

1 Revision of XGE-2018-0797R1 as invited by the action editor, Andrew R. A. Conway

2

3 **Abstract**

4

5 In their daily decisions, humans and animals are often confronted with the conflicting choice of

6 opting either for a rewarding familiar option (i.e., exploitation) or of opting for a novel, uncertain option

7 that may, however, yield a better reward in the near future (i.e., exploration). Despite extensive research,

8 the cognitive mechanisms that subtend the manner in which humans solve this *exploration-exploitation*

9 *dilemma* are still poorly understood. In this study, we challenge the popular assumption that exploitation

10 is a global default strategy that must be suppressed by means of cognitive control mechanisms so as to

11 enable exploratory strategies. To do so, we asked participants to engage in a challenging working-

12 memory task while performing repeated choices in a gambling task. Results showed that manipulating

13 cognitive control resources exclusively hindered participants' ability to explore the environment in a

14 directed, intentional manner. Moreover, under certain scenarios, adopting exploitative strategies was also

15 dependent on the availability of cognitive control resources. Additional analyses using a recent

16 computational model of information integration suggests that increasing cognitive load specifically

17 interferes with the ability to combine reward and information in order to inform choices. Our results shed

18 light on the cognitive mechanisms that underpin the resolution of the dilemma, and provide a formal

19 foundation through which to explore pathologies of goal-directed behavior.

20

21 keywords: exploration-exploitation dilemma, informative value, reinforcement learning, cognitive

22 control, adaptive behaviors

23

24

25

26

27

28

29 **Introduction**

30 Understanding the exploration-exploitation dilemma is widely taken to be one of the main
31 challenges in the domain of adaptive control and behavior (Cohen, McClure, & Yu, 2007). The dilemma
32 refers to the fact that when facing a choice, one may either choose to stick with what we know (familiar
33 rewarding outcomes) or engage in the risky exploration of unknown regions of the decision space. To
34 better picture this phenomenon, imagine that it is a nice day in your city. You are walking around
35 downtown in search of a pleasant place to eat. A good strategy would be to choose your favorite
36 restaurant, because the likelihood that you will find it satisfying is very high. However, new dining rooms
37 have recently opened in town. Do you select the restaurant that you know you will enjoy, or do you select
38 another restaurant you have never tried before, potentially finding either a new favorite, or profound
39 disappointment? Thus, the exploration-exploitation trade-off is a dilemma precisely because it involves
40 addressing a challenging conflict between maximizing reward and maximizing information. Solving it is
41 necessary in order to flexibly adapt to environments that are often both uncertain and dynamic. Because
42 all cognitive agents have to somehow address this challenge, the exploitation-exploration is ubiquitous
43 and has relevance for many organisms and for many types of decisions.

44 Although extensive research on the exploration-exploitation dilemma has been conducted over
45 the last decades in different scientific domains (e.g., artificial intelligence, animal foraging and
46 neuroscience), a complete understanding of the underlying mechanisms involved in the resolution of the
47 dilemma is still lacking. In the most popular framework (Daw, O'Doherty, Dayan, Seymour, & Dolan,
48 2006; Cohen et al., 2007), the dilemma is considered as a dual-process where exploitation is a default
49 strategy, and it appears to dominate choice behavior because of its association with stronger reward
50 histories. Following this framework, modifying behavior in an adaptive manner through exploration thus
51 requires overriding the exploitative strategies that tend to dominate the decision process by its stronger

52 association with rewards. To overcome this dominance, behavioral/cognitive control processes might play
53 a central role (i.e., inhibition) in enabling the switch to exploratory strategies (Daw et al., 2006; Cohen et
54 al., 2007). Cognitive control is the ability to coordinate sensory information and actions so as to align
55 them to internal states or intentions (Koechlin, Ody, & Kouneiher, 2003), and is required when the
56 mapping between sensory inputs and actions is rapidly changing or weakly established relative to other
57 existing stimulus-response associations (Miller & Cohen, 2001). Top-down control mechanisms could
58 therefore be the core process that underpins exploratory behavior by enabling the continuous monitoring
59 of the need for behavioral adjustments and by implementing new goal-directed behaviors (Daw et al.,
60 2006; Cohen et al., 2007). The “behavioral control” framework was introduced to explain activity in
61 prefrontal regions (i.e., frontopolar cortex), known to be involved in cognitive control (Mars, Sallet,
62 Rushworth, & Yeung, 2011) during exploratory decisions (Daw et al., 2006). Subsequent evidence has
63 confirmed the core involvement of higher cognitive-control functions in exploration (Badre, Doll, Long,
64 & Frank, 2012; Cavanagh, Figueroa, Cohen, & Frank, 2012; Frank, Doll, Oas-Terpstra, & Moreno,
65 2009).

66 To understand precisely how cognitive control is related to choice behavior in the exploration-
67 exploitation dilemma, it is important to note that, under the above framework, exploitation specifically
68 refers to “choosing the option that maximizes a (reward) prediction”. Exploration, on the other hand, is an
69 umbrella term that encompasses different type of strategies, essentially *random* and *directed* exploration
70 (Wilson, Geana, White, Ludvig, & Cohen, 2014). The concept of random exploration derives from
71 reinforcement learning (RL) theory (Sutton & Barto, 1998), wherein exploration is merely the product of
72 noise in the response-generation process. Under this scenario, a decision-maker who learns to maximize a
73 numerical reward signal may nevertheless make choices associated with lower reward values
74 (exploration) due to a noisy response. In contrast, the concept of directed exploration derives from
75 optimal decision-making theories, which take exploration to be an explicit, goal-directed strategy (Gittins
76 & Jones, 1974). In *directed exploration*, an animal ‘directs’ exploration toward uncertain options, thus

77 increasing its understanding of the surrounding environment through gaining new information. Thus, the
78 absence of information is the main driving factor in this subtype of exploration behavior.

79 Whether humans use information to direct their exploratory behaviors has been a matter of
80 intense discussion over the last decade, and a number of findings have suggested that this was not the case
81 (Daw et al., 2006; Payzan-LeNestour & Bossaerts, 2011). However, this view has been challenged
82 recently in studies using alternative paradigms which controlled for the availability of information in the
83 environment, suggesting that humans may adopt both random and directed exploration (Wilson et al.,
84 2014) the hidden mechanisms of which relate to the integration of reward and information into choice
85 values (Cogliati Dezza, Yu, Cleeremans, & Alexander, 2017). Although based on a common exploratory
86 drive, the two exploratory strategies showed different neural substrates (Warren et al., 2017; Zajkowski,
87 Kossut, & Wilson, 2017), different age-related development (Somerville et al., 2017), and they react
88 differently to changes in reward contingencies (Cogliati Dezza et al., 2017). Thus, the dilemma does not
89 seem to be a unitary binary process but instead a class of problems spanning different scales (Cohen et al.,
90 2007). Following this recent perspective, the dilemma is represented as a continuum (Mehlhorn et al.,
91 2015) where many behaviors fall in the extremes (e.g., choosing the highest valuable option or the most
92 uncertainty option) whereas others might fall somewhere in between (choosing a moderately valuable
93 option associated with some uncertainty). Behavior at these intermediate points on the continuum is less
94 amenable to interpretation, and controlled behavioral paradigms are required (Wilson et al., 2014).
95 Different cognitive mechanisms may therefore underlie the resolution of the dilemma, and the ability of a
96 decision-maker to deploy different exploratory strategies may depend on the availability of sufficient
97 cognitive control resources (Otto, Knox, Markman, & Love, 2014). However, a new framework that
98 attempts to integrate these new advances in understanding the exploration-exploitation dilemma and its
99 underlying cognitive mechanisms is still lacking.

100 Motivated by the behavioral control hypothesis of exploratory behavior and by recent
101 understanding over the resolution of the exploration-exploitation dilemma in humans, we consider
102 whether cognitive control processes might modulate the resolution of the exploration-exploitation

103 dilemma using a mixture of exploratory strategies (i.e., random and directed exploration). We
104 investigated this hypothesis using a variant of bandit tasks that has previously been used to disentangle
105 both random and directed exploratory strategies (Wilson et al., 2014; Cogliati Dezza et al., 2017). Bandit
106 tasks are a family of RL problems where, on each trial, participants must choose among a set of slot
107 machines (or “bandits”) with the goal of maximizing the total reward over a sequence of trials (Robbins,
108 1952). This new version of the bandit task used a two-phase gambling task where, on each game,
109 participants were initially instructed as to which options to choose (*forced-choice task*), after which they
110 were free to choose between options (*free-choice task*) so as to maximize their final gain. By adding a
111 forced-choice task on the top of the standard bandit task, the information participants had about the
112 payoffs of each option was controlled, thereby enabling the identification of the two exploratory strategies
113 in the first free-choice trial of each game (Wilson et al., 2014). In the current study, we additionally
114 manipulated cognitive control resources by asking participants to engage in a challenging working
115 memory task (Konstantinou & Lavie, 2013) while performing the sequential decision-making task. Under
116 the behavioral control hypothesis, depletion of cognitive control resources should lead to a more
117 pronounced expression of processes that operate independently of control such as exploitation, while
118 behaviors that require control - such as exploration - should be attenuated. In order to investigate the
119 effect of cognitive load manipulation on the learning and decision-making components of the dilemma,
120 we developed a computational model that is capable of capturing participants’ behavior on the new
121 version of the bandit task by associating a value with information on top of the standard reward-based
122 reinforcement learning formulation (Cogliati Dezza et al., 2017). Applying a computational model in this
123 context will help in understanding the underlined mechanisms affected by cognitive control manipulation
124 which might be not accessible with a “pure” behavioral analysis

125

126 **Methods**

127

128 *Participants*

129 Twenty-five young adults participated in this study (20 women; aged 18 - 24 years, mean age =
130 19.6). Based on a previous study (Cogliati Dezza et al., 2017), a power analysis suggested a sample size
131 of N=24, power = 0.999. Participants were students at the Faculty of Psychology (Université libre de
132 Bruxelles) and received credits for their participation to the study. The entire group belonged to the
133 Belgian French-speaking community. The experiment was approved by Faculty of Psychology Ethics
134 Committee, and was conducted according to the principles of the Declaration of Helsinki. Informed
135 consent was obtained from all participants prior to the experiment.

136 ***Procedure***

137 *Bandit task*

138 To investigate the effect of cognitive control on the exploration-exploitation dilemma, we asked
139 participants to perform 128 independent games of a new version of the multi-armed bandit task (Figure 1)
140 that has already been shown to elicit both random and directed exploratory strategies (Cogliati Dezza et
141 al., 2017; Wilson et al., 2014). As in standard bandit tasks, in this version, participants chose among
142 options with the goal of maximizing the total reward over a sequence of trials. When selected, each option
143 provides a reward (generated from a hidden distribution) that informs participants about the desirability of
144 each alternative. Contrary to standard bandit tasks, on each game participants performed a *forced-choice*
145 *task* followed by a *free-choice task* (Wilson et al., 2014; Figure 1a). During the forced-choice task,
146 participants were only allowed to select options that had been pre-selected by the computer (Figure 1c),
147 whereas during the free-choice task participants were able to make their own choices in view of
148 maximizing their final score (i.e., the amount of points earned throughout the game) (Figure 1d). Contrary
149 to the first version of this paradigm (Wilson et al., 2014), information regarding the points earned
150 following a choice did not remain visible following a feedback in order to allow learning to influence
151 participants' choices (Cogliati Dezza et al., 2017; Zajkowski et al., 2017). Each game was composed of 6
152 consecutive forced-choice trials and from 1 to 6 free-choice trials (Figure 1a). The number of free-choice

153 trials was manipulated so that participants were unable to predict the length of the free-choice task
154 (Cogliati Dezza et al., 2017) and to adjust their choices accordingly (Wilson et al., 2014).

155 In this version, options were represented as decks of cards and were placed on the left (blue deck),
156 right (green deck) and central (red deck) side of the computer screen (Figure 1b). The use of three options
157 allows us to discern between the strategic use of random and directed exploration (Cogliati Dezza et al.,
158 2017) without manipulating the prior knowledge participants had about horizon (i.e., the total number of
159 trials participants will experience in a game), as in previous versions (Wilson et al., 2014; Krueger, 2017).
160 In particular, if choice probabilities for the two non-exploitative options are equal, then exploratory
161 behavior is entirely driven by random exploration. On the contrary, if the choice probability is different
162 from chance, then choices are partially driven by directed exploration. Participants indicated their choices
163 using the buttons ‘c’, ‘v’ and ‘b’ of the computer keyboard (Figure 1b). After each choice, the card was
164 turned to reveal the points earned by the participant for selecting that deck. Participants could obtain
165 between 1 and 100 points on each trial, and the number of points earned for selecting a deck was sampled
166 from a truncated Gaussian distribution with standard deviation of 8 points (the standard deviation was
167 equal for the 3 decks). The generative mean of each deck was set to 30 and 50 points and adjusted by +/- 0,
168 4, 12, & 20 points to avoid the possibility that participants might be able to distinguish the generative mean
169 for a deck after a single observation (i.e., the generative means ranged from 10 to 70 points). As in our
170 previous study (Cogliati Dezza et al., 2017), the 3 decks of cards had the same generative means in 50% of
171 the games (*equal reward condition*) and different means in the rest of the games (*unequal reward*
172 *condition*); the intent of the different reward conditions in our previous study was to examine the influence
173 of reward context on exploration and exploitation. Although not the primary focus of this study, reward
174 context effects reported in our previous study were also observed here ($p < 10^{-3}$) replicating our previous
175 work. However, in the present study, the effect of reward context was not modulated by the cognitive
176 control manipulation. For this reason, the results concerning reward context will be not discussed any
177 further. The means of the generative Gaussian function were stable within a game and varied between
178 games. Participants were informed that the decks of cards did not change during the same game, but were

179 replaced by new decks at the beginning of each game. However, they were not informed of the reward
180 manipulation and the underlying generative distribution we adopted.

181 As in previous versions of this paradigm, during the forced-choice task we manipulated the
182 information about the decks of cards acquired by participants (i.e., the number of times each deck of cards
183 was played). On each game, participants were forced to either choose each deck 2 times (*equal information*
184 *condition*), or to choose one deck 4 times, another 2 times, and never for the remaining deck (*unequal*
185 *information condition*). The information manipulation guarantees the orthogonalization of reward and
186 information thus allowing the distinction of random and directed exploration in the first free-choice trial of
187 each game (Wilson et al., 2014). In 50% of the games, participants played with the equal information
188 condition. The order of card selection was randomized in both information conditions, as well as the
189 appearance of equal and unequal information condition.

