1    **Title:**

2    A computational model for individual differences in non-reinforced learning for
3    individual items
4

5    **Running Title**

6    Computational model of non-reinforced learning
7

8    **List of Authors**

9    Tom Salomon[1], Alon Itzkovitch[1], Nathaniel D. Daw[2], Tom Schonberg*[1,3]

10       [1] School of Neurobiology, Biochemistry and Biophysics, Faculty of Life Sciences, Tel

11       Aviv University.

12       [2] Neuroscience Institute, Department of Psychology, Princeton University.

13       [3] Sagol School of Neuroscience, Tel Aviv University.

14       * Correspondence concerning this article should be addressed to Tom Schonberg,

15       Department of Neurobiology, Faculty of Life Sciences and Sagol School of

16       Neuroscience, Tel Aviv University, Ramat Aviv 6997801, Tel Aviv, Israel. Email:

17       Schonberg@tauex.tau.ac.il.

18                                **Abstract**

19       Cue-Approach Training (CAT) is a paradigm that enhances preferences without external

20    reinforcmeents, suggesting a potential role for internal learning processes. Here, we developed a

21    novel Bayesian computational model to quantify anticipatory response patterns during the

22    training phase of CAT. This phase includes individual items and thus this marker is potentially of

23    internal learning signals at the item level. Our model, fitted to meta-analysis data from 29 prior

24    CAT experiments, was able to predict individual differences in non-reinforced preference

25    changes using a key computational marker. Crucially, two new experiments manipulated the

26    training procedure to influence the model's predicted learning marker. As predicted and

27    preregistered, the manipulation successfully induced differential preference changes, supporting

28    a causal role of our model. These findings demonstrate powerful potential of our computational

29    framework for investigating intrinsic learning processes. This framework could be used to

30    predict preference changes and opens new avenues for understanding intrinsic motivation and

31    decision-making.

32

33

34

35 **Teaser**

36 Bayesian modeling of response time predicts individual differences in non reinforced preference
37 change.

## Introduction

38

39     Value construction and modification are commonly understood as the result of learning

40     processes that rely on external reinforcements [1–4]. While reinforcement-based learning

41     dominated the research field of value-modifications, a recent review[5] highlighted the various

42     means to influence preferences and choices without external reinforcements covered in the

43     literature, dating back to the *mere exposure effect*. In those studies preferences for stimuli could

44     be enhanced merely by repeated exposure to stimuli [6]. Similarly, another non-reinforced

45     preference change paradigm that did not rely on external reinforcements was found in work

46     showing preference modification following previous choices[7]. These studies provide a strong

47     theoretical framing for efficacy of interventions such as advertisements on preferences and how

48     life choices affect future decisions [8,9], above and beyond external reinforcement-based

49     manipulations [10,11].

50     A decade ago, a unique procedure named the Cue-Approach Training (CAT) [12] task has

51     been developed as a reliable means to change preferences without external reinforcements. The

52     CAT procedure is a multiphase task with an initial training session on individual items and a

53     subsequent binary choice phase. The first studies with the CAT procedure, showed that

54     preferences for snack food stimuli could be modified via an association between images and a

55     neutral cue and rapid motor response[12–14]. The task included a Go/NoGo training phase, during

56     which a set of snack food stimuli were presented individually on the screen. Some snacks were

57     associated with a Go cue, to which participants were required to respond with a rapid button

58     press response (Go stimuli). Another (larger) set of snack food stimuli were passively presented

59     for the same exposure period, without a cue or response (NoGo stimuli). The training phase

60     included several repetitions (runs) in which participants could learn the association of individual

61     Go stimuli with the cue and response. In the subsequent post-training probe phase Participants

62     demonstrated enhanced preferences for the Go stimuli over NoGo stimuli. This phase included

63     binary choices between two snacks for actual consumption and these two snacks had a similar

64     initial subjective value [12–14].

65     Multiple studies with the CAT procedure exemplified its efficacy in modifying

66     preferences for a wide range of stimuli beyond snack food items, including healthy food items,

67     faces, positive affective stimuli and fractal art stimuli[15–18]. Overall, across dozens of studies with

68     different types of stimuli, participants consistently demonstrated enhanced preference for Go

3

69    stimuli over NoGo stimuli of similar initial value[12–21]. Studies that examined the long-term
70    maintenance of the CAT effect in additional follow-up sessions, found that the preference
71    modification effect persisted for months without any additional training sessions[12,15,18,21,22]. The
72    fact that a preference modification effect was observed following a training procedure with no
73    external reinforcement or feedback, on individually presented items, inspired the idea that CAT
74    relies on a non-reinforced valuation pathway, putatively involving attention and motor neural
75    circuits[5].

76        Several mechanisms have been suggested to underlie non-externally reinforced
77    paradigms. First, the preferences-as-memory theoretical framework suggests that value could
78    actually be framed as retrieval of relevant knowledge about alternatives from memory[23].
79    Relatedly, studies showed that Go stimuli that were better remembered were also more preferred
80    and that memory was positively correlated with likelihood of choosing Go stimuli[19,21]. Attention
81    has also been shown to play a key role in value-based decision making both on its own and in
82    synergy with memory[24,25]. Some support for the involvement of attention was found in eye-
83    tracking studies with CAT, which found that during the probe phase, participants view chosen
84    Go stimuli more [12,13,18,26], suggesting that increased gaze time during the probe phase indicated
85    attentional evidence-gathering process, which increased the likelihood of choosing the Go
86    stimuli over NoGo stimuli [27–29]. However, most of these studies [13,18,26] did not find enhanced
87    gaze for Go stimuli when they were not chosen, undermining the hypothesis that enhanced
88    attention could drive the CAT preference modification effect on its own.

89        However, these previous mechanistic views have focused on the probe phase where
90    preferences are being expressed, whereas preference is changed during the training phase.
91    During training all stimuli are presented individually for the same duration, thus mere exposure
92    or viewing time alone could not account for the preference modification effect. Studies that
93    examined the necessary features of the training task discovered that simply responding slowly to
94    the Go stimuli was not sufficient to induce preference changes, as preference modification did
95    not occur in a training version in which the Go cue was presented immediately with the Go
96    stimulus allowing a full 1 second to respond. Similarly, training with all Go cue items presented
97    in a block of Go stimuli[13] did not yield a preference change effect. These findings suggested that
98    simple attention or motor response are not sufficient, and that the rapid motor response during

99  CAT and the motor preparation learning are crucial factors for non-reinforced preference
100 modification with CAT.

101         The critical question that persists is what the mechanisms are, that operate at the
102 individual item level without external reinforcements to induce the subsequent observed
103 behavioral change in the binary choice phase. A prominent clue as to the underlying mechanisms
104 during training were identified in a recent imaging study with CAT[18], which found that
105 individual differences in the preference modification effect (measured during the binary choice
106 probe phase) correlated with increased neural activity during training within the supplementary-
107 motor cortex and the striatum. These regions are associated with motor-planning[30,31] and
108 reinforcement-based learning[32–35], respectively. The striatum has been suggested to be key neural
109 hub for learning thanks to its unique dual-role in both reward processing and motor regulation
110 networks [36,37]. The involvement of the striatum linking non-reinforced training and subsequent
111 preference change, could putatively indicate the presence of internal reinforcement [38] during the
112 training phase, particularly in the absence of external feedback.

113         These neuroimaging results, showing that enhance striatal activity during training
114 correlated with preference modification were the key inspired the current work's hypothesis. We
115 utilized the unique structure of CAT that trains individual items, which allows us to shed light on
116 valuation mechanisms at the individual item level. An important aspect of the CAT task is that
117 training is performed at the individual item basis, and preference change is tested in a subsequent
118 binary choice phase. Therefore, we surmised that identifying individual differences in motor
119 learning during the training phase could serve as the key factor to explain differential non-
120 reinforced behavioral change effect following CAT.

121         We tested our hypothesis by developing a novel Bayesian computational model with an
122 individualized learning marker based on training motor-response data on single items and testing
123 its association with the preference modification effect observed in the subsequent probe phase.
124 For the development of the new computational framework, we utilized 29 different experiments
125 that used the CAT procedure. Then, to examine the causal impact of our newly proposed
126 mechanism we designed a new non-reinforced training procedure which aimed to directly impact
127 learning and establish a framework for understanding the intrinsic mechanisms underlying non-
128 reinforced preference modification through motor-learning.

129    The current work thus consists of two main parts. The first part includes a meta-analysis
130    study, whereby using data of 864 participants from 29 previous CAT experiments, we devised a
131    novel marker for learning. By combining multiple samples with the CAT task into one meta-
132    analysis, we were able to create a unique dataset of a standardized non-reinforced preference
133    modification task. To identify an individualized marker for learning, we first examined RT
134    patterns during the training phase in an exploratory analysis and built a computational model of
135    these patterns using a Bayesian computational framework. The Bayesian model included an
136    individualized parameter, which modeled a transition from cue-dependent responses to
137    anticipatory responses and was used as a foreseeing marker of individualized learning. We
138    examined the association of the individualized learning marker with individual differences in
139    preference modification, as measured in the subsequent binary-choice probe phase. We further
140    examined the hypothesis that the mechanism of non-reinforced preference modification is
141    operating both at the participants-level (i.e., some participants learn better than others), as well as
142    on a more granular item-level (i.e., within participants, some stimuli are better learned than
143    other). We expanded our computational model to capture variability in within-participant item-
144    level learning and examined its association with subsequent choice pattern.

145    In the second part of this work, we implemented the findings from the Bayesian model to
146    develop a novel CAT design and directly manipulate the computational marker in order to
147    demonstrate a causal direction between the new model and preference change. In the new design,
148    we used two different cue-contingency conditions to manipulate learning difficulty. We
149    hypothesized that the cue contingency manipulation will affect the motor response and will be
150    captured by the individualized learning marker devised in the first part, as well as predict
151    subsequent preference modification in the probe phase. The second part included one
152    preliminary experiment ($n = 20$), which was followed by a larger pre-registered replication study
153    ($n = 59$). All the hypotheses, experimental design and analyses plans were preregistered before
154    data collection of the replication study began (https://osf.io/nwr4v). By manipulating the task
155    design to induce differential behavioral modification effect, we aimed to establish a causal
156    relationship between the proposed cognitive mechanism and the preference modification effect.

