

## 1 **Wagers for work: Decomposing the costs of cognitive effort**

2

3 Sarah L. Master<sup>1\*</sup>, Clayton E. Curtis<sup>1,2</sup>, and Peter Dayan<sup>3,4</sup>

4 1. Department of Psychology, New York University, New York NY

5 2. Center for Neural Science, New York University, New York NY

6 3. Max Planck Institute for Biological Cybernetics, Tuebingen DE

7 4. University of Tuebingen, Tuebingen DE

8

9 \* Corresponding author

10 email: sm4937@nyu.edu

11

### 12 **Abstract**

13 Some aspects of cognition are more taxing than others. Accordingly, many people will avoid  
14 cognitively demanding tasks in favor of simpler alternatives. Which components of these tasks  
15 are costly, and how much, remains unknown. Here, we use a novel task design in which subjects  
16 request wages for completing cognitive tasks and a computational modeling procedure that  
17 decomposes their wages into the costs driving them. Using working memory as a test case, our  
18 approach revealed that gating new information into memory and protecting against interference  
19 are costly. Critically, other factors, like memory load, appeared less costly. Other key factors which  
20 may drive effort costs, such as error avoidance, had minimal influence on wage requests. Our  
21 approach is sensitive to individual differences, and could be used in psychiatric populations to  
22 understand the true underlying nature of apparent cognitive deficits.

23

### 24 **Author Summary**

25 Anyone who has tried to mentally calculate how much to tip at a restaurant knows that cognitive  
26 effort can feel aversive. Doing math in your head, like most high-level cognitive abilities, depends  
27 critically on working memory (WM). We know that WM is sometimes effortful to use, but we don't  
28 know which aspects of WM use drive these effort costs. To address this question, we had  
29 participants request wages in exchange for performing various tasks that differed in their specific  
30 WM demands. Using computational models of their wage demands, we demonstrated that some  
31 aspects of WM are costly, such as bringing new information into memory and preventing  
32 interference. Other factors, like the amount of information in memory and attempts to avoid  
33 mistakes, were less costly. Our approach identified which specific subcomponents of WM are  
34 aversive. Future research could use these methods to test theories about how motivational  
35 problems might be masquerading as cognitive deficits in psychiatric populations.

36

## 37 **Introduction**

38

39           Some activities (e.g., getting dinner with friends) are more enjoyable than others (e.g.,  
40 calculating how to split the bill). Doing tasks which require greater cognitive effort, colloquially  
41 called “brain power,” can feel uniquely aversive, though to different degrees for different people  
42 (1–3). Indeed, despite tangible benefits, people often avoid cognitively demanding work (4,5).  
43 Such resistance suggests that we weigh the effort of mental activity, perhaps as a cost to be offset  
44 with reward.

45           Previous research has identified the experimental tasks which are more costly to perform  
46 by giving subjects control over which tasks they complete. Tasks which subjects demanded the  
47 most incentives to complete (5–7) or which subjects tended to avoid in favor of other tasks with  
48 equivalent rewards (4,8,9) are considered most effortful. Some costly aspects of these tasks are  
49 external, like time on task (10–12) or the complexity of the cognitive model required by the task  
50 (13–16), but other costs arise from the internal operations necessary to realize external actions.  
51 In general, cognitive resistance increases when tasks place substantial demands on working  
52 memory and cognitive control (5,17–19). However, it remains unknown which particular aspects  
53 of working memory and cognitive control may be most costly. For example, perhaps the sustained  
54 effort required during working memory maintenance is more costly than the transient effort  
55 required to inhibit a prepotent response.

56           Here, we decomposed simple and complex attention (i.e., detection) and working memory  
57 (N-back) tasks into putative elemental processes such as maintaining information in memory of  
58 different loads and resisting interference from task irrelevant lures. We assumed that the  
59 subjective costs of these operations are internally felt and consciously accessible, and that the  
60 total cost of completing a task is learned by experiencing these costs. We assessed these total  
61 costs using a modified auction procedure. Previous work has used such auctions to infer the  
62 subjective values of items on a menu (20,21); our modifications allowed us to infer the total effort  
63 costs associated with completing various cognitive tasks by asking subjects what a “fair wage” for  
64 task completion would be. Given the evidence that the allocation of cognitive resources is subject  
65 to a cost-benefit tradeoff (12,22–26), we hypothesized that subjects’ trial-by-trial fair wage  
66 demands would, at least to a first approximation, reflect the sum of the individual costs associated  
67 with task completion, as the amount of reward necessary to offset them. To assess the extent to  
68 which the costs we measured were related to the self-reported tendency to engage in effortful

69 cognitive tasks, albeit varying across trials, we collected Need For Cognition scores from each of  
70 our subjects (NFC; 2).

71 As our subjects were likely to experience costs other than those deriving from cognitive  
72 effort, we designed our experiment to try to limit the effects of these other factors. First, to  
73 minimize the influence of time on task on fair wage ratings, we gave subjects an easy task to  
74 complete when they wished to skip a harder one. We also ensured that every task round took the  
75 same amount of time. Second, cognitively effortful tasks often also elicit errors. This may be  
76 experienced as a cost, particularly in perfectionist subjects (27,28). While we could not as easily  
77 control for error avoidance costs as for time costs, we designed our task to minimize error  
78 avoidance behavior by not giving trial-by-trial feedback, not informing subjects of their accuracy  
79 round-by-round, and not reducing their compensation unless errors became overly prevalent. We  
80 also collected subject scores on the Short Almost Perfect Scale (SAPS; 29), to assess the degree  
81 to which subjects' fair wages were driven by the tendency to avoid making errors (i.e.  
82 perfectionism). Lastly, we included the costliness of errors alongside the costs of cognition in our  
83 computational analyses.

84 We found three non-zero cognitive effort costs: the cost of adding new information into  
85 working memory (WM), the cost of filtering out irrelevant information, and the cost of maintenance.  
86 More subtly, we found evidence that subjects learned the total costs of each task through task  
87 experience, and that the costs of cognition did not increase or decrease over the duration of the  
88 experiment. We found that the self-reported tendency to avoid effort was related to explicit ratings  
89 of task costliness and difficulty, as well as more implicit costs of cognition. This implies that effort  
90 avoidance may be driven both by the explicit, stable preference to avoid effort and by the implicit  
91 subjective experiences of the costs of cognition.

92  
93

## 94 **Results**

95 100 subjects completed the experiment online through Amazon Mechanical Turk. Subjects  
96 completed 32 task rounds and performed four different tasks in random order: an attentional  
97 vigilance task (1-detect), a vigilance task requiring more WM maintenance (3-detect), and the 1-  
98 and 2-back WM task (18,30). Before each task round began, subjects were shown the task they  
99 were to complete (an associated fractal), and were able to request a fair wage for that round of  
100 that task. A Becker-DeGroot-Marschak (BDM) auction mechanism then determined whether they  
101 completed 15 trials of the task they rated for the wage they requested, or 15 trials of the default,  
102 non-demanding task (the 1-detect) for a lower wage. We analyzed their performance and fair  
103 wage ratings across tasks. We used computational modeling to examine how fair wages were

104 influenced by the putative cognitive operations used to complete the previous task rounds, like  
105 WM maintenance or updating. We also related fair wage ratings to previous task behavior,  
106 including the number and types of errors they made.

107  
108 **Model-Agnostic Results:** There was a main effect of task identity on accuracy (Figure 2A;  $F =$   
109  $44.1$ ;  $p < 0.0001$ ), mean reaction time (RT;  $F = 31.5$ ,  $p < 0.0001$ ), and difficulty rating ( $F = 26$ ;  $p$   
110  $< 0.0001$ ). Post-hoc comparisons confirmed that subjects had lower accuracy, higher RTs, and  
111 higher difficulty ratings on the 2-back task than on all of the other tasks (Table 1; Accuracy: 2-  
112 back versus 1-detect  $p < 0.001$ ; 2-back versus 1-back  $p < 0.0001$ ; 2-back versus 3-detect  $p$   
113  $< 0.0001$ ; Mean RT: 2-back versus 1-detect  $p < 0.001$ ; 2-back versus 1-back  $p < 0.0001$ ; 2-back  
114 versus 3-detect  $p < 0.0001$ ; Difficulty ratings: 2-back versus 1-detect  $p < 0.001$ ; 2-back versus 1-  
115 back  $p < 0.0001$ ; 2-back versus 3-detect  $p < 0.0001$ ). Accuracy was highest on the 1-detect when  
116 compared with all the other tasks (all  $p$ 's  $< 0.0001$ ). Mean RTs on the 1-detect were lower than on  
117 the 1-back ( $p < 0.0001$ ), and 2-back ( $p < 0.0001$ ), but not on the 3-detect ( $p > 0.05$ ). Difficulty ratings  
118 were also lowest on the 1-detect compared to the 1-back ( $p < 0.0001$ ), 3-detect ( $p < 0.0001$ ), and  
119 2-back ( $p < 0.0001$ ). Mean accuracy was lower and mean RT was higher on the 1-back than on  
120 the 3-detect ( $p < 0.001$ ;  $p < 0.0001$ ). The mean difficulty rating was no different between the 1-back  
121 and 3-detect ( $p > 0.05$ ).

122 Subjects rated only the 1-back, 2-back, and 3-detect tasks, as the 1-detect task was the  
123 default task. A 2-way ANOVA on fair wage ratings showed a main effect of task identity (Figure  
124 2B; Table 1;  $F = 29.7$ ,  $p < 0.0001$ ) and a main effect of task iteration (Figure 2D; Supplementary  
125 Figure 1;  $F = 5.2$ ,  $p < 0.0001$ ). Subjects' mean fair wage ratings on the 2-back task were  
126 significantly higher than for the 1-back ( $p < 0.0001$ ). Comparing fair wage ratings for the 1- and  
127 2-back allows us to directly measure the costs of maintaining one more item in working memory,  
128 though the 1- and 2-back tasks also differ in the degree of interference present in WM and the  
129 number of errors made. Mean fair wage ratings on the 2-back were also higher than on the 3-  
130 detect ( $p < 0.0001$ ). Comparing fair wages from the 2-back and 3-detect, which both require the  
131 maintenance of 2 items, allows us to measure the cost of interfering stimuli in WM storage or the  
132 increased errors made on the 2-back task. Fair wages were not significantly different between the  
133 1-back and the 3-detect ( $p > 0.1$ ). Though the 1-back and 3-detect differ in their load on WM,  
134 subjects tended to rate them equivalently. These results suggest that increasing WM interference  
135 may be more subjectively costly than increasing WM load. We investigate further in our model-  
136 based analyses below.

137 While accuracy was significantly lower and fair wages significantly higher for the 2-back  
138 task, there was no relationship between mean 2-back accuracy and mean fair wage on the 2-  
139 back across subjects ( $r = -0.08$ ,  $p > 0.1$ ). There was also no relationship between mean accuracy  
140 and fair wage on the 3-detect task ( $r = -0.16$ ,  $p > 0.1$ ). However, there was a significant relationship  
141 of mean 1-back accuracy and mean fair wage ( $r = -0.36$ ,  $p < 0.01$ ). We further assess the influence  
142 of errors on fair wages below by using computational modeling (Model-based results).

143 Task accuracy was broadly stable across task iterations (Figure 2C; main effect of task  
144 iteration  $F = 1.3$ ,  $p > 0.05$ ). This indicates that performance did not improve with task experience.  
145 Across all subjects there was also no relationship between round number (out of 32) and mean  
146 task accuracy or mean RT (Pearson  $r = -0.02$ ,  $p > 0.1$ ; Pearson  $r = 0.003$ ,  $p > 0.1$ ). This is likely  
147 because subjects trained to 80% accuracy during practice and were already at their maximum  
148 performance levels by the start of the main task.

149 Fair wage ratings seem to decrease with task iteration ( $F = 5.2$ ,  $p < 0.0001$ ; Figure 2D),  
150 but potentially as a byproduct of the experimental design. That is, subjects who asked for lower  
151 wages completed the non-default tasks a higher number of times; therefore the lower mean fair  
152 wage on later task iterations may primarily come from subjects who had lower fair wage ratings  
153 overall (Supplementary Figure 2). Another possibility is that subjects ask for lower fair wages over  
154 time because they find that the tasks become less effortful with practice. If that were the case,  
155 then you might expect their accuracy to improve over the course of the experiment. However, the  
156 ANOVA on task accuracy by task iteration reported above found no main effect of task iteration.  
157 We investigated this further by averaging fair wages over each subject's first and last half of task  
158 completions, and comparing them via t-test to see whether their wage requests changed as their  
159 task experience increased. We did the same analysis for task accuracy. There was a significant  
160 decrease of fair wage ratings from the first to the second half of task completions for the 1-back  
161 task ( $p < 0.01$ ) and 3-detect task ( $p < 0.01$ ). There was no change in fair wage ratings across the  
162 first and second halves of experience with the 2-back task ( $p > 0.05$ ). There was no change in  
163 accuracy in the first and second halves of task completions on the 1-back task ( $p > 0.05$ ), 3-detect  
164 task ( $p > 0.05$ ), or 2-back task ( $p > 0.05$ ). Taken together, these results suggest that any decrease  
165 of fair wage ratings over task iterations stems from the experimental design, and not from learning  
166 or practice effects. We investigate this further with computational modeling below.

167  
168 **Analysis of Self-Report Measures:** We ran regressions on task behavior with linear and  
169 quadratic NFC and SAPS terms, using a model selection procedure which trimmed each  
170 regression down to an intercept term, and the self-report terms which were necessary for model

171 significance ( $p < 0.05$ ). NFC scores were linearly and quadratically related to mean 3-detect  
172 accuracy ( $\beta = -11.59$ ,  $\beta = 1.82$ ). NFC was quadratically related to difficulty ratings for the 1-detect  
173 ( $\beta = -0.06$ ). SAPS scores were linearly and quadratically related to mean 1-back accuracy (linear  
174  $\beta = 12.94$ , quadratic  $\beta = -1.60$ ), mean 3-detect accuracy (linear  $\beta = 7.58$ , quadratic  $\beta = -0.87$ ), and  
175 difficulty ratings for the 2-back task (linear  $\beta = -1.15$ , quadratic  $\beta = -0.13$ ). SAPS scores were also  
176 quadratically related to 2-back accuracy ( $\beta = -1.11$ ). Neither NFC nor SAPS score were linearly  
177 or quadratically related to mean RTs.

