

1 Why Do Irrelevant Alternatives Matter? An
2 fMRI-TMS Study of Context-Dependent
3 Preferences
4

5 Chen-Ying Huang^{1,*,**}, Hui-Kuan Chung^{2,*}, Hsin-Ju Lee³, Yi-Ta Lu¹,
6 Fu-Yun Tsuo⁴, Tzai-Shuen Chen⁵, Chi-Fu Chang⁶, Chi-Hung Juan⁶,
7 Wen-Jui Kuo^{3,**}, Tomas Sjöström^{7,**}
8

1 Department of Economics, National Taiwan University, Taipei, Taiwan

2 Department of Psychology, New York University, New York, NY 10003, U.S.A.

3 Institute of Neuroscience, National Yang-Ming University, Taipei, Taiwan

4 Department of Electrical Engineering, National Taiwan University, Taipei, Taiwan

5 Department of Economics, Washington University in St. Louis, St. Louis, MO 63130, U.S.A.

6 Institute of Cognitive Neuroscience, National Central University, Taoyuan, Taiwan

7 Department of Economics, Rutgers University, New Brunswick, NJ 08901, U.S.A.

* These authors contributed equally to this work.

** To whom correspondence should be addressed:

chenying@ntu.edu.tw, wjkuo@ym.edu.tw or tsjostrom@economics.rutgers.edu

9 August 2, 2015

10 **Abstract**

11 Both humans and animals are known to exhibit a violation of ratio-
12 nality known as "decoy effect": introducing an irrelevant alternative (a
13 decoy) can influence choices among other (relevant) alternatives^{1,2,3}. Ex-
14 actly how and why decoys trigger this effect is not known. It may be
15 an example of fast heuristic decision-making which is adaptive in natural
16 environments^{4,5}, but may lead to biased choices in certain markets or ex-
17 periments. We used functional magnetic resonance imaging (fMRI) and
18 transcranial magnetic stimulation (TMS) to investigate the neural under-
19 pinning of the decoy effect. The left ventral striatum was more active
20 when the chosen alternative dominated the decoy (compared with when
21 the same chosen alternative did not). This suggests that the decoy may
22 influence the valuation of other alternatives, making valuation context-
23 dependent. Consistent with the idea that control is recruited to prevent
24 heuristics from producing biased choices⁶, the right inferior frontal gyrus
25 (rIFG), often implicated in inhibiting prepotent responses^{7,8}, connected
26 more strongly with the striatum when participants successfully overrode
27 the decoy effect and made unbiased choices (compared with when they
28 were not successful). This is further supported by our TMS experiment:
29 participants whose rIFG was temporarily disrupted made biased choices
30 more often than a control group. Our results suggest that the decoy effect
31 is due to context-dependent activation of the reward area. But the dif-
32 ferential connectivity from the frontal area may indicate how deliberate
33 control monitors and corrects errors in heuristic decision making.

34 **Summary** We used functional magnetic resonance imaging and transcranial
35 magnetic stimulation to study the neural context effect caused by a decoy, its
36 trial-by-trial strength, and the possible control mechanism recruited to overcome
37 a potential decision-making bias.

38 In standard economic theory, decision-making is represented as assigning
39 value to each feasible alternative and choosing the alternative with the highest
40 value. In this theory, the value of an alternative depends only on its intrinsic
41 properties, not on the context in which it appears. Neuroscientific studies sug-
42 gest that the brain does encode values of available alternatives⁹, but there is
43 mixed evidence on the extent to which these values depend on the context^{10,11}.
44 Counterfactual outcomes of a lottery (i.e., outcomes that could have, but did
45 not occur) have been shown to be a source of context-dependence for human
46 subjects^{12,13}. Intuitively, winning five dollars in a lottery where this is the
47 smallest possible prize is not as rewarding as winning five dollars in a lottery
48 where this is the largest possible prize. Adjusting the responses of neurons in
49 the reward area to the range of possible outcomes in this way may be an efficient
50 way for the brain to utilize the neurons' limited firing range.

51 Since the counterfactual outcomes of a lottery had a chance of occurring,
52 they may have influenced anticipations, which may affect the evaluation of the
53 final outcome, perhaps by triggering disappointment or elation. In a sense,
54 any outcome with a positive probability of occurring is a relevant part of the
55 lottery. Suppose, however, that a participant chooses from a set of options, one
56 of which should be completely irrelevant to rational decision-making (in a sense
57 to be made precise below). Will the brain manage to disregard the irrelevant
58 alternative, and save its limited firing range for the viable options? Or will
59 the valuations of the viable (relevant) options be influenced by the irrelevant
60 option? We study this question in the context of the decoy problem, well known
61 from the literature on animal behavior¹ and consumer marketing².

62 The decoy problem can be illustrated by the pricing of *The Economist*
63 magazine³. A print-only subscription, and a print-and-digital subscription which
64 includes both print and digital access, are priced the same. The print-only sub-
65 scription is a decoy option which nobody is expected to choose, because it is
66 clearly worse to get only print access than to get *both* print and digital accesses
67 at the same price. However, controlled laboratory experiments² have shown a

68 "decoy effect": the existence of the decoy tends to increase the demand for the
69 print-and-digital subscription. But does the presence of a seemingly "irrelevant"
70 decoy actually influence the valuation of other options ? Our results indicate
71 that it does, even when choices are not influenced by the decoy (so that there is
72 no decoy effect in the usual, choice behavior, sense). Thus, valuation does seem
73 to depend on "irrelevant" aspects of the context. Our results hint at ways the
74 brain may try to prevent this from leading to biased decisions.

75 Using functional magnetic resonance imaging (fMRI), we scanned 32 par-
76 ticipants when they made choices from a series of two- and three-item menus.
77 The items on the menus were restaurant meals with specified prices and quality
78 levels. Formally, alternative X *dominates* alternative X' if and only if X is both
79 cheaper and has a higher quality than X' . (In our notation, an alternative that
80 is dominated by X will always be denoted X' .) Some of the three-item menus
81 contained a decoy item that was dominated by another item on the menu and
82 hence, according to economic theory, should be irrelevant to decision-making.
83 Domination, when it occurred, was always *asymmetric* in the sense that the
84 decoy was dominated only by one item, called the *target*, and not by the other,
85 called the *competitor* (supplementary information (SI), S1.4).

86 The trials that included a decoy (called the decoy trials) were paired, so
87 that each undominated alternative would be the target in one trial and the
88 competitor in the other. Such a pair of decoy trials has the form $\{A, B, A'\}$ and
89 $\{A, B, B'\}$, and this will be referred to as a *matching pair* (of decoy trials). The
90 undominated alternatives A and B are the same in both trials of the matching
91 pair. The only difference is that, in $\{A, B, A'\}$ the decoy A' is dominated by
92 A but not by B , while in $\{A, B, B'\}$ the decoy B' is dominated by B but not
93 by A ¹⁴. Matching pairs of decoy trials are the main focus of our analysis. We
94 do not compare three-item decoy trials with two-item trials without decoys,
95 because choosing from a larger menu may be different from choosing from a
96 smaller menu, and this may confound the analysis. Four trials are illustrated in
97 Figure 1a.