190 *Cognitive control manipulation*

191 Cognitive control resources were manipulated by asking participants to carry out a concurrent
192 working-memory task during the free-choice task. Specifically, we adopted Konstantinou and Lavie's
193 procedure (Konstantinou & Lavie, 2013), which has been shown to selectively interfere with cognitive
194 control processes (Baddeley, Emslie, Kolodny, & Duncan, 1998; D'Esposito, Postle, Ballard, & Lease,
195 1999). Prior to the beginning of the free-choice task, a sequence of 9 digits appeared on the screen (Figure
196 1e). Participants were asked to memorize and retain the sequence until the end of the game. After each
197 free-choice trial, a single memory probe digit was presented at fixation until a response was given. The
198 probe was equally likely to be any of the first 8 digits of the memory set. The participants' task was to
199 report the digit following the probe in the memory sample (e.g., if the memory set was '123456789' and
200 the probe was '3', the correct response would be '4'). The probe was displayed on the screen and
201 participants pressed the key corresponding to the selected digit.

202 In order to investigate the role of cognitive control resources on the exploration-exploitation
203 dilemma, participants were exposed to two different conditions: High Load vs. Low Load. In the High
204 Load condition, the digits were presented in random order (e.g., '371586249') for 2000 ms, and a new

205 sequence was generated on each game. In the Low Load condition, the digits were presented in fixed
206 numerical order (i.e., ‘123456789’) for 500 ms. Participants performed the two conditions on two
207 different days, with order randomized and counterbalanced (half of participants performed the High Load
208 condition on the first day and the Low Load condition the second day, and vice-versa). Performance on
209 the memory task was adopted as an inclusion criterion for the statistical analysis (see Results Section).
210 Due to technical problems, two participants failed to complete the entire 128 games in either the High
211 Load or the Low Load condition, but their data were included anyway since only a few games were
212 lacking (one participant played 124 games of High Load condition and the other 123 of the Low Load
213 condition) and removing those participants did not affect the main results.

214

215

216

217

Insert Figure 1

218

219 ***Computational models***

220 To investigate the hidden mechanisms involved in the resolution of the exploration and
221 exploitation dilemma under cognitive load, we adopted a previously implemented version of a
222 reinforcement learning model that learns reward values on each trial and incorporates a mechanism
223 reflecting the knowledge gained about each deck during previous experience - the gamma-knowledge
224 Reinforcement Learning model (gkRL). The gkRL model is able to reproduce participants’ behavior on
225 the above behavioral paradigm (Cogliati Dezza et al., 2017). Specifically, compared to a standard
226 reinforcement learning model is able to reproduce participants’ directed exploratory strategies in
227 scenarios where options are not sampled at the same rate.

228

229 On each trial, a simple δ learning rule (Rescorla & Wagner, 1972) is used to compute the
 230 expected reward value $Q(c)$ for each deck of cards c (= Left, Central or Right), using the following
 231 equation:

$$232 \quad Q_{t+1,j}(c) = Q_{t,j}(c) + \alpha \times \delta_{t,j} \quad (1)$$

233 where $Q_{t,j}(c)$ is the expected reward value for trial t and game j . $\delta_{t,j} = R_{t,j}(c) - Q_{t,j}(c)$ is the
 234 *prediction error*, which quantifies the discrepancy between the predicted outcome and the actual outcome
 235 obtained at trial t and game j . The expected reward $Q_{t,j}(c)$ is updated using the above rule only if an
 236 outcome from the deck c is observed, otherwise $Q_{t+1,j}(c) = Q_{t,j}(c)$. Considering participants were told
 237 that each game was independent from the others, Q_0 is initialized at the beginning of each game
 238 (Khamassi, Enel, Dominey, & Procyk, 2013) and set to the global estimate of Q (~ 40 points) (Cogliati
 239 Dezza et al., 2017).

240 Additionally, gkRL tracks information gained from each deck based on how often it is selected,
 241 as follows:

$$242 \quad I_{t,j}(c) = \left(\sum_1^t i_{t,j}(c) \right)^\gamma$$

$$243 \quad \text{where, } i_{t,j}(c) = \begin{cases} 0, & \text{choice} \neq c \\ 1, & \text{choice} = c \end{cases} \quad (2)$$

244 $I_{t,j}(c)$, is the amount of information associated with the deck c at trial t and game j . $I_{t,j}(c)$, is computed
 245 by including an exponential term γ that defines the degree of non-linearity in the amount of observations
 246 obtained from options after each observation. γ is constrained to be > 0 . Each time deck c is selected,
 247 $i_{t,j}(c)$ takes value of 1, and 0 otherwise. On each trial, the new value of $i_{t,j}(c)$ is summed to the previous
 248 $i_{t-1,j}(c)$ values and the resulting value is raised to γ , resulting in $I_{t,j}(c)$. For example, after six trials of
 249 the forced choice task, if one option has never been selected, $I_{t,j}(c)$ has value zero; whereas in the case
 250 that one option is selected 4 times, $I_{t,j}(c)$ has the value 4^γ . The parameter γ adds non-linearity to the
 251 information term (Cogliati Dezza et al., 2017), the intuition being that additional samples do not

252 contribute equally to the amount of information a subject has about an option (e.g., sampling an option
253 you have never observed is far more informative than sampling an option you have observed 100 times
254 previously).

255 Before selecting the appropriate option, gkRL subtracts the information gained $I_{t,j}(c)$ from the
256 expected reward value $Q_{t,j}(c)$:

257
$$V_{t,j}(c) = Q_{t,j}(c) - I_{t,j}(c) * \omega \quad (3)$$

258 $V_{t,j}(c)$ is the final value associated with deck c . Here, information accumulated during the past trials
259 scales values $V_{t,j}(c)$ so that increasing the number of observations of one option decreases its final value.
260 In other words, when one option is over-selected, $I_{t,j}(c)$ becomes larger resulting in lower $V_{t,j}(c)$. On the
261 contrary, if one option is never-selected $I_{t,j}(c)$ is zero, and $V_{t,j}(c) = Q_{t+1,j}(c)$. ω is the information
262 weight and determines the degree by which the model integrates information into choice values. In order
263 to generate choice probabilities based on expected reward values, the model uses a softmax choice
264 function (Daw et al., 2006; Humphries, Khamassi, & Gurney, 2012; Wilson & Niv, 2011). The softmax
265 rule is expressed as:

266
$$P(c/V_{t,j}(c_i)) = \frac{\exp(\beta \times V_{t,j}(c))}{\sum_i \exp(\beta \times V_{t,j}(c_i))} \quad (4)$$

267 where β is the inverse temperature that determines the degree to which choices are directed toward the
268 highest rewarded option. With higher β , the model mainly selects options associated with higher choice
269 value, whereas with lower β , the model's choices are more random.

270 The gkRL model can be informative concerning the effect of cognitive load on the dilemma in 2
271 ways. First, it can help to distinguish whether cognitive load effects on exploration are driven by
272 information computation (ω and γ), or whether they are instead driven by changes in choice variability
273 (β). Second, if changes are driven by alterations in information computation, the model can help to
274 distinguish whether these are driven by changes in information integration (ω) or by changes in the way
275 information availability decays with time (γ).

276

277 ***Model fitting and model comparison***

278 To estimate the model's parameters α , β , and ω , γ , we collected trial-by-trial participants' choices
279 in both High and Low Load condition (Table 1, Table2). During the fitting procedure, the objective
280 function - the negative log likelihood - $\sum_{j=1}^{128} \log (P_j(c))$ - for each participant under both load
281 conditions was computed and then minimized using MATLAB and Statistics Toolbox Release 2015b
282 function *fminsearchbnd* (which is exactly as *fminsearch* but does not search outside the fixed boundaries).
283 The boundaries adopted were as follows: α]0,1[, β]0, 10], ω [-300, 300], γ]0, 12]. To increase the
284 likelihood of finding a global rather than a local optimum, *fminsearchbnd* was iterated with 15 randomly
285 chosen starting points. The fitting procedure was validated by running a recovery analysis: the gkRL
286 model was simulated on the task using the retrieved parameter estimates to generate synthetic behavioral
287 data and then the fitting procedure was applied to the synthetic data in order to check whether previously
288 estimated parameters were indeed recovered ($r^2 > 0.4$). Likewise, we checked the model comparison
289 outcome by computing a confusion matrix and checking whether data generated from a model was indeed
290 best explained by that model.

291

292 ***Statistical analysis***

293 Statistical analysis was performed using RStudio (<https://www.rstudio.com/>), functions and
294 packages adopted are reported in the results section. To determine whether and how manipulating
295 cognitive control affected participants' decision strategies, we conducted repeated measure ANOVA
296 analyses. When violations of parametric tests were indicated, non-parametric tests were performed. *P*-
297 values of < 0.05 were considered significant.

298

299 **Results**

300 In this section, we first report the results concerning the cognitive load manipulation we adopted
301 and its effects on participants' performance. Subsequently, we examine the interaction between cognitive
302 load manipulation and decision strategies. Lastly, we investigate the possible hidden mechanisms affected
303 by manipulating cognitive/behavioral control mechanisms.

304 *Working memory Task*

305 First, we explored the effect of the cognitive load manipulation on memory accuracy. To do so,
306 trial-by-trial correct memory responses were collected. A Wilcoxon Signed Rank Test on the average
307 value of subjects' overall correct memory responses revealed a significant difference between High Load
308 ($M = 0.494$, $SD = 0.12$) and Low Load ($M = 0.986$, $SD = 0.012$), $p < 10^{-8}$, $r = .874$, indicating that, as
309 expected, increasing memory load affected participants' performance on the working memory task
310 (Figure 2a). Because it can be assumed that participants who scored at chance level on memory
311 performance were not reliably engaged in the memory task, accuracy on the memory task was used as an
312 inclusion criterion for further statistical analysis. A one-sample T-test on correct memory responses
313 revealed a significant difference between the High Load condition and chance level (12.5 %), $t(24) =$
314 15.29 , $p < 10^{-14}$, $d = 4.33$, suggesting that participants on averaged were actively engaged in the working-
315 memory task. Additionally, we investigated whether each participant performed above chance-level by
316 applying a one-sample sign test on participants' correct memory responses in the High Load condition.
317 Results revealed that each participant scored above chance level $p < 10^{-6}$. Following this result every
318 participant was included in the subsequent analysis.

319 *Cognitive load manipulation*

320 To check whether the cognitive load manipulation affected cognitive control processes by
321 increasing dual-task interference, we measured choice reaction times (RTs) during the free choice-task of
322 both High and Low Load condition (Figure 2b). A paired T-test on RTs revealed slower reaction times
323 during High Load condition ($M = 1005$ ms, $SD = 468$ ms) compared to Low Load condition ($M = 483$
324 ms, $SD = 145$ ms), $t(24) = 6.19$, $p < 10^{-6}$, $d = 1.24$, suggesting that less cognitive control resources were
325 available to participants during the High Load manipulation.

326

327

328

329

Insert Figure 2

330

331 *Performance*

332 We also examined whether the cognitive load manipulation affected the way participants
333 performed the gambling task. Here, performance refers to the ability to play strategically during the task
334 in order to maximize the total gain. To do so, we computed the probability of choosing the deck with the
335 highest average of points obtained in the previous trials (overall performance) during the entire free-
336 choice task under all reward conditions in both High Load and Low Load conditions. A Wilcoxon Signed
337 Rank Test on the average values of overall performance revealed a decrease in the High Load condition
338 ($M = 0.586$, $SD = 0.109$) compared to the Low Load condition ($M = 0.617$, $SD = 0.098$), $Z = 2.08$ $p =$
339 $.036$, $r = .417$, suggesting that loading cognitive control resources made it more difficult for participants
340 to retrieve previously learned information and act strategically. However, in both conditions all
341 participants scored above chance level set at 33%. A Wilcoxon Signed Rank Test on the average value of
342 participants' overall performance revealed a significant difference between choosing the deck associated
343 with the highest average points during the High Load condition and chance level, $p < 10^{-7}$, and between
344 choosing the deck associated with the highest averaged points during the Low Load condition and chance
345 level, $p < 10^{-7}$, indicating that participants played strategically during both load conditions.

346 *Cognitive control and decision strategies*

347 To investigate whether cognitive control plays a role in the resolution of the exploration-
348 exploitation dilemma, we first measured decision strategies when participants selected options unequally
349 during the forced-choice task (unequal information condition) in both the High and the Low Load
350 conditions (Figure 3a). We conducted the analysis on the first free-choice trial, being the only trial where
351 a clear distinction between random and directed exploration can be obtained (Wilson et al., 2014). Trials
352 were classified as “directed exploratory” when participants chose the option that had never been sampled

353 during forced-choice trials, as “exploitative” when participants chose the experienced deck with the
354 highest average of points (regardless of the number of times that deck had been selected during the
355 forced-choice task) and as “random exploratory” when the classification did not meet the previous
356 criteria. The sum of the 3 strategies defined the total choice probability equal 1 (choice probability=
357 probability to exploit + probability to random explore + probability to directed explore =1). We
358 conducted a 2 (condition: High Load, Low Load) by 3 (strategies: exploitation, random exploration,
359 directed exploration) non-parametric ANOVA. The test allows the use of two-way repeated measure
360 ANOVA in a non-parametric setting using aligned rank transformation (e.g., ART package in R; Conover
361 & Iman, 1981). Results showed an effect of strategy $F(2,120) = 44.83, p < 10^{-15}$, partial eta-squared =
362 0.428, and a condition X strategy interaction $F(2,120) = 5.87, p = .004$, partial eta-squared = 0.089 (Figure
363 3a). The effect of condition did not reach the significant threshold, $p = .974$. Post-hoc comparisons
364 showed an increase in random exploration in the High load condition ($M = 0.202, SD = 0.123$) compared
365 to the Low Load condition ($M = 0.13, SD = 0.09$), $p = .006$; a decrease in directed exploration in the High
366 Load condition ($M = 0.338, SD = 0.177$) compared to the Low Load condition ($M = 0.473, SD = 0.198$),
367 $p = .0012$; and an increase in exploitation in the High Load condition ($M = 0.459, SD = 0.149$) compared
368 to the Low Load condition ($M = 0.397, SD = 0.151$), $p = .031$.