178 We ran the same regression analysis on mean fair wage ratings, collapsed across all  
179 tasks. There was a significant quadratic relationship of NFC and mean fair wage ratings ( $\beta = -$   
180  $0.03$ ). We split subjects up into self-report tertiles to further investigate the significant quadratic  
181 relationships between task and self-report variables. The tertile split resulted in 25 low, 37 mid,  
182 and 37 high NFC subjects, and 34 low, 37 mid, and 28 high SAPS subjects. Post-hoc t-tests  
183 confirmed that the significant quadratic effect of NFC is driven by the difference in mean fair wages  
184 between the high and mid NFC subjects. Mid NFC subjects had higher fair wage ratings than high  
185 NFC subjects ( $p < 0.01$ ; Supplementary Figure 3). However, there were no differences between  
186 the low and high NFC groups ( $p > 0.05$ ), or the low and mid NFC groups ( $p > 0.05$ ). We supposed  
187 that high NFC subjects would ask for the lowest fair wages, but we did not find such a pattern in  
188 explicit fair wage ratings. We next investigated how NFC was related to the implicit costs of  
189 cognition captured by our computational model.

190  
191 **Model-based results:** Based on the model-agnostic results, we designed and tested a series of  
192 computational models to isolate the costs of distinct cognitive processes from fair wage ratings.  
193 These models allowed us to test the hypothesis that subjects have some internal awareness of  
194 the costs of certain cognitive operations, and to estimate the magnitude of these costs. We also  
195 measured the costs associated with all types of behavioral responses, including making errors.  
196 In doing so, we assessed whether fair wage ratings also captured costs stemming from physical  
197 effort (making key presses) or error avoidance, which are not cognitive process costs but are still  
198 potential modifiers of fair wages. Error avoidance in particular could explain, to some extent, effort  
199 avoidance in behavior; we fit error costs separately to assess this possibility (28).

200 We fit subjects' behavior with a series of models using the Computational Behavioral  
201 Modeling (CBM) toolbox (31). All models included a noise parameter ( $\sigma$ ), and an initial rating  
202 parameter for each task (*init*) as free parameters. One class of models assumed that the cost  
203 parameters were fixed across trials, but that the subjects learned about the total cost they  
204 experienced with a learning rate ( $\alpha$ ). A separate class of models assumed that subjects' demands

205 reflected the cost just on the previous iteration of the task, but that the cost parameters changed  
206 linearly with trial number at a rate given by parameter ( $\delta$ ). Within these model classes, we tested  
207 several combinations of cost parameters. The maintenance cost ( $C_{\text{maintenance}}$ ) captured the effect  
208 of maintaining more information in WM. The interference cost ( $C_{\text{interference}}$ ) captured the effect of  
209 “lure” trials in the 2-back task. The update cost ( $C_{\text{update}}$ ) captured the effect of updating WM with  
210 new information. The response cost parameter ( $C_{\text{response}}$ ) captured the influence of perceived  
211 matches (button presses) on subsequent BDM ratings. The miss cost ( $C_{\text{miss}}$ ) captured the effect  
212 of missed matches. The false alarm cost ( $C_{\text{fa}}$ ) captured the effect of making responses when there  
213 was no match.

214 The models with the highest frequencies in our subject pool included learning rate  $\alpha$ , rating  
215 noise  $\sigma$ , three initial rating parameters (one per task), update costs, interference costs, and  
216 maintenance costs (Figure 3A). Two subjects were best fit by a model including the cost changing  
217 parameter  $\delta$  and a fixed learning rate  $\alpha = 1$  but most subjects’ (98/100) experienced costs of  
218 cognition were stable across 32 task rounds. Most changes in fair wage ratings were likely driven  
219 by cost learning ( $\alpha$ ), differences in the cognitive operations required in different task rounds, or  
220 reporting noise ( $\sigma$ ).

221 The model with the highest model frequency included only update costs, and was the  
222 winning model overall with a protected exceedance probability of  $>0.99$  and a model frequency  
223 of 78.1%. The second most frequent model included only interference costs and had a model  
224 frequency of 10.3%. The third most frequent model included update, interference, and  
225 maintenance costs, and had a model frequency of 6.3%. The remaining five recovered models  
226 contained the rest of the cost components (including false alarm, miss, and response costs) in  
227 various combinations and accounted for the last 5.3% of model frequency (Figure 3A). They also  
228 contained two models with  $\delta$  cost-changing parameters instead of a cost-learning parameters.

229 Although most of our subjects were best fit by the winning model, one quarter of our  
230 subjects were best fit by other models. Subject fair wages were better fit by simulating data for  
231 each subject using their best-fitting model (mean  $r^2 = 0.516$ ; Figure 3C; Supplementary Figure 3),  
232 than by simulating data for all subjects with just the winning model (mean  $r^2 = 0.466$ ). In addition,  
233 10 subjects’ data were best explained by models containing multiple costs of cognition. Thus,  
234 subjects’ fair wages were influenced by more than just update costs.

235 There was scant evidence that button presses or errors were costly, as all models  
236 including response, false alarm, or miss cost parameters had a total model frequency less than  
237 3%. Models including response and miss costs each accounted for model frequencies less than  
238 1%, so these costs are not explored further below.

239 The mean update cost was 0.615 (Figure 3B), making it the highest magnitude cost  
240 parameter. The next highest mean parameter value was the interference cost, at 0.60, followed  
241 by the maintenance cost at 0.2, and the false alarm cost, at -0.65. Despite the near equivalence  
242 of the mean update and interference costs, lures in WM were much less frequent than updates to  
243 WM. Because of this, subjects lost more monetary bonuses due to the avoidance of update costs,  
244 resulting in their forfeiting an average of 0.87 cents extra per round. They were willing to forfeit  
245 0.26 cents and 0.38 cents per round to avoid maintenance and interference costs, respectively.  
246 While subjects did not know the exact mapping between BDM points and the monetary bonus at  
247 the conclusion of the experiment (1 point = 1 cent), this speaks to the true costliness of each  
248 component process, in terms of the overall monetary amounts subjects forfeited.

249 As we hypothesized, the mean difference between fair wage ratings on the 2-back and 3-  
250 detect tasks was predicted by the magnitude of the interference costs ( $r = 0.42$ ,  $p < 0.0001$ ). The  
251 mean difference between ratings on the 2-back and 1-back was predicted by the magnitude of  
252 the maintenance costs ( $r = 0.41$ ,  $p < 0.0001$ ). These correlations confirm that the tasks differ in  
253 their subjective costliness at least partially because of the differences in WM operations required  
254 by them.

255 We tested whether any self-report measures of effort avoidance or perfectionism related  
256 to fit cost parameters. Specifically, we wondered whether the need for cognition (NFC) or  
257 perfectionism (Short Almost Perfect Scale; SAPS) scales were predictive of any cost parameter  
258 values. For simplicity, we analyzed just parameter values from subjects best fit by the winning  
259 (update costs) model ( $N = 79$ ). We ran a regression including both linear and quadratic terms for  
260 the effect of NFC and SAPS scores on fit update cost parameters from the winning model. We  
261 found no significant linear ( $\beta = 0.218$ ,  $p > 0.1$ ) or quadratic ( $\beta = -0.044$ ,  $p > 0.1$ ) relationship  
262 between update cost and NFC. There was also no significant linear ( $\beta = -0.210$ ,  $p > 0.1$ ) or  
263 quadratic ( $\beta = 0.027$ ,  $p > 0.1$ ) relationship between update cost and SAPS score. NFC and SAPS  
264 scores were well-sampled across our sample of 100 subjects (Supplementary Figure 2).

265 We then examined whether there were parameter differences across NFC and SAPS  
266 tertiles. Within the subjects best fit by the winning model, high NFC subjects had significantly  
267 lower update costs than both low ( $p < 0.05$ ) and mid-NFC subjects ( $p < 0.05$ ). There were also  
268 differences in initial fair wage ratings across NFC groups (Figure 4), the general pattern being  
269 that mid NFC subjects asked for the highest initial fair wages. High NFC subjects had significantly  
270 lower initial fair wage ratings than mid NFC subjects for all three tasks (1-back  $p < 0.01$ ; 2-back  $p$   
271  $< 0.01$ ; 3-detect  $p < 0.05$ ). There were no significant differences between low and high NFC  
272 subjects' initial rating parameters. Mid NFC subjects had higher initial ratings for the 2-back task



273 than low NFC subjects ( $p < 0.05$ ). Mid NFC subjects had higher variance ( $\sigma$ ) around their fair  
274 wage ratings than high NFC subjects ( $p < 0.01$ ) and low NFC subjects ( $p < 0.05$ ). There were no  
275 significant differences in learning rates between subjects split into NFC tertiles ( $p$ 's  $> 0.1$ ). Taken  
276 together, these results suggest that both explicit reports about task costliness (initial fair wage  
277 ratings for each task), and more implicit experiences of the costs of cognitive operations (update  
278 costs) change with individual differences in NFC across subjects.

279 There were no significant differences in cost parameter magnitudes between subjects split  
280 into SAPS tertiles.

281

## 282 **Discussion**

283 Deploying working memory or paying attention can feel costly (32,33). In this work, we  
284 quantified the subjective costs of the cognitive operations demanded by commonly studied  
285 working memory and attention tasks, in a way sensitive to both the dynamics of cognitive effort  
286 exertion and individual differences in effort avoidance. Using a novel experimental paradigm  
287 which leverages an inverted Becker-DeGroot-Marschak auction procedure (20), we obtained  
288 subject ratings of the total cost of completing a working memory or attention task, one round at a  
289 time. We then used a computational model to decompose these ratings into the costs of the  
290 individual cognitive operations putatively used during that round, as well as aspects of subject  
291 behavior, like errors. Our computational models quantify the subjective costs of individual  
292 cognitive operations and allow us to test several hypotheses about how cognitive effort costs may  
293 change with time or task experience.

294 We found evidence that updating WM, interference from within WM storage, and WM load  
295 are subjectively costly. Most subjects tracked a single cost. The largest percentage of subjects  
296 tracked just update costs, and the next highest proportion tracked just interference costs.  
297 Although effortful cognition can be rewarding (3,34), we find that the costs, not the intrinsic  
298 rewards, of cognitive effort drove fair wages. Updating WM cost the most. Subjects forfeit on  
299 average 0.87 cents extra per round as a result of avoiding frequent WM updating demands.  
300 Interference costs (lure stimuli inside of WM) were similarly high, but because lures were  
301 somewhat infrequent, rating them highly (and thereby avoiding them) led subjects to lose less  
302 money per round. The third highest cost was that of maintaining more information in WM.

303 Increasing WM load (the N in N-back) has often been assumed to be the primary driver of  
304 increases in subjective difficulty. However, we show that WM load was only minimally costly and  
305 that updating and interference had a greater influence on subjective cognitive effort. Lure stimuli  
306 in WM storage demand an accurate maintenance of both stimulus identity and stimulus order.

307 The interference cost captures the confusability of stimuli in WM storage and the cost of  
308 disambiguating them by their temporal order. WM updating is similarly complex, as information  
309 must be gated in, gated out, and temporally re-ordered. WM updating has been compared to  
310 switching between WM attractor states, which could be an energetically costly process (35,36).  
311 Perhaps, the magnitude of the update cost parameter captures the complexity or energetic costs  
312 associated with this operation.

313 We find that subjects quickly learned the costs of completing each task through internal  
314 cost feedback signals, then exhibited stable fair wage ratings. Our models provided two surprising  
315 new insights into how the costs of cognition may figure into deciding between several paths of  
316 action. First, only 10 subjects were best fit by models which contain multiple cost parameters.  
317 Tracking multiple costs of cognition may be in itself costly, so subjects may have selected just  
318 one cost component to base their fair wage ratings on to minimize overall experimental demands,  
319 consciously or otherwise. Second and seemingly at odds with previous work (37–39), we found  
320 no evidence that fatigue impacted fair wage ratings as cost parameters did not increase or  
321 decrease over rounds. However, cognitive fatigue may only emerge after longer durations of  
322 cognitive work (40).

323 Our task design directly controlled for one possible confound of the costs of cognitive  
324 effort, time on task (10–12). Another key confound in cognitive effort avoidance work is error  
325 avoidance (27,28), which is harder to directly control for, as tasks which are cognitively effortful  
326 often also elicit more errors. Instead, we measured several potential markers of error avoidance  
327 and found that it had minimal influence over subjects' fair wage ratings. First, there was no  
328 relationship between round-by-round accuracy and fair wage ratings in two out of three tasks.  
329 Second, highly perfectionistic subjects, as measured by the Short Almost Perfect Scale (29), did  
330 not have higher fair wages overall, though they would be expected to have been particularly error  
331 avoidant. Third, while 2/100 subjects' fair wage ratings were responsive to false alarm errors, the  
332 fit cost of making false alarms was of the smallest magnitude, and in fact, numerically negative  
333 (Figure 3B). No subject was affected by the cost of making omission errors (misses). These  
334 results suggest that while error avoidance is a small factor in the overall costs of cognitive effort,  
335 it is not the most important component driving them.

336 The Need for Cognition (NFC) scale measures the self-reported tendency to engage in  
337 challenging cognitive work (2). Our task and modeling approach are sensitive to self-report NFC,  
338 as the cost of updating WM is lowest in subjects with high NFC. This establishes that our paradigm  
339 is sensitive to individual differences, and validates that what we measure with it is indeed related  
340 to the trait tendency to avoid cognitive effort. Interestingly, self-report NFC scores exhibited an

341 inverted U-shaped relationship with initial task ratings, where mid NFC subjects requested the  
342 highest fair wages on the 2-back task. This suggests that NFC interacts differently with more  
343 explicit task ratings versus the more implicit costs of cognition, and warrants further investigation.  
344 Though one would suspect that high NFC subjects would provide the lowest explicit fair wage  
345 ratings, their fair wage ratings were not significantly different from low NFC subjects' ratings.  
346 Instead, they differed in how their wages responded to the dynamic costs of cognition (WM update  
347 costs). This suggests a dissociation of explicit self-report measures and task behavior, but an  
348 association between explicit self-reports and the implicit costs of cognition measured through  
349 computational modeling.