98 Behaviorally, 81.20% of the matching pairs were *consistent*, in the sense that
99 the choice was independent of the decoy. That is, the participant chose the same
100 alternative from $\{A, B, A'\}$ and $\{A, B, B'\}$. 16.41% exhibited *preference reversals*,
101 in the sense that the participant switched her preference from A to B when
102 the decoy changed from A' to B' (Figure 1b)¹⁵. That is, the participant chose A
103 from $\{A, B, A'\}$ and B from $\{A, B, B'\}$. In terms of choice behavior, context-
104 independence seems to be the norm (accounting for 81.20% of the matching
105 pairs). But valuations may be context-dependent even when choice is not, be-
106 cause a change in values does not necessarily lead to a change in the chosen
107 alternative.

108 In a decoy trial, the choice is said to be *Along* the decoy if the target is chosen
109 and *Against* the decoy if the competitor is chosen. In a consistent matching pair,
110 the same alternative, say A , is chosen both from $\{A, B, A'\}$ and $\{A, B, B'\}$. The
111 choice is hence *Along* the decoy when A is chosen from $\{A, B, A'\}$ and *Against*
112 the decoy when A is chosen from $\{A, B, B'\}$. The response time (RT) suggests
113 that it was easier to make a choice along the decoy than against the decoy
114 in consistent trials. Across all participants, the average RT was 8.36 seconds
115 in *Along* trials and 9.94 seconds in *Against* trials. The average percentage
116 decrease of RT in *Along* is 16.20% which is strongly significantly different from
117 zero (Figure 1c). Thus, in consistent trials, where by definition the chosen
118 alternative does not depend on the decoy, choosing along the decoy still seems
119 easier than choosing against it.

120 We will look for a neural manifestation of a decoy effect, in the sense of
121 context-dependent valuations, in consistent trials. If valuations are context-
122 independent, the valuation area should not differentiate between the *Along* and
123 *Against* trials of a consistent matching pair. But if valuations are context-
124 dependent, an alternative will tend to be more highly valued when it dominates
125 the decoy, so the *Along* trial will be more rewarding than the *Against* trial
126 even though the chosen alternative is the same. Context-dependence therefore
127 leads to the prediction that reward-sensitive areas will be more activated in

128 Along than in Against trials of consistent matching pairs. This prediction was
129 supported by our fMRI data. The left ventral striatum was significantly more
130 active in the Along than in the Against consistent trials (Figure 2a). The
131 striatum activity was observed at putamen and extended medially to caudate
132 and anteriorly to insula¹⁶. Notice that the Along-Against contrast involves a
133 pair of trials with three-item menus, the same undominated items appear on
134 both menus, and the same item is chosen from both menus. The only difference
135 to which differential activity can be attributed is that the chosen item dominates
136 the decoy in the Along but not in the Against trials¹⁷.

137 Two additional pieces of evidence support the interpretation that choosing
138 along the decoy is more rewarding than choosing against it. First, although pre-
139 vious research has indicated that the striatum does code reward^{18,19}, we provide
140 direct evidence for the connection between the striatum and reward by estimat-
141 ing the utility of each alternative A , denoted $u(A)$, using data from post-test
142 choices in two-item menus with no decoys (SI, S4.2). The utility $u(A)$ depends
143 on item A 's price as well as its quality. We refer to $u(A)$ as the intrinsic utility
144 of A , because it corresponds to the utility from A 's intrinsic properties in the
145 absence of a decoy. We found that the striatum activity correlates parametri-
146 cally with the estimated intrinsic utility of the chosen item positively (Figure
147 2b). Hence in our data, there is evidence that the striatum may be coding
148 reward. The active cluster overlaps with what we discovered when contrasting
149 Along to Against (SI, Figure S1), rendering our interpretation that it was more
150 rewarding choosing along the decoy than choosing against it plausible. Second,
151 across participants, the differential striatum activity in Along versus in Against
152 positively correlated with two behavioral measures of how “decoyable” a par-
153 ticipant is. The first such measure is simply the number of preference reversals
154 (Figure 2c). The second measure, denoted μ , will be introduced below (Figure
155 2d). Thus, more decoyable participants experienced a larger increase in reward
156 when the chosen alternative dominated the decoy. This match between the be-
157 havioral measures and the neural effect is consistent with the hypothesis that

158 context-dependent valuation causes the decoy effect.

159 How well an alternative will satisfy a person’s wants and needs depends on
160 the intrinsic properties of price and quality. If valuations are influenced by irrel-
161 evant aspects of the context, choices may become biased (i.e., the intrinsically
162 less valuable option may be chosen). Unbiased choice may require some cogni-
163 tive control. To look into how control could be recruited, we first construct a
164 measure of the trial-by-trial strength of the decoy effect, which will guide us to
165 the potential effect of cognitive control on valuation.

166 The number of preference reversals can be used to compare, between-participant,
167 who is more decoyable. But it cannot be used to determine within-participant
168 when she feels the decoy effect more strongly. Moreover, as argued above, in
169 consistent trials the decoy may have an effect on valuations which is not reflected
170 in choices. To construct a trial-by-trial estimate of the decoy effect, we modify
171 a simplified linear ballistic accumulator (LBA) model. The idea behind LBA
172 is that choice is made when the accumulated evidence in favor of the chosen
173 alternative has reached a threshold. The accumulation of evidence is faster if
174 the evidence is stronger, so on average there is an inverse relationship between
175 the strength of evidence and RT^{20} . Translating these ideas into our experiment,
176 if the menu in trial t is $\{A, B, A'\}$, because A' is dominated the race will be
177 between A and B . Intrinsic utilities $u(A)$ and $u(B)$ were estimated in the post-
178 test as mentioned above. Thus, the intrinsic utility difference between A and
179 B is $u(A) - u(B)$. We hypothesized that the decoy effect in trial t causes the
180 relative value assigned to the target to be shifted by some amount $d(t)$. Since A
181 is the target in $\{A, B, A'\}$, due to the presence of the decoy the *decision utility*
182 *difference* will be $u(A) - u(B) + d(t)$. Here $d(t)$ represents the amount by which
183 the decoy effect favors the target, A , relative to the competitor, B , in trial t .
184 While we might expect $d(t)$ to be positive on average, we make no assumption
185 regarding the sign of $d(t)$.

186 In terms of the LBA model, the strength of the “evidence” favoring A is

187 $u(A) - u(B) + d(t)$, the sum of the intrinsic utility difference and the decoy
 188 effect of this trial $d(t)$. Thus, if A is chosen from $\{A, B, A'\}$ in trial t , we
 189 assume RT is inversely proportional to $u(A) - u(B) + d(t)$,

$$RT = \frac{T}{u(A) - u(B) + d(t)} \quad (1)$$

190 for some threshold constant $T > 0$. If instead the menu in trial t is $\{A, B, B'\}$,
 191 then since B is the target, the decoy effect favors B . Therefore, if $d(t)$ is the
 192 strength of the decoy effect in trial t , the decision utility difference, i.e., the
 193 strength of the “evidence” favoring A , is $u(A) - u(B) - d(t)$. Thus, if A is
 194 chosen from $\{A, B, B'\}$ then RT is inversely proportional to $u(A) - u(B) - d(t)$,

