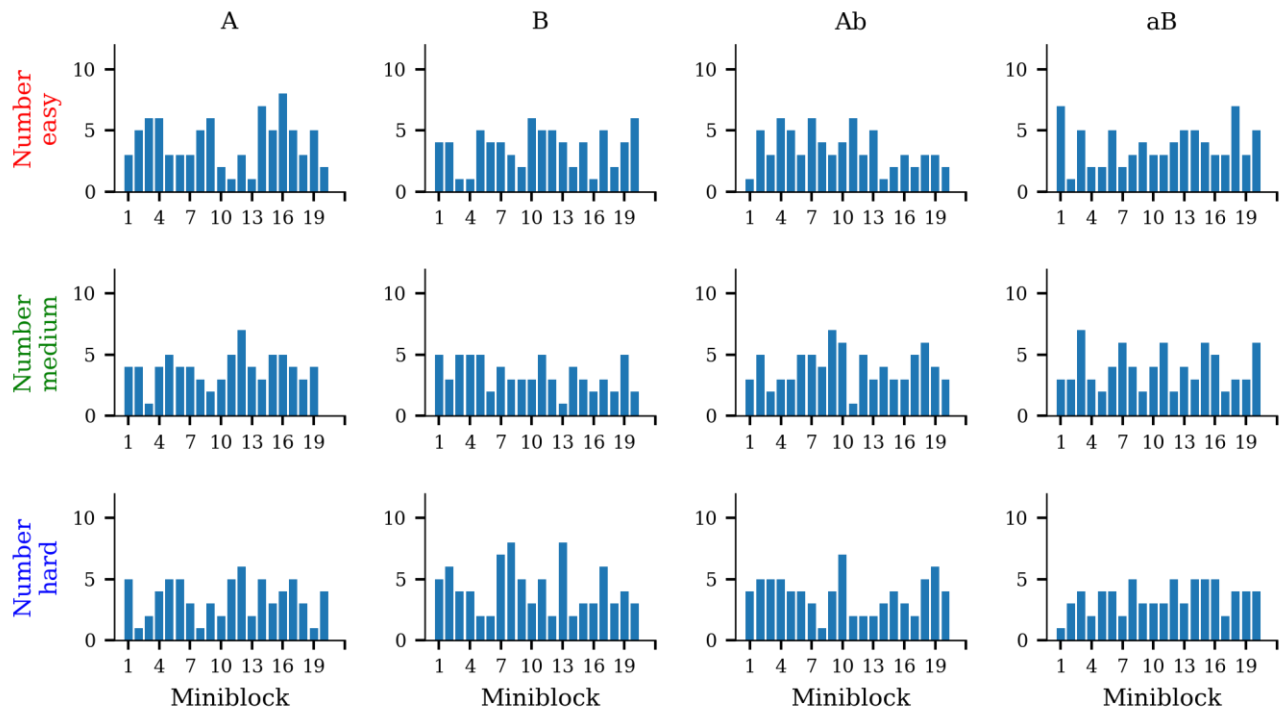
907

908 **S2 Fig. Occurrence of offer types binned with respect to trial.**



909

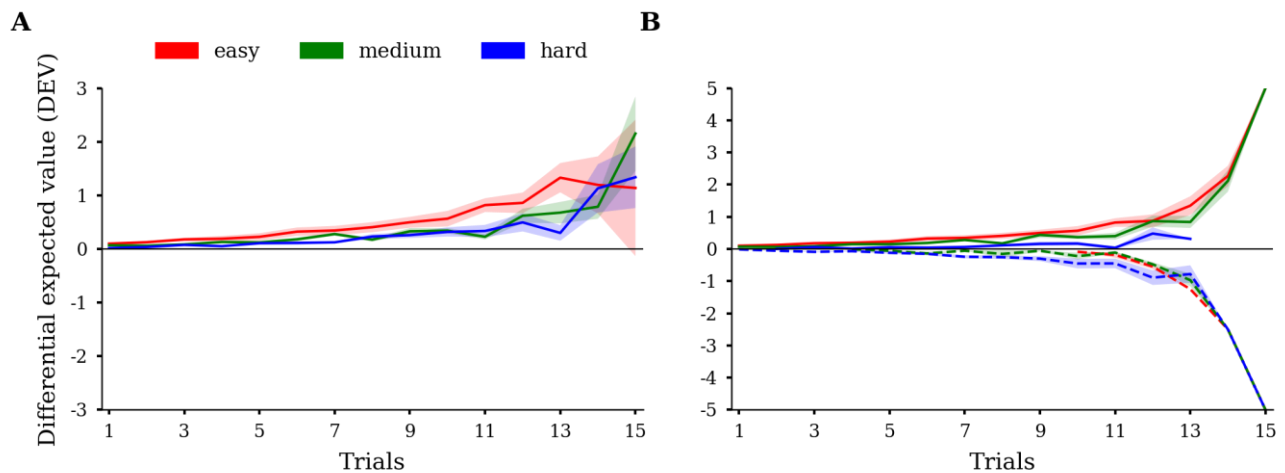
910 **S3 Fig. Occurrence of offer types binned with respect to miniblock.**