350 One limitation of our task design was the high degree of correlation between cost  
351 components, which may have impacted cost parameter recovery during model fitting. While  
352 maintenance demands were constant across the 2-back and 3-detect tasks, the 2-back was the  
353 only task which required subjects to filter out interference from lures stored in WM. In addition, as  
354 the 2-back was the most difficult task, it was associated with the most errors. Thus the total cost  
355 components increased from the 1-back to the 2-back, and to some extent from the 3-detect to the  
356 2-back. This resulted in high correlations between cost components within subjects. Despite this  
357 consequence of the experimental design, there remained a high degree of fidelity in parameter  
358 recovery (Supplementary Figure 4), and a low degree of tradeoff between fit parameter values  
359 (Supplementary Results). It remains an open question as to what extent these cognitive  
360 operations (i.e. WM updating, resistance to interference, and maintenance) depend on  
361 overlapping or independent mechanisms, and indeed whether the costs of these operations are  
362 related.

363 This work directly quantifies the costs associated with the cognitive operations required in  
364 working memory and attention tasks, not just how subjects avoid or approach each task. The N-  
365 back, a classic WM task, is useful in the study of working memory because it requires the use of  
366 many diverse WM operations (30). Here, we reveal that the N-back's strength may also be its  
367 weakness, in that the number of WM operations required to complete it is also what makes it so  
368 aversive (32).

369 There are many avenues for future work using this experimental and modeling approach.  
370 Here, we adopted one specific process model to decompose each round of each task into the  
371 component cognitive operations necessary to complete it, though there are many possible models  
372 to use. The use of a different process model could have resulted in a different cost component  
373 structure. Slight modifications to the tasks could also have given rise to different cost magnitudes.

374 For example, if we had provided explicit trial-by-trial feedback, we may have observed higher  
375 error costs.

376 These results have potential implications for treating cognitive dysfunction in psychiatric  
377 disorders. For one, the N-back task may not be suitable for use as a benchmark for WM ability in  
378 psychiatric populations, as many have comorbid cognitive and motivational deficits. Dopaminergic  
379 cortico-striatal loops, which are highly sensitive to reward, are thought to be a driver of WM  
380 performance (41–43). Our novel paradigm may be clinically useful, as cognitive dysfunction could  
381 be partially treated by comparing the costs of cognition across groups, then offsetting those costs  
382 with rewards (44–46).

383 In summary, along with a novel experimental approach in which subjects request wages  
384 for completing one round of one task, we implemented a modeling procedure that decomposes  
385 their wages into the costs driving them. We found that updating WM, interference among items in  
386 WM, and WM load are costly, independent of any error, time, or fatigue costs. This suggests that  
387 certain cognitive operations are inherently costly to perform, in alignment with the idea that human  
388 cognition is subject to cost-benefit analyses which can result in the use of less costly, less effective  
389 cognitive strategies (47). Surprisingly, the highest subjective cost of N-back performance was not  
390 WM load, but WM updating. We find a direct relationship between self-report individual differences  
391 in cognitive effort avoidance and the implicit costs associated with specific WM operations. That  
392 our task captures these individual differences, where others have not (48), suggests it could be  
393 implemented to capture other individual differences, perhaps in psychiatric or developmental  
394 populations.

395

396

397

398

399

400

## 401 **Methods**

402 100 subjects (35 female, 14 unspecified sex, mean(std) age: 39(12), 11 unspecified age)  
403 completed our online task in full. 281 unique workers opened our experiment on Amazon  
404 Mechanical Turk (AMT). Of the 270 subjects who consented to participate, 218 of them made it  
405 through the practice blocks, 142 successfully finished the quiz, 125 made it to the 16th block of  
406 the experiment, and then 100 completed the experiment in its entirety. Our final sample, which  
407 we analyze below, consisted of these 100 subjects who finished the experiment. We did not

408 include any data from any of the subjects who did not finish the experiment in our analyses. Given  
409 the strict accuracy and attention cutoffs we imposed, and the overall length of our task  
410 (mean(median) total time on task: 37(36) minutes) versus the typical length of tasks on AMT (one  
411 study reported that the mean time spent on submitted HITs was less than 2 minutes (49)), we  
412 considered a 37% completion rate to be acceptable.

413         Subjects were asked to complete 32 task rounds, alternating between 4 different tasks: a  
414 1-detect task (oddball detection), a 1-back task, a 3-detect task (detect 3 of the same stimulus in  
415 order), and a 2-back task. We chose these four tasks because they rely on many of the same  
416 cognitive processes, whilst also differing in important ways in the operations they require from  
417 those processes.

418

419 **Experimental procedure:** In a novel experimental paradigm, we leveraged the Becker-DeGroot-  
420 Marschak (BDM) auction procedure to measure the evolving subjective value of choice options  
421 (20). The experiment was coded using a pre-built Javascript framework for online Psychology  
422 experiments (JsPsych; 50) and custom Javascript functions. Subjects were introduced to 4 tasks,  
423 each of which was associated with a fractal image (a “task label”; see Figure 1): the 1- and 2-back  
424 working memory tasks, and two types of attentional vigilance task, which we refer to as the 1-  
425 detect (the default task) and 3-detect (4,5,30,51). Subjects completed a total of 32 rounds, using  
426 the BDM procedure before each round to report the wages they considered fair for performing the  
427 particular non-default task that was offered instead of the default 1-detect task.

428         In all tasks, subjects saw a sequence of 15 letters, one after the other. Subjects had to  
429 respond to the letters that matched a rule by pressing the “K” key on the keyboard. Stimuli  
430 remained on the screen for 1.5s; any response had to be made before they disappeared. If a  
431 subject responded late to a match, that trial was marked incorrect. The inter-stimulus interval was  
432 300ms. Time on task was standardized such that the time spent on each task could not influence  
433 subjective effort cost differences across tasks; each task round took approximately 24 seconds.

434         The 1-detect task was the default task, intended to involve minimal effort. Subjects had to  
435 respond only if they saw a “T” on screen. In the 3-detect task, subjects had to respond when any  
436 letter was presented 3 trials in a row. In the 1- and 2-back tasks, subjects had to respond when  
437 the letter on screen matched the one displayed 1 or 2 trials back, respectively. Letter sequences  
438 were standardized such that subjects were required to respond to 3 to 5 matches per round,  
439 regardless of task identity. We chose to run these four tasks because they involved similar  
440 cognitive processes, but differed in their rule structure and thus the number and complexity of the

441 operations they required. In particular, we sought to measure the costs of increased WM load and  
442 the information manipulation required by the N-back tasks.

443 Comparing the subjects' fair wage demands for the 1- and 2-back tasks allowed us to  
444 measure the cost of maintaining one more item in working memory ("maintenance"). Comparing  
445 the demands for the 2-back and 3-detect tasks, which both require the maintenance of 2 items,  
446 allowed us to measure the cost of protecting against interference in the contents of WM  
447 ("interference"). In the 3-detect task, subjects had to remember the 2 previous stimuli and  
448 compare them to the current stimulus. Detecting a match was simple as long as one recalled  
449 whether the previous 2 stimuli matched the current one. In the 2-back task it remained essential  
450 to recall the previous 2 stimuli. However, the stimulus from 1 trial ago was never relevant for the  
451 trial at hand; all that mattered was the identity of the stimulus 2 trials ago. Because both must be  
452 stored, however, it is possible that the stimulus from 1 trial ago was distracting, and if it matched  
453 the current stimulus, it may have served as a lure to respond. Identifying and filtering out this  
454 distraction may require significant attention and effort. Thus a "lure trial" was any trial where the  
455 irrelevant stimulus from 1 trial ago matched the stimulus on the current trial, in the 2-back task.  
456 The interference cost in our model captured the cost of these lure trials. This idea is also described  
457 in (52).

458 Stimuli were presented in pseudo-random order such that the use of other WM operations,  
459 like WM updating, also differed slightly across rounds. Additionally, forcing 3-5 matches per round  
460 allowed us to measure the costs of responding to perceived matches, not responding when  
461 matches occur (misses), or responding erroneously (false alarms).

462 To obtain fair wage ratings for each round in each task, we employed an inversion of the  
463 typical Becker-DeGroot-Marschak (BDM) auction procedure, in which subjects bid for items with  
464 points. In our procedure, subjects are asked to do some cognitive work in exchange for a fair  
465 wage. Before each round, subjects were shown a fractal image associated with one of the tasks,  
466 and were asked to use a slider to specify their "fair wage" for completing one round of that task.  
467 Possible fair wages ranged from 1 to 5 points. They were then shown a random computer offer,  
468 also from 1 to 5 points. If the computer offer was above their requested wage, they were given  
469 the computer offer for completing one round of the task associated with the fractal. If the computer  
470 offer was below their requested wage, they completed the default task for 1 point. All task rounds  
471 consisted of 15 trials.

472 We used the BDM procedure in this work because, via mechanism design, it motivated  
473 subjects to report the true subjective value of the effort they expected to expend on each instance  
474 of a task. If subjects were effort avoidant and wanted to earn higher wages or not complete

475 effortful tasks at all, they would ask for high wages. If subjects were effort seeking, or at least not  
476 effort avoidant, then their fair wages would be low as they should be satisfied with any number of  
477 points above the minimum. If one task was substantially more effortful, then our subjects should  
478 ask for higher wages on that task so that they would not have to complete that task without proper  
479 compensation. In an attempt to prevent subjects from being overly avoidant of making errors, we  
480 did not impose an accuracy cutoff for the receipt of points on individual rounds. Further, after the  
481 initial practice phase, subjects were not informed of their accuracy each round. However, subjects  
482 were aware that if they were inattentive to the task, or their overall accuracy fell below some cutoff,  
483 that the task would conclude early and they would receive less compensation (see exclusion  
484 criteria below). At the end of the task, subjects' points were tallied and converted into a monetary  
485 bonus.

486         At the end of the main experiment, subjects completed a basic demographic inventory,  
487 the Need For Cognition Scale (NFC; 2), and the Short Almost Perfect Scale (SAPS; 29). They  
488 then rated the difficulty of each of the tasks (signaled by its associated fractal) using the same  
489 slider that they used to provide their fair wage ratings. Subjects were also able to provide  
490 comments on their experiences completing the experiment. Subjects were given one hour and 15  
491 minutes to complete the entire experiment.

492

493 **Recruitment and exclusion criteria:** The subject pool was limited to Amazon Mechanical Turk  
494 workers based in the United States, to ensure English reading comprehension. We limited our  
495 recruitment to workers ages 18 and up with at least 100 completed Human Intelligence Tasks,  
496 and with at least 85% acceptance rates. We also ensured that subjects had not completed the  
497 task before using their Worker ID. To ensure that subjects understood the task and were able to  
498 maintain a high level of accuracy, we excluded subjects who did not demonstrate task proficiency  
499 or an understanding of the fair wage procedure after the practice phase. We implemented two  
500 tests that subjects had to pass to make it into the main experiment. First, subjects had to reach  
501 80% task accuracy on 15 trials of our most difficult task, the 2-back. They had up to 10 rounds to  
502 do so. 52 subjects failed to reach this criterion. Following that, subjects had to correctly answer 4  
503 out of 6 questions about the BDM procedure. 76 subjects did not pass this quiz. If subjects passed  
504 both those checks, then they proceeded to the main experiment. After these exclusions, 142  
505 subjects started the main experiment.

506         During the main experiment, subjects' performance was assessed 3 times (every 8  
507 rounds). If in 8 rounds, subjects missed the response deadline for 4 fair wage ratings or their  
508 overall accuracy went below 60%, the task ended early and their data were not used in the final

509 analyses. This eliminated another 42 subjects, resulting in a sample size of 100 subjects total.  
510 Subjects were given a 30 second rest between task rounds and no other breaks.

511  
512 **Model-agnostic analyses:** All model-agnostic and model-based analyses were run in MATLAB  
513 (53). Subject accuracy was calculated online as a weighted function of correct responses (hits)  
514 and correct withholding of responses (correct rejections), where hits were given three times more  
515 weight than correct rejections. We chose to emphasize hits over correct rejections in order to  
516 encourage participant engagement in the tasks, though subjects were not aware of the exact  
517 scoring procedure. In this way, subject accuracy was tracked while they completed the  
518 experiment, so that subjects who were not engaging with the task could be removed from the  
519 experiment early. Once subjects completed the experiment, we examined their behavior on each  
520 task by running ANOVAs on accuracy, response time, fair wages, and difficulty ratings, looking  
521 for an effect of task identity. We examined significant main effects of task identity with post-hoc t-  
522 tests. We assessed the linear relationships of mean accuracy on each task and mean fair wage  
523 for that task across subjects. Additionally, we ran linear analyses of accuracy versus task iteration  
524 and overall experimental round, to examine potential learning or fatigue effects on accuracy. We  
525 ran these same analyses on fair wage demands to determine whether subjects' fair wages may  
526 have changed with time or task practice. Additionally, we ran a comparison of reaction times on  
527 fair wage ratings at the start and end of the experiment.

528 We scored subjects Short Almost Perfect Scale (SAPS) and Need for Cognition Scale  
529 (NFC) responses by summing the numerical values of all their answers, reversing some values  
530 as indicated by published scoring guidelines, then dividing by the number of questions answered.  
531 We used this normalization to ensure that any questions that were not responded to would not  
532 artificially lower questionnaire scores. We excluded questionnaire data from subjects who  
533 incorrectly answered one or both of our screener questions (i.e. "Please select 'Strongly Agree'  
534 for this question"). This type of attention check has been shown to be a reliable way of removing  
535 subjects who are randomly responding to questionnaires, especially when administered more  
536 than once during an experiment (54).

537 The mean(std) normalized NFC score was 3.4(0.9) and the mean(std) normalized SAPS  
538 score was 4.4(1.3). 1 subject chose not to finish those questionnaires and as such has no NFC  
539 or SAPS score.

540 We correlated these questionnaire scores with each other and with participant age. NFC  
541 and SAPS scores were positively correlated ( $r = 0.24$ ;  $p < 0.05$ ). There was no relationship  
542 between participant age and NFC (NFC/age  $r = -0.17$ ;  $p > 0.1$ ) or SAPS score (SAPS/age  $r = -$



543 0.16;  $p > 0.1$ ). We also regressed NFC and SAPS scores, and their squares, against mean fair  
544 wage ratings, average accuracy, and average response time. We used a model selection  
545 procedure which trimmed each regression down to an intercept term, and the self-report terms  
546 which were necessary for model significance ( $p < 0.05$ ). If two reduced models were significant,  
547 and they included different terms, we selected the model with the lower mean squared error (MSE)  
548 in predicting each task variable. We did this to assess both linear and quadratic relationships  
549 between individual difference scores and task performance measures.