369 The above analysis appears to suggest that the effect of cognitive load manipulation affected
370 directed and random exploration in an opposite fashion: directed exploration decreased, whereas random
371 exploration increased under High Load compared to Low Load condition. However, in the unequal
372 information condition, trials labelled as random exploration correspond to the deck of cards that is
373 sampled either twice or 4 times during the forced-choice task. Therefore, in this condition trials labelled
374 as random exploration might be confounded with information-based processing (i.e., when subjects select
375 the option observed twice during the forced-choice task). In order to gain insight into this issue we
376 conducted two additional analyses: 1) In the unequal information condition we repeated the above 2X3
377 ANOVA, but only for trials where random exploratory trials were those associated with the deck of
378 cards sampled 4 times during the forced-choice task; 2) we investigated participants’ behavior in the

379 equal information condition where random exploration was not confounded with the number of
380 observations of each option (being the outcomes of the 3 options equally experienced; Wilson et al.,
381 2014). In the first analysis, we labelled trials as exploitative when the option was associated with highest
382 reward and selected twice during the forced-choice task, random-exploratory when the option was
383 associated with lowest reward and selected 4 times during the forced-choice task and directed exploratory
384 as previously described. Next, we conducted a 2 (condition: High Load, Low Load) by 3 (strategies:
385 exploitation(2seen), random exploration(4seen), directed exploration) non-parametric ANOVA. Results
386 showed an effect of strategy $F(2,120) = 79.8, p < 10^{-15}$, partial eta-squared = 0.57, and a condition X
387 strategy interaction $F(2,120) = 7.48, p < 10^{-3}$, partial eta-squared = 0.111, whereas the effect of condition
388 was not significant, $p = .137$. Post-hoc comparison revealed an increase in random exploration in the
389 High Load condition ($M = 0.119, SD = 0.073$) compared to Low Load ($M = 0.08, SD = 0.053$), $p = .025$,
390 whereas exploitation did not differ. Results concerning directed exploration are already reported in the
391 previous paragraph. In the second analysis, we investigated the effect of cognitive load manipulation on
392 decision strategies when participants were forced to equally select options (equal information condition;
393 Figure 3b). We classified choices as “exploitative” when participants chose the experienced deck with the
394 highest average of points and “random explorative” otherwise. A 2 (condition: High-load, Low-load) by 2
395 (strategy: exploitation, random exploration) non-parametric ANOVA on participants’ choices showed an
396 effect of strategy $F(1,45) = 64.06, p < 10^{-10}$, partial eta-squared = 0.587, and a condition X strategy
397 interaction $F(1,45) = 5.9, p = .019$, partial eta-squared = 0.116. Post-hoc comparisons revealed an increase
398 in random exploration in the High Load condition ($M = 0.366, SD = 0.155$) compared to the Low Load
399 condition ($M = 0.273, SD = 0.129$), $p = .0009$; and a decrease in exploitation in the High Load condition
400 ($M = 0.633, SD = 0.155$) compared to the Low Load condition ($M = 0.725, SD = 0.129$), $p = .001$. Taken
401 together, these analyses confirm that cognitive control manipulation affected the two exploratory
402 strategies in a different fashion.

403

404

405 Insert Figure 3

406 -----

407

408 Subsequently, we investigated whether the above results could have been driven by an ineffective
 409 High Load manipulation in trials where participants incorrectly performed the working-memory task. To
 410 do so, we compared RTs from correct and incorrect memory trials during the High Load condition. If the
 411 behavioral pattern observed above was driven by an ineffective load manipulation during incorrect
 412 memory trials, participants should have shown differences in their RTs as a function of memory accuracy.
 413 We computed participants' RTs during correct and incorrect memory trials and compared the average
 414 values. A Wilcoxon Signed Rank Test on choice RTs showed no differences between correct ($M = 1023$
 415 ms , $SD = 562 ms$) and incorrect memory trials ($M = 993 ms$, $SD = 423 ms$) in all free-choice trials, Z
 416 $= 0.04$, $p = .979$, $r = .008$, and a marginal difference in the first free-choice trials (correct: $M = 2036.8 ms$,
 417 $SD = 1831.6 ms$; incorrect: $M = 2116.3 ms$, $SD = 1425.2 ms$), $Z = -1.95$, $p = .051$, $r = -.39$. However, this
 418 marginal difference was in the direction of higher RTs for incorrect trials as participants were taking more
 419 time to retrieve incorrectly memorized sequence. Overall, these results suggest that even if participants
 420 were not correctly performing the memory task, they were still in a “loaded state” during the High Load
 421 condition suggesting that the observed effects on the decision strategies were a direct consequence of
 422 lowering cognitive control resources.

423 *Randomness vs. information integration under cognitive load*

424 Our previous analysis showed that manipulating cognitive control resources affected how
 425 participants balanced the exploration-exploitation dilemma, exploring more randomly overall during high
 426 working memory load conditions. In this section, we asked whether this effect was due to an increase in
 427 the randomness in participants' choices, or whether this effect was due to alterations in reward and
 428 information processing that subtend the resolution of the dilemma through directed exploration (Cogliati
 429 Dezza et al., 2017). To better investigate the mechanisms affected by the load manipulation, we fit the
 430 gkRL model to all participants' first free choices during both the High and Low Load condition to obtain

Running head: cognitive control and information integration in the dilemma

431 the estimates of the values of the following parameters: learning rate α , inverse of the temperature β , the
432 non-linear parameter γ and the information parameter ω (Table 1). We then compared the estimated
433 parameters for the High Load condition with the parameters of the Low Load condition to investigate the
434 effect of the cognitive control manipulation. As expected, because the learning processes during the
435 forced-choice task were not affected, a Wilcoxon Signed Rank Test on the learning rate α showed no
436 difference between the Low Load condition ($M = 0.426$, $SD = 0.249$) and the High Load condition ($M =$
437 0.456 , $SD = 0.3$), $Z = -0.4171$, $p = .691$, $r = -.083$. Furthermore, a Wilcoxon Signed Rank Test on the
438 inverse temperature parameter β showed no difference between the Low Load ($M = 0.606$, $SD = 1.383$)
439 and the High Load condition ($M = 0.865$, $SD = 1.8$), $Z = 0.094$, $p = .937$, $r = .019$. Additionally, a
440 Wilcoxon Signed Rank Test on the parameter γ showed no difference between the Low Load condition
441 ($M = 1.66$, $SD = 3.44$) and the High Load condition ($M = 1.44$, $SD = 2.62$), $Z = -0.444$, $p = .672$, $r = -$
442 $.089$. On the contrary, the information parameter ω showed a decrease in the High Load ($M = -2.81$, SD
443 $= 51.76$) compared to the Low Load condition ($M = 4.82$, $SD = 9.89$), $Z = -2.058$, $p = .039$, $r = -.412$,
444 suggesting that the increase in random exploration was due to an inability to integrate the learned
445 information into a choice value, rather than an increase in the randomness of participants' choices or by
446 alteration in how information is decay with time (Figure 4a).

447 Furthermore, we fitted the model to all free-choice trials so as to have a more comprehensive
448 view over the underlying process as well as to obtain a better estimate of the parameter values due to the
449 higher number of data points (Table 2). As before, a Wilcoxon Signed Rank Test showed no difference
450 between the Low Load and the High Load conditions, neither for the inverse temperature parameter β (M
451 $= 0.344$, $SD = 0.737$; $M = 0.424$, $SD = 0.81$), $Z = 0.202$, $p = .853$, $r = .04$, nor for the γ parameter ($M =$
452 2.247 , $SD = 3.201$; $M = 2.771$, $SD = 3.448$), $Z = -0.336$, $p = .751$, $r = -.067$. Again, a Wilcoxon Signed
453 Rank Test on the information parameter ω did reveal a decrease in information integration from Low
454 Load ($M = 5.19$, $SD = 9$) to High Load ($M = -1.993$, $SD = 9.3$), $Z = -3.35$, $p = .0003$, $r = -.67$. However,
455 the same test applied to the learning rate α revealed a decrease in the speed of integration of new reward

456 information from Low Load ($M = 0.495$, $SD = 0.198$) to High Load ($M = 0.312$, $SD = 0.186$), $Z = 3.108$,
457 $p = .001$, $r = -.621$ (Figure 4b). The effect on learning rate in this analysis is explained by the fact that we
458 considered all free-choice trials during which participants were performing the memory task while
459 repeatedly selecting options. As a consequence, the ability to integrate new reward information
460 (expressed by the learning rate) was also affected.

461 As an additional check, we fit the gkRL model exclusively on the free-choices trials where
462 memory responses were correct in both Low and High Load conditions. Wilcoxon Signed Rank Tests
463 confirmed our previous results: no differences in parameter β , $Z = -0.525$, $p = .615$, $r = -.105$, and
464 parameter γ between Low and High Load condition, $Z = -1.0$, $p = .325$, $r = -.202$, whereas a higher α was
465 observed in the Low Load condition ($M = 0.493$, $SD = 0.229$) compared to the High Load condition ($M =$
466 0.321 , $SD = 0.248$), $Z = 2.516$, $p = 0.001$, $r = .503$. A higher information parameter ω was also obtained
467 in the Low Load condition ($M = 6.02$, $SD = 8.9$) compared to the High Load condition ($M = -0.827$, $SD =$
468 7.63), $Z = 3.4$, $p = .0002$, $r = .686$ (Figure 4c).

469 -----
470 Insert Figure 4
471 -----

472 *Cognitive control and information integration*

473 Following the above results, cognitive load seems to affect participants' ability to integrate
474 learned information into choice values in order to solve exploration-exploitation problems. As a further
475 investigation, we asked whether a standard reinforcement learning (sRL) model that learned reward
476 values following equation (1) and entered directly in equation (4) without any integration of information,
477 could better explain this 'inability' in integrating information during cognitive control manipulation. To
478 do so, we compared fits of both the gkRL model and sRL model. During the fitting procedure, we
479 computed negative-log likelihoods of both models and their model evidence (or the log model evidence -
480 the probability of obtaining the observed data given a particular model). We adopted an approximation to
481 the (log) model evidence, namely the Bayesian Information Criterion (BIC; (Schwarz, 1978)). We

482 conducted a frequentist analysis with BIC values of the two models (fitted to the first free-choice trials)
483 entered into a t test. Results showed that during the Low Load condition the gkRL model ($BIC_{gkRL} = 184$)
484 best represented participants' data compared to sRL ($BIC_{sRL} = 203$), $t(24) = -3.034$, $p = .005$, $d = 0.455$,
485 replicating our previous findings on reward and information integration during this new version of the
486 bandit task (Cogliati Dezza et al., 2017). However, in the High Load condition neither the gkRL model
487 ($BIC_{gkRL} = 218$) nor the sRL model ($BIC_{sRL} = 214$) better represented participants' data, $t(24) = 0.437$, $p =$
488 $.666$, $d = -0.076$. To better understand this point, we individually investigated the BIC values of each
489 model (Figure 5). While in the Low Load condition the performance of the majority of participant was
490 better explained by the gkRL model (Figure 5a), in the High Load condition approximately 60 % of
491 participants were better represented by the sRL model (whereas the behavior of 20% were better
492 explained by the gkRL model, and 20% were equally explained by both models, Figure 5b), confirming
493 that during the High Load condition information processing was heavily compromised and that for the
494 majority of subjects the computation of information was nullified. Furthermore, we extended the
495 comparison of the two computational models to all free-choice trials to have an exhaustive understanding
496 of the hidden processes. Contrary to our previous model comparison in the High Load condition, results
497 showed that, when fit to all free-choice trials, the gkRL model ($BIC_{gkRL} = 802$) best represented
498 participants' data compared to sRL ($BIC_{sRL} = 849$), $t(24) = -3.4$, $p = .002$, $d = 0.258$ (we obtained the same
499 result in the Low Load condition so the results are not reported here).

500 A possible reason behind this apparently incoherent result could be related to the working
501 memory process itself. The memory sequence was presented to participants at the beginning of the free-
502 choice task only: cognitive load may be reduced during later free-choice trials either as a consequence of
503 inability to maintain the complete sequence over the course of the free-choice task (and thus freeing
504 cognitive resources for making choices), or because cognitive demands related to maintaining the
505 sequence are higher immediately following the presentation of the sequence. We therefore investigated
506 participants' behavior during all free-choice trials to better clarify this point. However, after the first free-
507 choice trial it is not possible to distinguish between random and directed exploration due to a confound

508 between reward and information (Wilson et al., 2014). For this reason, in order to investigate participants'
509 behavior during the all free-choice trials, we focused on information-based processes only. To do so, we
510 computed the probability of selecting the least-seen deck (the option visited the least number of times in
511 the previous trials), the most-seen deck (the option visited the most number of times in the previous trials)
512 and the middle-seen deck (when previous criteria did not match) during both load conditions. When we
513 investigated the behavior globally, the analysis gave similar results observed in the previous behavioral
514 analysis (section *Cognitive control and decision strategies*), where the probability of selecting the least-
515 seen option was reduced during the High Load condition ($M = 0.255$, $SD = 0.099$) compared to the Low
516 Load condition ($M = 0.304$, $SD = 0.071$, $Z = -2.652$, $p = .006$, $r = -.53$), whereas the most-seen showed
517 the opposite pattern an increase in the High Load ($M = 0.573$, $SD = 0.121$) compared to the Low Load
518 condition ($M = 0.531$, $SD = 0.092$), $Z = 2.18$, $p = .028$, $r = .436$ (the probability of choosing the middle
519 seen option did not differ and so we will not consider this strategy any further- Figure 6a). However,
520 investigating the trial-by-trial probability revealed a more exhaustive view. Indeed, the above result was
521 true only for the first 3 free-choice trials (all p values $< 10^{-2}$), whereas we did not observe differences in
522 terms of the most-seen and least-seen options during the last 3 trials (all $p > 0.05$) (Figure 6b). These
523 results suggest that the effect of cognitive load was greatest during the first free-choice trials and vanished
524 during the last trials suggesting that the reason behind the better performance of gKRL compared to sRL
525 in explaining all participants' free choices was due to a decrease in cognitive load in the last trials of each
526 game. Considering that the above analyses focused on information only, it is possible that additional
527 factors may inform choice behavior in free choice trials. To examine this, we also computed switch/stay
528 probabilities for free choice trials. Switch/stay behavior changed in the High Load ($M_{\text{switch}} = 0.416$ SD_{switch}
529 $= 0.165$; $M_{\text{stay}} = 0.584$ $SD_{\text{stay}} = 0.165$) compared to Low Load Condition ($M_{\text{switch}} = 0.476$ $SD_{\text{switch}} = 0.118$;
530 $M_{\text{stay}} = 0.533$ $SD_{\text{stay}} = 0.118$), both $p = .042$. However, differences in switch/stay behavior were most
531 apparent on the first free-choice trial – subjects tended to switch choices, but did so more often in the low-
532 load condition (Figure 7). Results showed that in the last trials of each game stay (switch) probability did

533 not change between High Load and Low Load condition (all $p > 0.05$) confirming that a decrease in
534 cognitive load occurred in the last trials of each game.