157    By identifying an individualized computational marker for learning in CAT, the current
158    work provides empirical evidence that non-externally reinforced preference modification occurs
159    at the individual item level via motor-learning cognitive mechanisms. We establish a reliable

160  quantification of this individualized learning, without the need to explicitly ask participants to

161  reveal their preference overtly. The ability to exemplify a marker for non-externally reinforced

162  preference changes at the individual item level, this work will shed light on the putative

163  mechanism of this effect. Furthermore, it will allow us to identify learning patterns both

164  between-participants and within-participant, at the individual item level. A computational marker

165  for learning holds the potential to passively monitor learning in real-time and support the

166  development of novel closed-loop interventions.

167

168                           **Study 1: Meta-analysis of CAT studies**

169          In the first part of the current work we performed a meta-analysis of 29 CAT experiments

170  with a total of $N = 864$ participants, which had been conducted by our research group and

171  colleagues (see methods for sample size, publication origin and demographics of the individual

172  experiments). All experiments included the three main phases of the CAT procedure: initial

173  preference evaluation, cue-approach training, and a preference modification probe task (Figure 1;

174  see detailed description in the Methods section). In the first phase of the task, initial baseline

175  preferences for a set of stimuli were evaluated in an auction procedure [39], when snacks were

176  used, or a forced choice task [15] for all other stimulus types . Following initial preference

177  evaluation, in the cue-approach training phase, approximately 30% of the stimuli were

178  consistently associated with a neutral Go cue (Go stimuli), to which participants needed to

179  respond with a rapid button press, while the rest of the stimuli were presented passively without

180  cue or response (NoGo stimuli). Participants were instructed to wait for the Go cue and respond

181  as fast as they could, once they perceived the cue. In the final probe phase, preference

182  modification following CAT was evaluated. Participants were asked to indicate their preferred

183  stimulus out of pairs of stimuli with similar initial value, in which one of the two stimuli was a

184  Go stimulus and the other was a NoGo stimulus.

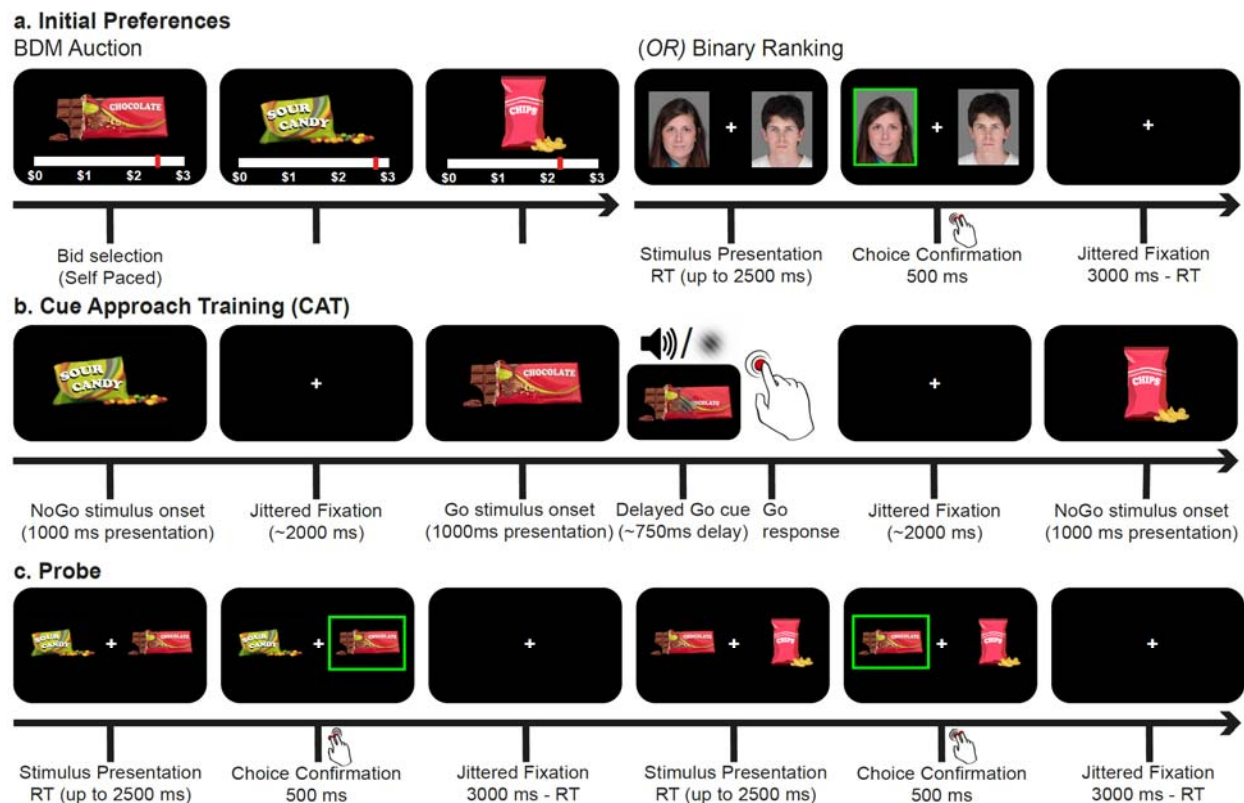Computational marker of non-reinforced learning



Figure 1. General outline of the three main procedural components of the Cue-approach training (CAT) paradigm. (a) Initial preference evaluation task. Baseline preferences for all stimuli were evaluated either with using a Becker-DeGroot-Marschak (BDM) auction (for consumable stimuli) or using a forced choice binary ranking task (for non-consumable stimuli). (b) In the cue approach training (CAT) task, approximately 30% of stimuli were presented in association with a delayed cue (auditory or visual) to which participants responded with a rapid button press (Go stimuli). All other stimuli were presented without cue and response (NoGo stimuli). (c) In the probe phase, preference modification was evaluated using a binary forced choice between pairs of stimuli of similar initial value, where one was a Go stimulus and the other a NoGo stimulus. Face images are included with permission from the copyright holder [40].

The meta-analysis study focused on identifying response patterns in the training task which could indicate learning efficacy and predict the behavioral change measured in the subsequent probe task. We hypothesized that faster Reaction Time (RT) could indicate improvement in learning.

201

**Results**

**RT analysis of CAT.** The training task in the different experiments consisted of 8-20 training runs (the number of training runs varied between experiments and identical within each experiment). During each run, all Go and NoGo stimuli were presented once individually on-

8

206  screen. During NoGo trials, stimuli were presented without a Go cue and required no response.

207  Whenever Go stimulus appeared (approximately 30% of trials), a Go cue was followed, to which

208  participants were asked to respond with a rapid button press, before stimuli offset. Participants

209  were instructed to respond *only* when they hear or see the cue. Each trial started with stimulus

210  onset, which was presented on screen for 1000ms. During Go trials, a Go cue appeared

211  approximately 750ms following stimulus onset (thus leaving participants approximately 250ms

212  to respond). The cue onset changed according to participants' performance, so that each

213  successful Go response was followed with an increase in Go onset time (thus leaving less time to

214  respond and making the next trial more challenging; see detailed description in the methods

215  section).

216  To identify unique response patterns, we performed an exploratory analysis of the RT

217  distribution during Go trials. Before analyzing RT in the CAT task, invalid trial trials were

218  excluded. Resulting in the exclusion of missed trials, in which participants failed to make a

219  response within the allocated 1500ms from stimulus onset timeframe (0.90% of trials), and trials

220  in which the response was shorter than a threshold of 100ms (0.02% of trials), as these trials are

221  likely to be indicative of inattentional response. In total, 1731 trial out of 188,224 Go trials were

222  excluded (0.92% of trials).

223  Examining the RT density distribution in the training task as a function of training run,

224  revealed a distinct pattern, wherein as training progressed, participants' RTs were less

225  homogeneous and started to form a growing peak of early RTs (Figure 2a). In all 29 experiments,

226  the Go cue changed in each trial, according to participants' performance (see methods), thus, a

227  more informative measurement of RT was the time from cue onset (RT minus cue-onset),

228  referred to here as *effective RT*. In earlier training runs, when the Go stimuli were associated for

229  the first time with the Go cue, most effective RTs were clustered around a unimodal center,

230  approximately 300ms following the Go signal onset (effective RT $M$ = 293ms) with 99% of

231  effective RTs larger than 145ms. We use here this empirical quantile of 145ms as threshold to

232  define rapid *anticipatory responses*, as these responses can be elicited when the participants

233  predicted a go cue would appear. As training progressed, a consistent pattern appeared in the

234  effective RT data - the main peak of the RT distribution was reduced, while a growing portion of

235  participants' responses consisted of faster anticipatory RTs, in many cases preceding the Go cue

236  onset (Figure 2b). For example, while the proportion of anticipatory responses comprised only

237    1% of RTs in the first training run, they comprised 2.7% of RTs by the 5th run, 12.2% of the

238    trials by the 10th run, 20.5% of trials by the 15th run, and finally 29.4% of trials in the final 20th

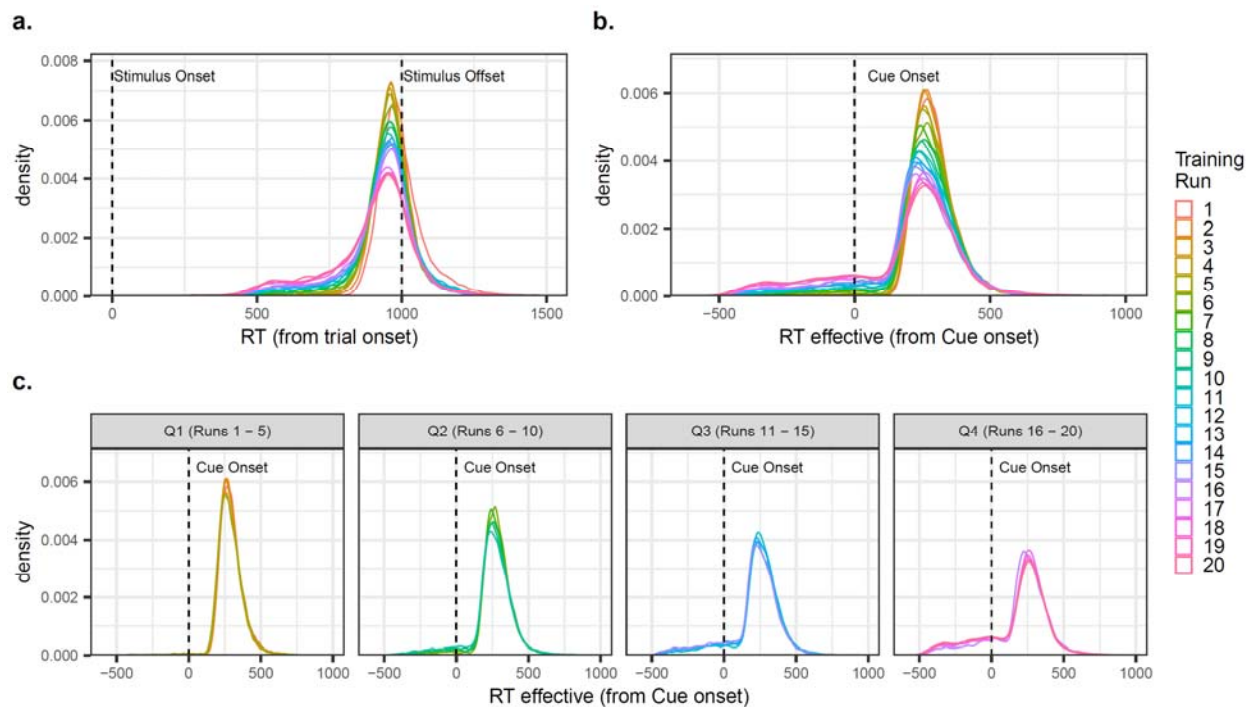239    run (aggregated across all samples).