550 To build upon the quadratic relationships observed, we also split subjects into tertiles  
551 based on their questionnaire scores. Because both scales administered were short-form, many  
552 subjects have the same score. Thus after splitting subjects into low, mid-, or high scoring groups  
553 based on these scores, the resulting tertiles did not have the same number of subjects in them.  
554 Nevertheless, we ran a series of ANOVAs and post-hoc t-tests to examine whether these groups  
555 differed in their task accuracy, or fair wages.

556

557 **Computational methods and model-based analyses:** We use a computational model to  
558 quantify the putative cognitive processes used in task completion and their influence on fair wage  
559 ratings. We use a process model to decompose each task into the cognitive operations putatively  
560 involved in its completion. We fit 84 candidate models to subject data. We fit this many models in  
561 an attempt to account for most possible combinations of cost parameters, while also limiting  
562 model fitting to those models with high individual parameter recoverability. This number is also  
563 elevated by our use of two different functions of how fair wages change with time. We modeled  
564 subjects' fair wage ratings as a dynamic process driven by subject learning or by the changing  
565 costs of cognitive effort. The first class of models tests the hypothesis that the total cost associated  
566 with each task is learnt through experience with the task and the number of costly components  
567 required to complete it ( $\alpha$  class of models). The second class tests the hypothesis that cognitive  
568 effort costs may themselves change over time, as costly processes become either less costly with  
569 practice or more costly as subjects grow fatigued ( $\delta_j$  class of models).

570 The fair wage ratings for each task were initialized in the model by fitting initial rating  
571 parameters for each task and each subject, thus capturing each subjects' initial ratings with very  
572 high fidelity (Supplementary Figure 4). Each subject's initial fair wage ratings for each task were  
573 captured using a free parameter  $init_i$ .

574

575  $rating^0(\text{task} = \text{"1-back"}) = init_{1\text{-back},i}$ ;

576  $rating^0(\text{task} = \text{"2-back"}) = init_{2\text{-back},i}$ ;

577 rating<sup>0</sup>(task = “3-detect”) = *init*<sub>3-detect</sub>

578

579 While the inclusion of three extra free parameters to determine initial fair wages may seem  
580 unnecessary, correctly capturing each subject’s starting point allows us to fit most accurately how  
581 subjects’ fair wage ratings evolve over course of the experiment, as well as how they respond to  
582 individual cost components. However, because there are already extra parameters in the  $\delta_j$  class  
583 of models (+1  $\delta_j$  for each cost parameter, so that they can change independently), we did not fit  
584 individual *init* parameters to each task in this class of model, to avoid overfitting. After the initial  
585 fair wage ratings, the total cost on task round  $r_k$  of task  $k$  was then used to determine the fair wage  
586 rating on the next round of that same task (round  $r_k + 1$ ). This round may arise some trials later;  
587 we denote trials by  $t$ .

588 We approached cost decomposition with a simple program which was capable of  
589 accurately completing each task with the same “cognitive” functions, but switched between rule  
590 structures depending on the task at hand. We tallied each operation that the model had to use to  
591 complete each task round with 100% accuracy, including how many items had to be maintained  
592 in WM, how many times WM storage had to be updated with new information, or how many times  
593 there were interfering “lure” stimuli in WM storage. In addition, we tallied the mistakes (misses  
594 and false alarms) and button press responses made by each subject in each round. All those  
595 components were then scaled by their associated costs (which might change over trials  $t$ , and  
596 were fit through the modeling), and were summed to produce the total cost incurred on that round  
597 of that task. For round  $r = r_k$  of task  $k$ :

$$(1) \quad \text{cost}^r(k) = \sum_{j \in C_{\text{params}}} \text{components}_j^r c_j^t$$

598

600 The most complex model included six cost parameters (set  $C_{\text{params}}$ ): the cost of  
601 responding to a perceived match ( $C_{\text{response}}$ ), the cost of maintaining information in WM ( $C_{\text{maintenance}}$ ),  
602 the cost of protecting against interference in the contents of WM ( $C_{\text{interference}}$ ), the cost of updating  
603 WM with new information ( $C_{\text{update}}$ ), the cost of false alarm responding when there was no match  
604 ( $C_{\text{fa}}$ ), and the cost of missing a match ( $C_{\text{miss}}$ ). Other than the interference cost, which was only  
605 present in the 2-back task, each cost was fit from ratings of all 3 rated tasks. However, we tested  
606 models containing all combinations of 6 different possible costs. All cost parameters were  
607 unbounded such that they could be positive, or negative. If any components were perceived to be  
608 rewarding, instead of costly, then our model would capture that with a negative cost magnitude.

609 It is important to note that a different choice of process model could result in a different cost  
610 structure.

611 We tested two possible fair wage rating updating mechanisms: a class of model which  
612 assumed that the costs themselves changed over trials and that subjects directly reported their  
613 experienced costs as they changed, and a class of model which learned the total cost of  
614 completing each task following task experience. These updating mechanisms are subtly different,  
615 and involve two different free parameters:  $\delta$ , the scalar with which costs are changed trial-by-trial,  
616 and  $\alpha$ , the cost learning rate. It should be noted, however, that it is theoretically possible that both  
617 mechanisms contribute to cost ratings simultaneously. For simplicity and for robustness of model  
618 recovery, we chose to fit these updating mechanisms as separate model classes.

619 In the cost-changing class of models,  $\delta_j$  ( $j \in C_{params}$ ) is the cost-specific change  
620 parameter which captures how costs linearly change over time (trial number  $t$ ), i.e. with task  
621 experience or fatigue:

$$622 \quad (2) \quad c_j^t = c_j^0 * (1 + (\frac{\delta_j * t}{T}))$$

623

624  $T$  is the total number of trials, and  $\delta$  can be positive or negative. The flexibility of  $\delta$  allows  
625 the cost of each cognitive operation  $c_j$  to increase or decrease linearly. Note that because the cost  
626 components are shared over tasks, and fatigue is supposed to generally increase with time on  
627 task, in this model class each cost is changed according to overall trial number ( $t$ ), instead of task  
628 round number ( $r_k$  for task  $k$ ). In this class of models, fair wage ratings on round  $r_k + 1$  are a direct  
629 function of the cost parameters and task components involved to complete the previous task  
630 round  $r = r_k$  (which is equivalent to having a cost learning rate  $\alpha = 1$ ):

$$631 \quad (3) \quad \text{rating}^{r+1}(k) = \text{cost}^{r_k}(k)$$

632

633 In the cost-learning version of the model, the costs do not change with trial number, as  
634 they do in the other class of models, so:  $c_j^1 = c_j^2 = \dots = c_j^T$ . This class of models learns  
635 incrementally about the total cost of completing each task.  $\alpha$  is the subject-specific cost learning  
636 rate which captures how much each subject adjusts their ratings for an individual task  $k$  based on  
637 the most recent round  $r = r_k$  of that task:

$$638 \quad (4) \quad \text{rating}^{r+1}(k) = \text{rating}^r(k) + \alpha (\text{cost}^r(k) - \text{rating}^r(k))$$

639

640

641           Lastly, we modeled noise in the fair wage rating process with a Gaussian noise process  
642 centered on 0 with standard deviation  $\sigma$ , also a free parameter, and by applying this noise to each  
643 fair wage rating independently. This makes the generated rating follow:

$$644 \quad (5) \text{ rating}^{r+1}(k) = \bar{\text{rating}}^{r+1}(k) + \mathcal{N}(0, \sigma^2)$$

645

646           Given the modest number of ratings provided by each subject (32 in total, split amongst 3  
647 tasks), and the overall similarity of ratings between subjects, we fit our models using a hierarchical  
648 Bayesian inference (HBI) for computational behavioral modeling (CBM) package (31). Employing  
649 a hierarchical parameter estimation procedure allows for similarity across subjects to be  
650 leveraged to fit individual parameter values accurately, especially when fitting few individual data  
651 points. The package leverages estimations of group parameter means and variances in the  
652 individual parameter estimation process. In addition, this package allows for the possibility that  
653 not every subject is best fit by one model. Model responsibilities are calculated subject-by-subject  
654 such that subjects who are not well-described by a model do not influence the overall parameter  
655 probability distributions from that model. In our case, this allows for individual differences in what  
656 processes are perceived as costly. If the ratings of some subjects are not affected by a certain  
657 cost term, then the group-level estimate of this cost is not driven down by their inclusion in the  
658 pool. The Bayesian model fitting procedure constrains the group parameters to have Gaussian  
659 distributions, and so, as is common, we transformed the parameter associated with the learning  
660 rate  $\alpha$  using a logistic sigmoid (so it lies between 0 and 1) and the parameter associated with the  
661 rating noise  $\sigma$  using an exponential, so that it is positive (with a log normal distribution).

662           To validate the winning models further by assessing their ability to produce the behavioral  
663 effects of interest, we simulated fair wage ratings using each of the winning models. We then  
664 compared these model simulations to real subject behavior via visual inspection, and by  
665 computing mean r-squared values for each model. Because stochasticity is one feature of model  
666 behavior (via the standard deviation parameter  $\sigma$ ), we simulated each subject's data using their  
667 fit parameter values 10 different times to control for the stochasticity of these simulations. Each  
668 time, we correlated the true fair wage ratings of all subjects with the set of simulated fair wages,  
669 and then squared the r-value obtained. We ran this over 1000 iterations, and then took the  
670 average r-squared value to produce a mean r-squared value for each model. This was then used  
671 to validate that the models could reproduce subject behavior.

672           In the CBM toolbox, the group-level mean for each parameter is calculated separately for  
673 each model. This allows group-level cost parameter magnitudes to be compared within-model,  
674 but not across-model. In order to compare the magnitudes of the cost parameters across all our

675 models, we constructed posterior probability distributions over the magnitude of each cost. We  
676 used parameter estimates from every subject and every model, weighing the contribution of each  
677 subject  $s$  and model  $m$  by their fit responsibility  $\rho$ :

$$(6) \quad P(\theta|\mathcal{D}^s) = \sum_m \rho_m^s P(\theta_m^s|\mathcal{D}^s)P(\theta_{\tilde{m}})$$

678  
679  
680 where  $P(\theta_{\tilde{m}})$  is the group-level prior distribution over the costs not in model  $m$ . This prior is a  
681 weighted average over the group-level parameter distributions derived from each model, where  
682 the weights are again derived from the model responsibilities  $\rho_m^s$ . We assumed that these prior  
683 distributions were Gaussian within-model, then averaged them across models to produce non-  
684 Gaussian mixture models of across-model priors. Here, the responsibility  $\rho_m^s$  reflects the degree  
685 to which each subject's fair wage data were explained by that model.

686  
687 Using equation 6, we constructed a 4D distribution over the four cost parameters included in  
688 models with at least 1% fit model frequency. We summed over the 4D joint distribution to produce  
689 the marginal distributions of each cost. Additionally, we subdivided our subjects into tertiles based  
690 on self-report scores (NFC and SAPS), and calculated the degree to which the posterior  
691 parameter distributions overlapped across these score groups  $g$ :

$$(7) \quad P(\theta; g) = \prod_{s=1}^{S_g} \left( \sum_m \rho_m^s P(\theta_m^s|\mathcal{D}^s)P(\theta_{\tilde{m}}) \right)$$

692  
693 where subjects 1 through  $S_g$  belong to the group of interest.

694 We obtained the means and standard deviations of the marginal posterior distributions  
695 over individual cost magnitudes. In this way, we assessed the degree to which the cost  
696 magnitudes were separable within- and across-subjects, and across models which did not share  
697 all the same parameters.

698 We confirmed the validity of our models and model fitting procedure by running a generate  
699 and recover procedure. For each model, we simulated a data set of 30-100 subjects with known  
700 parameter values. We used trial-by-trial cost components taken directly from subject behavior to  
701 ensure that real responses, including errors, and task characteristics were compatible with our  
702 modeling procedure. To determine which models were sufficiently robust in parameter recovery,  
703 we ran a generate and recover of all 126 possible models (combining different costs and using  
704 an  $\alpha$  or  $\delta$  update mechanism). In this way, we selected 84 models to test that showed reliable  
705 parameter recovery and minimal cost parameter tradeoff. We wanted to test a broad array of

706 models since we had limited a priori knowledge of which cost components would drive fair wages,  
707 or what form cost updates would take. At the same time, we wanted to fit real subjects' data only  
708 with models that had recoverable free parameters and minimal tradeoff between costs, despite  
709 possible correlations of cost components, as individual differences were of particular interest.

710         Supplementary Figure 4 shows the results of this generate and recover procedure for one  
711 example model, which includes update, maintenance, and false alarm costs (N = 50 simulated  
712 subjects). All fit and real parameters were highly correlated ( $\sigma$   $r = 0.84$ ,  $p < 0.001$ ;  $\alpha$   $r = 0.94$ ,  $p <$   
713  $0.001$ ; update costs  $r = 0.55$ ,  $p < 0.001$ ; maintenance costs  $r = 0.66$ ,  $p < 0.001$ ; false alarm costs  
714  $r = 0.78$ ,  $p < 0.001$ ;  $\text{init}_{1\text{-back}}$   $r = 0.88$ ,  $p < 0.001$ ;  $\text{init}_{2\text{-back}}$   $r = 0.96$ ,  $p < 0.001$ ;  $\text{init}_{3\text{-detect}}$   $r = 0.92$ ,  $p$   
715  $< 0.001$ ). This indicates that our models supported the reliable recovery of individual parameters,  
716 despite the modest number of trials that were fit per subject.

717

718

## 719 **Acknowledgements**

720 We thank our subjects on Amazon Mechanical Turk for their participation in this experiment, as  
721 well as our earliest pilot subjects, the members of Arbeitsgruppe Peter Dayan (AGPD) at the Max  
722 Planck Institute for Biological Cybernetics. We also thank the members of AGPD for many helpful  
723 comments on the initial design of the experiment, and specifically thank Dr. Franziska Broeker for  
724 her assistance in establishing the online data collection pipeline, and Drs. Andrew Webb and  
725 Aenne Brielmann for thoroughly reviewing the modeling and analysis code for this project. We  
726 thank Dr. Martin Breidt, Dr. Holger Dinkel, Mihai Vintiloiu, and the entire IT Core Facility at the  
727 Max Planck Campus in Tuebingen for their technological and research administrative support.  
728 We thank Finn van Krieken for his assistance with front-end online experiment development. We  
729 also thank Dagmar Maier for administrative support.