$$RT = \frac{T}{u(A) - u(B) - d(t)}. \quad (2)$$

195 For any decoy trial, $d(t)$ satisfies equation (1) in Along trials (where A
 196 is chosen from $\{A, B, A'\}$) and (2) in Against trials (where A is chosen from
 197 $\{A, B, B'\}$). We assume that $d(t)$ is drawn from a normal distribution with
 198 participant-specific mean μ and standard deviation σ . Note that μ can be inter-
 199 preted as a participant’s average decoy effect. Even though we expect it to be
 200 positive, whether it is positive, zero or negative is left open. We use maximum
 201 likelihood to estimate T , μ and σ for each participant. The trial-by-trial decoy
 202 effect $d(t)$ is then backed out from (1) in Along trials and from (2) in Against
 203 trials. Notice that the sign in front of $d(t)$ is different in (1) and in (2), because
 204 the decoy effect favors a different alternative in the two cases. A random com-
 205 ponent of the intrinsic utility difference would not flip sign in this way, so $d(t)$
 206 cannot be interpreted as an estimate of random utility²¹. Similarly, $d(t)$ is not
 207 a monotonic function of RT. In Along trials, the shorter RT is the *stronger* the
 208 decoy effect is (because the decoy is “helping” the target to be chosen quickly),
 209 but in Against trials, the shorter RT is the *weaker* the decoy effect is (because
 210 the decoy is working against the choice of the competitor).

211 The results of the estimation suggest that the model is appropriate. A partic-
 212 ipant’s estimated average decoy effect, μ , is a measure of how “decoyable”

213 she is. If the decoy indeed makes the target more valuable, μ is expected to
214 be positive. For all but two participants, it is indeed positive. Furthermore, μ
215 strongly positively correlated with the number of preference reversals (Figure
216 3a) suggesting it effectively captures how decoyable each participant is. As men-
217 tioned above, μ also strongly positively correlated with the differential striatum
218 activity in Along versus in Against. This strengthens the case for the higher
219 activity in Along than in Against being evidence that participants find choos-
220 ing along the decoy more rewarding than choosing against it. Since T can be
221 interpreted as the threshold for making a decision, a larger T should imply a
222 longer RT, and indeed the positive correlation between T and the average RT
223 of decoy trials is very strong (SI, S4.5).

224 For each participant, we divided the decoy trials into two halves depending
225 on the size of $d(t)$. The half with large $d(t)$ was classified as Strong trials, sug-
226 gesting a strong decoy effect. The half with small $d(t)$ was classified as Weak
227 trials. Neurally, no region was found contrasting Strong to Weak. The left in-
228 ferior parietal lobule (IPL) was the only region more active contrasting Weak
229 to Strong (Figure 3b). Its activity in fact parametrically tracks $d(t)$ negatively
230 (Figure 3c)²². IPL has previously been implicated in goal-directed preparation
231 of attention^{23,24,25}. In particular, IPL close to our activation cluster is activated
232 more when there is a switch than when a task is repeated, when trials are in-
233 congruent than when they are congruent in the Stroop or flanker tasks, when
234 stimulus and response are incompatible than when they are compatible, and
235 when working memory is more heavily taxed^{26,27,28,29}. The common denomi-
236 nator in these experiments is that IPL may support the allocation of attention
237 to facilitate task-relevant representations. As part of a general network sub-
238 serving voluntary attention, this may explain why there is higher IPL activity,
239 suggesting heightened attention, when the decoy effect $d(t)$ is smaller.

240 If the decoy effect increases the valuation of the target, it could potentially
241 lead to the wrong decision in trials where the target is intrinsically worse than
242 the competitor. Presumably, it is impossible to inhibit processing of the decoy

243 at the perceptual level. But at a higher level, control may prevent interference
244 from the decoy and help refocus on the relevant aspects of the choice situation
245 (the intrinsic properties of the undominated alternatives). We do not find a
246 direct modulation from the IPL in our data. We do find a plausible indirect
247 modulation which we will further verify by a transcranial magnetic stimulation
248 (TMS) experiment.

249 We looked into how the putative control was recruited by the following two
250 psychophysiological interaction analyses. We divided the decoy trials into two
251 halves. The half where the target has a lower intrinsic utility than the com-
252 petitor is categorized as Conflict. The other half, where the target has higher
253 intrinsic utility, is categorized as NoConflict³⁰. In a Conflict trial, the target has
254 lower intrinsic utility, and thus control may be required to override the decoy
255 effect and prevent the target from being chosen. In a NoConflict trial, the target
256 has higher intrinsic utility, so the decoy effect is less problematic and control
257 presumably less critical. When attention is heightened and control is more likely
258 to make a difference, we might expect a stronger connectivity between the area
259 supporting attention and that implementing control.

260 We took the average activity of a 4-mm sphere surrounding a peak voxel of
261 IPL as the seed to examine whole-brain whether any area exhibits stronger func-
262 tional connectivity with IPL in Conflict than in NoConflict. The right inferior
263 frontal gyrus (rIFG) is the only region that has this stronger task-related con-
264 nectivity (Figure 4a). The rIFG is postulated to be a site where goal-directed
265 and stimulus-driven attention converge^{31,23}. Neuroimaging studies implicate this
266 area in supporting inhibition to implement control^{32,33,34}. It is more active in
267 No-Go/Stop-Signal trials than in Go trials, in invalid Posner cueing trials than
268 in valid ones, in trials where recent history may interfere than when it may not –
269 all possibly reflecting processes related to overriding prepotent responses^{7,8}. The
270 rIFG has been implicated in a number of other exertions requiring self-control,
271 such as inhibiting an incorrect answer to a problem of logic³⁵ and focusing on
272 the facts rather than the framing of a question³⁶. It has even been implicated

273 in the decision to quit smoking³⁷. The presence of the irrelevant decoy is more
274 critical to decision-making in a Conflict trial than in a NoConflict trial³⁸. The
275 stronger functional connectivity between IPL and rIFG could possibly be due to
276 the posterior area signaling the conflicting representations to the frontal cortex
277 for control²⁸. Anatomical connections between IPL and IFG have been demon-
278 strated using diffusion-weighted imaging^{39,40}. As previous studies indicate that
279 rIFG may play a role in inhibiting irrelevant responses, our connectivity result
280 caused us to explore further the role of rIFG in our data.

281 We have argued that the striatum seems to code the decision utility of the
282 chosen item. The observation that the striatum is more active in Along than
283 in Against suggests that the decision utility of a chosen item tends to increase
284 when it is the target, which may lead to biased choices in Conflict trials. If
285 the role of control is to reduce this bias, it would tend to offset this differential
286 activity in Along versus in Against by either increasing the decision utility of
287 the chosen option in Against trials, or reducing it in Along trials, or both. This
288 suggests that the control area may correlate more strongly with the reward
289 area in Against than in Along Conflict trials. In Conflict trials the target is
290 intrinsically worse than the competitor, so if the choice is along the decoy then
291 control was not successfully applied, as the intrinsically best alternative is not
292 chosen; if the choice is against the decoy then the intrinsically best alternative
293 is chosen⁴¹. Hence, by considering the differential connectivity with a control
294 area in Against versus in Along Conflict trials, we are contrasting trials where
295 control is successful with trials where control is unsuccessful.