911

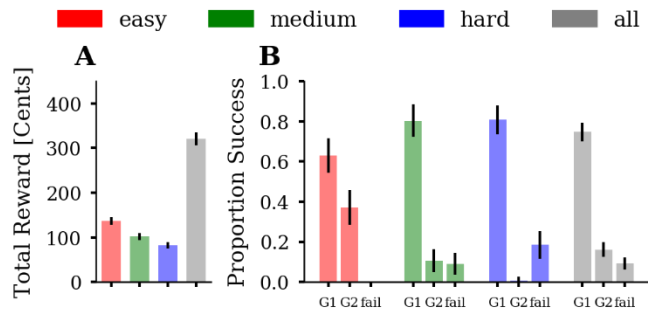
912 **S4 Fig. Occurrence of offer types binned with respect to miniblock and difficulty.**

913



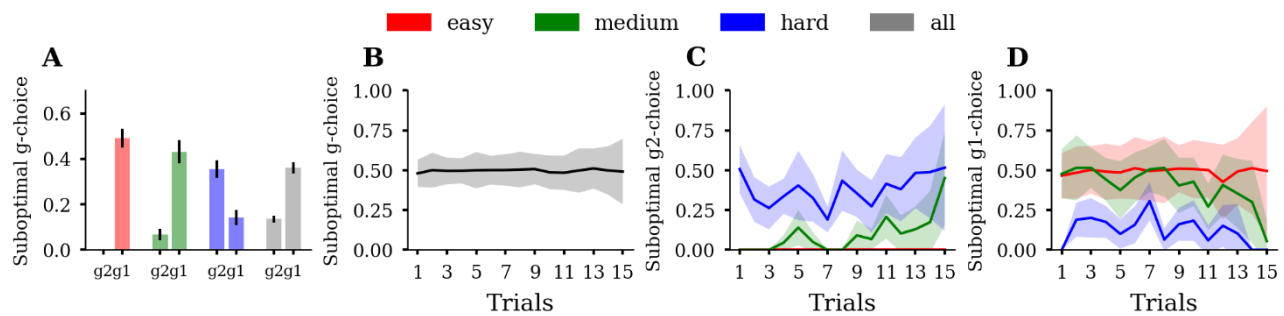
914

915 **S5 Fig. Average absolute (A) and signed (B) differential expected value (DEV) per trial and**
 916 **condition.** Discount and reward ratio had been fixed ($\gamma = 1$, $\kappa = 1$). Average absolute DEVs at the
 917 beginning of the miniblock are smaller than in the end, indicating the relative importance of decisions
 918 close to the final trial of miniblocks. Conditions are colour coded. The shaded areas represent SD.