730

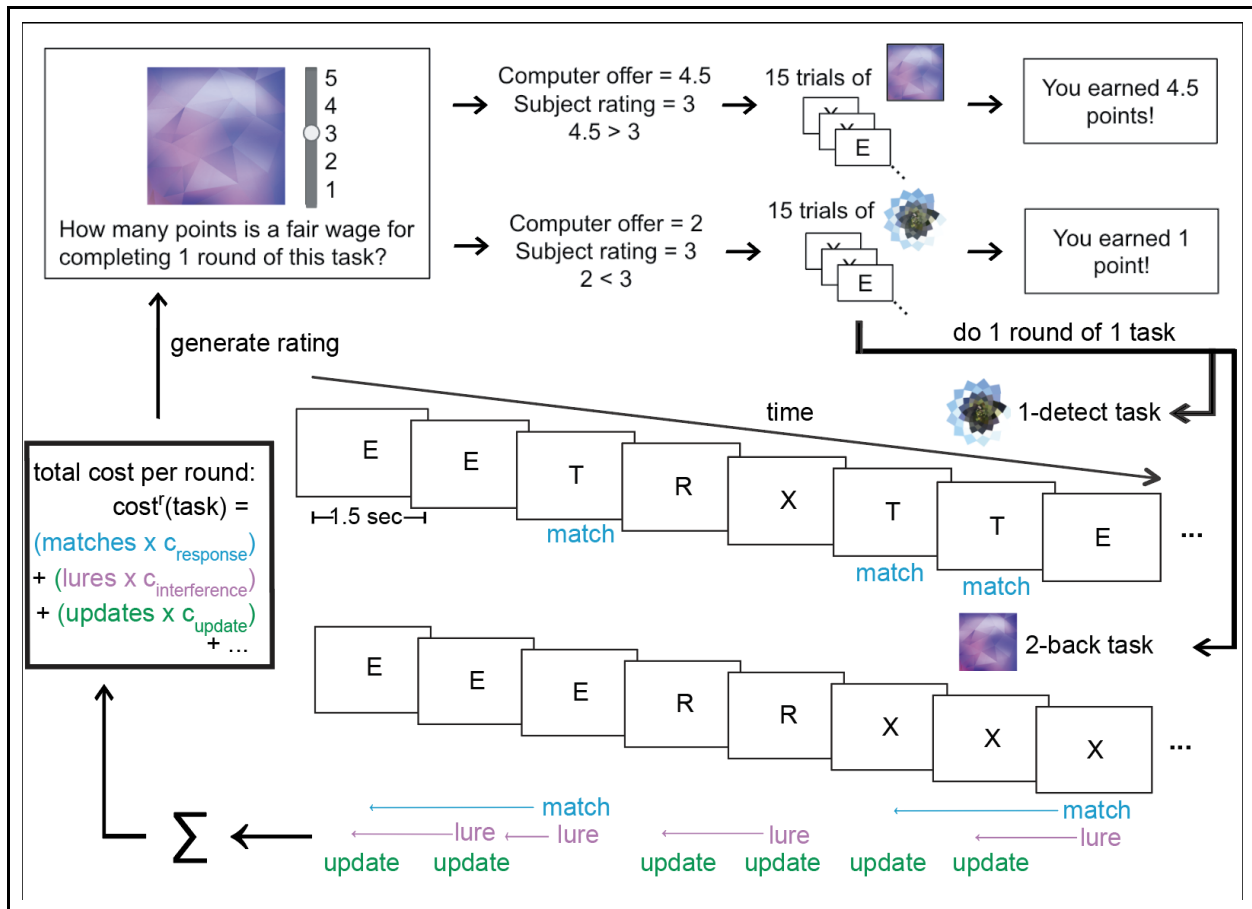
731

732

733

734

735 FIGURES  
736



**Figure 1. The behavioral paradigm & computational modeling approach.** Before each round of the experiment, subjects were shown an image which was associated with one of three possible tasks. They then indicated the wages (in points) that they would like to receive for completing 1 round of that task. If their fair wage rating was below a random computer offer, then they would complete that task and receive the computer's offer. If their fair wage was above a random computer offer, then they would complete a different, easier task instead. We employed this inversion of the Becker-DeGroot-Marschak auction procedure to incentivize subjects to be truthful in their fair wage ratings.

737  
738  
739

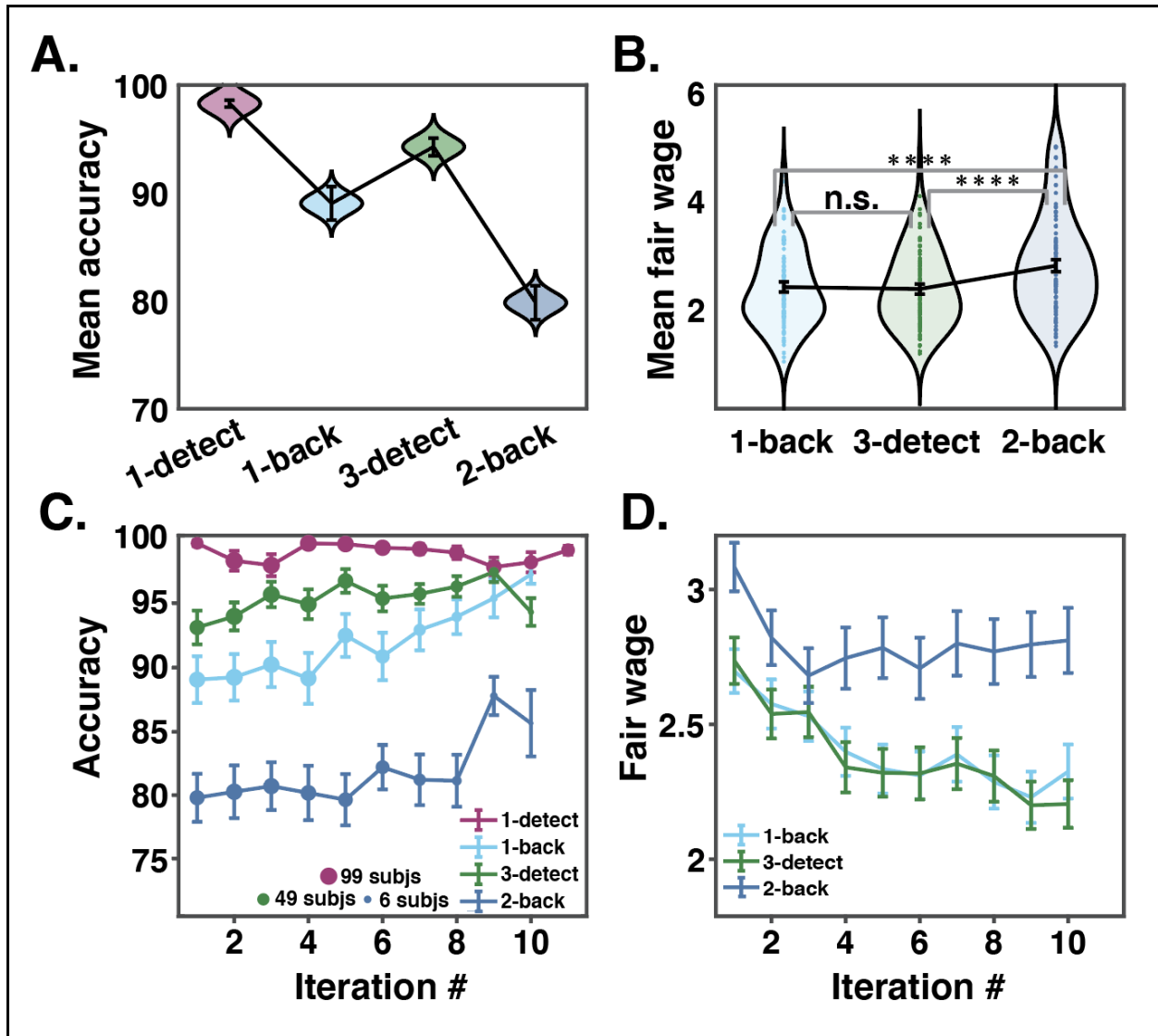
740 **Table 1:** Mean accuracy, reaction time (RT) in milliseconds, and difficulty ratings across all  
741 subjects for the default task, the 1-detect, and for the three rated tasks, the 1-back, 3-detect, and  
742 2-back tasks. The maximum RT was 1500 milliseconds. The minimum fair wage and difficulty  
743 rating was a 1, and the maximum was a 5.  
744

Group means	<b>1-detect</b>	<b>1-back</b>	<b>3-detect</b>	<b>2-back</b>
<b>Percent accuracy</b>	98.29	89.03	94.28	79.83
<b>RT (msec)</b>	548.33	610.97	532.59	717.37
<b>Difficulty rating</b>	1.90	2.44	2.40	3.32
<b>Fair wage</b>	NA	2.41	2.37	2.80

745



746



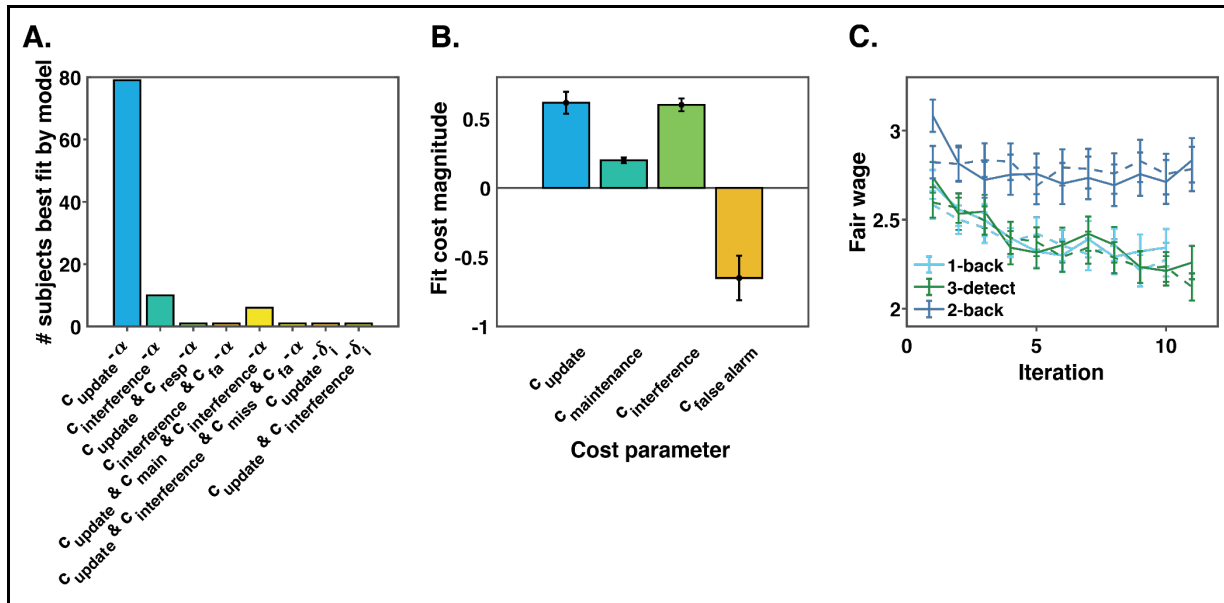
**Figure 2. Model-agnostic behavioral results.** **A.** Distributions of mean accuracies across all subjects for the default task (1-detect), and the three rated tasks (1-back, 3-detect, and 2-back). The black bars depict the means and standard errors of the mean (SEMs) of each distribution. The distribution of all subjects' mean accuracies was plotted using a Gaussian kernel via violin.m. All mean accuracies for each task were significantly different from each other (all  $p$ 's < 0.001). **B.** Distributions of mean fair wages across all subjects for the three rated tasks. The lowest possible rating was 1, and the highest possible rating was 5. The black bars depict the means and SEMs of each distribution. The distribution of ratings was plotted using violin.m. \*\*\*\* indicates significance at the  $p < 0.0001$  level. **C.** Mean accuracy across all subjects on each iteration of each task. Due to the stochasticity inherent to the BDM auction procedure, individual subjects completed the 1-back, 3-detect, and 2-back tasks a variable number of times, but a maximum of 11 times each. The relative number of subjects who completed each iteration is depicted by the size of the dot plotted at the mean. Error bars are plotted with standard error of the mean. A two-way ANOVA of task and task iteration revealed a main effect of task identity ( $F = 15$ ,  $p < 0.0001$ ) but no effect of task iteration ( $F = 1.3$ ,  $p > 0.05$ ). Thus mean accuracy was

different across tasks but did not change with task experience. **D.** Mean fair wage rating by rating number, where the maximum is 11 ratings of one task. A 2-way ANOVA on BDM ratings showed a main effect of task identity (Table 1;  $F = 33$ ;  $p < 0.0001$ ) and a main effect of task iteration (Figure 1;  $F = 21$ ;  $p < 0.0001$ ).