240
241    Figure 2. Reaction time distribution of data pooled from meta-analysis with 29 studies.
242    Density plots of (a) RTs time-locked to the trial onset and (b) RTs time-locked to the varying Go
243    signal onset (*effective RT*). Color indicates training run (repetition). As training progressed, the
244    RT distribution shifted from a cue-dependent unimodal distribution to left-tailed distribution
245    with increasing proportion of early anticipatory responses. (c) a clear differentiation in effective
246    RT is apparent when splitting the data (from Figure 2b) to four quartiles, according to the
247    training run (each quartile indicating five consecutive training runs).
248
249        Thus, RT data suggest a temporal-dependent response pattern. As training progress,

250    participants tended to rely less on the Go cue (late cue-dependent RTs) and generated more early

251    anticipatory responses. Interestingly, the mode of the RT distribution seemed to remain stable,

252    and only the relative proportion of the two distributions changes. We hypothesized that this

253    transition pattern from the cue-dependent response to a stimulus-triggered anticipatory response

254    reflected a process of learning about the stimuli that could be associated with behavioral change

255    in the later preference probe. Thus, we next sought to develop a model to characterize the

256    strength of this effect within-individual.

257        ***Individualized-learning computational model.*** Based on the unique RT pattern observed

258    in the exploratory analysis of CAT task data, we developed a novel computational Bayesian

259  model. To capture the form of RT distribution with rapid anticipatory responses and slower cue-

260  dependent responses, we modeled the effective RT as a mixture of two Gaussian distributions,

261  with two different means ($\mu_1$ and $\mu_2$ free parameters), two standard deviations ($\sigma_{\varepsilon_1}$ and $\sigma_{\varepsilon_2}$), and

262  a mixture proportion ($\Theta_{i,t}$) indicating the relative mixture proportion for participant$_i$ at trial$_t$. The

263  model was implemented with the Stan probabilistic programming language [39], which optimized

264  the model parameter estimates using a Markov-chain Monte Carlo (MCMC) approach. See

265  methods section for a detailed description of the model and Supplementary Code S1. The first

266  Gaussian was restricted to have lower mean than the second Gaussian to maintain their

267  correspondence across different MCMC sampling chains (i.e., $\mu_1, \sigma_{\varepsilon_1}$ would correspond with the

268  early anticipatory RT distribution). To capture the gradual change in anticipatory responses over

269  training, the mixture proportion was modeled as a time-dependent variable, expressed by the

270  following formula:

$$\Theta_{i,t} = \phi\big(\theta_0 + \theta_{slope_i} Run_{i,t}\big)$$

271  Where the mixture probability $\Theta_{i,t}$ for participant$_i$ at trial$_t$ was modeled as a linear

272  function (with a fixed $\theta_0$ intercept and random participant-level $\theta_{slope_i}$ free paramter) of training

273  run (where the $Run_{i,t}$ ranged between [0,1]). The linear trend was scaled non-linearly to

274  probability-appropriate values of [0, 1] range using a normal distribution cumulative distribution

275  function (normal CDF; denoted with $\phi$), resulting in a sigmoid-like link function. See further

276  details in the methods section.

277  To associate all learning with a single parameter and a simple interpretation of $\theta_{slope_i}$

278  free parameter, we chose to fix the $\theta_0$ intercept term, which represents the mixture proportion at

279  the very first run ($Run_{i,t} = 0$). The $\theta_0$ intercept term was fixed at -3.1, corresponding with

280  $\Theta_{i,t} = \phi(-3.1) \approx 0.1\%$ mixture proportion at the very first run for all participants. Thus, higher

281  $\theta_{slope_i}$ parameter estimate indicates an individual participant with faster transition from slow

282  cue-dependent RTs to fast anticipatory responses, which corresponds with larger $\Theta_{i,t}$ mixture

283  proportion at the very last run. For example, parameter estimate of $\theta_{slope_i} = 3.1$, indicated that at

284  the last run ($Run_{i,t} = 1$), the mixture proportion of anticipatory responses would be $\phi(0) = 50\%$,

285  while $\theta_{slope_i} = 4.745$ would correspond with mixture proportion of $\phi(1.645) \approx 95\%$, at the last

286  run. We hypothesized that this individualized parameter would be indicative of improved

11

287    learning in the training phase of CAT and would also be positively associated with stronger

288    preference modification effect, measured in the subsequent probe phase.

289         Fixing the expected anticipatory responses to very low number at the first run, was an

290    assumption we chose to adapt in the current work to be able to quantify different participants'

291    learning using $\theta_{slope_i}$ as a single comparable parameter which captured all aspects related to

292    learning. If $\theta_0$ intercept term would also have been left as a free parameter, we would have

293    needed to account for both parameters and their interaction to define a computational learning

294    marker.

295         After running four independent chains, the Bayesian model converged to a stable solution

296    (see Supplementary Fig. S1). The converged model estimated effective RTs as a time-dependent

297    mixture of two RT distributions: one of early anticipatory responses ($\mu_1$ = 107.50ms, 95%CI

298    [103.77, 111.21], $\sigma_{\varepsilon_1}$ = 276.30ms, 95%CI [274.13, 278.42]) and late cue-dependent responses

299    ($\mu_2$ = 286.12ms, 95%CI [285.67, 286.58], $\sigma_{\varepsilon_2}$ = 72.61ms, 95%CI [72.23, 73.00]). Overall,

300    participants demonstrated an increase in proportion of anticipatory responses as training

301    progressed, as manifested in the positive group-level parameter ($\theta_{slope}$ = 4.19, 95%CI [3.84,

302    4.54]), with variation between participants ($\sigma_{\theta_{slope}}$ = 5.14, 95% CI [4.85, 5.46]), indicating some

303    participants made faster transitions and some generated nearly only cue-dependent responses

304    (non-positive $\theta_{slope_i}$ estimates).

305         Posterior predictive checks (simulated distributions of RT based on estimated model

306    parameters) revealed a good fit of the model to the actual data. Simulated posterior distributions

307    recreated the patterns observed empirically, showing more rapid transition to anticipatory

308    responses in participants with higher $\theta_{slope_i}$ parameter estimate (Figure 3; Supplementary Fig.

309    S2).

12

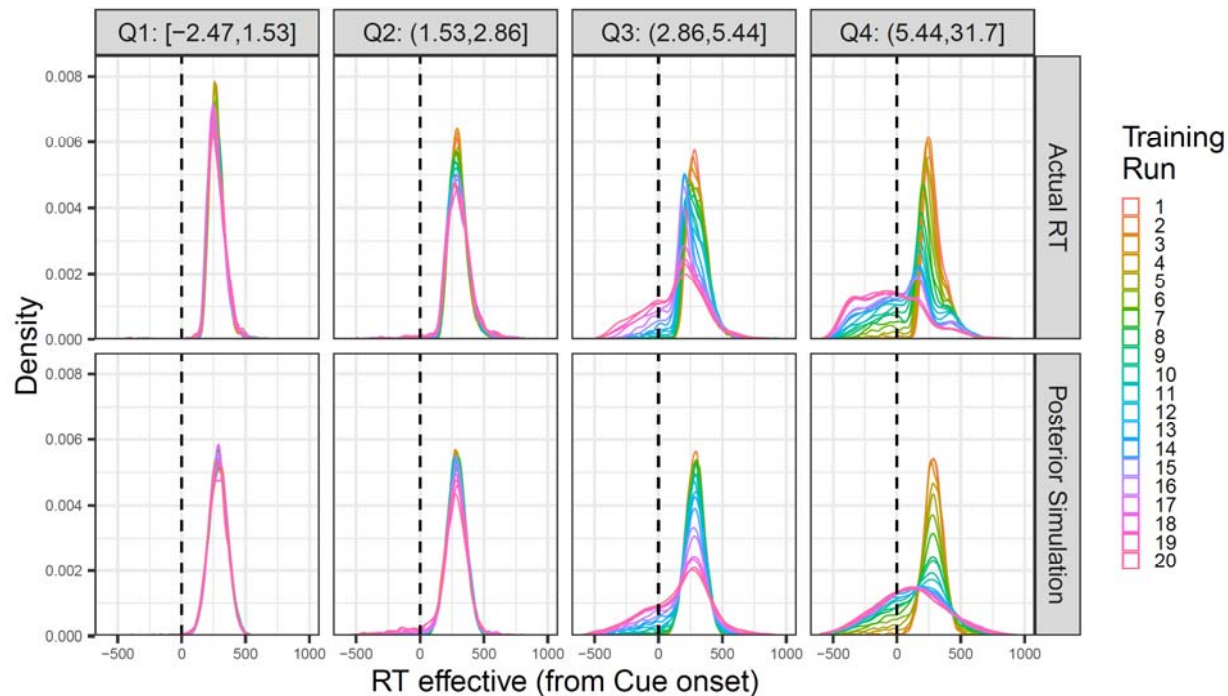Computational marker of non-reinforced learning



310
311 Figure 3. Actual RT distributions versus simulated posterior distributions, by quantile
312 group. Participants were categorized into four equal quantile groups, according to their parameter
313 estimates (denoted here as Q1-Q4; columns). Participants with higher parameter estimates were
314 characterized with faster transition to anticipatory responses (top row). Posterior simulated RT
315 distributions using mixture of Gaussians (bottom row) recreated relatively well this transition
316 pattern. Vertical dashed line represents cue-onset. See also Supplementary Fig. S2 for more
317 detailed comparison.
318
319 ***Learning parameter association with choices.*** Previous findings with CAT consistently

320 showed that preferences for Go stimuli were enhanced following CAT, as manifested in choice

321 behavior during the probe phase. When presented with a choice between a Go stimulus and a

322 NoGo stimulus of similar initial value, participants consistently chose the Go stimulus [12,14–20,22].

323 To test the hypothesis that this probe performance is associated with RT patterns during CAT,

324 we examined whether variation in the slope parameter fit to RTs correlated with probe

325 performance. Note that, under the model, the slope parameter controls the final proportion of

326 anticipatory responses, which we understood as a measure of the extent of learning about the

327 stimulus achieved by the endpoint of CAT. To evaluate the association of the parameter

328 with preference-change effect following CAT, we analyzed the proportion of probe trials in

329 which participants chose the Go over the NoGo stimulus, as a per-participant linear function of

330 (mixed logistic regression model, including random intercept and slope terms for the 29

331 experiments; see methods).