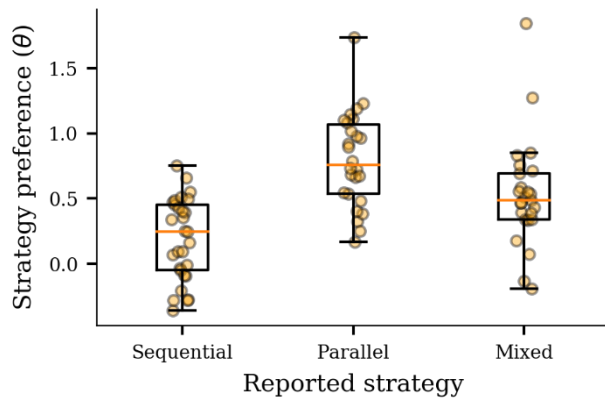
919

920 **S6 Fig. Simulated goal success and total reward of a random agent that always accepts basic**
 921 **offers but guesses for mixed offers ($\theta = 0, \beta \rightarrow 0, \gamma = 1, \kappa = 1$).** (A) Average total reward across
 922 agent instances ($n = 1000$). (B) Proportion of successful goal-reaching, averaged across agent
 923 instances, for each of the three conditions. We plot the proportion of reaching, at the end of a
 924 miniblock, a single goal (G1), both goals (G2), or no goal (fail). The random agent achieves fewer
 925 G2-successes in easy and medium than the participants but fails more often in medium and hard. The
 926 three conditions are colour-coded (easy = red, medium = green, blue = hard) and the average over
 927 conditions is shown in grey. Error bars depict SD.



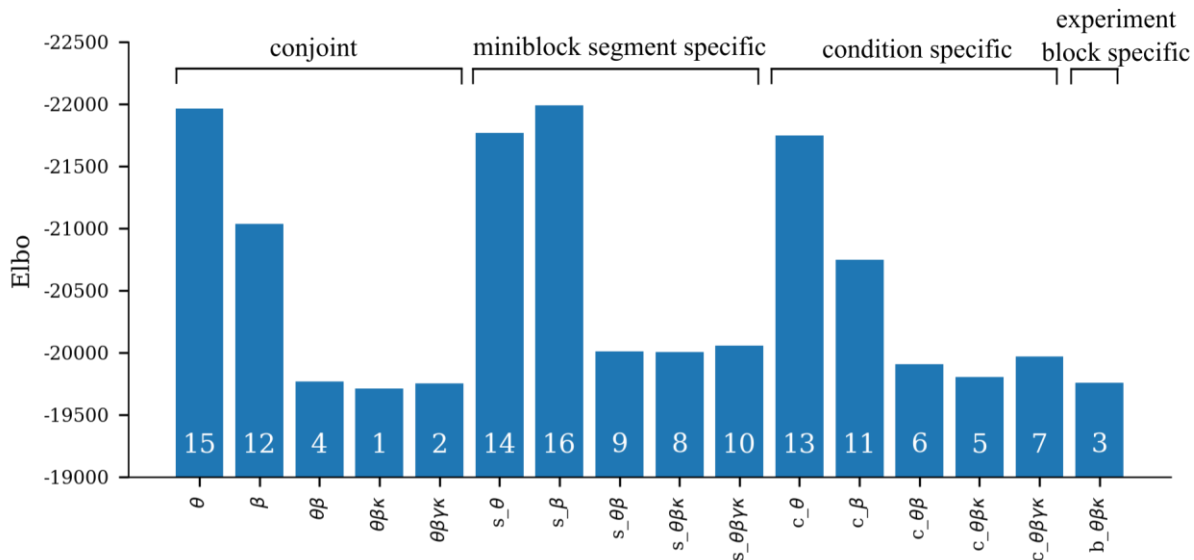
928

929 **S7 Fig. Simulated suboptimal g-choices of a random agent that always accepts basic offers but**
 930 **guesses for mixed offers ($\theta = 0, \beta \rightarrow 0, \gamma = 1, \kappa = 1$).** (A) Proportions of suboptimal g1-choices
 931 (g1) and suboptimal g2-choices (g2), averaged over agent instances ($n = 1000$). The random agent
 932 makes many suboptimal g1-choices in the easy and medium and many suboptimal g2-choices in the
 933 hard conditions. Summing together g1 and g2 yields approximately 50% suboptimal g-choices. (B)
 934 Suboptimal g-choices as a function of trial averaged over agent instances. The random agent makes
 935 approximately 50% suboptimal g-choices across all trials in the miniblock. If participants use non-
 936 random response strategies, i.e. planning or heuristics, their pattern of suboptimality across trials
 937 should deviate from the straight-line pattern of the random agent. (C) Suboptimal g2-choices as a
 938 function of trial averaged over agent instances. (D) Suboptimal g1-choices as a function of trial
 939 averaged over agent instances. Summing together g1 (D) and g2 (C) yields approximately 50%
 940 suboptimal g-choices across trials. Error bars and shaded areas depict SD. Conditions are colour
 941 coded.