747

748

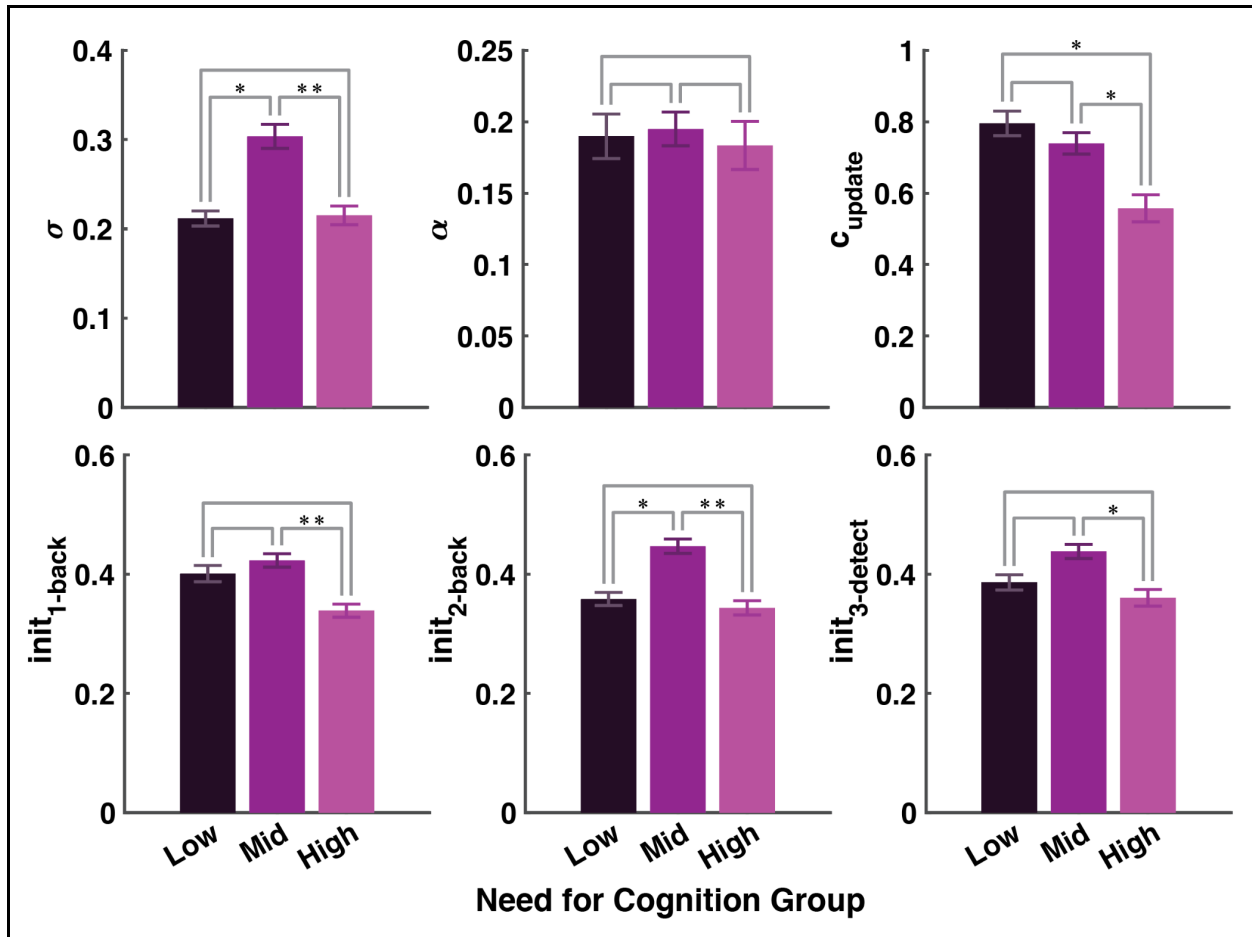
749



**Figure 3. Computational modeling results.** **A.** The number of subjects best fit by each model with a non-zero model frequency. Of the 84 computational models fit to subjects' fair wages, the winning models were alpha cost-learning models containing update costs ( $C_{\text{update}}$ ), interference costs ( $C_{\text{interference}}$ ), and maintenance costs ( $C_{\text{maintenance}}$ ), and false alarm costs ( $C_{\text{fa}}$ ), in various combinations. The model with the highest model frequency was the model including update costs alone. **B.** The mean of the posterior distribution of each cost parameter from the models that best fit at least 1 subject's fair wages. These posterior distributions were calculated by combining inferred parameter distributions across subjects and across models. Inference was performed over joint 4D distributions to capture co-variance between update, interference, maintenance, and false alarm costs. For plotting purposes we summed over the three irrelevant dimensions for each parameter to construct its marginal distribution, and then calculated the means and variances of the marginals. Error bars reflect the hierarchical standard error of the mean; they were calculated not with the square root of the total number of subjects in the denominator, but with the square root of the number of subjects' data explained by models containing that parameter. Note that the error bars describe the spread of the marginal parameter distributions, not variance in the fitting process, and thus are not suitable for estimating the statistical significance of the effects plotted. **C.** Real (solid lines) versus simulated (dashed lines) fair wage ratings on each rating iteration for each task. Data simulated using each subjects' best model faithfully reproduces real subject data ( $r^2 = 0.51$ ).

750

751



**Figure 4. Winning model parameter values by Need for Cognition (NFC) Group.** Mean parameter magnitudes from the winning 6-parameter update cost model.  $\sigma$  is the standard deviation parameter which dictates how noisy each subject's fair wage ratings are, on average.  $\alpha$  is the subject-specific task cost-learning rate. The update cost is the magnitude of the influence of WM updates on each subject's fair wage ratings. The init parameters dictate each subject's initial fair wage for each task. Subjects were split into NFC tertiles resulting in low (N = 25), mid (N = 37), and high (N = 37) NFC groups. Fit parameter values were then averaged within-group to produce each bar. Error bars are standard error of the mean. \* indicates significant difference as assessed with a t-test at  $p < 0.05$  level. \*\*  $p < 0.01$

752

753 **SUPPLEMENTARY INFORMATION**

754

755 **SUPPLEMENTARY RESULTS**

756

757           Because the 2-back was associated with the most WM updating, most errors, and most  
758 interference in WM, many of the cognitive components of task completion that came out of our  
759 process model were highly correlated. For example, across all subjects & task rounds, updates,  
760 maintenance, lures, and false alarms were all correlated (interference vs. false alarms  $r = 0.48$ ,  $p$   
761  $< 0.001$ ; interference vs. maintenance  $r = 0.53$ ,  $p < 0.001$ ; maintenance vs. false alarms  $r = 0.31$ ,  
762  $p < 0.001$ ; updates vs. maintenance  $r = 0.86$ ; updates vs. false alarms  $r = 0.43$ ; updates vs.  
763 interference  $r = 0.62$ ). In addition, when we ran within-subject correlations of these components  
764 across rounds, they were significantly correlated within 91% (maintenance vs. interference), 44%  
765 (maintenance vs. false alarms), 74% (interference vs. false alarms), 100% (updates vs.  
766 maintenance), 66% (updates vs. false alarms), and 90% (updates vs. interference) of subjects.  
767 One consequence of this may be that the cost parameter values associated with these  
768 components may trade off with one another in model fitting, artificially raising or lowering each  
769 other. In addition, model selection may have been impacted, resulting in a low number of subjects  
770 who were best fit by models including multiple costs. Because most subjects' data are best  
771 captured by a model including only one cost of cognitive effort, one might wonder whether the  
772 cost parameters obtained from our models are capturing one cost only, but incorrectly assigning  
773 them to three different components due to the relatedness of the components. To compare cost  
774 parameter magnitudes across models including only one or two cost parameters each and to  
775 ensure their separability, we constructed posterior distributions over parameter values using the  
776 outputs of the CBM toolbox(31). We also examined whether these parameter values traded off  
777 during model fitting by examining their covariances, which are derived from the inverse Hessian  
778 of the search gradient within the multidimensional parameter space.

779           The most frequent model with multiple cost parameters contained update, maintenance,  
780 and interference costs. The covariances between these parameters, which is influenced both by  
781 their empirical covariances, and their covariance during parameter fitting, were all within an  
782 acceptable range. The covariance between the update and maintenance costs was largest, at -  
783 0.22. Between the update and interference costs, the covariance was 0.0492. Between the  
784 interference and maintenance costs, it was -0.0763.

785           We verified in a pre-model fitting generate and recover procedure that individual cost  
786 parameters were being accurately fit even in models with multiple costs (Supplementary Figure  
787 4). In addition, there was no evidence that these single cost parameters somehow capture just

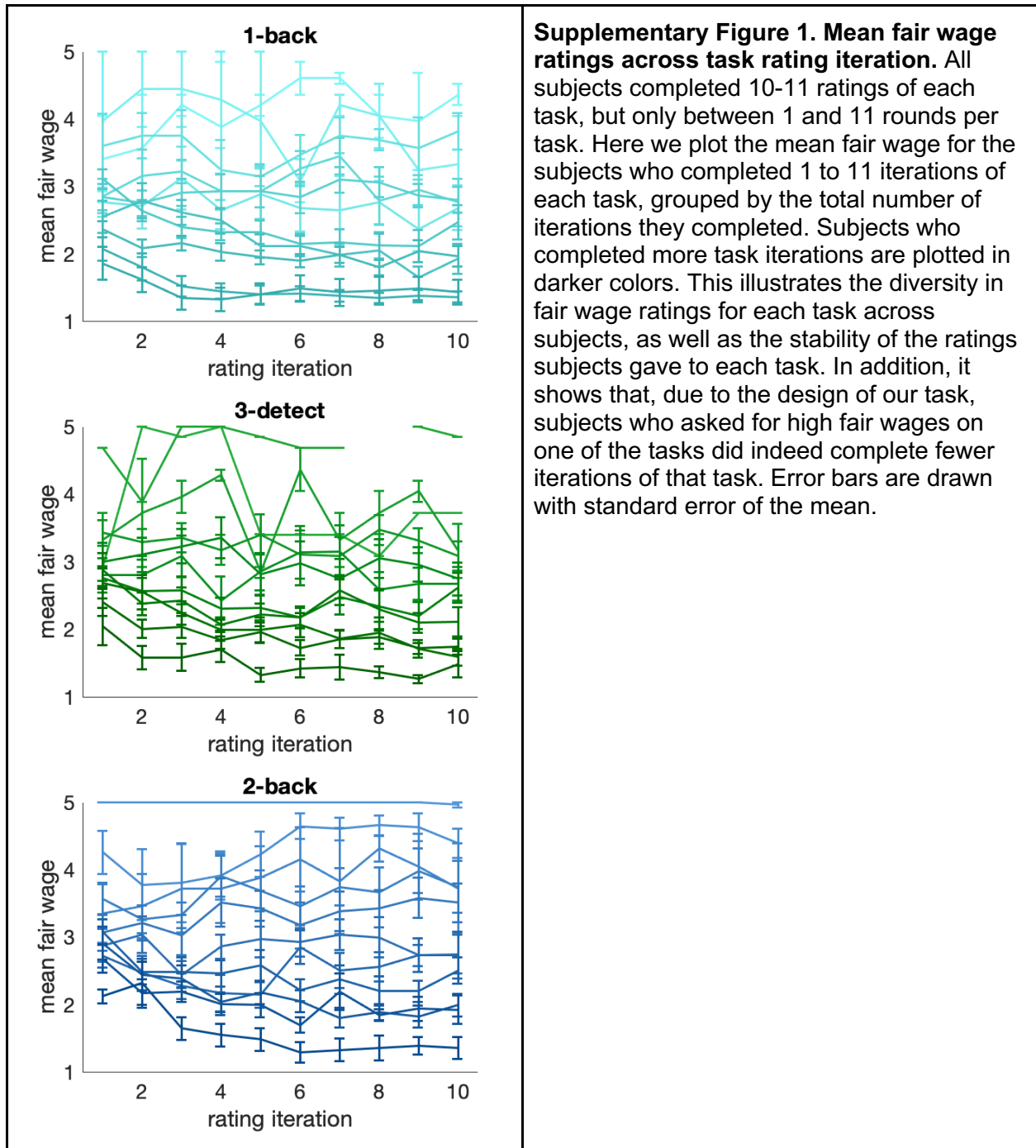
788 one underlying component, rather than 3 separate ones, as the posterior distributions over their  
789 magnitudes are mostly non-overlapping on the group level (Figure 3B). While the update and  
790 interference costs are of similar magnitudes, and therefore overlapping, the negligible covariance  
791 between update and interference costs suggests that they did not trade off in model fitting.  
792