13

332    The meta-analysis showed           was positively associated with the preference

333    modification effect – i.e., participants with higher         parameter estimate, also demonstrated

334    greater odds of choosing Go stimuli (OR = 1.04, 95% CI = [1.02, 1.06], $Z$ = 5.17, $p$ = 4.6E$^{-7}$;

335    two-sided mixed model logistic regression; Figure 4). The model's intercept was significantly

336    greater than zero, i.e., even when extrapolated to very low         value, the model forecasted

337    enhanced preference for Go stimuli (intercept odds = 1.26, 95% CI = [1.18, 1.34], $Z$ = 7.03, $p$ =

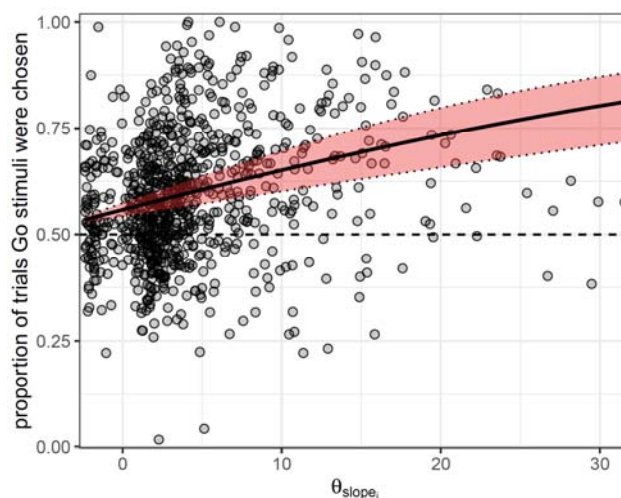338    2.0E$^{-12}$).



339
340    Figure 4. Meta-analysis results - computational marker and preference modification effect.
341    Participants who transitioned faster to anticipatory response during cue-approach training (larger
342         estimates) also demonstrated stronger preference modification effect (proportion of trials
343    Go stimuli were chosen). Trend line and surrounding red margins represent estimated preference
344    modification effect and 95% CI, respectively (mixed model logistic regression). Dots represent
345    individual participants.
346
347    The estimated variances associated with the fixed-effect terms, random-effect terms and

348    residuals were used to evaluate two      scores for the generalized (logistic) linear mixed mode: a

349    marginal              (representing the relative proportion of variance associated with the fixed-

350    effects) and a conditional              score (representing the relative proportion of variance

351    associated with the fixed-effects and random-effects; see more details in the methods section).

352    The         fixed effect accounted for              = 0.311 of the variance. With the random

353    intercept and slopes, the fixed and random effects combined accounted for              = 0.828 of

354    the overall variance.

14

355         *Additional model validation.* In a post-hoc analysis (see Supplementary materials), we

356      examined the explanatory power of $\theta_{slope_i}$ above and beyond a simpler RT-based marker. Using

357      the proportion of anticipatory responses (RTs which were faster than the top 1% of RTs in the

358      first run) each participant made as an alternative marker, we found the simpler marker was also

359      associated with subsequent probe choices. Furthermore, $\theta_{slope_i}$ showed no significant

360      contribution above and beyond this simpler RT-based marker.

361         In an additional post-hoc analysis, we used an alternative Bayesian model in which RT

362      were modeled using log-normal distributions instead of Gaussians (see Supplementary). This

363      new model better captured the right-tail shape of RTs (Supplementary Fig. S13) and replicated

364      similar correlation pattern between $\theta_{slope_i}$ parameter estimate and Go stimuli choices during

365      probe (OR = 1.05, 95% CI = [1.03, 1.07], $Z$ = 6.05, $p$ = 1.4E$^{-9}$; see Supplementary Fig. S14).

366         *Stimulus-specific learning parameter.* While the previous model aggregated the data

367      within participants, in fact, each participant has learned the Go cue association with several

368      different Go stimuli. Thus, it is possible that learning was not uniform across the entire training

369      stimuli, meaning every participant had a slightly different learning rate for each of the Go stimuli

370      she or he encountered. In an additional model, we aimed to expend our model by examining

371      behavior on a more granular scale – focusing on the variability in RT and choice for individual

372      Go stimuli within participant. Accordingly, the computational model was elaborated by

373      decomposing the per-subject learning parameter $\theta_{slope_i}$ into a set of subject- and stimulus-

374      specific parameters $\theta_{slope_{i,s}}$, indicating the transition speed to anticipatory RTs of each Go

375      stimulus$_s$ encountered by participant$_i$ during training. The model was otherwise identical in

376      design (see methods). This reanalysis of the data was conducted after the collection and analysis

377      of the following two studies which are reported in the current work.

378         The reanalyzed model converged around slightly different parameter estimates: estimated

379      early anticipatory RT distribution was characterized with lower mean ($\mu_1$ = 88.30ms, 95% CI

380      [84.1, 92.34], $\sigma_{\varepsilon_1}$ = 278.2ms, 95% CI [275.84, 280.65]), late cue-dependent responses were

381      similar ($\mu_2$ = 289.31ms, 95% CI [288.82, 289.79], $\sigma_{\varepsilon_2}$ = 75.64ms, 95% CI [75.23, 76.05]), and

382      the $\theta_{slope_{i,s}}$ parameters tended to be lower but with greater variability ($\theta_{slope}$ = 2.18, 95% CI

383      [2.01, 2.34], $\sigma_{\theta_{slope}}$ = 5.95, 95% CI [5.76, 6.16]; see Supplementary Fig. S3).

384    Similarly to the previous logistic regression analysis with per-participant parameter,

385    using the current analysis' per-stimulus parameters also significantly predicted choices of the

386    corresponding stimuli in the subsequent probe phase (Intercept: Odds = 1.45, 95%CI [1.36,

387    1.54], $p$ = 1.4E$^{-29}$; $\theta_{slope_{i,s}}$ OR = 1.02, 95%CI [1.01, 1.03], Z = 3.65, $p$ = 2.6E$^{-4}$; two-sided mixed

388    model logistic regression; see Supplementary Fig. S4). Thus, Go stimuli which were characterize

389    with larger proportion of anticipatory responses were also chosen more during the probe phase.

390    The models' fixed effects accounted for $R^2_{GLMM(m)}$ = 0.011 of the total variance and the fixed

391    effect with mixed effects accounted for $R^2_{GLMM(c)}$ = 0.775 of the total variance.

392    However, this overall effect comprises both between-participant and within-subject,

393    between-stimulus contributions. To examine the unique contribution of $\theta_{slope_{i,s}}$ above and

394    beyond participant-level $\theta_{slope_i}$ parameter (which was evaluated in the previous analysis), in a

395    third logistic regression model we included both per-participant and per-stimulus parameters

396    from the two different Stan models, as independent variable (and their random effects; see

397    methods) to predict choices during probe. Our results indicate that each parameter provided some

398    unique contribution above and beyond the other parameter (Intercept: Odds = 1.37, 95%CI [1.27,

399    1.47], $p$ = 2.38E$^{-18}$; stimulus-level $\theta_{slope_{i,s}}$: OR = 1.02, 95%CI [1.01, 1.03], Z = 3.34, $p$ = 4.2E$^{-4}$;

400    participant-level $\theta_{slope_i}$: OR = 1.01, 95%CI [1.00, 1.03], Z = 2.07, $p$ = 0.019). Thus, stronger

401    learning of anticipatory RTs, both per-individual and per-stimulus within individual, were

402    associated with increased odds of choosing Go stimuli.

403    **Interim discussion**

404    Examining RT patterns during CAT revealed for the first time a distinct time-dependent

405    RT pattern, in which participants gradually transition from slow cue-dependent responses to

406    rapid anticipatory RTs. Using Bayesian modeling, we were able to quantify this RT transition

407    pattern with stable parameter estimates, of which the $\theta_{slope_i}$ parameter provided a promising

408    computational marker to evaluated individualized differences in learning during training.

409    Furthermore, as expected, this computational marker was also found to be associated with the

410    preference modification effect in the subsequent probe phase, as participants with a more robust

411    learning marker also demonstrated stronger preference modification effect. Examining the data

412    on a stimulus-level learning parameter, demonstrated an improved explanatory power for the

413    variability in probe choices, above and beyond the wider participant-level parameter. Since the

414 computational marker was calculate based on measurements that preceded the probe phase, it

415 holds a potential to provide predictive marker for future preference modification.

416     Several challenges were noted in modeling the original experimental data with the new

417 computational model, largely owing to the fact that the original experimental procedures were

418 not designed with the present analyses in mind. First, in the training task, Go cue onsets were

419 derived based on performance using a staircase procedure. Participants who performed well and

420 responded rapidly were presented in the next trials with a more challenging cue, appearing later

421 into the trial, while participants who performed poorly were exposed to earlier occurring (less

422 difficult) cues. Consequently, anticipatory responses were not time-locked to the same event as

423 the late cue-dependent responses. While a model that uses RTs time-locked to the stimulus onset

424 could potentially model anticipatory responses shape well, using a model that is agnostic of the

425 cue-onset time, would result in difficulty to differentiate between cue-dependent and anticipatory

426 responses due to the changing cue onset. For example, a response at 750ms could be attributed to

427 cue-dependent response of a participant with very poor performance which was presented with a

428 cue at 500ms, or an anticipatory response of a well-performing participant that encountered a

429 challenging 800ms cue onset delay. Thus, making on the surface the two response patterns hard

430 to distinguish and introducing a potentially confounding source of variation across participants.

431 Therefore, in the current study we use RTs time-locked to the cue onset, which improved the

432 distinction between cue-dependent and anticipatory responses, at the expense of less accurate

433 tracking of the anticipatory response distribution shape.