296 We took the average activity of a 4-mm sphere surrounding a peak voxel of
297 rIFG identified above as the seed to examine whether the striatum has a dif-
298 ferential functional connectivity with rIFG in Against than in Along of Conflict
299 trials. In a whole-brain search, the left striatum shows a stronger task-related
300 functional connectivity with rIFG (Figure 4b). The active cluster overlaps with
301 what we discovered when contrasting Along to Against of consistent trials (SI,
302 Figure S5). This is consistent with a possible role of rIFG in reducing the decoy-

303 induced bias in the decision utility. Anatomical connections between IFG and
304 striatum have been demonstrated. Restricted frontostriatal diffusion seems to
305 correlate with greater control, hinting at its contribution to the recruitment of
306 control⁴².

307 The two connectivity results provide evidence on the role of rIFG. rIFG
308 exhibits differential functional connectivity with the striatum, hinting at its
309 possible influence on choices. As it is often implicated in overriding irrelevant
310 responses, rIFG may play a role in implementing control to overcome a potential
311 decoy-induced decision-bias. We used TMS to investigate this further (SI, S5).
312 We applied theta burst TMS to temporarily interfere with the region of interest
313 before participants started the choice task^{43,44}. In one group of participants,
314 the IFG group, the site of stimulation was at the peak voxel of rIFG identified
315 from the fMRI experiment. In the other group, the vertex group, the site of
316 stimulation was the vertex. Each group had 32 participants, the same number
317 as in the fMRI experiment. If rIFG plays a role in inhibiting the decoy-induced
318 bias, then when it is temporarily disrupted, the inhibitory control is expected to
319 be weaker, and the decoy effect is expected to be stronger⁴⁵. In terms of choices,
320 we thus expected the IFG group to exhibit more preference reversals than the
321 vertex group. In terms of RT, because a strong decoy effect is expected to make
322 it much easier to choose along the decoy than against it, we expected the IFG
323 group to have a larger percentage decrease of RT in Along trials, compared with
324 Against trials, than the vertex group.

325 These expectations were born out in the TMS experiment. The IFG group
326 on average had a preference reversal rate of 25.30% whereas the vertex group
327 had 18.72%, an increase of 6.58 percentage points. The t-test for comparing the
328 preference reversal rate of the IFG group with that of the vertex group has a
329 p-value (one-tailed) of 0.048 (Figure 4c). For the IFG group, on average, the RT
330 was 7.23 seconds in Along trials and 9.08 seconds in Against trials. The average
331 percentage decrease of RT in Along trials, compared with Against trials, was
332 19.22%. For the vertex group, on average, the RT was 6.68 seconds in Along

333 trials and 7.47 seconds in Against trials. The average percentage decrease of
334 RT in Along trials, compared with Against trials, was 9.05%. The t-test for
335 comparing the average percentage decrease of RT in Along trials, compared
336 with Against trials, of the IFG group with that of the vertex group has a p-
337 value (one-tailed) of 0.023 (Figure 4d).

338 In some naturally occurring choice situations, context-dependent valuation
339 may be adaptive⁴. In sequential decision problems, optimal decisions depend
340 on a comparison of the current option with its background⁵. Discovering an
341 option that dominates other alternatives in the background can make it opti-
342 mal to choose it right away. Thus, a context-dependent valuation that gives a
343 positive connotation to dominance may be adaptive. There are, however, sev-
344 eral reasons why such heuristics may sometimes lead to suboptimal decisions⁴⁶.
345 First, someone may be trying to manipulate the decision maker. For example,
346 a product’s intrinsic properties of price and quality will determine how well
347 it will satisfy the consumer’s wants and needs, i.e., will determine the expe-
348 rienced (or intrinsic) utility – but by introducing artificial decoy alternatives,
349 marketers make the decision utility exceed the experienced utility. Second, the
350 assumptions underlying the normal operation of the system may be violated.
351 For example, the heuristics may be well adapted to sequential tasks, but not
352 to the rather artificial static tasks in our experiment. Third, some decision
353 problems are too complex for purely affective judgments, and require more de-
354 liberative and “rational” cognitive processes. All three reasons are relevant to
355 our experiment. The decoy effect may be an evolved adaptation that illumi-
356 nates the on-line judgments made by the decision-making system, rather than
357 a “design flaw.” It may be analogous to visual illusions, e.g., the Ebbinghaus
358 illusion where the perceived size of an object depends on the sizes of neighboring
359 objects^{47,48}, which is adaptive in ecologically relevant scenarios.

References and Notes

1. Sasaki, T. & Pratt, S. C. Emergence of group rationality from irrational individuals. *Behav Ecol* **22**, 276-281, doi:DOI 10.1093/beheco/arq198 (2011).
2. Huber, J., Payne, J. W. & Puto, C. Adding Asymmetrically Dominated Alternatives - Violations of Regularity and the Similarity Hypothesis. *J Consum Res* **9**, 90-98, doi:Doi 10.1086/208899 (1982).
3. Ariely, D. *Predictably irrational : the hidden forces that shape our decisions*. 1st edn, (Harper, 2008), chap. 1.
4. Houston, A. I. Natural selection and context-dependent values. *P Roy Soc B-Biol Sci* **264**, 1539-1541, doi:DOI 10.1098/rspb.1997.0213 (1997).
5. Freidin, E. & Kacelnik, A. Rational choice, context dependence, and the value of information in European starlings (*Sturnus vulgaris*). *Science* **334**, 1000-1002, doi:10.1126/science.1209626 (2011).
6. Kahneman, D. & Frederick, S. Frames and brains: elicitation and control of response tendencies. *Trends in cognitive sciences* **11**, 45-46, doi:10.1016/j.tics.2006.11.007 (2007).
7. Levy, B. J. & Wagner, A. D. Cognitive control and right ventrolateral prefrontal cortex: reflexive reorienting, motor inhibition, and action updating. *Annals of the New York Academy of Sciences* **1224**, 40-62, doi:10.1111/j.1749-6632.2011.05958.x (2011).
8. Bunge, S. A., Ochsner, K. N., Desmond, J. E., Glover, G. H. & Gabrieli, J. D. Prefrontal regions involved in keeping information in and out of mind. *Brain : a journal of neurology* **124**, 2074-2086 (2001).
9. Padoa-Schioppa, C. & Assad, J. A. Neurons in the orbitofrontal cortex encode economic value. *Nature* **441**, 223-226, doi:10.1038/nature04676 (2006).