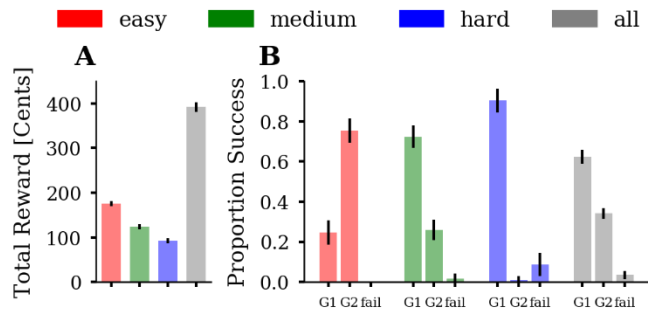
942

943 **S8 Fig. Qualitative comparison of participants' reported strategy use and fitted strategy**
 944 **preference parameter.** Participants who reported the use of a sequential strategy had lower
 945 estimated strategy preference, including the most negative values, than participants who reported the
 946 use of a parallel strategy. Participants who reported mixed use of a parallel and sequential strategy
 947 had greater strategy preference than the sequential group but lower estimates than the parallel group.
 948 The plot shows 80 of 89 participants whose verbal reports matched with one of the three strategy
 949 categories.



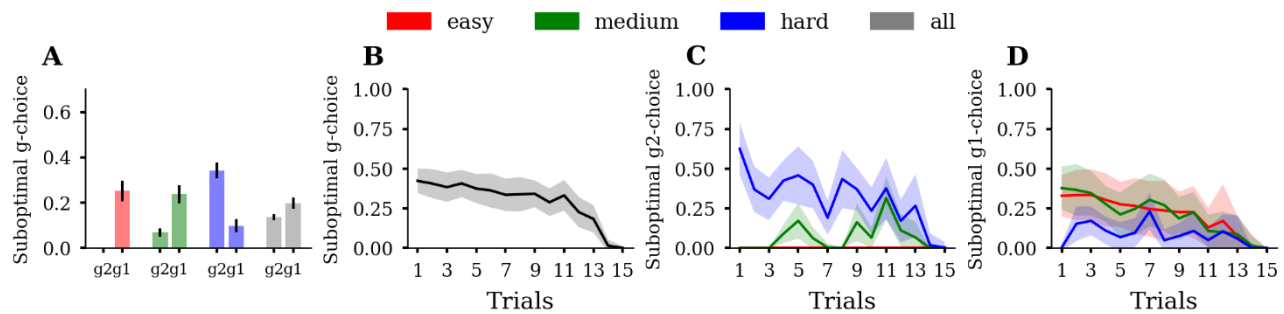
950

951 **S9 Fig. Comparing Elbo (evidence lower bound) between different model variants.** White
 952 numbers represent the rank from highest to lowest Elbo. Model comparisons showed that the three
 953 parameter model (θ, β, κ) had the highest model evidence. Adding γ did not increase model evidence
 954 ($elbo_{\theta\beta\kappa} - elbo_{\theta\beta\gamma\kappa} = -44$). Estimating model parameters separately for miniblock segments (trial
 955 1-5, trial 6-10, trial 11-15; prefix 's_' in the figure) had lower model evidence compared to the
 956 winning model ($elbo_{\theta\beta\kappa} - elbo_{s_\theta\beta\kappa} = -294$). Estimating model parameters separately for
 957 conditions (easy, medium, hard; prefix 'c' in the figure) had lower model evidence compared to the
 958 winning model ($elbo_{\theta\beta\kappa} - elbo_{c_\theta\beta\kappa} = -94$). Estimating model parameters separately for
 959 experiment blocks (miniblock 1-20, miniblock 21-40, miniblock 41-60; prefix 'b' in the figure) had
 960 also lower model evidence compared to the winning model ($elbo_{\theta\beta\kappa} - elbo_{s_\theta\beta\kappa} = -48$). Bars in
 961 the plot depict Elbo averaged over the last 20 posterior samples.