434     Some lack of fit of the predicted data could also be indicative that the model's simplistic

435 assumption that anticipatory RTs are normally distributed around the same mean for all

436 participants is unlikely. This challenging effect may have led to some counter-intuitive results;

437 for instance, due to high variability in anticipatory responses, the fitted model predicted that very

438 slow responses (RT effective > 500) would be more likely in the fastest learners' late runs,

439 compared with earlier runs or with slower learners. In an additional post-hoc analysis (performed

440 after the initial write-up of this manuscript), we modeled RT using a log-normal distribution,

441 which better captured the skewed right-tail shape of RTs, which replicated the conclusions of the

442 current design. Nonetheless, the large heterogenous data structure (with over 800 participants

443 from 29 different experiments) provided a computational challenge when trying to fit more

444 complex models (e.g., with random effects per participant on this parameter or using non-normal

445  distributions). The current model was selected as a reasonable compromise that provided

446  sensible results. The fact that using an additional alternative modeling approach resulted in

447  similar conclusions as those gathered with the preregistered Gaussian model, provided vital

448  evidence that the conclusions of the current study are stable above and beyond implementation

449  choices.

450       Furthermore, the task instructions of CAT specifically mentioned that participants are

451  required to respond *after* cue-onset. An additional crucial factor which could not be controlled in

452  the current study was how rigorously participants complied with this instruction. It is possible

453  that some participants might have learned the stimulus-cue association well and could have

454  produced anticipatory responses but were reluctant to deviate from the task's guidelines. If so,

455  their $\theta_{slope_i}$ would not reliably reflect their actual learning of the task. It is important to note that

456  although the logistic regression results showed a significant positive association of choices with

457  participant-level and even stimulus-level individualized learning parameters, significant

458  preference change remains even when the model extrapolated for no learning effect. The effect

459  sizes reported above are quite modest and in a post-hoc analysis they were not found

460  significantly better compared to predicting choices with a simpler RT-based marker. Thus,

461  suggesting of additional contributions to the choice pattern which are not otherwise accounted

462  for by the computational marker. The small effect sizes might be due to the inaccurate

463  measurement of additional constructs within the computation parameter, such as the instruction

464  adherence and interaction with performance-based dynamic cue onset, which were proposed

465  here.

466       Despite these drawbacks, our unique model provided a prospective computational marker

467  for learning and subsequent preference changes. Based on the promising results of Study 1 meta-

468  analysis and considering the challenges raised above in fitting the computational model to the

469  CAT task data, in Study 2, we devised a novel design of the CAT procedure, which was tested

470  with two independent experiments, a smaller preliminary study ($n = 25$) and a larger pre-

471  registered replication study ($n = 59$). In these new experiments we addressed the limitations of

472  the previous study design. We also modified the training design to answer a more directional

473  hypothesis which does not only examine correlational association, but also try to provide

474  evidence for more consequential directionality, by inducing a differential behavioral change.

475

476        **Study 2: Novel CAT design – preliminary and replication experiments**

477        A novel experimental design was tested in two pre-registered experiments - a preliminary

478    experiment and a larger replication experiment. Based on the conclusions from Study 1, the CAT

479    procedure and instructions were altered to optimize the task for the Bayesian computational

480    framework and manipulate behavior in accordance with the hypothesized cognitive mechanism.

481        In the training task, the Go cue onset time was fixed at 850ms from stimulus onset and

482    participants were clearly instructed that they were permitted to make anticipatory responses

483    before cue onset. Furthermore, based on the conclusions from the meta-analysis study, we

484    introduced within the training task a manipulation of the predictability of the stimulus-Go cue

485    contingency, which we expected would affect the $\theta_{slope_i}$ parameter and subsequently also

486    manipulate the preference modification effect. Of the Go stimuli, half of the stimuli were always

487    associated with the Go cue (100% contingency condition), while the rest of the Go stimuli were

488    only associated in half of the presentations with the Go cue (50% contingency condition). The

489    rest of the stimuli were never associated with a Go cue (NoGo stimuli; See methods section for

490    detailed description of the new design).

491        We hypothesized that in the 50% contingency condition, learning the association of Go

492    stimuli with the Go cue would be more challenging, and would therefore be identified by a lower

493    $\theta_{slope_i}$ learning parameter. We also hypothesized that manipulating learning efficacy during

494    training would therefore induce a differential preference modification effect in the subsequent

495    probe phase, with more robust preference modification for Go stimuli trained in the 100%

496    contingency condition, compared with the more challenging 50% contingency condition.

497        The new task design was first tested in a preliminary (not preregistered) study with $n = 20$

498    valid participants, aimed to evaluate the efficacy of the new task and the replicability of the

499    meta-analysis conclusions. Following the preliminary experiment, we run an additional pre-

500    registered direct replication experiment with a larger sample size ($n = 59$), a size chosen based on

501    the preliminary experiment results. As the two experiments were identical in design, they are

502    reported together in all sections.

503

504    **Results**

505        **Reaction times in the training task.** As in the meta-analysis study, RT patterns of the

506    training task were examined. Before examining RT patterns, we removed Go trials where

19

507  participants did not respond at all (preliminary exp. = 6.33%. replication exp. 5.97% of Go trials)

508  or trials with an unlikely fast RT (RT< 250ms, preliminary exp. = 0.31%, replication experiment

509  0.21% of trials). In total, 636 (of 9600; 6.62%) and 1747 (of 28,320; 6.17%) of trials were

510  excluded from the preliminary and replication experiments, respectively.

511  In both experiments, we were able to replicate the RT pattern observed in the meta-

512  analysis study – while in early training runs, participants mostly relied on cue-dependent RTs, as

513  training progressed, an increased portion of RTs were of earlier anticipatory RTs (Figure 5).

514  Examining mean (*SD*) effective RT at Run 1, revealed similar RTs in all contingency conditions

515  - preliminary exp.: $M_{50\%Cont.} = 326.88$ms (85.13ms), $M_{100\%Cont.} = 336.60$ms (64.39ms);

516  replication exp.: $M_{50\%Cont.} = 322.49$ms (63.68ms), $M_{100\%Cont.} = 319.90$ms (78.88ms). As

517  training progressed, participants were able to anticipate the cue and respond before cue-onset.

518  Examining the differential pattern between the 100% contingency condition and the 50%

519  contingency condition showed, as expected, that participants initiated more anticipatory

520  responses in the 100% contingency condition, compared with the 50% contingency condition,

521  resulting in faster mean-RTs at Run 20 - preliminary exp.: $M_{50\%Cont.} = 253.27$ms (187.65ms),

522  $M_{100\%Cont.} = 175.55$ms (245.5ms); replication exp.: $M_{50\%Cont.} = 256.57$ms (188.94ms),

523  $M_{100\%Cont.} = 81.39$ms (271.99ms); see Figure 5.

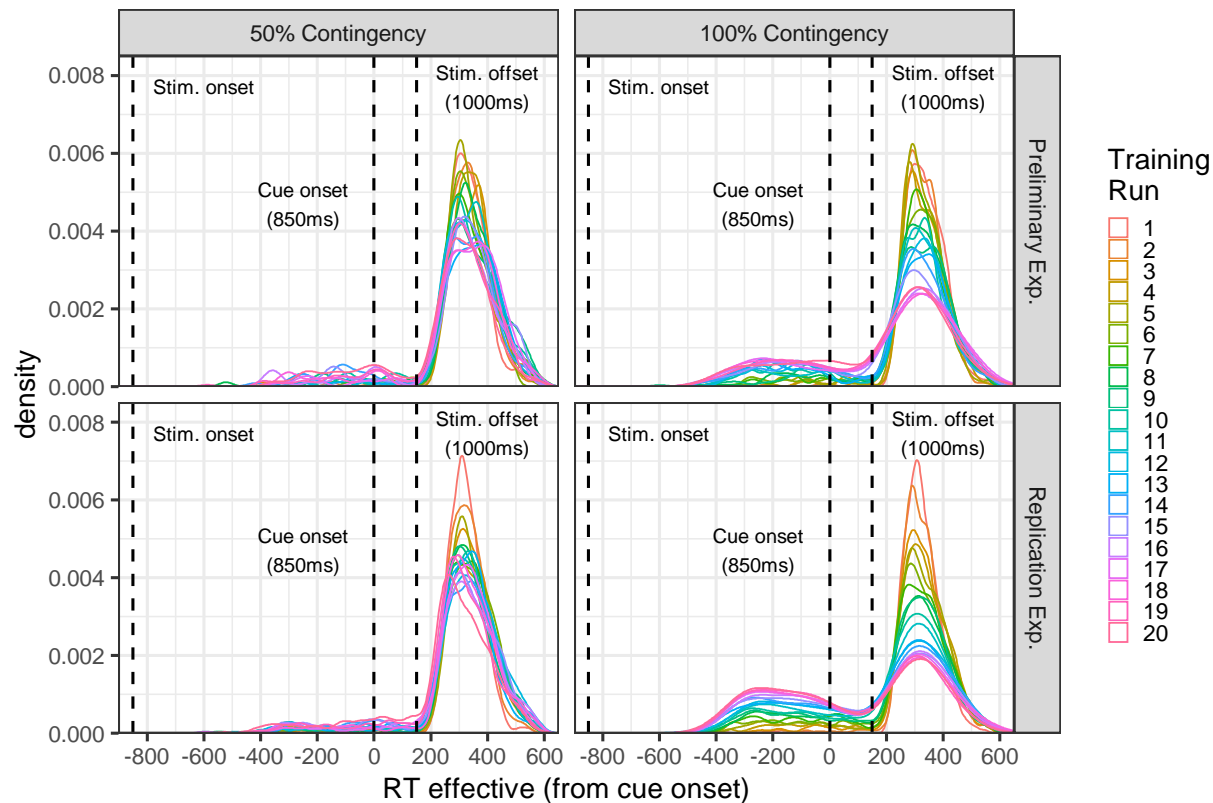Computational marker of non-reinforced learning



Figure 5. Density plots of effective RT distribution (time-locked to cue onset) in Study 2 (preliminary experiment on top row and replication experiment on bottom row). The Go cue onset was fixed at 850ms from stimulus onset - vertical dashed lines represents the trial onset (time 0), cue-onset, and trial offset (850ms and 1000ms from trial onset, respectively). As training progressed (indicated by the training run in color), responses were less homogeneous following cue onset, with increased proportion of anticipatory responses. The transition to anticipatory responses was more robust in the 100% contingency condition (right column) compared with the 50% contingency condition (left column).