535 -----

536 Insert Figure 5

537 -----

538 -----

539 Insert Figure 6

540 -----

541 -----

542 Insert Figure 7

543 -----

544 *From the model to behavior*

545 The above results demonstrate that the gkRL model better accounts for our behavioral data
546 relative to sRL. In order to demonstrate that the gkRL model parameters are behaviorally-relevant, we
547 correlated the differences observed between the two load conditions in the information integration
548 parameter ω with the differences in exploitation in the unequal information condition. If the model
549 captures behavioral dynamics, we should expect increased differences between the estimate of parameter
550 ω in the two load conditions as well as increased differences in exploitation between the two load
551 conditions. We observed a positive correlation between the difference in ω and the difference in
552 exploitation (Pearson correlation $r(23) = .413, p = .039$) suggesting that reduction of the integration of new
553 information was associated with increased exploitative behaviors. Additionally, simulations of the model
554 are also able to reproduce the condition-dependent behavioral results we observe in our data. We
555 simulated the gkRL model 80 times under the two loading conditions. In the High Load condition
556 ω values were randomly drawn from a uniform distribution with mean -2, whereas for the Low Load
557 condition the mean was set to 5. The other parameters did not change between the two conditions and
558 their values were randomly chosen from a uniform distribution with mean set around the mean values
559 observed in participants' data. We then labeled model's choices in the unequal information as directed
560 exploratory, random exploratory and exploitative. We conducted a 2 (condition: High Load, Low Load)

561 by 3 (strategies: exploitation, random exploration, directed exploration) non-parametric ANOVA. Results
562 showed an effect of strategy $F(2,395) = 223.04$, $p < 10^{-15}$, partial eta-squared = 0.53, and a condition X
563 strategy interaction $F(2,395) = 240.52$, $p < 10^{-15}$, partial eta-squared = 0.549. The effect of condition did
564 not reach the significant threshold, $p = .5$. The results mimicked the same behavioral pattern observed in
565 participants' data (Figure 8a). Additionally, we computed random exploration and exploitation in the
566 equal information condition. We conducted a 2 (condition: High Load, Low Load) by 2 (strategies:
567 exploitation, random exploration) non-parametric ANOVA. Results showed an effect of strategy $F(2,237)$
568 = 382.89, $p < 10^{-15}$, partial eta-squared = 0.617, however neither an effect of condition X strategy and an
569 effect of condition was observed (all $p > 0.05$) (Figure 8b). The behavior of the model in the equal
570 information condition, however, did not replicate the findings observed in participants' data. We better
571 discuss this result in the next section.

572 -----
573 Insert Figure 8
574 -----

575

576 *Cognitive control and value degradation*

577 In order to understand the underlying mechanisms affected in the equal information condition that
578 are not captured by the information integration account expressed by the gkRL model, we implemented a
579 new version of the gkRL model- the value gamma knowledge RL (vgkRL). The rationale behind this
580 additional model implementation is that cognitive load might have affected processes concerning both
581 information integration (as captured by the gkRL model) as well as reward information. Indeed, the gkRL
582 model was developed primarily in order to capture participants' behavior in the unequal sampling
583 scenario where differences in information are expected to have a large influence on exploration-
584 exploitation decisions (Cogliati Dezza et al., 2017). However, model simulations in the equal information
585 condition appears to suggest that cognitive load may additionally degrade the integration of reward
586 information into an overall choice value. In order to investigate this reward degradation account, the

587 vgkRL adds an integration of reward values on top of the information integration expressed in gkRL.

588 Equation (3) thus becomes:

589

$$590 \quad V_{t,j}(c) = (Q_{t+1,j}(c) * \rho) - (I_{t,j}(c) * \omega) \quad (5)$$

591

592 ρ indicates the degree by which reward values are integrated in choice values. We fitted vgkRL to

593 participants' data and simulated the model using the retrieved parameters. We then analyzed model

594 behavior in both unequal and equal information condition. In the unequal condition, we conducted a 2

595 (condition: High Load, Low Load) by 3 (strategies: exploitation, random exploration, directed

596 exploration) non-parametric ANOVA. Results showed an effect of strategy $F(2,110) = 21$, $p < 10^{-7}$,

597 partial eta-squared = 0.144, and a condition X strategy interaction $F(2,110) = 4.79$, $p = .01$, partial eta-

598 squared = 0.743 (Figure 9a). The effect of condition did not reach the significant threshold, $p = .9$. Post-

599 hoc comparisons revealed the same pattern observed in participants' behavior where directed exploration

600 decreases whereas random exploration increases in the High Load condition compared to the Low Load

601 condition (all $p < .05$). On the contrary, exploitation did not differ between the two conditions ($p > .05$).

602 Subsequently, we conducted a 2 (condition: High Load, Low Load) by 2 (strategies: exploitation, random

603 exploration) non-parametric ANOVA in the equal information condition. Results showed an effect of

604 strategy $F(2,72) = 23.87$, $p < 10^{-5}$, partial eta-squared = 0.249, and a condition X strategy interaction

605 $F(2,72) = 6.16$, $p = .015$, partial eta-squared = 0.079 (Figure 9b). The effect of condition did not reach the

606 significant threshold, $p = .986$. Post-hoc comparisons revealed the same pattern observed in participants'

607 behavior where exploitation decreases whereas random exploration increases in the High Load condition

608 compared to the Low Load condition (all $p < .05$). These results seem to suggest that on the top of the

609 information degradation process occurring in the unequal information condition, cognitive load also

610 affected reward value degradation captured by the ρ parameter in vgkRL model. Therefore, cognitive load

611 appears to specifically interferes with the ability to combine reward and information in order to inform

612 choices. To better test this hypothesis, we compared the estimated parameters of the model between the
613 two conditions. Unfortunately, the analysis did not reveal any differences in the estimated parameters
614 between Low Load and High Load condition (all $p > .05$). The reason behind this counterintuitive result
615 might be that when adding parameters to the model higher number of data points are necessary in order to
616 obtain a reliable estimate within the same statistical power. Thus, the fitting procedure was less powerful
617 and less able to recover the accurate estimates.

618 -----
619 Insert Figure 9
620 -----

621

622 **Discussion**

623 The results of this study challenge a popular view concerning the cognitive mechanisms underlying
624 the resolution of the exploration-exploitation dilemma. Specifically, following this perspective the dilemma
625 is considered as a binary process and cognitive control as the main underlying mechanism which is
626 required in order to override default exploitative strategies in favor of exploration of the surrounding
627 environment (Daw et al., 2006; Cohen et al., 2007). Our results showed that indeed the need for cognitive
628 control seems necessary when resolving the dilemma. However, increased cognitive load appears to affect
629 only one aspect of exploration, namely directed exploration, and the effect of cognitive load on exploration
630 seems to mostly be driven by information degradation. Additionally, our results unveiled a different facet
631 of exploitative behaviors that moves away from the traditional view of exploitation as a ‘default strategy’.
632 Together, these findings shed additional light on the mechanisms underlying adaptive control and behavior
633 and suggest new approaches for interpreting the exploration-exploitation dilemma. In the following, we
634 discuss the implications of our main results.

635 In line with what could be expected due to dual-task interference (Herath, Klingberg, Young,
636 Amunts, & Roland, 2001), participants’ choice RTs were affected by high cognitive load, suggesting that
637 participants cognitive control resources were effectively reduced in this condition. Further analyses

638 showed that high cognitive load affected participants' performance on the gambling task in terms of
639 choosing the option associated with highest reward (i.e., overall performance). Under both load
640 conditions, overall performance was above chance-level. However, during the High Load condition
641 participants were slower in integrating new evidence, as shown by the decrease in the learning rate during
642 the free-choice task, which, in turn, might explain the decrease in overall participants' performance.

643 One of the main results of this study concerns the antagonist effects of cognitive load on the two
644 exploratory strategies. Specifically, increased cognitive load resulted in a decrease in directed exploration
645 and in an increase in random exploration, suggesting that directed exploration depends on the availability
646 of sufficient control resources, and that depletion of such resources promotes random exploration. This
647 result presents a different picture concerning the involvement of cognitive control in the resolution of the
648 exploration-exploitation dilemma compared to that suggested by the behavioral control hypothesis (Daw
649 et al., 2006; Cohen et al., 2007). Resolving the dilemma through exploration seems not to be a unitary
650 process that always requires cognitive resources to be mustered, independent of the type of exploratory
651 strategies adopted. On the contrary, the resolution of the dilemma through exploration is a multi-faceted
652 phenomenon (Wilson et al., 2014; Warren et al., 2017; Somerville et al., 2017; Zajkowski et al., 2017),
653 and cognitive control seems to intervene only when exploring the environment in a directed, intentional
654 manner. These results are in line with recent studies that suggest that random and directed exploration are
655 distinct strategies, even if based on a common exploratory drive (Cogliati Dezza et al., 2017; Zajkowski
656 et al., 2017).

657 Furthermore, when interfering with the resolution of the dilemma, cognitive control cooperates
658 with those aspects of exploration related with the integration of information into choice values. Under
659 cognitive load participants were more prone to stay with the same option (as shown by effects of
660 cognitive load in the switch/stay behavior) penalizing the search for new information. This result is in line
661 with several studies on information-based processes concerning the exploration-exploitation dilemma that
662 collectively highlight a tight association between information-based exploration (directed exploration)
663 with pre-frontal areas involved in higher-level cognitive processes (Badre, Doll, Long, & Frank, 2012;

664 Cavanagh, Figueroa, Cohen, & Frank, 2012) as well as the prefrontal dopamine network (Frank, Doll,
665 Oas-Terpstra, & Moreno, 2009; Kayser, Mitchell, Weinstein, & Frank, 2015). However, our results
666 appear to contrast with a study by Daw and colleagues that suggested a crucial role for top-down control
667 processes in random exploration (Daw et al., 2006). In their study, activity in brain regions associated
668 with higher-level cognitive functions (i.e. frontopolar cortex) was associated with the probability of
669 randomly exploring options. Frontopolar cortex was subsequently associated with switching between
670 strategies instead of targeting exploratory strategy itself (Boorman, Behrens, Woolrich, & Rushworth,
671 2009) and TMS studies of this region affected only directed exploration (Zajkowski et al., 2017). Clearly,
672 more research is needed to understand the neuronal and neurochemical mechanisms underlying
673 exploration in light of the new and recent evidence on directed and random exploration.

674 Our results are in line with a recent finding that showed higher cognitive costs associated with
675 those processes involved in reflexive exploration (Otto et al., 2014). Specifically, cognitive load seems to
676 affect participants' ability to use a model of the environment where environmental statistics (i.e., state
677 transition probabilities) and reward structure are integrated into choice values in order to guide
678 exploratory behaviors. However, our results suggest a more nuanced view concerning this phenomenon:
679 the results of our model fits suggested that reducing cognitive resources specifically affected those
680 processes involved in information-integration, while processes involved in transforming probability
681 distributions into action selection (decreasing or increasing the level of noise in the system through a
682 softmax function) were unaffected. Moreover, the effect of cognitive load on information is restricted to
683 integration and not to other aspects of the information processing, such as information decay (which
684 might be captured by differences in the gamma parameter). In our study, however, we approached the
685 computational problem using a model-free strategy where choices are only driven by past experience
686 (information and reward history) without a representational characterization of the environment (contrary
687 to a model-based strategy where choices are driven by the model of the world; Daw, Niv, & Dayan,
688 2005). It might be possible that in real life scenarios humans adopt model-based approaches when facing
689 exploration and exploitation problems, requiring more complex, and resource-intensive computations that

690 are only approximated by the manner in which information is integrated in the gkRL model. The relation
691 between model-based strategy and information integration should be addressed by future research.

692 Our results further question the interpretation of exploitation as a default strategy that requires
693 cognitive control to be inhibited (Daw et al., 2006; Cohen et al., 2007). Contrary to what might have been
694 expected by the behavioral control hypothesis (Cohen et al., 2007), exploitation was affected by the
695 cognitive control manipulation in such a way that when participants visited the options the same number of
696 times (i.e., equal information condition), they decreased exploitative choices during the High Load
697 condition. This finding seems to suggest that, under certain scenarios, cognitive control is necessary to
698 achieve exploitation, as in the other goal-directed behaviors. That is, choosing to exploit requires cognitive
699 resources in a fashion similar to choosing to explore. Our results are in line with recent findings on
700 cognitive foraging, where exploring other patches or exploiting familiar patches involved similar cognitive
701 mechanisms (Hills, 2010). Our results also provide support for the view that considers exploitation not only
702 as the strategy that selects best rewarded actions, but also as a strategy that relies on cognitive control
703 resources to maintain task demands (Hills, Todd, & Goldstone, 2010). Sticking with the same option can be
704 considered as a sub-goal of the higher goal of maximizing reward in the long run, and maintaining attention
705 between competing task demands required higher cognitive control functions (e.g., the cocktail party
706 phenomenon; Conway, Cowan, & Bunting, 2001; Hills et al., 2010). A drawback, however, is that our
707 model was unable to capture this phenomenon (Figure 8b). Indeed, gkRL was developed in order to capture
708 human behavior in unequal sampling scenarios (Cogliati Dezza et al., 2017). In order to understand the
709 underlying mechanisms of the effect of cognitive load on the exploitation, we presented an implementation
710 of the gkRL model where the integration of reward into choice value was also modulated. Simulations of
711 this model showed that the reward value degradation might be the underlying mechanism behind the
712 decrease in exploitation in the equal information condition. However, the limited number of trials available
713 in our paradigm precluded a definitive answer to this question. Further work is needed in order to
714 understand how cognitive control might influence choice value computation.