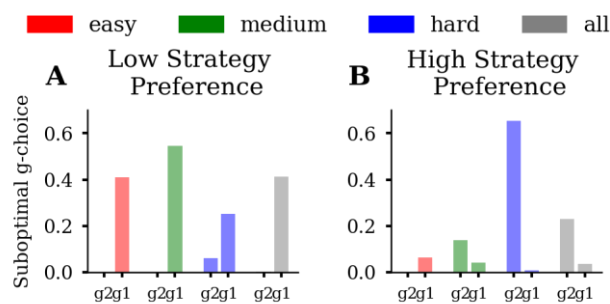
962

963 **S10 Fig. Posterior predictive checks: Simulated goal success and total reward closely resemble**
 964 **observed participant behaviour.** (A) Average total reward across samples ($n = 1,000$). (B)
 965 Proportion of successful goal-reaching, averaged across samples, for each of the three conditions. We
 966 plot the proportion of reaching, at the end of a miniblock, a single goal (G1), both goals (G2), or no
 967 goal (fail). The three conditions are colour-coded (easy = red, medium = green, blue = hard) and the
 968 average over conditions is shown in grey. Error bars depict SD. Data were generated using 1,000
 969 posterior samples from the group hyper parameters.



970

971 **S11 Fig. Posterior predictive checks: Simulated suboptimal g-choices closely resemble**
 972 **observed participant behaviour.** (A) Proportions of suboptimal g1-choices (g1) and suboptimal g2-
 973 choices (g2), averaged over samples ($n = 1,000$). (B) Suboptimal g-choices as a function of trial
 974 averaged over samples. (C) Suboptimal g2-choices as a function of trial averaged over samples. (D)
 975 Suboptimal g1-choices as a function of trial averaged over samples. Error bars and shaded areas
 976 depict SD. Conditions are colour coded. Data were generated using 1,000 posterior samples from the
 977 group hyper parameters.



978

979 **S12 Fig. Comparison of suboptimal g-choices between a low strategy preference and high**
 980 **strategy preference participant.** The plot shows proportions of suboptimal g1-choices (g1) and
 981 suboptimal g2-choices (g2) (A) of the participant with the lowest fitted strategy preference ($\theta =$

982 -0.36) and **(B)** of the participant with the highest fitted strategy preference ($\theta = 1.84$). The low
983 strategy preference participant prefers a sequential strategy leading to suboptimal g1-choices in the
984 easy and medium condition. The participant with a high strategy preference parameter prefers a
985 parallel strategy, resulting in a few suboptimal g1-choices in easy in and medium but a large number
986 of suboptimal g2-choices in the hard condition.

987 **S1 Text: Task instructions (translated from German)**

988 Dear participant,
989 your task in this experiment is to reach goals. Within a block, consisting of 15 trials, you can either
990 reach goal A, goal B or both goals at the same time. For one reached goal you will gain additional 5
991 Cents and for two reached goals additional 10 Cents. Your task is to obtain as much money as
992 possible.

993 To reach goals, you must collect points. You can get points by accepting an offer. Some offers
994 however, might have a negative effect on the state of a goal. Your task is to decide in every trial,
995 whether to accept an offer or wait for the next offer. Press “up arrow” to accept an offer and “down
996 arrow” to wait.

997 Important: Please decide deliberately but speedily. If you decide too slowly, you will get a
998 notification. After every 5 notifications, 50 Cents will be subtracted from your bonus-payout. (The
999 experiment starts with a training phase, in which no money can be lost.)

1000 More about the goals:

1001 Your goal progress will be represented by a bar, which is labelled with A or B. A goal counts as
1002 achieved, if one of the bars reaches or surpasses the white horizontal mark. The goal state will be
1003 evaluated after the end of the 15 trials.

1004 More about the offers:

1005 There are 4 different offers – A, B, Ab and aB. All offers have the same occurrence probability of
1006 25%. The offers differ with respect to their effect on the goal state. A increases the A-bar by one
1007 point. B increases the B-bar by one point. Ab increases the A-bar by one point and subtracts one
1008 point from the B-bar. aB increases the B-bar by one point and subtracts 1 point from the A-bar.

1009 Initial conditions:

1010 At the beginning of the block, you already have some A- and B-points. The amount of initial points
1011 varies from block to block.