*Computational marker of learning.* As in the meta-analysis study, RTs within each condition were modeled with a time-dependent mixture model, implemented with Stan Bayesian framework. RT data within each condition (100% contingency and 50% contingency) were fitted using two Gaussian distribution parameters (two means and two SDs). A $\theta_{slope_i}$ time dependent parameter was modeled for each condition and determined the rate of transition from cue-dependent responses to anticipatory responses of individual participants. Since the Go cue in Study 2 was fixed at 850ms, the effective RT (relative to stimulus rather than cue) measurement perfectly correlated with actual RT; thus, we used only the actual RT in the model.

In both the preliminary and replication experiment, the computational model converged around similar hyperparameter estimates (see MCMC trace plot and posterior simulations fit in

21

544  Supplementary Fig. S6 and Supplementary Fig. S7, respectively). Cue-dependent responses were

545  characterized with later mean RT and smaller variance, compared with earlier anticipatory

546  responses (preliminary experiment: $\mu_1 = 735.07$, 95%CI [677.14, 792.87], $\sigma_{\varepsilon_1} = 140.14$, 95% CI

547  [130.13, 153.49], $\mu_2 = 1190.23$, 95% CI [1188.46, 1192.03], $\sigma_{\varepsilon_2} = 78.11$, 95%CI [76.78, 79.47];

548  replication experiment: $\mu_1 = 740.44$, 95% CI [710.10, 771.76], $\sigma_{\varepsilon_1} = 140.35$, 95%CI [135.79,

549  144.74], $\mu_2 = 1187.46$, 95% CI [1186.25, 1188.68], $\sigma_{\varepsilon_2} = 80.30$, 95% CI [79.46, 81.15]). As the

550  cue onset was fixed in Study 2 at 850ms, the mean cue-dependent RTs of approximately 1190ms

551  corresponded with an effective RTs center of approximately 340ms in both experiments.

552  Anticipatory RT mean varied considerably between participants, as identified by the $\sigma_{\mu_1}$

553  parameter, which modeled between-participant variability in $\mu_{1_i}$ parameter estimates

554  (preliminary experiment: $\sigma_{\mu_1} = 110.81$, 95%CI [73.29, 166.69]; replication experiment: $\sigma_{\mu_1} =$

555  110.99, 95%CI [90.92, 136.34]).

556      Transition from cue-dependent to anticipatory responses, as captured by $\theta_{slope_i}$

557  parameter estimates of the preliminary experiment, indicated faster transition in the 100%

558  contingency condition, compared with the 50% contingency condition ($\theta_{slope100\%} = 1.95$,

559  95%CI [1.03, 2.75], $\theta_{slope50\%} = 0.90$, 95%CI [-0.01, 1.68]; mean difference in $\theta_{slopes} = 1.04$,

560  95%CI = [-0.17, 2.23]), with prominent variability between participants ($\sigma_{\theta_{slope100\%}} = 1.98$,

561  95%CI [1.45, 2.67], $\sigma_{\theta_{slope50\%}} = 1.79$, 95%CI [1.21, 2.57]). This effect was replicated even more

562  distinctly in the larger replication experiment ($\theta_{slope100\%} = 3.25$, 95%CI [2.54, 3.89], $\theta_{slope50\%}$

563  = 1.54, 95%CI [1.11, 1.92]; mean difference in $\theta_{slopes} = 1.70$, 95%CI = [0.90, 2.47]; $\sigma_{\theta_{slope100\%}}$

564  = 2.82, 95%CI [2.37, 3.42], $\sigma_{\theta_{slope50\%}} = 1.45$, 95%CI [1.16, 1.82]). The 95% credible interval of

565  the difference between the two conditions $\theta_{slope_i}$ parameters indicated a general trend in the

566  preliminary experiment which was distinct in the replication experiment.

567      ***Preference modification in probe task.*** Based on the results of the meta-analysis in Study

568  1, we hypothesized that manipulating the Go cue contingency would correspondingly manipulate

569  the preference modification effect, manifested as more pronounced preference modification for

570  stimuli in the 100% contingency condition, compared with the 50% contingency condition.

571      As expected, in the preliminary experiment, participants chose Go stimuli over NoGo

572  stimuli above chance level, both in the 50% contingency condition (prop. = 57.11%, $Z = 2.85$, *p*

573    = 0.002, odds = 1.33, 95% CI [1.09, 1.62]; one-sided logistic mixed model), as well as in the

574    100% contingency condition (prop. = 64.45%, $Z = 3.40$, $p = 3.4E^{-4}$, odds = 1.81, 95% CI [1.29,

575    2.56]). The preference modification effect was more robust in the 100% contingency condition,

576    compared with the 50% contingency condition ($Z = 1.92$, $p = 0.028$, OR = 1.36, 95% CI = [0.99,

577    1.87]). These effects were prominently replicated in the larger replication experiment (50%

578    contingency condition: prop. = 58.41%, $Z = 5.50$, $p = 1.9E^{-8}$, odds = 1.40, 95% CI [1.24, 1.59];

579    100% contingency condition: prop. = 66.31%, $Z = 6.41$, $p = 7.1E^{-11}$, odds = 1.97, 95% CI [1.60,

580    2.42]; condition difference-effect: $Z = 3.81$, $p = 7.0E^{-5}$, odds = 1.40, 95% CI [1.18, 1.67]; one-

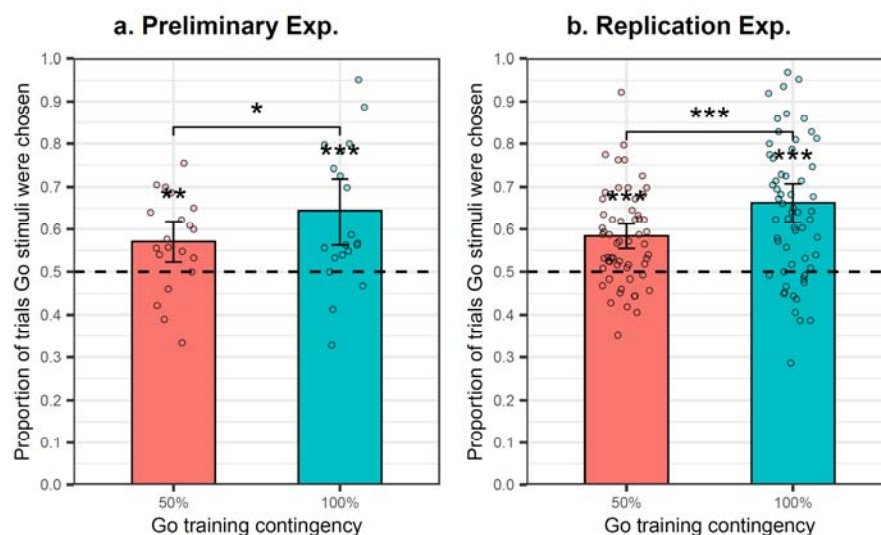581    sided logistic mixed model), see Figure 6.



582
583    Figure 6. Probe results. Proportion of trials participants chose Go stimuli over NoGo stimuli of
584    similar initial value, in the preliminary experiment (a) and replication experiment (b). Dots
585    represent individual participants, error-bars represents 95% CI based on a mixed model logistic
586    regression. Participants demonstrated enhance preference both for Go stimuli in the 50%
587    contingency condition and to a larger extent in the 100% condition. Statistical significance is
588    denoted with asterisks (* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$; one-sided mixed model logistic
589    regression). Dashed line represents 50% chance level.
590
591        We also examined the impact of initial subjective value on the reported effects by

592    controlling for the initial value difference within each probe choice, as well as by examining the

593    interaction of contingency effect with value categories. In all of the analyses, our effects of

594    interest remained consisted, i.e., we observed overall enhanced preferences for Go over NoGo

595    stimuli, and a more robust preference for Go stimuli within the 100% contingency condition (see

596    supplementary materials).

597   ***Prediction of probe using computational model.*** Most importantly, we aimed to examine

598   whether the preference modification effect could be predicted using the individually fitted

599   $\theta_{slope_i}$ parameter estimated for each participant. Adding the individual learning parameter as an

600   additional independent variable to the logistic regression model revealed that indeed $\theta_{slope_i}$

601   parameter estimated from the training task, could be used to predict preference modification in

602   the subsequent probe task (Figure 7). In the preliminary experiment a significant contribution

603   was found for $\theta_{slope_i}$ in the 50% contingency condition (log-OR = 0.12, $Z$ = 2.30, $p$ = 0.011, OR

604   = 1.13, 95% CI [1.02, 1.26]; one-sided mixed model logistic regression), as well as in the 100%

605   contingency condition (log-OR = 0.26, $Z$ = 3.59, $p$ = 1.7E$^{-4}$, OR = 1.29, 95% CI [1.12, 1.49]). No

606   significant difference was found between the slopes under the two conditions (log-OR = 0.13, $Z$

607   = 1.52, $p$ = 0.13, OR = 1.14, 95% CI [0.96, 1.36]; two-sided mixed model logistic regression).

608   These results were replicated also in the larger replication experiment, where $\theta_{slope_i}$ parameter

609   estimates predicted the preference modification effect, both in the 50% contingency condition

610   choices (log-OR = 0.10, $Z$ = 2.39, $p$ = 0.008, OR = 1.11, 95% CI [1.02, 1.20]; one-sided mixed

611   model logistic regression), as well as in the 100% contingency condition (log-OR = 0.11, $Z$ =

612   3.49, $p$ = 2.4E$^{-4}$, OR = 1.12, 95% CI [1.05, 1.19]). No differential slope effect was found

613   between the two conditions (log-OR = 0.01, $Z$ = 0.26, $p$ = 0.79, OR = 1.01, 95% CI [0.93, 1.10];

614   two-sided mixed model logistic regression).