715 Taken together, our results suggest a new perspective on the exploration-exploitation dilemma as
716 the product of multiple competing control modes that jointly promote adaptive behavior through
717 increased emphasis on stability or flexibility. Similar to cognitive search modes (Hommel, 2012), the
718 differences between these control modes might be in the control-style they call for: a divergent decision-
719 making style- one goal representation that diverges to different action selections (or perceptual
720 representations in the case of cognitive search)- and a convergent style where a potential number of
721 possible actions (or a number of representations) converges toward an optimal solution (Hommel, 2012).
722 At the neural level, these different modes may be represented by tonic and phasic activity in the Locus
723 Coeruleus expressed by the release of norepinephrine (NE) (Aston-Jones & Cohen, 2005). LC-NE is the
724 target of projections from cortical regions implicated in cognitive control and adaptive behavior,
725 including regions involved in processing information regarding behaviorally salient changes in the
726 environment (e.g., Anterior Cingulate Cortex, Anterior Insula, and Orbitofrontal Cortex). Following
727 unexpected changes in the environment, tonic LC-NE activity may favor adaptive exploration by allowing
728 disengagement from current task demands (Yu & Dayan, 2005). On the other hand, in stable
729 environments, phasic LC-NE activity may promote exploitative behavior by increasing attention toward
730 task-relevant stimuli and maintenance of the current goal (Aston-Jones & Cohen, 2005; Jepma &
731 Nieuwenhuis, 2011). This perspective, however leaves many questions unanswered. For example, the
732 interaction between these control modes and the regions previously associated with exploration (i.e.,
733 frontopolar cortex) is still unknown and needs to be addressed by future research. Moreover, random
734 exploration, but not directed exploration, was affected by pharmacological manipulation of baseline NE
735 levels (Warren et al., 2017) questioning how the LC-NE system may control the two exploratory
736 strategies and which is the role of random exploration in this mode-based trade-off. So far, random
737 exploration seems to be a low-level (Warren et al., 2017) or automatic action control process (Humphries
738 et al., 2012) that might be necessary when a less engaging or faster behavioral adaptation is required.
739 However, the exact manner in which low-level control interacts with higher cognitive control remains an
740 open question and should be the subject of future research.

741 Although our study adds additional perspective on the cognitive mechanisms underlying the
742 resolution of the exploration-exploitation dilemma by humans, there are nonetheless limitations that may
743 influence the scope of our results. Besides the limitation of the computational model discussed above, the
744 absence of horizon manipulation in our paradigm makes impossible to distinguish whether the increase in
745 random exploration in the High Load condition was due to random exploration itself (changes in
746 randomness in long horizon) or by overall increase in randomness (Krueger, 2017). On the same line, the
747 information integration parameter was not horizon-dependent. Thus, we cannot explain the effect of
748 cognitive load on the information integration on a trial-basis. Additionally, although ambiguity appears to
749 modulate the tension between exploration and exploitation (Wilson et al., 2014; Krueger, 2017), we did not
750 specifically investigate this aspect in this study. Lastly, we did not compute participants' memory span,
751 preventing us from delineating individual profiles concerning the efficacy of our experimental
752 manipulation.

753 Regardless of these limitations, using a recently developed behavioral paradigm (Wilson et al.,
754 2014), we disentangle the role of cognitive control in the resolution of the exploration-exploitation
755 dilemma. Our results emphasized the multifaceted nature of the resolution of the dilemma and suggest
756 that multiple-cognitive control modes are the underlying cognitive mechanisms. This study is in line with
757 a new perspective on how to look at the exploration-exploitation dilemma, and provides a formal
758 foundation within which to explore pathologies of goal-directed behavior such as manifest in addiction,
759 obsessive-compulsive disorders and attentional deficits.

760

761 **Data availability**

762 Data will be provided on Open Science Framework after publication of the manuscript. The DOI
763 will be provided in this manuscript after the email of acceptance.

764

765 **Context of the Research**

766 This study is part on a broader research project that aims to investigate the neurobehavioral and
767 neurocognitive mechanisms underlying the resolution of the exploration-exploitation dilemma in humans
768 in order to develop a solid framework within which to explore decision-making alterations in psychiatry
769 disorders. The exploration-exploitation dilemma provides, indeed, a powerful tool to investigate
770 motivation, outcome evaluation, effort as well as risk-taking and impulsivity which are the main decision-
771 making components disrupted in psychiatry disorders (Addicott, Pearson, Sweitzer, Barack, & Platt,
772 2017).

773

774 **References**

- 775 Addicott, M. A., Pearson, J. M., Sweitzer, M. M., Barack, D. L., & Platt, M. L. (2017). A Primer
776 on Foraging and the Explore/Exploit Trade-Off for Psychiatry Research.
777 *Neuropsychopharmacology*, 42(10), 1931-1939. doi:10.1038/npp.2017.108
- 778 Aston-Jones, G., & Cohen, J. D. (2005). Adaptive gain and the role of the locus coeruleus-
779 norepinephrine system in optimal performance. *J Comp Neurol*, 493(1), 99-110.
780 doi:10.1002/cne.20723
- 781 Baddeley, A., Emslie, H., Kolodny, J., & Duncan, J. (1998). Random generation and the
782 executive control of working memory. *Q J Exp Psychol A*, 51(4), 819-852.
783 doi:10.1080/713755788
- 784 Badre, D., Doll, B. B., Long, N. M., & Frank, M. J. (2012). Rostrolateral prefrontal cortex and
785 individual differences in uncertainty-driven exploration. *Neuron*, 73(3), 595-607.
786 doi:10.1016/j.neuron.2011.12.025
- 787 Boorman, E. D., Behrens, T. E., Woolrich, M. W., & Rushworth, M. F. (2009). How green is the
788 grass on the other side? Frontopolar cortex and the evidence in favor of alternative
789 courses of action. *Neuron*, 62(5), 733-743. doi:10.1016/j.neuron.2009.05.014
- 790 Cavanagh, J. F., Figueroa, C. M., Cohen, M. X., & Frank, M. J. (2012). Frontal theta reflects
791 uncertainty and unexpectedness during exploration and exploitation. *Cereb Cortex*,
792 22(11), 2575-2586. doi:10.1093/cercor/bhr332
- 793 Cogliati Dezza, I., Yu, A. J., Cleeremans, A., & Alexander, W. (2017). Learning the value of
794 information and reward over time when solving exploration-exploitation problems.
795 *Sci Rep*, 7(1), 16919. doi:10.1038/s41598-017-17237-w
- 796 Cohen, J. D., McClure, S. M., & Yu, A. J. (2007). Should I stay or should I go? How the human
797 brain manages the trade-off between exploitation and exploration. *Philos Trans R Soc
798 Lond B Biol Sci*, 362(1481), 933-942. doi:10.1098/rstb.2007.2098
- 799 Conover, W. J., & Iman, R. L. (1981). Rank transformations as a bridge between parametric
800 and nonparametric statistics. *American Statistician*, 35, 124-129.
- 801 Conway, A. R., Cowan, N., & Bunting, M. F. (2001). The cocktail party phenomenon revisited:
802 the importance of working memory capacity. *Psychon Bull Rev*, 8(2), 331-335.

- 803 D'Esposito, M., Postle, B. R., Ballard, D., & Lease, J. (1999). Maintenance versus manipulation
804 of information held in working memory: an event-related fMRI study. *Brain Cogn*,
805 *41*(1), 66-86. doi:10.1006/brcg.1999.1096
- 806 Daw, N. D., Niv, Y., & Dayan, P. (2005). Uncertainty-based competition between prefrontal
807 and dorsolateral striatal systems for behavioral control. *Nat Neurosci*, *8*(12), 1704-
808 1711. doi:10.1038/nn1560
- 809 Daw, N. D., O'Doherty, J. P., Dayan, P., Seymour, B., & Dolan, R. J. (2006). Cortical substrates
810 for exploratory decisions in humans. *Nature*, *441*(7095), 876-879.
811 doi:10.1038/nature04766
- 812 Frank, M. J., Doll, B. B., Oas-Terpstra, J., & Moreno, F. (2009). Prefrontal and striatal
813 dopaminergic genes predict individual differences in exploration and exploitation.
814 *Nat Neurosci*, *12*(8), 1062-1068. doi:10.1038/nn.2342
- 815 Gittins, J., & Jones, D. (1974). A dynamic allocation index for the sequential design of
816 experiments. In J. Gans (Ed.), *Progress in statistics* (pp. 241-266). Amsterdam: **The**
817 **Netherlands: North-Holland.**
- 818 Herath, P., Klingberg, T., Young, J., Amunts, K., & Roland, P. (2001). Neural correlates of dual
819 task interference can be dissociated from those of divided attention: an fMRI study.
820 *Cereb Cortex*, *11*(9), 796-805.
- 821 Hills, T. T., Todd, P. M., & Goldstone, R. L. (2010). The central executive as a search process:
822 priming exploration and exploitation across domains. *J Exp Psychol Gen*, *139*(4), 590-
823 609. doi:10.1037/a0020666
- 824 Hommel, B. (2012). Convergent and divergent operations in cognitive search. In P. M. Todd,
825 T. T. Hills, & H. Robbins (Eds.), *Cognitive search: evolution, algorithms, and the brain.*
- 826 Humphries, M. D., Khamassi, M., & Gurney, K. (2012). Dopaminergic Control of the
827 Exploration-Exploitation Trade-Off via the Basal Ganglia. *Front Neurosci*, *6*, 9.
828 doi:10.3389/fnins.2012.00009
- 829 Jepma, M., & Nieuwenhuis, S. (2011). Pupil diameter predicts changes in the exploration-
830 exploitation trade-off: evidence for the adaptive gain theory. *J Cogn Neurosci*, *23*(7),
831 1587-1596. doi:10.1162/jocn.2010.21548
- 832 Kayser, A. S., Mitchell, J. M., Weinstein, D., & Frank, M. J. (2015). Dopamine, locus of control,
833 and the exploration-exploitation tradeoff. *Neuropsychopharmacology*, *40*(2), 454-
834 462. doi:10.1038/npp.2014.193
- 835 Khamassi, M., Enel, P., Dominey, P. F., & Procyk, E. (2013). Medial prefrontal cortex and the
836 adaptive regulation of reinforcement learning parameters. *Prog Brain Res*, *202*, 441-
837 464. doi:10.1016/B978-0-444-62604-2.00022-8
- 838 Koechlin, E., Ody, C., & Kouneiher, F. (2003). The architecture of cognitive control in the
839 human prefrontal cortex. *Science*, *302*(5648), 1181-1185.
840 doi:10.1126/science.1088545
- 841 Konstantinou, N., & Lavie, N. (2013). Dissociable roles of different types of working memory
842 load in visual detection. *J Exp Psychol Hum Percept Perform*, *39*(4), 919-924.
843 doi:10.1037/a0033037
- 844 Krueger, P. M. (2017). Strategies for exploration in the domain of losses. *Judgement and*
845 *Decision Making*, *12*(2), 104-117.
- 846 Mars, R. B., Sallet, J., Rushworth, M. F., & Yeung, N. (2011). *Neural basis of motivational and*
847 *cognitive control*. Cambridge, MA: MIT Press.

- 848 Mehlhorn, K., Newell, B. R., Todd, P. M., Lee, M. D., Morgan, K., Braithwaite, V. A., ... , &
849 Gonzalez, C. (2015). Unpacking the exploration–exploitation tradeoff: A synthesis of
850 human and animal literatures. *Decision*, 2(3), 191-215.
- 851 Miller, E. K., & Cohen, J. D. (2001). An integrative theory of prefrontal cortex function. *Annu*
852 *Rev Neurosci*, 24, 167-202. doi:10.1146/annurev.neuro.24.1.167
- 853 Otto, A. R., Knox, W. B., Markman, A. B., & Love, B. C. (2014). Physiological and behavioral
854 signatures of reflective exploratory choice. *Cogn Affect Behav Neurosci*, 14(4), 1167-
855 1183. doi:10.3758/s13415-014-0260-4
- 856 Payzan-LeNestour, E., & Bossaerts, P. (2011). Risk, unexpected uncertainty, and estimation
857 uncertainty: Bayesian learning in unstable settings. *PLoS Computational Biology*(7),
858 e1001048. doi:<http://dx.doi.org/10.1371/journal.pcbi.1001048>
- 859 Rescorla, R. A., & Wagner, A. R. (1972). A theory of Pavlovian conditioning: Variations in the
860 effectiveness of reinforcement and nonreinforcement. *Classical conditioning: Current*
861 *research and theory*, 64-99.
- 862 Robbins, H. (1952). Some aspects of the sequential design of experiments. *Bulletin of the*
863 *American Mathematical Society*, 58 527-535.
- 864 Schwarz, G. (1978). Estimating the dimension of a model. *Ann. Stat*, 6, 461-464.
- 865 Somerville, L. H., Sasse, S. F., Garrad, M. C., Drysdale, A. T., Abi Akar, N., Insel, C., & Wilson, R.
866 C. (2017). Charting the expansion of strategic exploratory behavior during
867 adolescence. *J Exp Psychol Gen*, 146(2), 155-164. doi:10.1037/xge0000250
- 868 Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. Cambridge, MA:
869 MIT Press.
- 870 Warren, C. M., Wilson, R. C., van der Wee, N. J., Giltay, E. J., van Noorden, M. S., Cohen, J. D., &
871 Nieuwenhuis, S. (2017). The effect of atomoxetine on random and directed
872 exploration in humans. *PLoS One*, 12(4), e0176034.
873 doi:10.1371/journal.pone.0176034
- 874 Wilson, R. C., Geana, A., White, J. M., Ludvig, E. A., & Cohen, J. D. (2014). Humans use directed
875 and random exploration to solve the explore-exploit dilemma. *J Exp Psychol Gen*,
876 143(6), 2074-2081. doi:10.1037/a0038199
- 877 Wilson, R. C., & Niv, Y. (2011). Inferring relevance in a changing world. *Front Hum Neurosci*,
878 5, 189. doi:10.3389/fnhum.2011.00189
- 879 Yu, A. J., & Dayan, P. (2005). Uncertainty, neuromodulation, and attention. *Neuron*, 46(4),
880 681-692. doi:10.1016/j.neuron.2005.04.026
- 881 Zajkowski, W. K., Kossut, M., & Wilson, R. C. (2017). A causal role for right frontopolar cortex
882 in directed, but not random, exploration. *Elife*, 6. doi:10.7554/eLife.27430

884 Acknowledgments

885 ICD is a researcher supported by F.R.S.-FNRS grant (Belgium). AC is a research director with the F.R.S.-
886 FNRS (Belgium). WA was supported in part by FWO-Flanders Odysseus II Award #G.OC44.13N.