- 387 10. Padoa-Schioppa, C. & Assad, J. A. The representation of economic value
388 in the orbitofrontal cortex is invariant for changes of menu. *Nature neu-*
389 *roscience* **11**, 95-102, doi:10.1038/nn2020 (2008).
- 390 11. Padoa-Schioppa, C. Range-adapting representation of economic value in
391 the orbitofrontal cortex. *The Journal of neuroscience : the official journal*
392 *of the Society for Neuroscience* **29**, 14004-14014, doi:10.1523/JNEUROSCI.3751-
393 09.2009 (2009).
- 394 12. Nieuwenhuis, S. *et al.* Activity in human reward-sensitive brain areas is
395 strongly context dependent. *NeuroImage* **25**, 1302-1309 (2005).
- 396 13. Breiter, H. C., Aharon, I., Kahneman, D., Dale, A. & Shizgal, P. Func-
397 tional imaging of neural responses to expectancy and experience of mon-
398 etary gains and losses. *Neuron* **30**, 619-639 (2001).
- 399 14. The order of the trials and the ordering of the items in a trial were ran-
400 domized, except that in the matching pairs $\{A, B, A'\}$ and $\{A, B, B'\}$, the
401 ordering of target, competitor, and decoy was controlled to be the same
402 to make the two trials in a matching pair as similar as possible.
- 403 15. We denote a participant by her.
- 404 16. Two participants were removed from the fMRI analysis because of exces-
405 sive head motions.
- 406 17. A previous fMRI study⁴⁹ contrasted three-item trials of the form $\{A, B, A'\}$
407 with two-item trials of the form $\{A, B\}$. In the three-item trials, there was
408 decreased activation in the amygdala, MPFC, and right IPL, and increased
409 activation in DLPFC and ACC. Based on this, it was argued that the de-
410 coy effect may be caused by a shift toward more reason-based (heuristic)
411 choice processes. Contrasting menus of different sizes introduces a pos-
412 sible confound, as it may be more cognitively demanding to choose from
413 a larger menu. Our fMRI analysis was designed to contrast very similar
414 three-item trials.

- 415 18. Delgado, M. R. Reward-related responses in the human striatum. *Annals*
416 *of the New York Academy of Sciences* **1104**, 70-88, doi:10.1196/annals.1390.002
417 (2007).
- 418 19. Kable, J. W. & Glimcher, P. W. The neural correlates of subjective
419 value during intertemporal choice. *Nature neuroscience* **10**, 1625-1633,
420 doi:10.1038/nn2007 (2007).
- 421 20. Brown, S. D. & Heathcote, A. The simplest complete model of choice
422 response time: linear ballistic accumulation. *Cognitive psychology* **57**,
423 153-178, doi:10.1016/j.cogpsych.2007.12.002 (2008).
- 424 21. If we define $d(t)$ by the *same* formula, say (1), in *both* Along and Against
425 trials, then $d(t)$ would correspond to a random component of utility, but
426 it would not be a measure of the decoy effect. In SI, S4.11 and S4.12, we
427 address the issue of random utility and show that if $d(t)$ were defined by
428 the same formula Along and Against then it would not be significantly
429 correlated with the inferior parietal lobule, whereas – as reported below –
430 there is significant correlation when we use (1) in Along trials and (2) in
431 Against trials.
- 432 22. We perform several robustness checks on IPL's negative correlation with
433 $d(t)$. See SI, S4.7.
- 434 23. Shomstein, S. Cognitive functions of the posterior parietal cortex: top-
435 down and bottom-up attentional control. *Frontiers in integrative neuro-*
436 *science* **6**, 38, doi:10.3389/fnint.2012.00038 (2012).
- 437 24. Hopfinger, J. B., Buonocore, M. H. & Mangun, G. R. The neural mech-
438 anisms of top-down attentional control. *Nature neuroscience* **3**, 284-291,
439 doi:10.1038/72999 (2000).
- 440 25. Raz, A. & Buhle, J. Typologies of attentional networks. *Nature reviews.*
441 *Neuroscience* **7**, 367-379, doi:10.1038/nrn1903 (2006).

- 442 26. Liston, C., Matalon, S., Hare, T. A., Davidson, M. C. & Casey, B. J.
443 Anterior cingulate and posterior parietal cortices are sensitive to dissocia-
444 ble forms of conflict in a task-switching paradigm. *Neuron* **50**, 643-653,
445 doi:10.1016/j.neuron.2006.04.015 (2006).
- 446 27. Ye, Z. & Zhou, X. Conflict control during sentence comprehension: fMRI
447 evidence. *NeuroImage* **48**, 280-290, doi:10.1016/j.neuroimage.2009.06.032
448 (2009).
- 449 28. Sylvester, C. Y. et al. Switching attention and resolving interference:
450 fMRI measures of executive functions. *Neuropsychologia* **41**, 357-370
451 (2003).
- 452 29. McNab, F. et al. Common and unique components of inhibition and work-
453 ing memory: an fMRI, within-subjects investigation. *Neuropsychologia*
454 **46**, 2668-2682, doi:10.1016/j.neuropsychologia.2008.04.023 (2008).
- 455 30. The decoy trial $\{A, B, A'\}$ is defined as a Conflict trial if $u(A) - u(B) < 0$,
456 and as a NoConflict trial if $u(A) - u(B) > 0$. As the decoy trials are
457 in matching pairs, exactly half of them are Conflict whereas the other
458 half are NoConflict. For instance, if $u(A) - u(B) < 0$, then $\{A, B, A'\}$ is
459 Conflict and $\{A, B, B'\}$ is NoConflict. See SI, S4.8.
- 460 31. Asplund, C. L., Todd, J. J., Snyder, A. P. & Marois, R. A central role for
461 the lateral prefrontal cortex in goal-directed and stimulus-driven attention.
462 *Nature neuroscience* **13**, 507-512, doi:10.1038/nn.2509 (2010).
- 463 32. Aron, A. R., Robbins, T. W. & Poldrack, R. A. Inhibition and the right
464 inferior frontal cortex: one decade on. *Trends in cognitive sciences* **18**,
465 177-185, doi:10.1016/j.tics.2013.12.003 (2014).
- 466 33. Nubert, F.-X., Mars, R. B. & Rushworth, M. F. S. Is there an inferior
467 frontal cortical network for cognitive control and inhibition? In *Princi-*
468 *ples of Frontal Lobe Function*, Stuss, D. T., Knight, R. T. Eds. (Oxford
469 University Press, New York, 2012), pp. 332-352.

- 470 34. Dillon, D. G. & Pizzagalli, D. A. Inhibition of Action, Thought, and Emo-
471 tion: A Selective Neurobiological Review. *Applied & preventive psychology*
472 : *journal of the American Association of Applied and Preventive Psychol-*
473 *ogy* **12**, 99-114, doi:10.1016/j.appsy.2007.09.004 (2007).
- 474 35. Goel, V. & Dolan, R. J. Explaining modulation of reasoning by belief.
475 *Cognition* **87**, B11-22 (2003).
- 476 36. De Martino, B., Kumaran, D., Seymour, B. & Dolan, R. J. Frames, biases,
477 and rational decision-making in the human brain. *Science* **313**, 684-687,
478 doi:10.1126/science.1128356 (2006).
- 479 37. Berkman, E. T., Falk, E. B. & Lieberman, M. D. In the trenches of real-
480 world self-control: neural correlates of breaking the link between craving
481 and smoking. *Psychological science* **22**, 498-506, doi:10.1177/0956797611400918
482 (2011).
- 483 38. Unlike the previous experiments that manipulate control explicitly, our
484 design is not aimed at studying control. The classification of Conflict
485 and NoConflict is based on the estimation of intrinsic utilities and so the
486 difference between them could be quite subtle.
- 487 39. Caspers, S. et al. Probabilistic fibre tract analysis of cytoarchitecton-
488 ically defined human inferior parietal lobule areas reveals similarities to
489 macaques. *NeuroImage* **58**, 362-380, doi:10.1016/j.neuroimage.2011.06.027
490 (2011).
- 491 40. Rushworth, M. F., Behrens, T. E. & Johansen-Berg, H. Connection pat-
492 terns distinguish 3 regions of human parietal cortex. *Cerebral cortex* **16**,
493 1418-1430, doi:10.1093/cercor/bhj079 (2006).
- 494 41. In Conflict trials, 66.30% of Along trials (together with their matching
495 trials) exhibit preference reversal (signifying unsuccessful control), while
496 96.96% of Against trials (together with their matching trials) are consis-
497 tent. In other words, the differential connectivity in Against than in Along