615



616
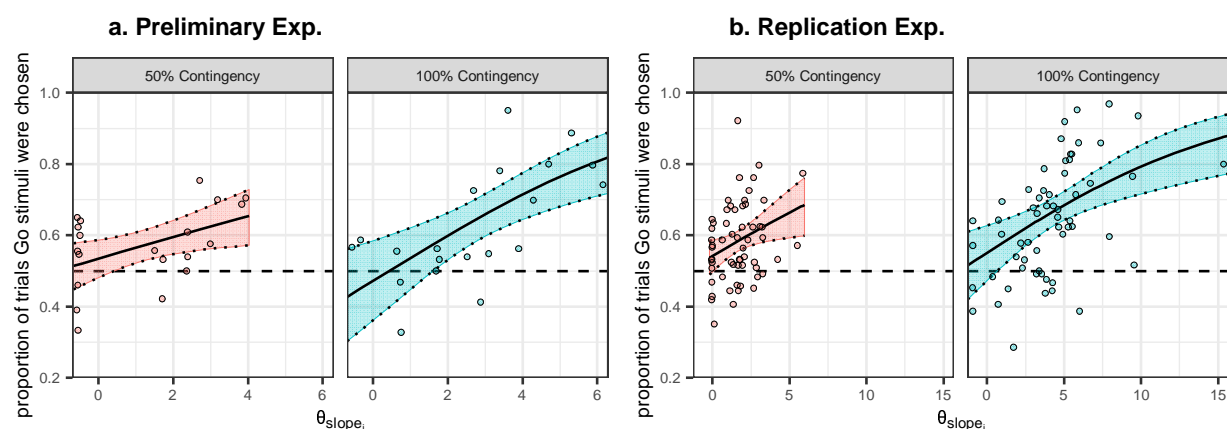617   Figure 7. Association of probe choices with training computational marker in study 2. Larger
618   $\theta_{slope_i}$ estimates parameter estimates were positively associated with subsequent preference
619   modification in probe phase, both in the preliminary experiment (a) and in the larger replication
620   experiment (b). The 100% contingency condition was characterized with larger $\theta_{slope_i}$ values,

621 and respectively also stronger preference modification effect. Trend line and surrounding color
622 margins represent estimated preference modification effect and 95% CI, respectively (mixed
623 model logistic regression). Dots represent individual participants. Horizontal dashed line
624 represents 50% chance level.
625
626       Examining the sizes of estimated fixed-effect, random-effect and residual variances,

627 revealed that the model accounted for a large portion of the variance in participants' choice

628 pattern in both experiments, both when examining the contribution of the fixed effects only

629 (preliminary exp.: $R^2_{GLMM(m)}$ = 0.421; replication exp.: $R^2_{GLMM(m)}$ = 0.223), and furthermore

630 when examining the joint contribution of fixed and random effects (preliminary exp.: $R^2_{GLMM(c)}$ =

631 0.869; replication exp.: $R^2_{GLMM(c)}$ = 0.872).

632       To examine the overall effect of contingency condition above and beyond $\theta_{slope_i}$, using a

633 likelihood ratio test, the logistic regression model was compared with a restricted (nested) model

634 that did not contain contingency fixed effects (contingency intercept and slope interaction;

635 restriction of 2 parameters). In both experiments, comparing the two models revealed no

636 significant contribution of contingency above and beyond $\theta_{slope_i}$ (preliminary experiment:

637 $\Delta AIC_{(rest.-full)}$ = -1.64, $\chi^2_{(2)}$ = 2.36, $p$ = 0.307; replication experiment: $\Delta AIC_{(rest.-full)}$ = -3.61,

638 $\chi^2_{(2)}$ = 0.39, $p$ = 0.824; two-sided likelihood ratio test).

639       *Post-hoc model validation.* Like in Study 1, in a post-hoc analysis we fitted a non-

640 Gaussian model (log-normal distributions) and found similar fit to actual RT data (see

641 Supplementary Fig. 16), as well as positive association of the $\theta_{slope_i}$ parameter estimates with

642 subsequent probe choice both in preliminary experiment (50% condition: log-OR = 0.03, $Z$ =

643 2.45, $p$ = 0.024, OR = 1.03, 95% CI [1.02, 1.04]; 100% condition: log-OR = 0.04, $Z$ = 3.75, $p$ =

644 0.001, OR = 1.04, 95% CI [1.03, 1.51]; one-sided mixed model logistic regression), and in the

645 replication study in 50% condition (50% condition: log-OR = 0.014, $Z$ = 2.12, $p$ = 0.038, OR =

646 1.01, 95% CI [1.008, 1.02]; one-sided mixed model logistic regression). We also performed a

647 post-hoc analysis in which we compared the explanatory power of $\theta_{slope_i}$ with a simpler marker

648 using the proportion of anticipatory responses each participant made (i.e. the proportion of RTs

649 which were faster than the top 1% of RTs in the first run; see Supplementary). In contrast to the

650 previous meta-analysis results, in the two Study 2 experiments, we found that $\theta_{slope_i}$ had a

651   significant predictive power, above and beyond a simpler marker which is based on the

652   proportion of anticipatory responses participants made during training.

653   Examining the interplay with choice RT. Finally, to examine the hypothesis that $\theta_{slope_i}$

654   impacts choices via faster automated response to Go stimuli, we introduced two new exploratory

655   analysis which examined whether $\theta_{slope_i}$ correlated with faster RT during probe. We found no

656   evidence supporting that probe choices were associated with choice RT during probe

657   (preliminary experiment - 50% contingency choices: $b = 4.56$, $t(18) = 0.18$, $p = 0.859$, 100%

658   contingency: $b = 9.27$, $t(18) = 0.42$, $p = 0.682$; replication experiment - 50% contingency

659   choices: $b = 1.18$, $t(57) = 0.09$, $p = 0.928$, 100% contingency: $b = 3.79$, $t(57) = 0.57$, $p = 0.568$;

660   two-sided linear mixed model).

661   We further examined whether the predictive power of $\theta_{slope_i}$ remained consistent above

662   and beyond the time it took participants to make their probe choices. We introduced to the

663   logistic regression analysis used to examine the predictive power of $\theta_{slope_i}$ an additional

664   regressor of the choice RT (in seconds). We found that choice RT often had a significant

665   explanatory power, wherein faster choices were associated with higher likelihood of choosing to

666   Go stimuli (preliminary experiment – 50% contingency: log-OR = 0.35, $Z = 1.31$, $p = 0.19$, OR =

667   1.42, 95% CI [0.84, 2.40], 100% contingency: log-OR = -0.72, $Z = -2.33$, $p = 0.019$, OR = 0.49,

668   95% CI [0.26, 0.89]; replication experiment – 50% contingency: log-OR = -0.58, $Z = -3.50$, $p =$

669   $4.6E^{-4}$, OR = 0.56, 95% CI [0.41, 0.78], 100% contingency: log-OR = -1.01, $Z = -5.59$, $p = 2.2^{-8}$,

670   OR = 0.36, 95% CI [0.26, 0.52]; two-sided mixed logistic regression). Nonetheless $\theta_{slope_i}$ was

671   consistently found to be predictive of choosing Go stimuli, above and beyond choice RT

672   (preliminary experiment – 50% contingency: log-OR = 0.12, $Z = 2.44$, $p = 0.007$, OR = 1.13,

673   95% CI [1.02, 1.25], 100% contingency: log-OR = 0.26, $Z = 3.36$, $p = 3.8^{-4}$, OR = 1.30, 95% CI

674   [1.12, 1.52]; replication experiment – 50% contingency: log-OR = 0.1, $Z = 2.42$, $p = 0.008$, OR =

675   1.11, 95% CI [1.02, 1.21], 100% contingency: log-OR = 0.12, $Z = 3.53$, $p = 2.0^{-4}$, OR = 1.12,

676   95% CI [1.05, 1.20]).

677   ***Stimulus-specific learning parameter.*** In an additional analysis (mentioned but not

678   detailed in the preregistration), we extended this analysis to stimulus-level effects. Each Go

679   stimulus$_s$ presented to participant$_i$ was fitted a $\theta_{slope_{i,s}}$ parameter, which was expected to capture

680   within-participant variability in learning - i.e., model for difference in learning of stimuli within

681 a contingency condition. We aimed to use these participant and stimulus-specific learning

682 parameters as independent variables in a mixed model logistic regression (with a random

683 intercept, random slope for contingency and random slope for $\theta_{slope_{i,s}}$ independent variables, see

684 methods for full description of the statistical models).

685      Fitting a full Stan model for both stimuli of 50% and 100% contingency simultaneously

686 did not converge. We hypothesized that the full model, with 2 conditions $\times$ 16 stimuli $\times$ n

687 participants inter-dependent $\theta_{slope_{i,s}}$ parameters was too complex to resolve using MCMC

688 process. Thus, we decided to reduce complexity by eliminating the participant-level effects on

689 the mean RT for the early anticipatory responses (estimating only group-level $\mu_1$, without

690 participant-level $\mu_{1_i}$ parameters; as was done in the previous meta-analysis study). We further

691 split the data based on contingency condition and analyzing each contingency as an independent

692 dataset. Using this approach resolved the convergence issues. This solution for an unexpected

693 issue was not preregistered or planned prior to data analysis.

694      Examining the association between the $\theta_{slope_{i,s}}$ parameter estimates and the participant-

695 level $\theta_{slope_i}$ parameter estimates of the previous model, revealed high similarity (preliminary

696 exp.: $r = 0.84$; replication exp.: $r = 0.80$). In addition, participants with negative $\theta_{slope_i}$

697 parameter estimates had very small variance (SD $< 0.1$) in their $\theta_{slope_{i,s}}$ estimates (see

698 Supplementary Fig. S8). When such $\theta_{slope_{i,s}}$ were introduced in the mixed-model logistic

699 regression, the low within-participant variability caused convergence warnings. Therefore, all

700 participants with $\theta_{slope_i} < 0.2$ (which is equivalent to SD $< 0.1$ of $\theta_{slope_{i,s}}$) in either contingency

701 condition, were excluded from the analysis. This resulted in the exclusion of nine participants in

702 the preliminary experiment and 14 participants in the replication experiment.