887

888 Author contributions

889 ICD designed and carried out the experiment. ICD performed the analysis of the data and the model
890 analysis. ICD and WA discussed and interpreted the data. ICD, AC and WA wrote the manuscript.

891

892 **Additional information**

893 *Competing financial interests*

894

895 The authors declare no competing financial interests.

896

897

898 **Figure Captions**

899 **Figure 1 Behavioral paradigm.** a) *Organization of games and trials.* On each game, participants faced 6
900 consecutive trials of the forced-choice task and between 1 and 6 trials of the free-choice task. In the first
901 free-choice trial (in yellow), reward and information are orthogonalized enabling the distinction between
902 random and directed exploration. The number of free-choice trials was exponentially distributed such that a
903 higher proportion of games allowed subjects to make 6 free choices. b) *Choices.* Participants indicated their
904 choices using the forefinger, middle finger and ring finger and pressing the keyboard keys ‘c’, ‘v’ and ‘b’,
905 respectively. c) *Forced-choice task.* Three decks of cards were displayed on the screen (a blue, a red and
906 green deck) and participants were forced to choose a preselected deck (outlined in blue in the figure). After
907 selecting the deck, the card turned and revealed the points associated with the selected option, between 1
908 and 100 points. At this stage, the points displayed to participants were not added to their total score. d)
909 *Free-choice task.* Participants made their own decisions among the same three decks of cards displayed
910 during the forced-choice task. After each trial, the points displayed on the screen were added to the
911 participants’ total score and participants were instructed to attempt to maximize the total points earned at
912 the end of the experiment. e) *Cognitive load manipulation.* Before the 1st trial of the free-choice task, a
913 sequence of 9-digits was displayed on the screen. During the Low Load condition, the digits were presented
914 in fixed numerical order (i.e., ‘123456789’) for 500 ms. On the contrary, during the High Load condition
915 the digits were presented in random order (i.e., ‘371586249’), for 2000 ms, and a new sequence was
916 generated on each game. After each free-choice trial a digit (randomly selected from the 9-digit sequence)
917 was displayed to participants who needed to report (‘R_m’- memory response) the number that followed the
918 presented number in the previous 9-digit sequence presented before the 1st free choice trial.

919

920 **Figure 2 Memory performance and Cognitive load manipulation.** a) Memory accuracy measured by
921 averaging trial-by-trial correct memory responses obtained by participants during both High and Low
922 Load condition. Error bars are also represented as the standard error from the mean (s.e.m). b) Cognitive

923 load manipulation increased participants' choice reaction time (RT in ms) during the High Load
924 compared to the Low Load condition. Error bars are also represented as s.e.m.

925

926 **Figure 3** *Cognitive load and decision strategies.* a) In the unequal information condition, directed
927 exploration decreased in the High Load condition compared to the Low Load condition, whereas random
928 exploration and exploitation showed the opposite trend. Error bars are also represented as s.e.m. b) In the
929 equal information condition, random exploration increased under High Load condition whereas
930 exploitation decreased. Error bars are also represented as s.e.m.

931

932 **Figure 4** *Information integration.* a) *First-free trials.* Model fit on the first free-choice only revealed a
933 decrease in the information weigh parameter ω (that modulates to integration of information into choice
934 values) during High Load condition compared to Low Load condition, whereas the inverse of temperature
935 β , the learning rate α and the γ parameter were not affected by the cognitive load. Error bars are also
936 represented as s.e.m. b) *All-free trials.* Model fit on all free-choices showed a decrease in information
937 parameter ω and the learning rate α in the High Load condition, whereas both β and γ were not affected
938 by the cognitive manipulation. Error bars are also represented as s.e.m. c) *Correct memory choices.*
939 Model fit on the trials where participants correctly performed the memory task. Error bars are also
940 represented as s.e.m. The results showed the same pattern observed when fitting all free-choices.

941

942 **Figure 5** *Comparative fit of the gkRL and sRL.* The comparison of the fit is based on the BIC values of
943 both models during the Low Load (a) and High Load condition (b). Each point is one participant. The
944 sRL fit better when the point is below the identity line. When a point lays on the identity line the models
945 equally explain participants' behavior.

946

947 **Figure 6** *Seen analysis*. a) Probability to choose the option seen least, middle and most of the time during
948 the free-choice task. Choices towards the least seen option decreased during the High Load condition
949 compared to the Low Load condition, whereas choices toward most seen options showed the opposite
950 pattern. Error bars are also represented as s.e.m. b) Probability to choose the option seen most (*Most-Seen*
951 *options*), least (*Least-Seen options*) and middle (*Middle-Seen options*) of the time during the free-choice
952 task split by trial. During the first three free-trials, the probability to choose the least seen option (and
953 most seen option) differs significantly, whereas in the last free-choice trials no difference was observed.
954 To avoid overloading the visualization, we reported only when the comparisons did not reach significance
955 threshold.

956
957 **Figure 7** *Switch/Stay strategy*. a) Probability to stay with the same option chosen at trial t-1 during the
958 free-choice task. During last free-choice trials, the probability to stay with the same option did not differ
959 between the two loading conditions. Error bars are also represented as s.e.m. b) Probability to switch from
960 the option chosen at trial t-1 during the free-choice task. During last free-choice trials, the probability to
961 switch did not differ between the two loading conditions. Error bars are also represented as s.e.m.

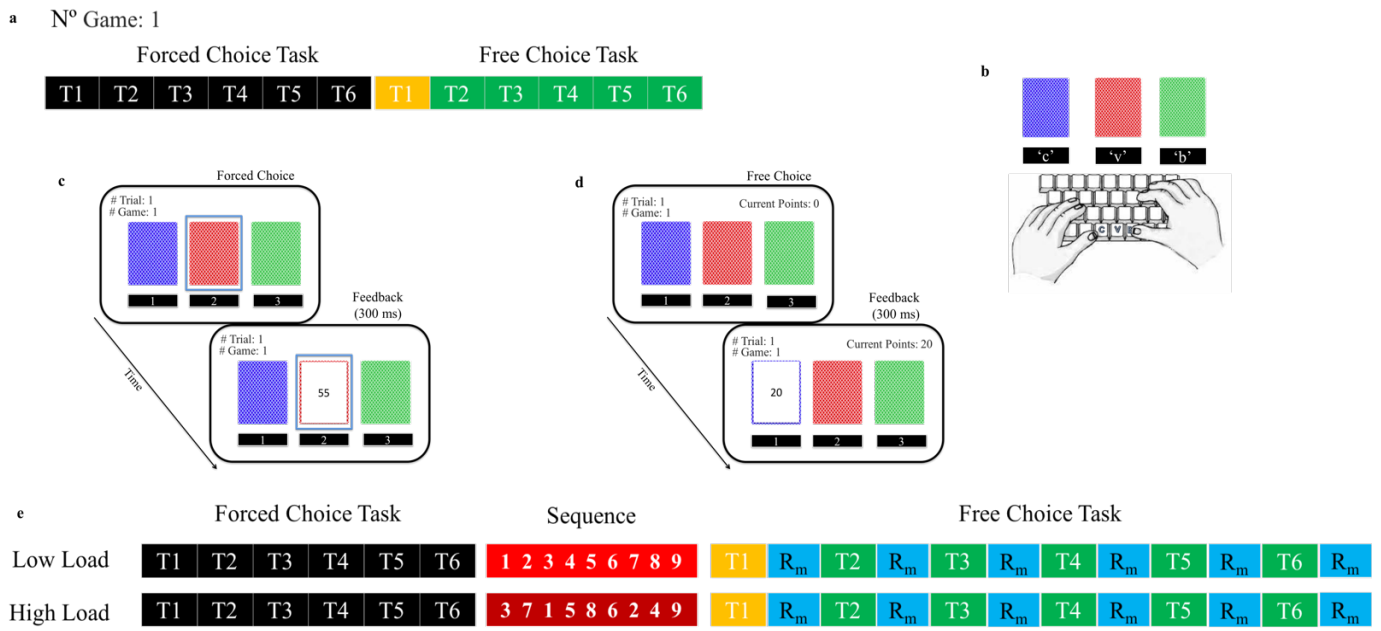
962
963 **Figure 8** *gkRL Simulation*. a) In the unequal information condition, the model simulated under the two
964 loading conditions reproduced the same behavioral pattern observed in participants: directed exploration
965 decreased in the High Load condition, whereas random exploration and exploitation increased in Low
966 Load condition. b) In the equal information condition, no behavioral differences in exploitation and
967 random exploration were observed between the two loading conditions. Only the comparisons that did not
968 reach significance threshold are reported.

969
970 **Figure 9** *vgkRL Simulation*. Model simulation reproduced a similar patter observed in participants' data
971 in both unequal (a) and equal (b) information condition.