1012 **S1 Movie. Simulated goal success and total reward where the precision parameter β varies**
1013 **between 0.25 and 3 with θ , γ , and κ sampled from their fitted population mean. (A)** Average
1014 total reward across agent instances ($n = 1,000$). An increase in β increases total reward obtained in the
1015 easy and medium but decreases total reward in the hard condition. **(B)** Proportion of successful goal-
1016 reaching, averaged across agent instances, for each of the three conditions. We plot the proportion of
1017 reaching, at the end of a miniblock, a single goal (G1), both goals (G2), or no goal (fail). An increase
1018 in β increases G2 success rate in easy and medium but also increases fail rate in medium and hard.
1019 The three conditions are colour-coded (easy = red, medium = green, blue = hard) and the average
1020 over conditions is shown in grey. Error bars depict SD.

1021

1022 **S2 Movie. Simulated suboptimal g-choices where the precision parameter β varies between**
1023 **0.25 and 3 with θ , γ , and κ sampled from their fitted population mean. (A)** Proportions of
1024 suboptimal g1-choices (g1) and suboptimal g2-choices (g2), averaged over agent instances (n
1025 =1000). An increase in β decreases suboptimal g1- and g2-choices. **(B)** Suboptimal g-choices as a
1026 function of trial averaged over agent instances. The influence of β and the associated decrease of
1027 suboptimal g-choices successively increases towards the end of the miniblock. Suboptimal g-choices
1028 in the first half of the miniblock are largely unaffected by the β parameter. **(C)** Suboptimal g2-
1029 choices as a function of trial averaged over agent instances. An increase in β decreases suboptimal
1030 g2-choices late in the miniblock in medium and hard but not in easy. **(D)** Suboptimal g1-choices as a
1031 function of trial averaged over agent instances. An increase in β decreases suboptimal g1-choices late
1032 in the miniblock in easy and medium but not in hard. Error bars and shaded areas depict SD.
1033 Conditions are colour coded.

1034 **S3 Movie. Simulated goal success and total reward where the strategy preference parameter θ**
1035 **varies between -1 and 1 with β , γ , and κ sampled from their fitted population mean. (A)**
1036 Average total reward across agent instances (n =1000). An increase in θ increases total reward
1037 obtained in easy and medium but decreases total reward in hard. **(B)** Proportion of successful goal-
1038 reaching, averaged across agent instances, for each of the three conditions. We plot the proportion of
1039 reaching, at the end of a miniblock, a single goal (G1), both goals (G2), or no goal (fail). An increase
1040 in θ increases G2 success rate in easy and medium but also increases fail rate in medium and hard.
1041 The three conditions are colour-coded (easy = red, medium = green, blue = hard) and the average
1042 over conditions is shown in grey. Error bars depict SD.

1043 **S4 Movie. Simulated suboptimal g-choices where the strategy preference parameter θ varies**
1044 **between -1 and 1 with β , γ , and κ sampled from their fitted population mean. (A)** Proportions of
1045 suboptimal g1-choices (g1) and suboptimal g2-choices (g2), averaged over agent instances (n
1046 =1000). An increase in θ decreases suboptimal g1- choices and increases suboptimal g2-choices.
1047 Suboptimal g1-choices decrease more in easy and medium than in hard. Suboptimal g2-choices
1048 decrease more in hard than in easy and medium. **(B)** Suboptimal g-choices as a function of trial
1049 averaged over agent instances. A change in θ affects the number of suboptimal g-choices made at the
1050 beginning but not at the end of the miniblock. For $\theta > 0$ suboptimal g-choices further decrease,
1051 because g2-choices are often optimal in easy and medium. **(C)** Suboptimal g2-choices as a function
1052 of trial averaged over agent instances. An increase in θ increases suboptimal g2-choices early in the
1053 miniblock, predominantly in the hard condition. **(D)** Suboptimal g1-choices as a function of trial
1054 averaged over agent instances. An increase in θ decreases suboptimal g1-choices early in the
1055 miniblock, predominately in easy and medium. Error bars and shaded areas depict SD. Conditions
1056 are colour coded.

1057 **S1 Notebook. Parameter recovery simulations.**