- 498 Conflict trials could also be phrased as the differential connectivity when
499 choices are consistent than when choices largely exhibit reversal.
- 500 42. Liston, C. et al. Frontostriatal microstructure modulates efficient recruit-
501 ment of cognitive control. *Cerebral cortex* **16**, 553-560, doi:10.1093/cercor/bhj003
502 (2006).
- 503 43. Huang, Y. Z., Edwards, M. J., Rounis, E., Bhatia, K. P. & Rothwell, J. C.
504 Theta burst stimulation of the human motor cortex. *Neuron* **45**, 201-206,
505 doi:10.1016/j.neuron.2004.12.033 (2005).
- 506 44. Chao, C. M. et al. Predictability of saccadic behaviors is modified by
507 transcranial magnetic stimulation over human posterior parietal cortex.
508 *Human brain mapping* **32**, 1961-1972, doi:10.1002/hbm.21162 (2011).
- 509 45. Jacobson, L., Javitt, D. C. & Lavidor, M. Activation of inhibition: di-
510 minishing impulsive behavior by direct current stimulation over the in-
511 ferior frontal gyrus. *Journal of cognitive neuroscience* **23**, 3380-3387,
512 doi:10.1162/jocn_a_00020 (2011).
- 513 46. Slovic, P., Finucane, M. L., Peters, E. & MacGregor, D. G. The affect
514 heuristic. *Eur J Oper Res* **177**, 1333-1352, doi:DOI 10.1016/j.ejor.2005.04.006
515 (2007).
- 516 47. Eagleman, D. M. Visual illusions and neurobiology. *Nature reviews. Neu-*
517 *roscience* **2**, 920-926, doi:10.1038/35104092 (2001).
- 518 48. Goodale, M. A. & Haffenden, A. Frames of reference for perception and
519 action in the human visual system. *Neuroscience and biobehavioral reviews*
520 **22**, 161-172 (1998).
- 521 49. Hedgcock, W. & Rao, A. R. Trade-Off Aversion as an Explanation for the
522 Attraction Effect: A Functional Magnetic Resonance Imaging Study. *J*
523 *Marketing Res* **46**, 1-13 (2009).

Figure 1

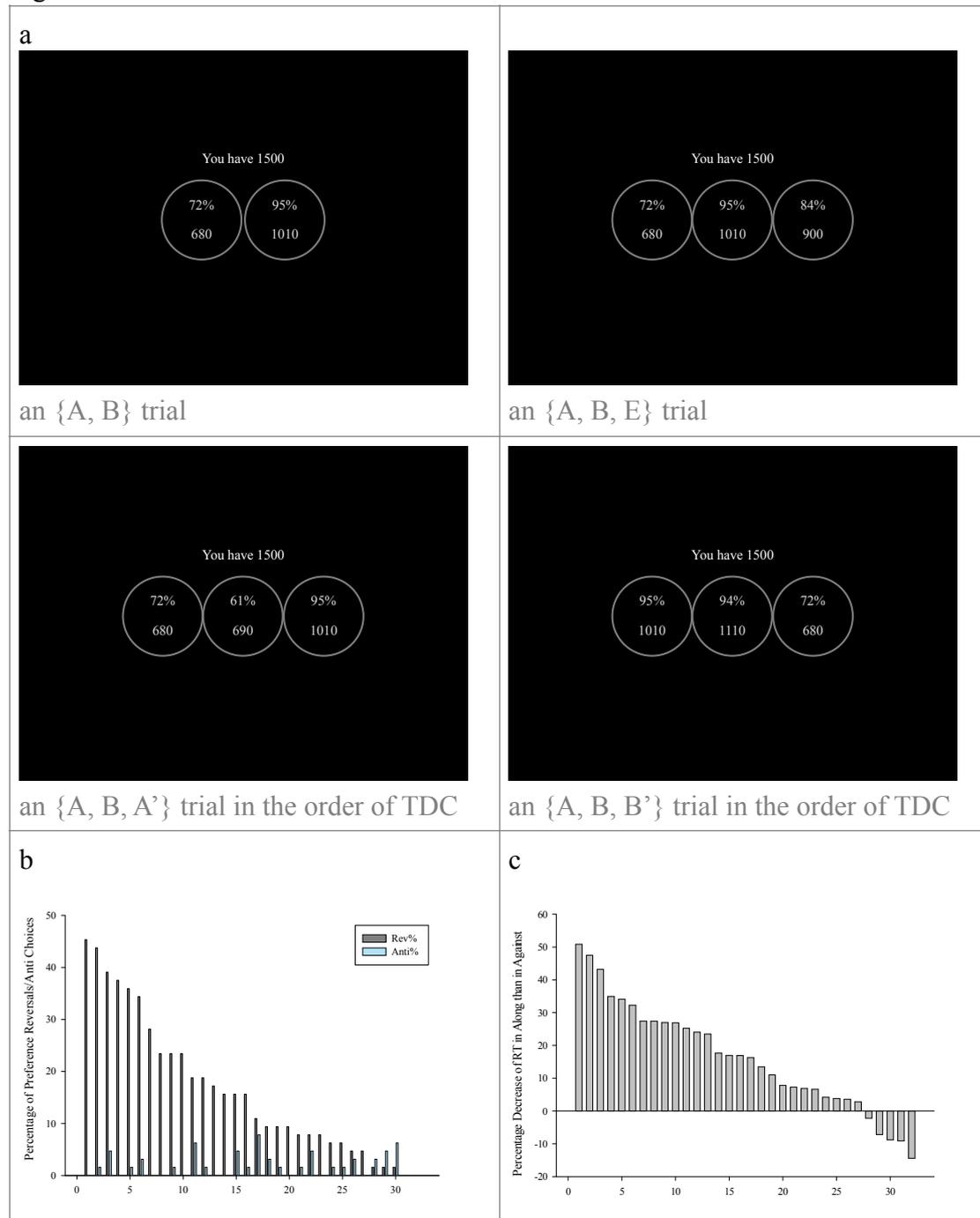


Figure 1a Sample screens from the experiment. Participants were endowed with 1500 New Taiwan Dollars (roughly 50 US dollars) to spend in each trial and made choices from a series of two- and three-item menus. The items on the menus were restaurant meals with specified prices and quality levels in which higher percentages imply better qualities (SI, S1.2 and S1.4). The trial in the top left corner is a two-item trial whereas that in the top right corner is a three-item trial. The two trials in the bottom panel are a matching pair of decoy trials. In this example, option A is the item with quality 72% and price \$680; option B has quality 95% and price \$1010; option E has

quality 84% and price \$900; option A' has quality 61% and price \$690; option B' has quality 94% and price \$1110. An item is dominated by another item if it has both lower quality and higher price. Thus, A' is dominated by A but not by B, whereas B' is dominated by B but not by A. The target, denoted T, is the item that dominates the decoy, denoted D; there is no domination relationship between the competitor, C, and the decoy. Hence in the bottom left trial, the target is A (quality 72% and price \$680), the decoy is A' (61%, \$690) and the competitor is B (95%, \$1010). In the bottom right trial, B (95%, \$1010) is the target, B' (94%, \$1110) the decoy and A (72%, \$680) the competitor.