972

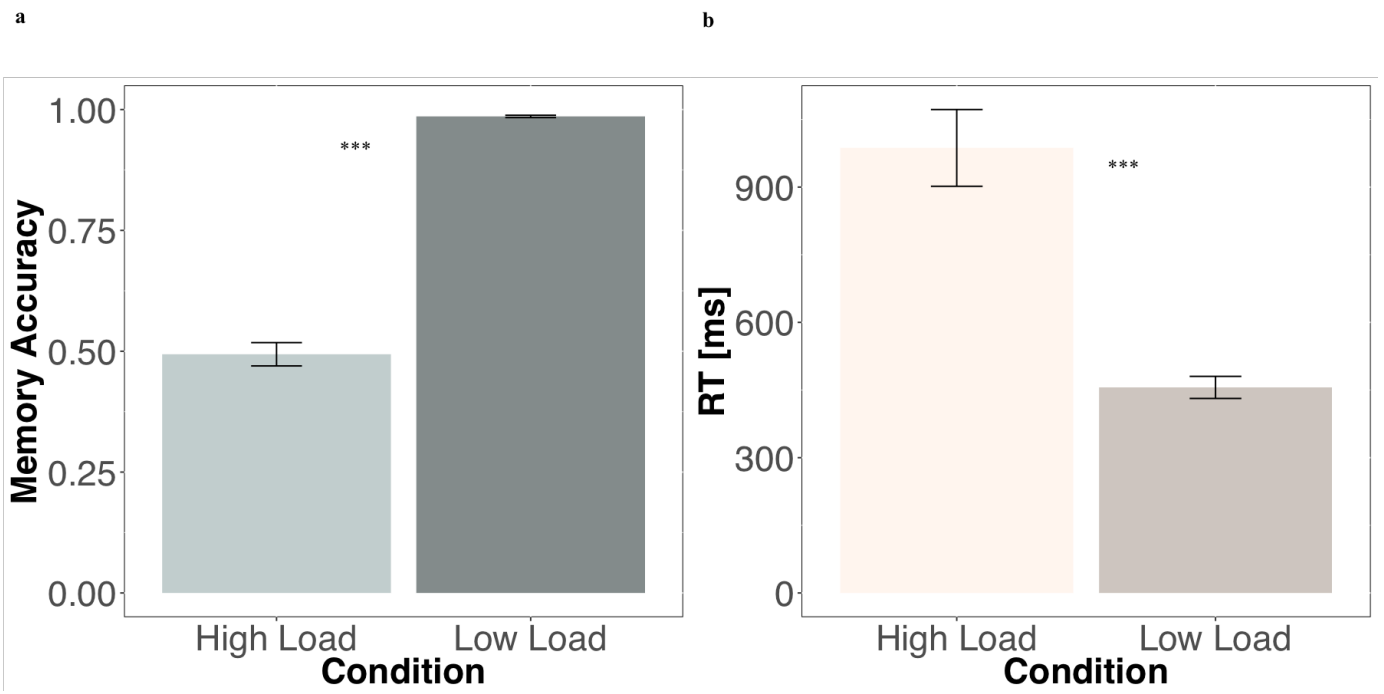
973 **Figure**

974 **Figure 1**



975
976

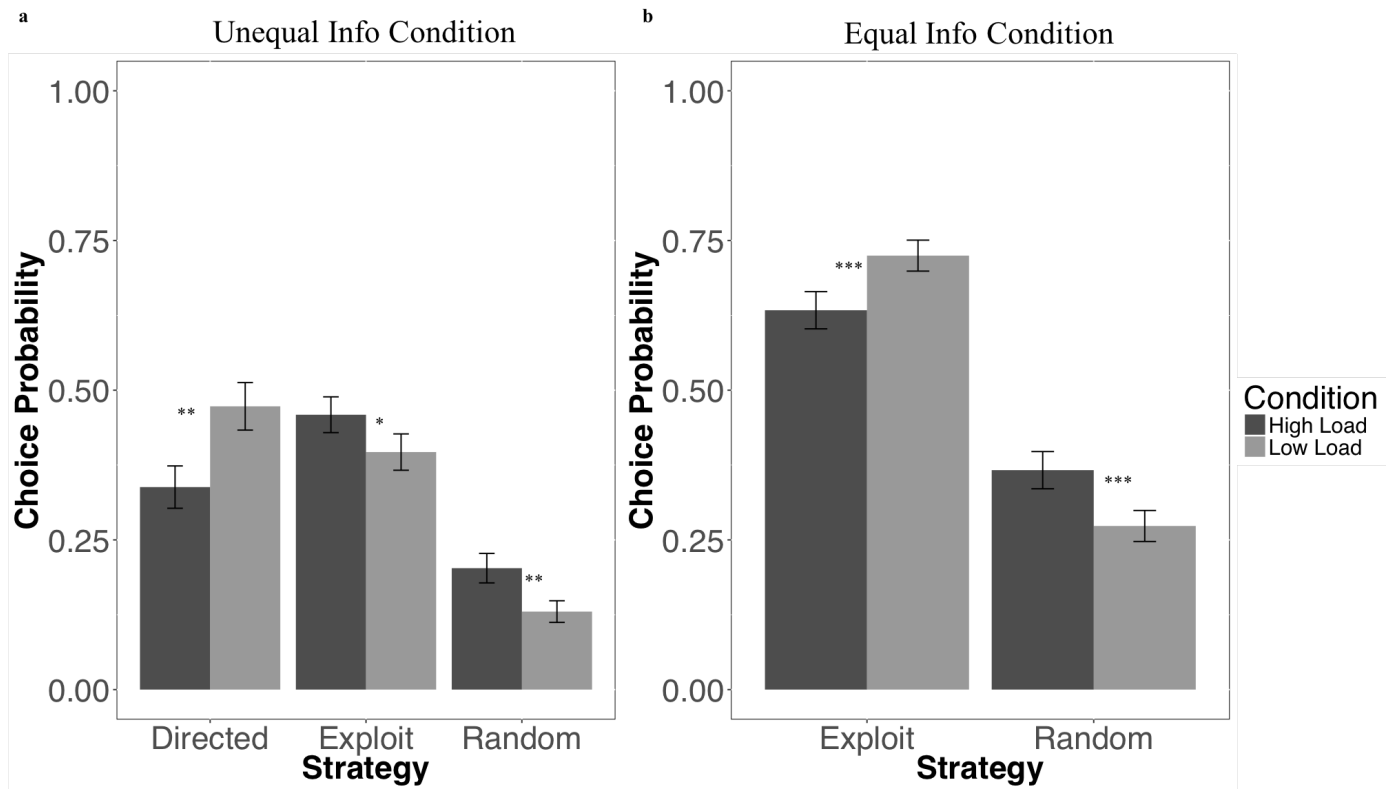
977 **Figure 2**



978

979

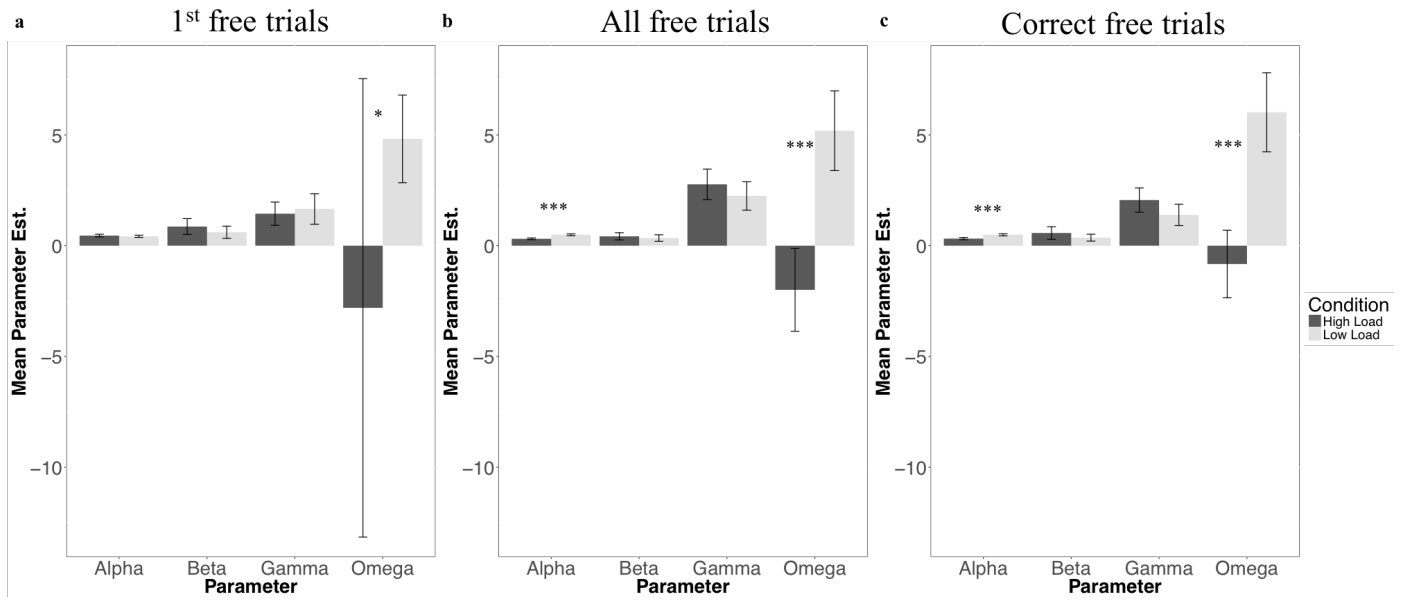
980 **Figure 3**



981

982

983 **Figure 4**

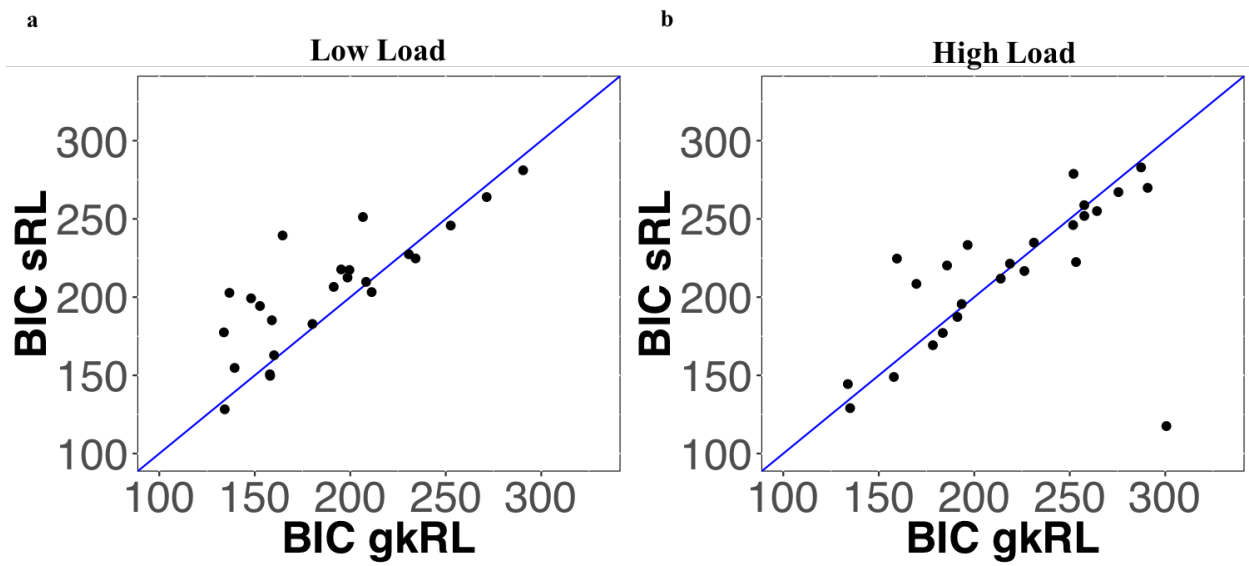


984

985

986 **Figure 5**

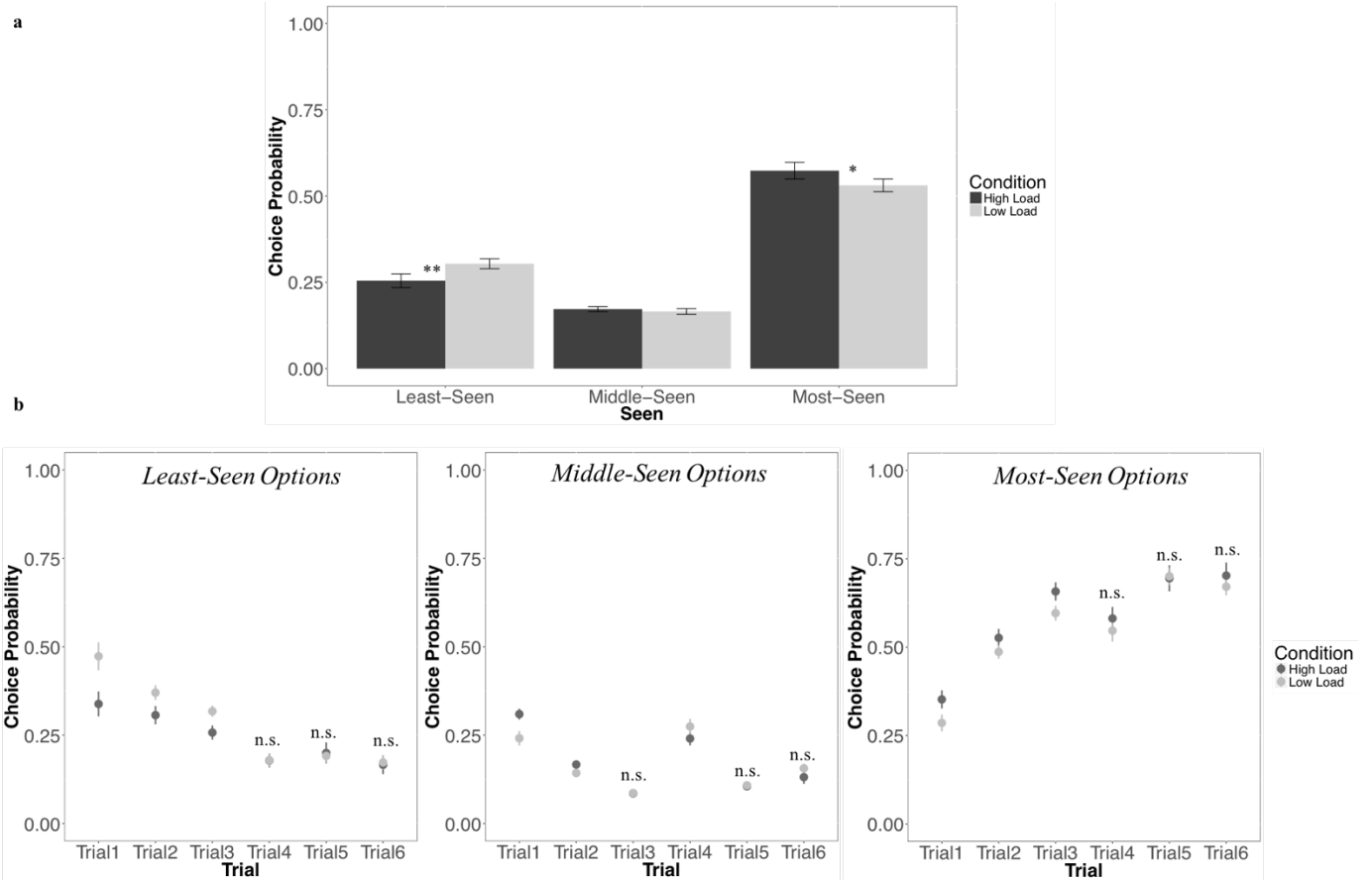
sRL vs. gkRL with the respect to the Identity Line