793

794 **SUPPLEMENTARY FIGURES**

795

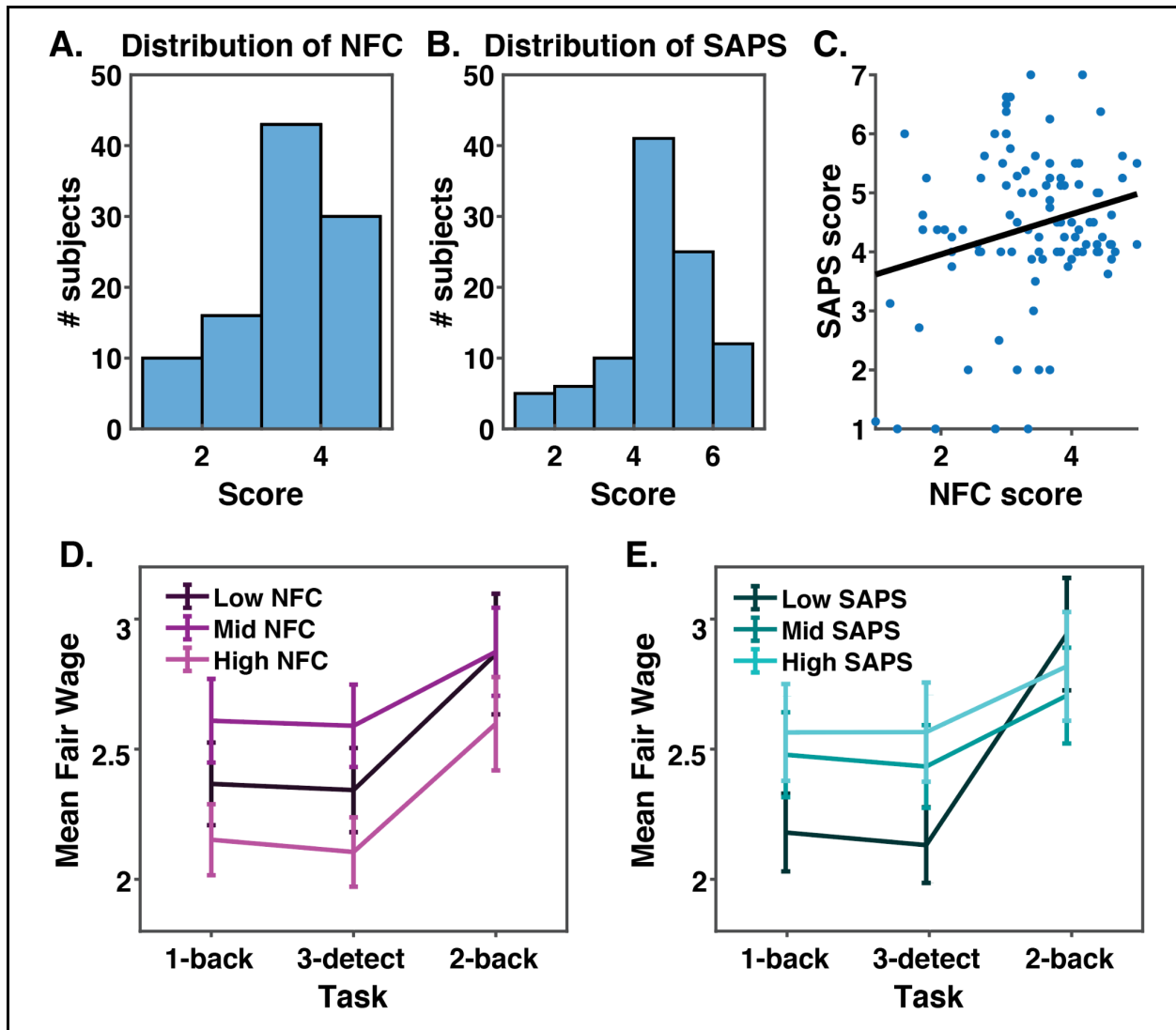
796



797

798

799



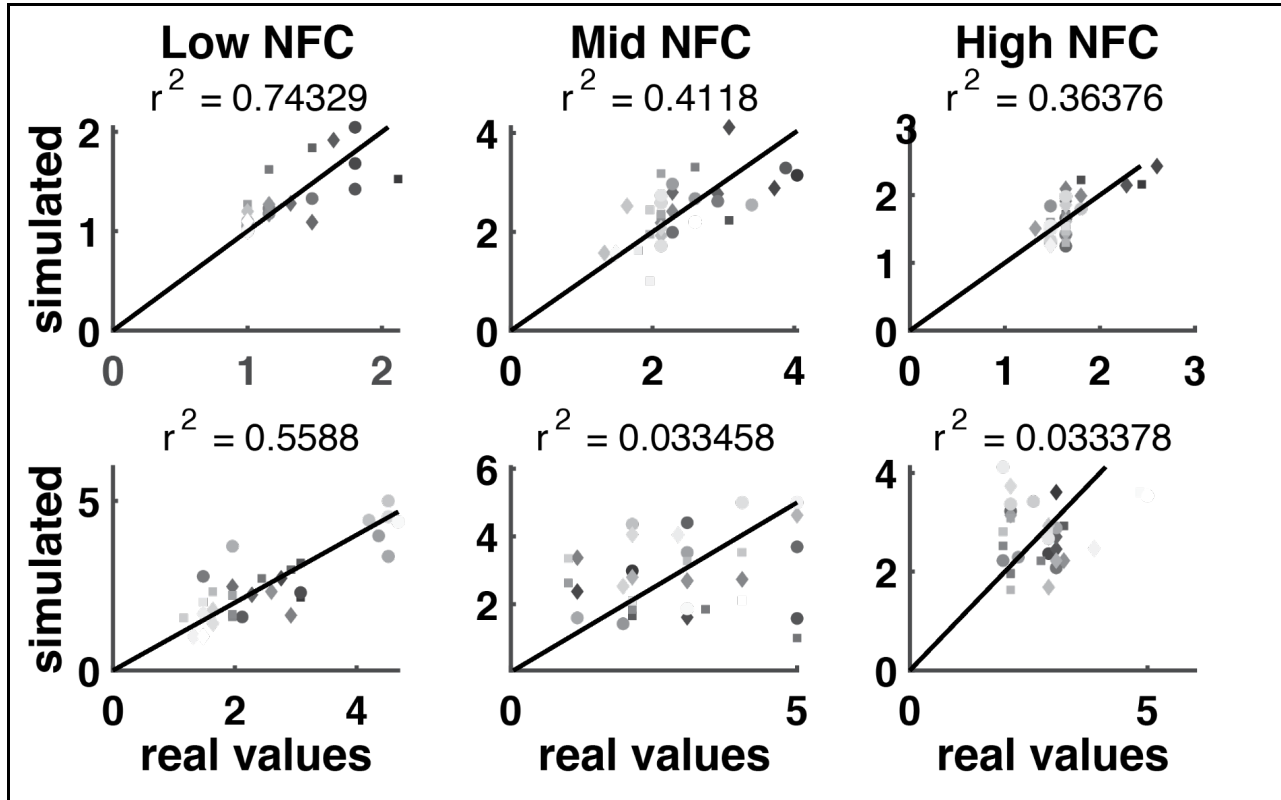
**Supplementary Figure 2. Self-report scores and their relationships to mean fair wages.**

**A.** Distribution of Need for Cognition (NFC) scores within the experimental sample. Scores have been normalized by the number of questions answered such as not to lower the mean of the distribution artificially. The distribution of NFC scores in our sample is right-skewed compared to the typical distribution of NFC scores. However, this is typical of subjects on MTurk (Berinsky, Huber & Lenz, 2012). **B.** Distribution of Short Almost Perfect Scale (SAPS) scores. Scores have been normalized by the number of questions answered such as not to artificially lower the mean of the distribution. The distribution of SAPS scores in our sample is typical of both in-person samples (Rice, Richardson, & Tueller, 2013) and other samples on MTurk, including one sample of 400 subjects (Stricker, Flett, Hewitt, & Pietrowsky). **C.** NFC scores versus SAPS scores. NFC and SAPS scores were positively correlated ( $r = 0.24$ ;  $p < 0.01$ ). **D.** Mean fair wage rating on the 1-back, 3-detect, and 2-back tasks by tertile split NFC groups. Error bars were drawn using the standard error of the mean (SEM). There was a significant quadratic relationship of NFC and mean fair wage ratings ( $\beta = -0.03$ ). Post-hoc t-tests confirmed that the significant quadratic effect of NFC was only driven by mid NFC subjects having significantly higher fair wage ratings than high NFC subjects ( $p < 0.01$ ). **E.** Mean fair wage rating on the 1-back, 3-detect, and 2-back tasks by tertile split SAPS groups. Error bars were drawn using the



SEM. A 3-way ANOVA, revealed no effect of SAPS group ( $F = 2.2, p > 0.05$ ) or of the interaction of SAPS group and task identity ( $F = 1.5, p > 0.05$ ) on fair wages.

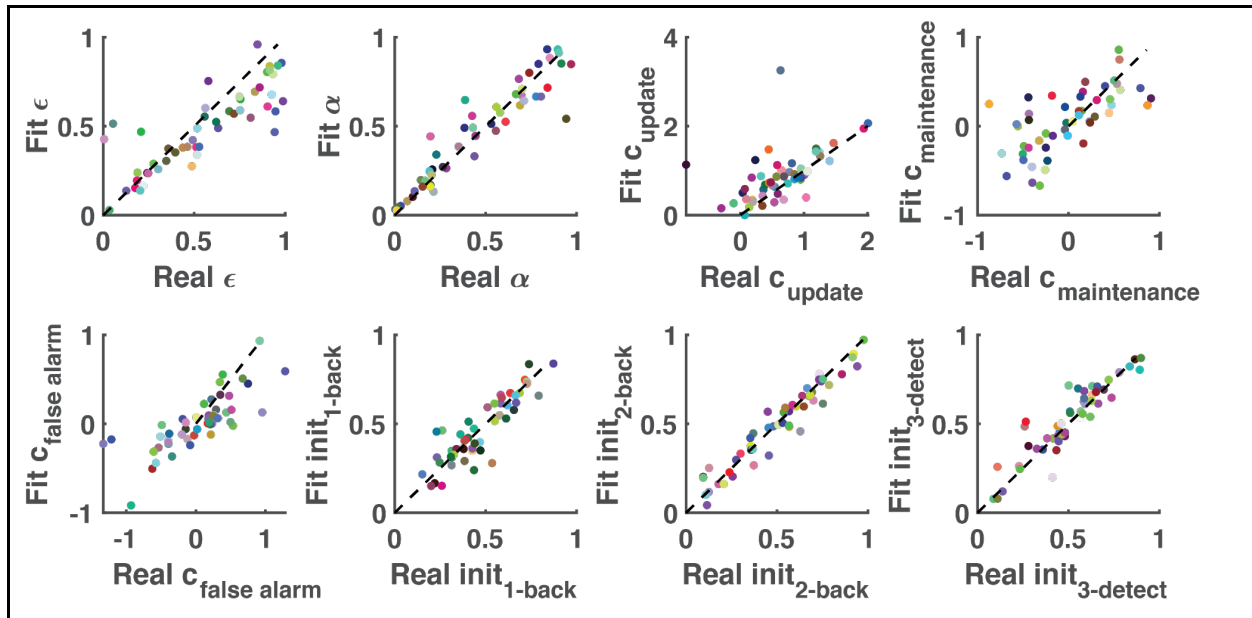
800  
801  
802



**Supplementary Figure 3. Real fair wage values versus simulated fair wage values.** For each subject, we simulated their data using the model which took the highest responsibility for their data, and their fit parameter values. Here we have selected 2 random subjects from each NFC tertile (left: low NFC, middle: middle NFC, right: high NFC) and plotted their real and fit fair wage values. The title of each plot is the mean  $r^2$  value after 100 simulations with the subject's best fit model and best fit parameter values. Markers are shaded such that later trials are displayed in darker colors, and the shape of the marker indicates which task was rated (squares are 1-back ratings, circles are 3-detect ratings, and diamonds are 2-back ratings).

803  
804

805



**Supplementary Figure 4. The results of a generate and recover procedure on a model including update, maintenance, and false alarm (FA) costs.** A dataset of simulated subjects was produced with random parameter values (constrained by the bounds of those parameters), and then fit with the same procedure as real subject data. Here we show the fits for 50 subjects out of 100, where each subject's fits are plotted in a unique color. The identity line is overlaid on each subplot in black. Comparing the fit parameter values to the real values reveals the high fidelity of the model fitting procedure. Models were fitted using the Computational Behavioral Modeling (cbm) toolbox of Piray et al (2019). All candidate models were visually inspected and verified as recoverable to avoid fitting models with parameter tradeoffs. Only models with parameter recoverability were fit to real subject data.

806

807

808

## 809 References

- 810 1. Cacioppo JT, Petty RE. The need for cognition. *J Pers Soc Psychol* [Internet].  
811 1982;42(1):116. Available from:  
812 [https://psycnet.apa.org/journals/psp/42/1/116/?casa\\_token=oYXHA6mcNkAAAAAA:Jph-](https://psycnet.apa.org/journals/psp/42/1/116/?casa_token=oYXHA6mcNkAAAAAA:Jph-9Cj-zdYqP301sDS3i56bZotUZ23_-z5vesiy792dee98du-E54y4RZj0vvgzZPNc5Mf9REHq5w-n2uaLr9Jc)  
813 [9Cj-zdYqP301sDS3i56bZotUZ23\\_-z5vesiy792dee98du-E54y4RZj0vvgzZPNc5Mf9REHq5w-](https://psycnet.apa.org/journals/psp/42/1/116/?casa_token=oYXHA6mcNkAAAAAA:Jph-9Cj-zdYqP301sDS3i56bZotUZ23_-z5vesiy792dee98du-E54y4RZj0vvgzZPNc5Mf9REHq5w-n2uaLr9Jc)  
814 [n2uaLr9Jc](https://psycnet.apa.org/journals/psp/42/1/116/?casa_token=oYXHA6mcNkAAAAAA:Jph-9Cj-zdYqP301sDS3i56bZotUZ23_-z5vesiy792dee98du-E54y4RZj0vvgzZPNc5Mf9REHq5w-n2uaLr9Jc)
- 815 2. Cacioppo JT, Petty RE, Feng Kao C. The efficient assessment of need for cognition. *J Pers*  
816 *Assess*. 1984;48(3):306–7.
- 817 3. Inzlicht M, Shenhav A, Olivola CY. The Effort Paradox: Effort Is Both Costly and Valued.  
818 *Trends Cogn Sci* [Internet]. 2018 Apr;22(4):337–49. Available from:  
819 <http://dx.doi.org/10.1016/j.tics.2018.01.007>
- 820 4. Kool W, McGuire JT, Rosen ZB, Botvinick MM. Decision making and the avoidance of  
821 cognitive demand. *J Exp Psychol Gen*. 2010;139(4):665.
- 822 5. Kool W, Botvinick M. A labor/leisure tradeoff in cognitive control. *Journal of experimental*  
823 *psychology General* [Internet]. 2014;143(1):131–41. Available from:  
824 <http://dx.doi.org/10.1037/a0031048>
- 825 6. Sandra DA, Otto AR. Cognitive capacity limitations and Need for Cognition differentially  
826 predict reward-induced cognitive effort expenditure. *Cognition* [Internet]. 2018  
827 Mar;172:101–6. Available from: <http://dx.doi.org/10.1016/j.cognition.2017.12.004>
- 828 7. Westbrook A, Kester D, Braver TS. What is the subjective cost of cognitive effort? Load,  
829 trait, and aging effects revealed by economic preference. *PLoS One* [Internet]. 2013 Jul  
830 22;8(7):e68210. Available from: <http://dx.doi.org/10.1371/journal.pone.0068210>
- 831 8. Sayalı C, Badre D. Neural systems of cognitive demand avoidance. *Neuropsychologia*  
832 [Internet]. 2019 Feb 4;123:41–54. Available from:  
833 <http://dx.doi.org/10.1016/j.neuropsychologia.2018.06.016>
- 834 9. Sayalı C, Badre D. Neural systems underlying the learning of cognitive effort costs  
835 [Internet]. Cold Spring Harbor Laboratory. 2020 [cited 2021 Mar 5]. p. 2020.06.08.139618.  
836 Available from: <https://www.biorxiv.org/content/10.1101/2020.06.08.139618v1.abstract>
- 837 10. Agrawal M, Mattar MG, Cohen JD, Daw ND. The temporal dynamics of opportunity costs: A  
838 normative account of cognitive fatigue and boredom. *Psychol Rev* [Internet]. 2022  
839 Apr;129(3):564–85. Available from: <http://dx.doi.org/10.1037/rev0000309>
- 840 11. Constantino SM, Daw ND. Learning the opportunity cost of time in a patch-foraging task.  
841 *Cogn Affect Behav Neurosci* [Internet]. 2015 Dec;15(4):837–53. Available from:  
842 <http://dx.doi.org/10.3758/s13415-015-0350-y>
- 843 12. Otto AR, Daw ND. The opportunity cost of time modulates cognitive effort.  
844 *Neuropsychologia* [Internet]. 2019 Feb 4;123:92–105. Available from:  
845 <http://dx.doi.org/10.1016/j.neuropsychologia.2018.05.006>
- 846 13. Callaway F, Jain YR, van Opheusden B, Das P, Iwama G, Gul S, et al. Leveraging artificial