Figure 1b Percentage of preference reversals and anti choices. The figure is a participant-by-participant breakdown of the percentage of preference reversals and anti choices of decoy trials (in descending order of the former). Preference reversals occur when A is chosen from {A, B, A'} and B from {A, B, B'} whereas anti choices occur when B is chosen from {A, B, A'} and A from {A, B, B'}. The average percentage of preference reversals was 16.41 and that of anti choices was 2.05 over all 32 participants. Two participants had neither preference reversals nor anti choices.

Figure 1c Response times were shorter in Along trials than in Against trials. In consistent trials, the response times were significantly shorter when choices are along the decoy than when they are against the decoy. The figure is a participant-by-participant breakdown of the percentage decrease of RT in Along trials compared with Against trials (in descending order). The average percentage decrease of RT over all 32 participants was 16.20. The percentage decrease of RT is the difference of RT in Against trials and Along trials divided by the average of them in percentage.

Figure 2

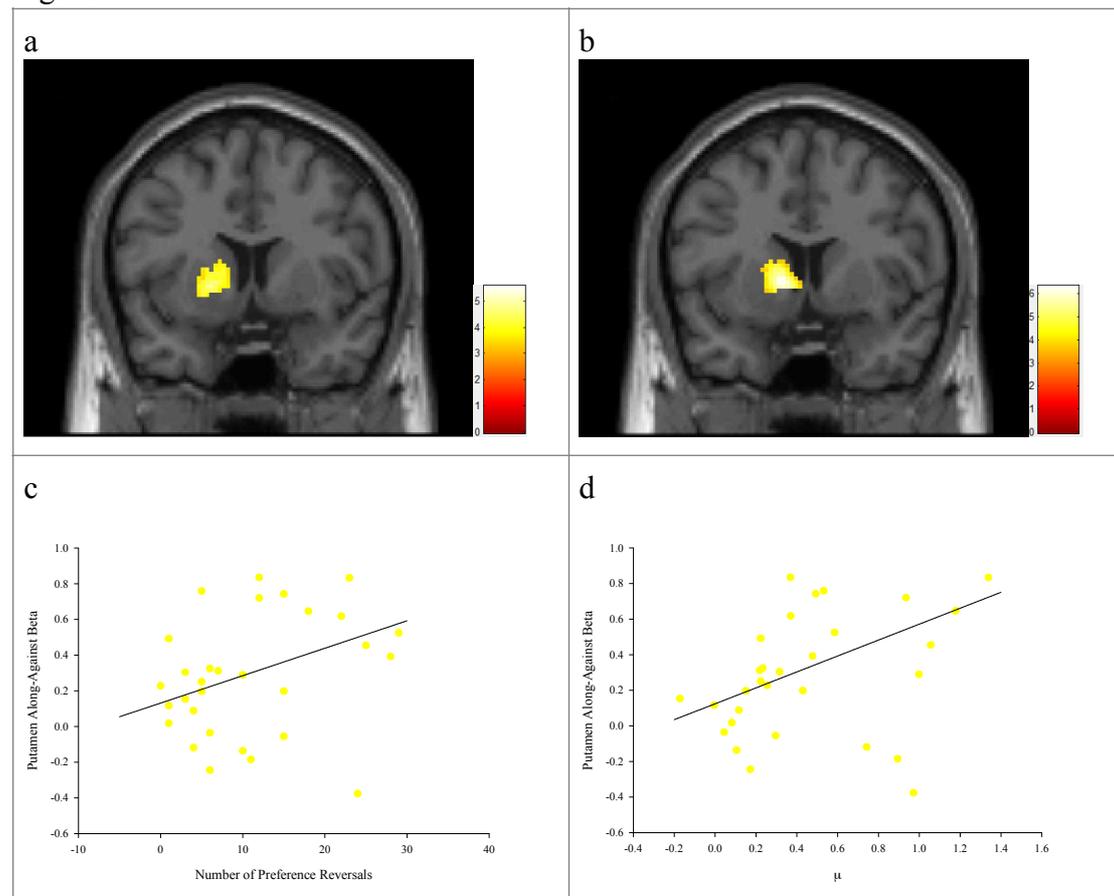


Figure 2a The left striatum was more active Along than Against in consistent trials ($y=8$) [$P<0.05$ at whole-brain cluster correction with a t threshold of 3.40 and an extent of 398 voxels] (SI, S4.1 and table S1).

Figure 2b The left striatum ($y=8$) correlated parametrically with the estimated intrinsic utility of the chosen option positively [$P<0.05$ at whole-brain cluster correction with a t threshold of 3.40 and an extent of 248 voxels] (SI, S4.2 and table S2).

Figure 2c Positive correlation between the number of preference reversals and the differential striatum activity Along versus Against. The robust regression slope is 0.015 and the p -value (two-tailed) is 0.057. The putamen activity is based on a 4-mm sphere centered at (-18, 8, 0). For other peak voxels of striatum, see SI, S4.3, table S5 and figure S2.

Figure 2d Positive correlation between the estimated μ and the differential striatum activity Along versus Against. The robust regression slope is 0.447 and the p -value (two-tailed) is 0.006. The putamen activity is based on a 4-mm sphere centered at (-18, 8, 0). For other peak voxels of striatum, see SI, S4.3, table S5 and figure S2.