987

988

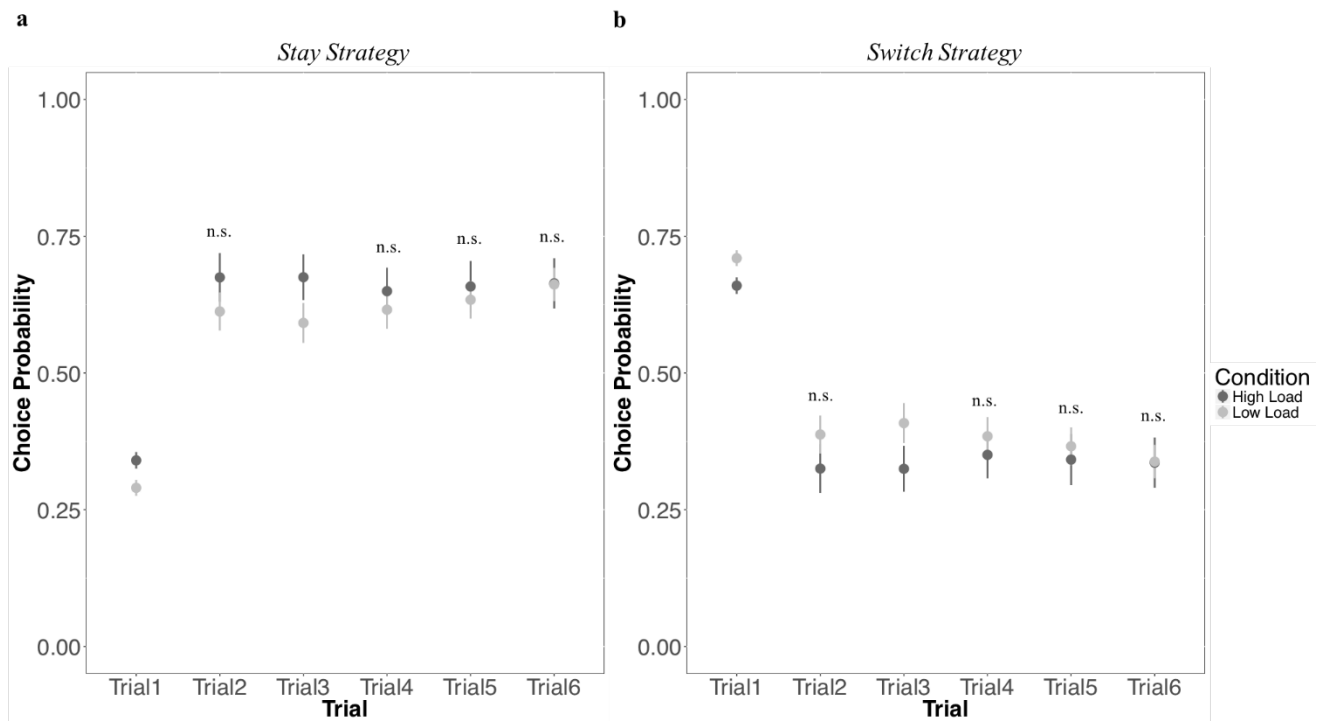
989 **Figure 6**



990

991

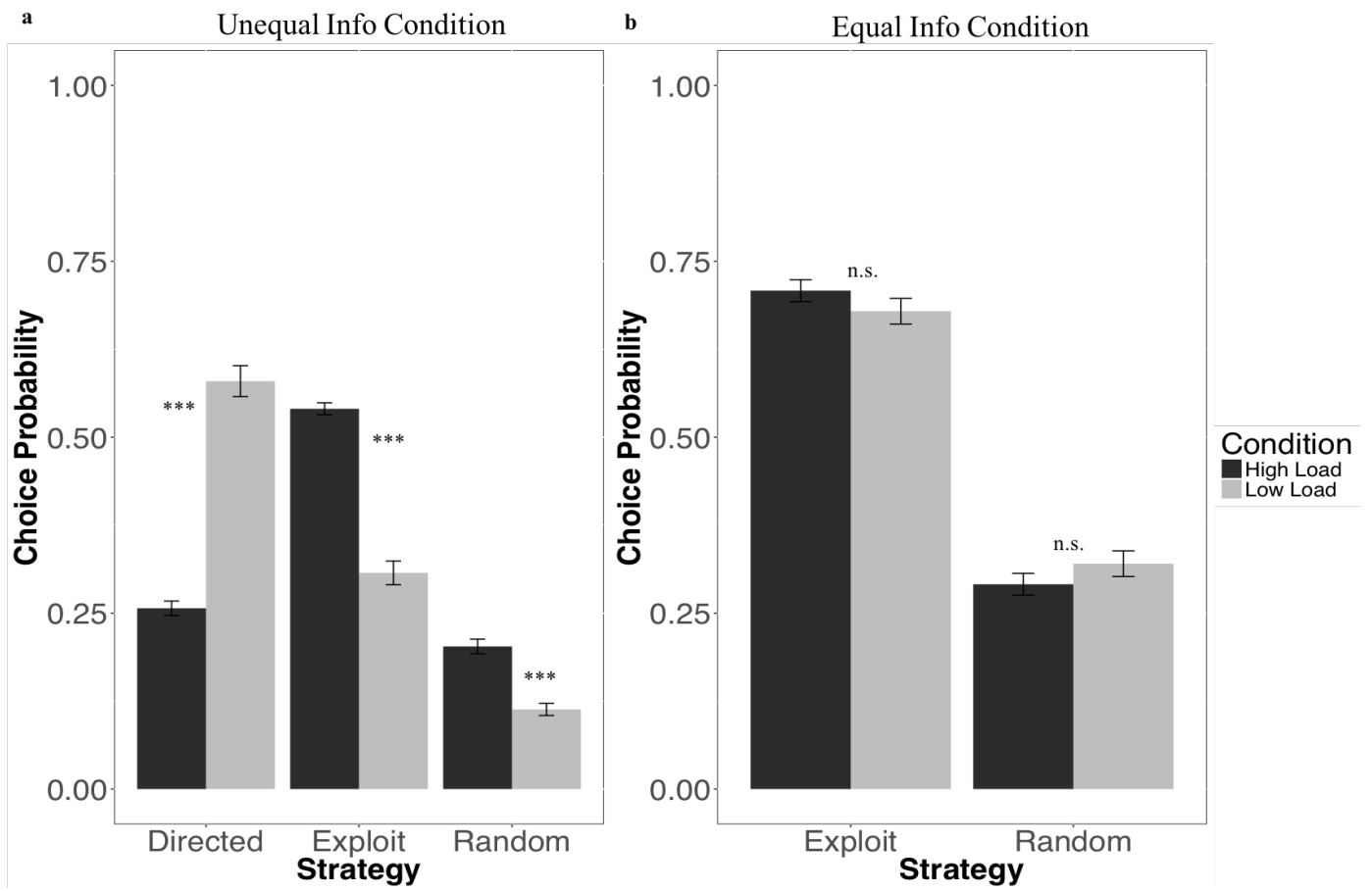
992 **Figure 7**



993

994

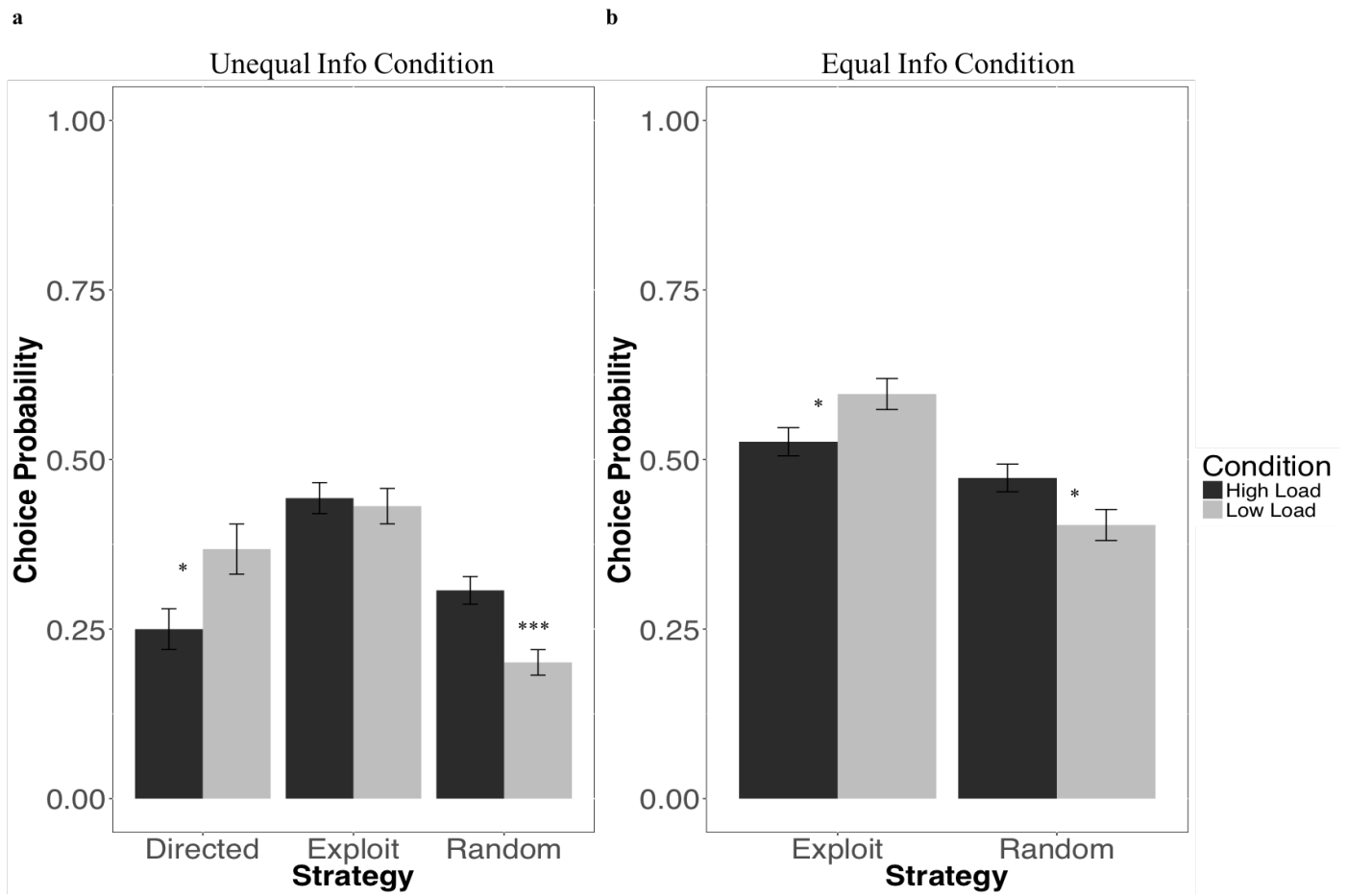
995 **Figure 8**



996

997

998 **Figure 9**



999

1000

1001 **Table Captions**

1002 **Table 1** *Model fit results: 1st free trials.* Estimated parameters for each subject using gkRL model during
1003 High Load and Low Load condition. Group average of the estimated parameters are also reported. Group
1004 standard deviation are reported within parenthesis.

1005 **Table 2** *Model fit results: all free trials.* Estimated parameters for each subject using gkRL model during
1006 High Load and Low Load condition. Group average of the estimated parameters are also reported. Group
1007 standard deviation are reported within parenthesis.

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023 **Table 1**

Participants	Low Load					High Load				
	α	β	ω	γ	Log(γ)	α	β	ω	γ	Log(γ)
SubjectNumber1	0.870	0.142	21.754	0.002	-6.458	0.562	0.145	12.740	0.520	-0.653
SubjectNumber2	0.464	0.283	-10.186	0.000	-19.612	0.566	0.094	1.017	0.202	-1.601
SubjectNumber3	0.371	0.337	1.979	0.289	-1.240	0.547	0.080	2.667	1.313	0.272
SubjectNumber4	0.587	0.186	0.000	10.386	2.340	0.482	0.297	-2.314	0.359	-1.024
SubjectNumber5	0.000	6.392	-0.067	0.595	-0.518	0.000	7.500	0.000	0.919	-0.084
SubjectNumber6	0.724	0.132	9.246	0.000	-19.089	0.599	0.132	-8.941	0.110	-2.209
SubjectNumber7	0.254	0.173	6.058	0.000	-25.509	0.109	0.201	0.000	7.882	2.065
SubjectNumber8	0.463	0.131	18.008	0.321	-1.137	0.011	1.367	0.230	0.000	-21.745
SubjectNumber9	0.258	0.042	-3.445	0.000	-19.636	0.455	0.000	148.72	2.125	0.754
SubjectNumber10	0.190	0.400	0.011	4.429	1.488	0.004	4.693	-0.048	0.811	-0.209
SubjectNumber11	0.455	0.120	12.844	0.000	-25.225	1.000	0.040	-13.591	0.558	-0.583
SubjectNumber12	0.270	0.283	6.120	0.310	-1.170	0.343	0.072	2.610	0.000	-21.931
SubjectNumber13	0.385	0.077	2.828	0.664	-0.410	0.661	0.084	9.516	0.251	-1.383
SubjectNumber14	0.502	0.076	28.386	0.000	-21.935	0.345	0.351	1.306	1.354	0.303
SubjectNumber15	0.510	0.248	9.855	0.326	-1.119	0.422	0.165	10.952	0.291	-1.233
SubjectNumber16	0.698	0.096	-4.486	0.000	-21.560	0.563	0.102	-26.370	0.145	-1.931
SubjectNumber17	0.413	0.134	-16.646	0.000	-20.527	0.432	0.150	-21.224	0.000	-22.509
SubjectNumber18	0.564	0.114	14.187	0.000	-24.509	0.004	2.181	-0.337	0.803	-0.219
SubjectNumber19	0.463	0.219	-1.686	0.506	-0.681	0.556	0.142	-1.872	0.742	-0.298
SubjectNumber20	0.536	0.462	14.166	0.298	-1.211	0.733	0.187	0.000	7.368	1.997
SubjectNumber21	0.143	0.530	2.390	1.090	0.086	0.534	0.140	5.736	0.000	-20.851
SubjectNumber22	0.004	3.518	0.000	11.000	2.398	1	0.004	-201	0.000	-25.761
SubjectNumber23	0.054	0.771	1.145	0.706	-0.348	0.569	0.140	0.000	9.404	2.241
SubjectNumber24	1.000	0.065	0.000	10.000	2.303	0.002	3.278	-0.057	0.976	-0.024
SubjectNumber25	0.468	0.215	8.107	0.519	-0.657	0.897	0.081	10.165	0.000	-20.222
Total	0.426	0.606	4.82	1.66	-8.16	0.456	0.865	-2.81	1.44	-5.47
	(0.249)	(1.38)	(9.89)	(1.40)	(10.79)	(0.3)	(1.8)	(51.76)	(2.62)	(9.68)

1024

1025

1026 **Table 2**

Participants	Low Load					High Load				
	α	β	ω	γ	Log(γ)	α	β	ω	γ	Log(γ)
SubjectNumber1	0.640	0.119	19.633	0.070	-2.659	0.570	0.114	15.735	0.072	-2.631
SubjectNumber2	0.562	0.112	-0.465	1.403	0.338	0.339	0.121	-0.018	3.718	1.313
SubjectNumber3	0.475	0.168	4.188	0.000	-21.921	0.590	0.070	5.611	0.000	-21.464
SubjectNumber4	0.446	0.254	0.958	2.943	1.079	0.412	0.238	-0.550	1.952	0.669
SubjectNumber5	0.000	2.543	-0.040	0.000	-19.081	0.306	0.010	-37.732	0.674	-0.395
SubjectNumber6	0.698	0.124	4.107	0.000	-21.476	0.580	0.072	-4.313	1.064	0.062
SubjectNumber7	0.397	0.103	9.887	0.000	-22.202	0.336	0.049	-0.075	3.504	1.254
SubjectNumber8	0.627	0.110	26.525	0.000	-27.211	0.006	1.500	0.041	12.000	2.485
SubjectNumber9	0.002	2.958	-0.005	11.847	2.472	0.096	0.009	6.579	1.546	0.435
SubjectNumber10	0.526	0.180	-0.001	4.443	1.491	0.210	0.090	-2.064	0.636	-0.452
SubjectNumber11	0.388	0.123	-0.033	3.179	1.157	0.132	0.131	-5.146	0.682	-0.383
SubjectNumber12	0.443	0.225	9.320	0.000	-22.414	0.451	0.067	-0.072	3.044	1.113
SubjectNumber13	0.513	0.067	-0.005	4.412	1.484	0.258	0.125	4.366	0.000	-21.627
SubjectNumber14	0.597	0.091	17.966	0.007	-4.918	0.296	0.234	-0.031	2.916	1.070
SubjectNumber15	0.467	0.171	0.000	9.762	2.279	0.279	0.192	0.000	8.301	2.116
SubjectNumber16	0.478	0.106	-0.237	2.153	0.767	0.433	0.096	-10.481	0.704	-0.350
SubjectNumber17	0.440	0.094	-15.048	0.330	-1.109	0.455	0.068	-15.817	0.541	-0.615
SubjectNumber18	0.675	0.085	-0.048	3.136	1.143	0.005	2.961	-0.224	0.695	-0.364
SubjectNumber19	0.609	0.130	-0.319	1.955	0.671	0.429	0.144	-4.790	0.601	-0.509
SubjectNumber20	0.430	0.268	15.801	0.000	-21.862	0.422	0.181	0.000	9.228	2.222
SubjectNumber21	0.545	0.186	15.156	0.000	-20.273	0.291	0.168	-0.010	4.014	1.390
SubjectNumber22	0.387	0.024	-0.147	3.334	1.204	0.000	0.771	-0.579	0.360	-1.022
SubjectNumber23	0.523	0.107	10.980	0.259	-1.352	0.523	0.134	0.000	9.999	2.303
SubjectNumber24	1.000	0.049	5.261	6.845	1.924	0.002	2.814	-0.189	0.057	-2.870
SubjectNumber25	0.398	0.188	12.575	0.085	-2.467	0.381	0.133	-0.018	2.961	1.085
Total	0.495 (0.198)	0.344 (0.737)	5.19 (9.0)	2.247 (3.201)	-6.917 (10.802)	0.312 (0.186)	0.424 (0.81)	-1.785 (9.425)	2.771 (3.448)	-1.407 (6.215)

1028
1029
1030
1031
1032
1033