- 847 intelligence to improve people's planning strategies. *Proc Natl Acad Sci U S A* [Internet].  
848 2022 Mar 22;119(12):e2117432119. Available from:  
849 <https://www.pnas.org/doi/abs/10.1073/pnas.2117432119>
- 850 14. Felso V, Lieder F. Measuring individual differences in the depth of planning [Internet]. 2022.  
851 Available from: [psyarxiv.com/xmf3y](https://psyarxiv.com/xmf3y)
- 852 15. Ho MK, Abel D, Correa CG, Littman ML, Cohen JD, Griffiths TL. People construct simplified  
853 mental representations to plan. *Nature* [Internet]. 2022 Jun;606(7912):129–36. Available  
854 from: <http://dx.doi.org/10.1038/s41586-022-04743-9>
- 855 16. Kool W, Gershman SJ, Cushman FA. Planning Complexity Registers as a Cost in  
856 Metacontrol. *J Cogn Neurosci* [Internet]. 2018 Oct;30(10):1391–404. Available from:  
857 [http://dx.doi.org/10.1162/jocn\\_a\\_01263](http://dx.doi.org/10.1162/jocn_a_01263)
- 858 17. Bustamante L, Lieder F, Musslick S, Shenhav A, Cohen J. Learning to Overexert Cognitive  
859 Control in a Stroop Task. *Cogn Affect Behav Neurosci* [Internet]. 2021 Jan 6; Available  
860 from: <http://dx.doi.org/10.3758/s13415-020-00845-x>
- 861 18. Cohen JD, Perlstein WM, Braver TS, Nystrom LE, Noll DC, Jonides J, et al. Temporal  
862 dynamics of brain activation during a working memory task. *Nature* [Internet]. 1997 Apr  
863 10;386(6625):604–8. Available from: <http://dx.doi.org/10.1038/386604a0>
- 864 19. MacLeod CM. The Stroop task: The “gold standard” of attentional measures. *J Exp Psychol*  
865 *Gen* [Internet]. 1992;121(1):12. Available from: [https://psycnet.apa.org/record/1992-22285-001?casa\\_token=un3vMHeRB6AAAAA:6eTn4jyRxAi75AStUH0QgySs3jxVwS\\_BhFeZ3vsRZ-Nyf51B41\\_JoyaTTbxvxhp0WFKUZpgpDAG1acynF5cfkOA](https://psycnet.apa.org/record/1992-22285-001?casa_token=un3vMHeRB6AAAAA:6eTn4jyRxAi75AStUH0QgySs3jxVwS_BhFeZ3vsRZ-Nyf51B41_JoyaTTbxvxhp0WFKUZpgpDAG1acynF5cfkOA)
- 868 20. Becker GM, DeGroot MH, Marschak J. Measuring utility by a single-response sequential  
869 method. *Behav Sci* [Internet]. 1964 Jul;9(3):226–32. Available from:  
870 <http://doi.wiley.com/10.1002/bs.3830090304>
- 871 21. De Martino B, Kumaran D, Holt B, Dolan RJ. The neurobiology of reference-dependent  
872 value computation. *J Neurosci* [Internet]. 2009 Mar 25;29(12):3833–42. Available from:  
873 <http://dx.doi.org/10.1523/JNEUROSCI.4832-08.2009>
- 874 22. Boureau YL, Sokol-Hessner P, Daw ND. Deciding How To Decide: Self-Control and Meta-  
875 Decision Making. *Trends Cogn Sci* [Internet]. 2015 Nov;19(11):700–10. Available from:  
876 <http://dx.doi.org/10.1016/j.tics.2015.08.013>
- 877 23. Frömer R, Lin H, Dean Wolf CK, Inzlicht M, Shenhav A. When effort matters: Expectations  
878 of reward and efficacy guide cognitive control allocation [Internet]. *bioRxiv*. 2020 [cited 2022  
879 Aug 11]. p. 2020.05.14.095935. Available from:  
880 <https://www.biorxiv.org/content/10.1101/2020.05.14.095935>
- 881 24. Kurzban R, Duckworth A, Kable JW, Myers J. An opportunity cost model of subjective effort  
882 and task performance. *Behav Brain Sci* [Internet]. 2013 Dec;36(6):661–79. Available from:  
883 [https://www.cambridge.org/core/product/identifier/S0140525X12003196/type/journal\\_article](https://www.cambridge.org/core/product/identifier/S0140525X12003196/type/journal_article)
- 884 25. Lieder F, Shenhav A, Musslick S, Griffiths TL. Rational metareasoning and the plasticity of  
885 cognitive control. *PLoS Comput Biol* [Internet]. 2018 Apr;14(4):e1006043. Available from:  
886 <http://dx.doi.org/10.1371/journal.pcbi.1006043>

- 887 26. Shenhav A, Botvinick MM, Cohen JD. The expected value of control: an integrative theory  
888 of anterior cingulate cortex function. *Neuron* [Internet]. 2013 Jul 24;79(2):217–40. Available  
889 from: <http://dx.doi.org/10.1016/j.neuron.2013.07.007>
- 890 27. Dunn TL, Inzlicht M, Risko EF. Anticipating cognitive effort: roles of perceived error-  
891 likelihood and time demands. *Psychol Res* [Internet]. 2019 Jul;83(5):1033–56. Available  
892 from: <http://dx.doi.org/10.1007/s00426-017-0943-x>
- 893 28. Shenhav A, Fahey MP, Grahek I. Decomposing the motivation to exert mental effort. *Curr*  
894 *Dir Psychol Sci* [Internet]. 2021 Aug 1;30(4):307–14. Available from:  
895 <http://dx.doi.org/10.1177/09637214211009510>
- 896 29. Rice KG, Richardson CME, Tueller S. The short form of the revised almost perfect scale. *J*  
897 *Pers Assess* [Internet]. 2014;96(3):368–79. Available from:  
898 <http://dx.doi.org/10.1080/00223891.2013.838172>
- 899 30. Braver TS, Cohen JD, Nystrom LE, Jonides J, Smith EE, Noll DC. A parametric study of  
900 prefrontal cortex involvement in human working memory. *Neuroimage* [Internet]. 1997  
901 Jan;5(1):49–62. Available from: <http://dx.doi.org/10.1006/nimg.1996.0247>
- 902 31. Piray P, Dezfouli A, Heskes T, Frank MJ, Daw ND. Hierarchical Bayesian inference for  
903 concurrent model fitting and comparison for group studies. *PLoS Comput Biol* [Internet].  
904 2019 Jun;15(6):e1007043. Available from: <http://dx.doi.org/10.1371/journal.pcbi.1007043>
- 905 32. Kool W, Botvinick M. A labor/leisure tradeoff in cognitive control. *J Exp Psychol Gen*  
906 [Internet]. 2014 Feb;143(1):131–41. Available from: <http://dx.doi.org/10.1037/a0031048>
- 907 33. Westbrook A, Braver TS. Cognitive effort: A neuroeconomic approach. *Cogn Affect Behav*  
908 *Neurosci* [Internet]. 2015 Jun;15(2):395–415. Available from:  
909 <http://dx.doi.org/10.3758/s13415-015-0334-y>
- 910 34. Milyavskaya M, Inzlicht M, Johnson T, Larson MJ. Reward sensitivity following boredom  
911 and cognitive effort: A high-powered neurophysiological investigation. *Neuropsychologia*  
912 [Internet]. 2019 Feb 4;123:159–68. Available from:  
913 <http://dx.doi.org/10.1016/j.neuropsychologia.2018.03.033>
- 914 35. Murray JD, Jaramillo J, Wang XJ. Working Memory and Decision-Making in a  
915 Frontoparietal Circuit Model. *J Neurosci* [Internet]. 2017 Dec 13;37(50):12167–86. Available  
916 from: <http://dx.doi.org/10.1523/JNEUROSCI.0343-17.2017>
- 917 36. Musslick S, Bizyaeva A, Agaron S, Leonard N, Cohen JD. Stability-flexibility dilemma in  
918 cognitive control: A dynamical system perspective. *Proceedings of the 41st Annual Meeting*  
919 *of the Cognitive Science Society* [Internet]. 2019 Jan [cited 2022 Dec 3]; Available from:  
920 <https://par.nsf.gov/biblio/10125021>
- 921 37. Baumeister RF, Muraven M, Tice DM. Ego Depletion: A Resource Model of Volition, Self-  
922 Regulation, and Controlled Processing. *Soc Cogn* [Internet]. 2000 Jun 1;18(2):130–50.  
923 Available from: <https://doi.org/10.1521/soco.2000.18.2.130>
- 924 38. Kurzban R, Duckworth A, Kable JW, Myers J. An opportunity cost model of subjective effort  
925 and task performance [Internet]. Vol. 36, *Behavioral and Brain Sciences*. 2013. p. 661–79.  
926 Available from: <http://dx.doi.org/10.1017/s0140525x12003196>

- 927 39. Wiehler A, Branzoli F, Adanyeguh I, Mochel F, Pessiglione M. A neuro-metabolic account  
928 of why daylong cognitive work alters the control of economic decisions. *Curr Biol* [Internet].  
929 2022 Aug 22;32(16):3564–75.e5. Available from:  
930 <http://dx.doi.org/10.1016/j.cub.2022.07.010>
- 931 40. Blain B, Hollard G, Pessiglione M. Neural mechanisms underlying the impact of daylong  
932 cognitive work on economic decisions. *Proc Natl Acad Sci U S A* [Internet].  
933 2016;113(25):6967–72. Available from: <http://dx.doi.org/10.1073/pnas.1520527113>
- 934 41. Cools R, D'Esposito M. Inverted-U–Shaped Dopamine Actions on Human Working Memory  
935 and Cognitive Control. *Biol Psychiatry* [Internet]. 2011 Jun 15;69(12):e113–25. Available  
936 from: <https://www.sciencedirect.com/science/article/pii/S0006322311002782>
- 937 42. Froudust-Walsh S, Bliss DP, Ding X, Rapan L, Niu M, Knoblauch K, et al. A dopamine  
938 gradient controls access to distributed working memory in the large-scale monkey cortex.  
939 *Neuron* [Internet]. 2021 Nov 3;109(21):3500–20.e13. Available from:  
940 <http://dx.doi.org/10.1016/j.neuron.2021.08.024>
- 941 43. O'Reilly RC, Frank MJ. Making working memory work: a computational model of learning in  
942 the prefrontal cortex and basal ganglia. *Neural Comput* [Internet]. 2006 Feb;18(2):283–328.  
943 Available from: <http://dx.doi.org/10.1162/089976606775093909>
- 944 44. Lieder F, Griffiths T. Helping people make better decisions using optimal gamification. In:  
945 *CogSci* [Internet]. 2016. Available from:  
946 [https://re.is.mpg.de/uploads\\_file/attachment/attachment/610/GamificationCogSciRevised.p](https://re.is.mpg.de/uploads_file/attachment/attachment/610/GamificationCogSciRevised.pdf)  
947 [df](https://re.is.mpg.de/uploads_file/attachment/attachment/610/GamificationCogSciRevised.pdf)
- 948 45. Lieder F, Krueger PM, Callaway F, Griffiths T. A reward shaping method for promoting  
949 metacognitive learning [Internet]. 2017. Available from: [psyarxiv.com/qj346](https://psyarxiv.com/qj346)
- 950 46. Lieder F, Chen OX, Krueger PM, Griffiths TL. Cognitive prostheses for goal achievement.  
951 *Nat Hum Behav* [Internet]. 2019 Oct;3(10):1096–106. Available from:  
952 <http://dx.doi.org/10.1038/s41562-019-0672-9>
- 953 47. Kool W, Botvinick M. Mental labour. *Nat Hum Behav* [Internet]. 2018 Dec;2(12):899–908.  
954 Available from: <http://dx.doi.org/10.1038/s41562-018-0401-9>
- 955 48. Froböse MI, Westbrook A, Bloemendaal M, Aarts E, Cools R. Catecholaminergic  
956 modulation of the cost of cognitive control in healthy older adults. *PLoS One* [Internet].  
957 2020 Feb 21;15(2):e0229294. Available from:  
958 <http://dx.doi.org/10.1371/journal.pone.0229294>
- 959 49. Hara K, Adams A, Milland K, Savage S, Callison-Burch C, Bigham JP. A Data-Driven  
960 Analysis of Workers' Earnings on Amazon Mechanical Turk. In: *Proceedings of the 2018*  
961 *CHI Conference on Human Factors in Computing Systems* [Internet]. New York, NY, USA:  
962 Association for Computing Machinery; 2018 [cited 2021 Mar 5]. p. 1–14. (CHI '18).  
963 Available from: <https://doi.org/10.1145/3173574.3174023>
- 964 50. de Leeuw JR. jsPsych: a JavaScript library for creating behavioral experiments in a Web  
965 browser. *Behav Res Methods* [Internet]. 2015 Mar;47(1):1–12. Available from:  
966 <http://dx.doi.org/10.3758/s13428-014-0458-y>

- 967 51. Cohen JD, Perlstein WM, Braver TS, Nystrom LE, Noll DC, Jonides J, et al. Temporal  
968 dynamics of brain activation during a working memory task. *Nature* [Internet]. 1997 Apr  
969 10;386(6625):604–8. Available from: <http://dx.doi.org/10.1038/386604a0>
- 970 52. Rac-Lubashevsky R, Kessler Y. Dissociating working memory updating and automatic  
971 updating: The reference-back paradigm. *J Exp Psychol Learn Mem Cogn* [Internet]. 2016  
972 Jun;42(6):951–69. Available from: <http://dx.doi.org/10.1037/xlm0000219>
- 973 53. Matlab S. Matlab. The MathWorks, Natick, MA [Internet]. 2012; Available from:  
974 [https://itb.biologie.hu-berlin.de/~kempter/Teaching/2003\\_SS/gettingstarted.pdf](https://itb.biologie.hu-berlin.de/~kempter/Teaching/2003_SS/gettingstarted.pdf)
- 975 54. Berinsky AJ, Margolis MF, Sances MW. Separating the shirkers from the workers? Making  
976 sure respondents pay attention on self-administered surveys: Separating the shirkers from  
977 the workers? *Am J Pol Sci* [Internet]. 2014 Jul;58(3):739–53. Available from:  
978 <http://doi.wiley.com/10.1111/ajps.12081>