Figure 3

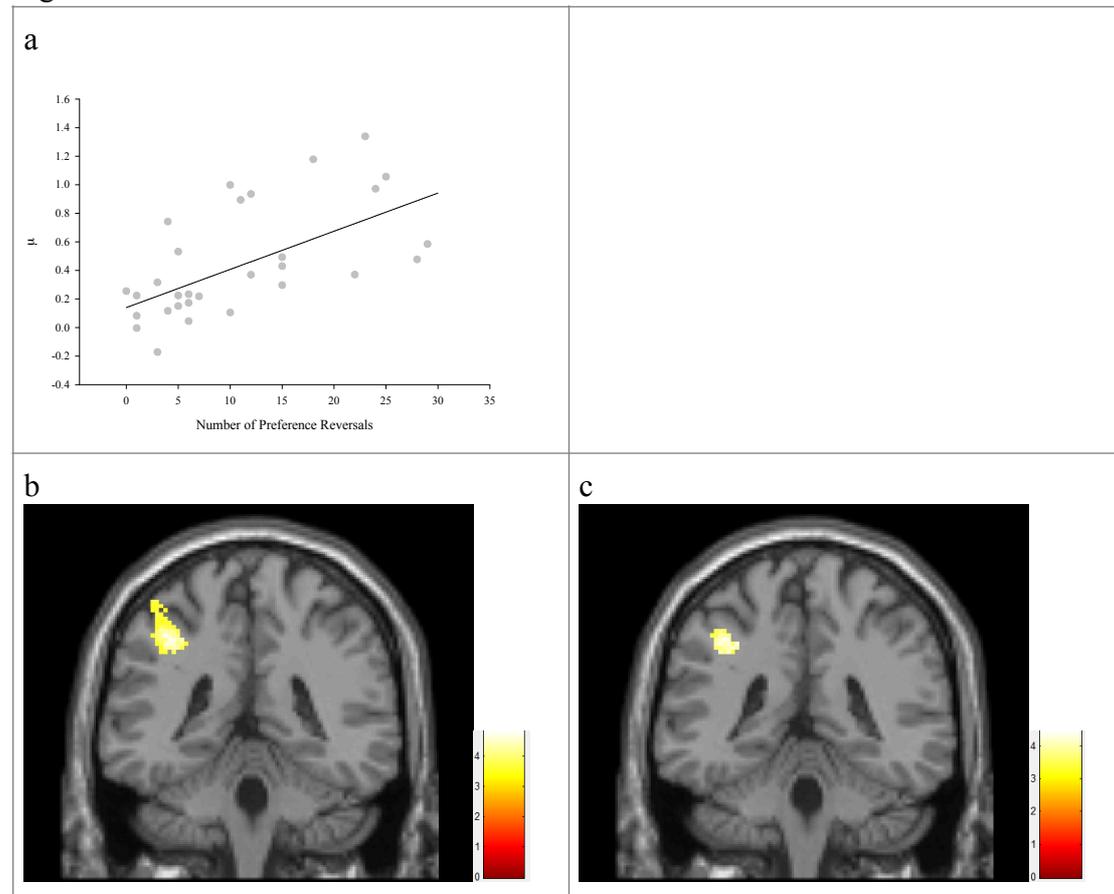


Figure 3a Strong positive correlation between the number of preference reversals and the estimated μ . The Pearson correlation coefficient is 0.610. The robust regression slope is 0.027 and the p-value (two-tailed) is 0.001. For other correlations, see SI, S4.5, table S6 and figure S3.

Figure 3b The left inferior parietal lobule (IPL) was more active in Weak trials than in Strong trials ($y=-40$) [$P<0.05$ at whole-brain cluster correction with a t threshold of 3.40 and an extent of 461 voxels] (SI, S4.6 and table S7).

Figure 3c The left IPL activity parametrically tracked the trial-by-trial decoy effect $d(t)$ negatively ($y=-40$) [$P<0.05$ at whole-brain cluster correction with a t threshold of 3.40 and an extent of 185 voxels] (SI, S4.6 and table S8).

Figure 4

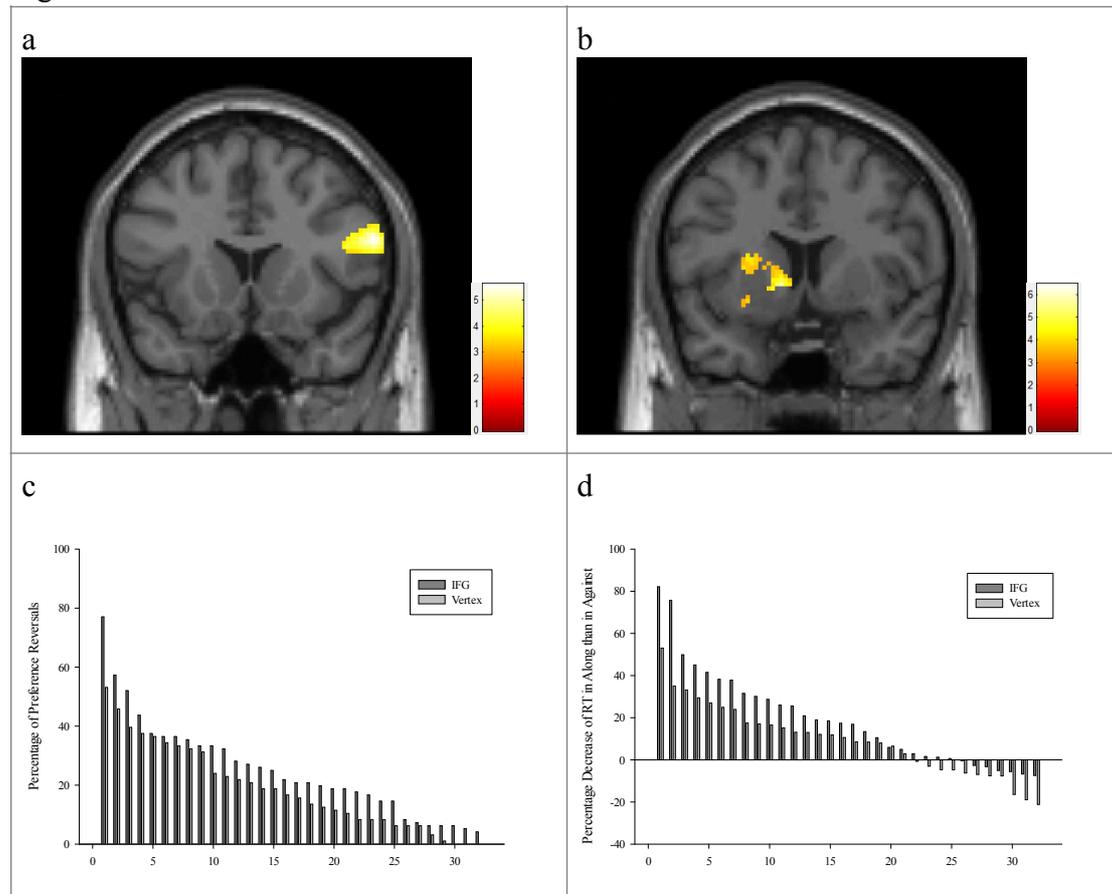


Figure 4a The right inferior frontal gyrus (rIFG) exhibited stronger functional connectivity with the IPL in Conflict trials than in NoConflict trials ($y=16$) [$P<0.05$ at whole-brain cluster correction with a t threshold of 3.40 and an extent of 342 voxels] (SI, S4.8 and table S13). The psychophysiological interaction is based on a 4-mm sphere of IPL centered at $(-40, -40, 46)$.

Figure 4b The left striatum exhibited stronger functional connectivity with the rIFG in Against than in Along of Conflict trials ($y=8$) [$P<0.05$ at whole-brain cluster correction with a t threshold of 3.41 and an extent of 188 voxels] (SI, S4.9 and table S14). The psychophysiological interaction is based on a 4-mm sphere of rIFG centered at $(60, 16, 22)$.

Figure 4c Percentage of preference reversals of the rIFG TMS group was larger than that of the vertex TMS group. The figure is a participant-by-participant breakdown of the percentage of preference reversals of decoy trials (in descending order), in the IFG group and in the vertex group, respectively. The average percentage of preference reversals was 25.30 for the IFG group and 18.72 for the vertex group. Three vertex participants had no preference reversals (SI, S5.7).

Figure 4d Percentage decrease of response time was larger in the rIFG TMS group than in the vertex TMS group. In consistent trials, the response times were significantly shorter when choices are along the decoy than when they are against the decoy, and this was even more so in the IFG group than in the vertex group. The figure is a participant-by-participant breakdown of the percentage decrease of RT in Along trials compared with Against trials (in descending order), of the IFG group and the vertex group, respectively. The average percentage decrease of RT over all IFG participants was 19.22 and that over all vertex participants was 9.05. The percentage decrease of RT is the difference of RT in Against trials and Along trials divided by the average of them in percentage (SI, S5.8).