703      Analyzing the remaining 11 and 45 participants in the preliminary and replication

704 experiment, resulted in similar conclusions, as found with the participant-level $\theta_{slope_i}$ parameter

705 estimate. Overall, in both contingency conditions, a significant positive association was found

706 between $\theta_{slope_{i,s}}$ and preference of Go stimuli over NoGo stimuli (preliminary experiment - 50%

707 contingency: $Z = 4.34$, $p = 7.0\text{E}^{-6}$, OR $= 1.33$, 95% CI [1.17, 1.52]; 100% contingency: $Z = 2.85$,

708 $p = 0.002$, OR $= 1.21$, 95% CI [1.06, 1.38], overall model $R^2_{GLMM(m)} = 0.269$, $R^2_{GLMM(c)} = 0.631$;

709 replication experiment - 50% contingency: $Z = 3.14$, $p = 8.4\text{E}^{-4}$, OR $= 1.09$, 95% CI [1.03, 1.15];

710 100% contingency: $Z = 5.09$, $p = 1.7\text{E}^{-7}$, OR $= 1.09$, 95% CI [1.05, 1.12], overall model

Computational marker of non-reinforced learning

711     $R^2_{GLMM(m)} = 0.147$, $R^2_{GLMM(c)} = 0.475$; one-sided mixed model logistic regression; Figure 8). No

712     significant differences were found between the slopes of the two conditions (preliminary

713     experiment: $Z = -1.35$, $p = 0.18$, $OR = 0.91$, 95% CI [0.79, 1.04]; replication experiment: $Z = -$

714     0.08, $p = 0.94$, $OR = 1.00$, 95% CI [0.95, 1.05]; two-sided mixed model logistic regression).
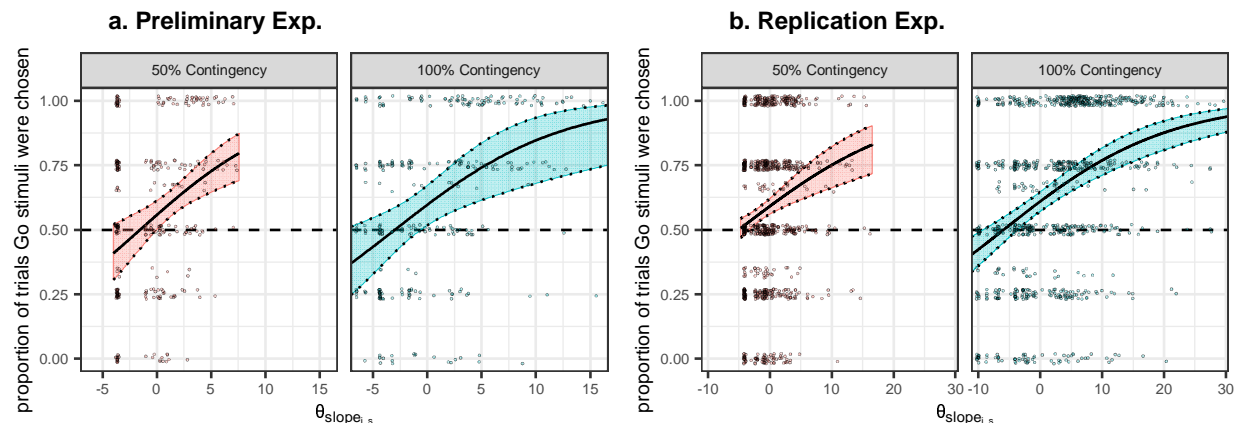
715



716
717     Figure 8. Association of probe choices with stimulus-level computational marker in study 2.
718     Larger $\theta_{slope_{i,s}}$ estimates parameter estimates were positively associated with subsequent
719     preference modification in probe phase, both in the preliminary experiment (a) and in the larger
720     replication experiment (b). Dots represent individual stimuli (small vertical jitter was added to
721     better visualize high density areas). Trend lines and surrounding color margins represent
722     estimated preference modification effect and 95% CI, respectively (mixed model logistic
723     regression). Horizontal dashed line represents 50% chance level.
724

725     The contribution of stimulus-level $\theta_{slope_{i,s}}$ parameter to predicting preferences was

726     examined by comparing a full model, containing both stimulus-level $\theta_{slope_{i,s}}$ and participant-

727     level $\theta_{slope_i}$ individualized learning parameters to a restricted model containing only $\theta_{slope_i}$ as

728     independent variable. The addition of $\theta_{slope_{i,s}}$ had a significant contribution both in the

729     preliminary experiment ($\Delta AIC = 59.65$, $\chi^2_{(9)} = 77.65$, $p = 4.7\text{E}^{-13}$; likelihood ratio test), as well as

730     in the replication experiment ($\Delta AIC = 93.97$, $\chi^2_{(9)} = 111.97$, $p = 5.8\text{E}^{-20}$).

731     **Interim discussion**

732     Based on the conclusions of the meta-analysis study, a novel preference modification

733     procedure was designed and tested in Study 2. In the new design, two cue contingency

734     conditions (50% and 100% contingency) were tested within-participant, which we hypothesize

735     would manipulate the difficulty of learning the stimulus-cue association during the training phase

736     of the CAT paradigm. The introduction of the two contingency conditions allowed us to

737  experimentally impact the $\theta_{slope_i}$ individualized computational marker of learning, which in turn

738  provided a reliable predictor for the subsequent preference modification effect in the probe

739  phase. As expected, we found a more robust preference modification effect for Go stimuli in the

740  100% contingency condition, compared with the 50% contingency condition. The fact that a

741  directed experimental intervention in the preceding training phase induced a differential effect in

742  the probe phase, provides support for a directional causal impact of training on preference

743  change. In light of these results, it is plausible to conclude that the association demonstrated in

744  Study 1 between the individualized learning parameter and behavioral change might represents a

745  causal relationship, in which improved learning in the training phase induces stronger behavioral

746  change in the subsequent probe phase.

747  Furthermore, between-participant variability in the $\theta_{slope_i}$ computational marker was also

748  predictive of between-participant difference in preference modification, as measured during the

749  subsequent probe phase. Participants with more robust $\theta_{slope_i}$ also demonstrated stronger

750  preference for Go stimuli over NoGo stimuli of similar initial value. This effect was observed

751  both in the 100% and 50% contingency conditions, with no significant difference in the effect

752  between the two conditions. This suggests that $\theta_{slope_i}$ had a similar predictive pattern for both

753  conditions, one unit increase in $\theta_{slope_i}$ resulted in similar increase in the odds of choosing Go

754  over NoGo, above and beyond the contingency condition. However, as the 100% contingency

755  condition was characterized with larger $\theta_{slope_i}$ parameter estimates, so was the probe phase

756  characterized with greater odds of choosing Go stimuli (i.e., stronger preference modification

757  effect). Taken together, these results provide evidence for novel means to both predict individual

758  differences in learning and preference modification between participants, as well as means to

759  manipulate this preference modification effect. All the non-exploratory results in the preliminary

760  experiment were carefully detailed in a preregistration, and were all replicated in a larger

761  replication experiment, which included the exact same procedure and analysis pipelines. It is

762  important to note that using a larger sample size not only replicated the significance of the

763  effects, but also the descriptive trends, resulting in outcomes of similar effect size.

764  In contrast to post-hoc analysis in Study 1, which found that the computational marker

765  performed similarly well as a simple RT-based marker, in a post-hoc analysis of the two new

766  experiments, we did find that $\theta_{slope_i}$ captured a unique predictive signal above and beyond a

29

767   simple proportion of anticipatory responses. These results reinforce our hypothesis, that by

768   optimizing the training procedure for our computational model (by fixing the Go cue onset time

769   and explicitly encouraging participant to make anticipatory responses), we were able to better

770   capture learning signals using $\theta_{slope_i}$ marker.

771        It is interesting to note that descriptively the new design induced a strong preference

772   modification effect in the 100% Go contingency (group estimate of approximately 65% in both

773   experiments), compared with other CAT and Go/NoGo experiments, were a more modest effect

774   of 55%-60% preference bias for Go stimuli was commonly observed [12–20]. Without carefully

775   controlling for all factors, it is hard to conclude whether this difference is significant, and which

776   one of the changes in the experimental design induced this enhanced effect, however, it would be

777   interesting to try and address this directly in future work.

778        In an exploratory analysis, we examined variation across stimuli as well as participants,

779   focusing both on the variability in learning patterns between different participants, as well as

780   within-participant variability in learning for different stimuli. Due to the technical requirements,

781   such an in-depth analysis was applicable only for a smaller subset of our data, however, the

782   partial results we found suggest that the prediction model is applicable in all levels of

783   granularity, both at the individual participant level, as well as in the individual stimulus level.

784

785                              **General discussion**

786        In the current work, we aimed to study and characterize the cognitive mechanisms

787   underlying non-reinforced preference modification with CAT[12], focusing on the training phase

788   where individual items are presented. While behavioral change is manifested in a binary probe

789   phase, actual preference change occurs at the individual item level. Based on previous studies

790   that alluded to the potential involvement of internal non-reinforced learning in CAT, the current

791   work aimed to use computational modeling to measure this internal process at the training phase

792   with individual items and establish its causal role on the subsequent preference modification

793   effect. The CAT task is a multiphase standardized procedure, which includes an initial

794   preference evaluation task, a non-reinforced training phase on individual items, and a binary

795   probe phase. Therefore, the task provided a unique opportunity to aggregate data from multiple

796   studies and build a large corpus of training data.

797        In the first part of the current work, we examined previously collected data of n=864

798    participants from 29 different CAT experiments in a meta-analysis and identified a distinct time-

799    dependent RT pattern during the training task. We found that while in early stages of training,

800    participants depended on the Go cue to initialize responses (homogeneous RTs following cue

801    onset), as training progressed, participants relied less on the Go cue and started to generate fast

802    anticipatory responses. This raised the hypothesis that a marker of transition from cue-dependent

803    responses to early anticipatory responses could be indicative of internal non-reinforced learning

804    and subsequent preference changes.

805        Using a Bayesian computational framework implemented with Stan, a statistical model of

806    RT patterns was formulated. RTs were postulated to derive from a mixture model of two

807    Gaussian distributions and a time-dependent mixture proportion, which accounted for the

808    transition speed from late cue-dependent responses to early anticipatory responses. The model

809    parameters were estimated with MCMC algorithm, with a key participant-level $\theta_{slope_i}$

810    parameter, which modeled the individualized difference between participants in transition to

811    anticipatory responses rate. This Bayesian parameter was used as a passive computational

812    marker of learning.

813        Examining the association between the individualized learning parameter and preference

814    modification effect in the subsequent probe task, showed a positive association in which

815    participants with more robust learning parameter estimates during training also demonstrated

816    stronger preference modification effect in the subsequent probe phase (i.e., enhanced preference

817    for the Go stimuli over NoGo stimuli of similar initial value). Thus the $\theta_{slope_i}$ parameter, which

818    was optimize to model RT pattern during training, was found to be good predictor of preference

819    modification effect in a future probe task. An additional analysis modeling RT using a more

820    finely tuned learning $\theta_{slope_{i,s}}$ parameter, fitted for the individual stimuli, further expended these

821    findings, and demonstrated positive predictive power of $\theta_{slope_{i,s}}$ in evaluating preference

822    modification effect on a stimulus-level basis, above and beyond the more general participant-

823    level parameter.

824        Based on the results of Study 1, we derived a theoretical hypothesis with implication for

825    shaping learning efficacy and consequently also the effect of training on preference behavior. If

826    faster transition to anticipatory responses is related to learning, then hindering the ease of cue

827 anticipation could also obstruct learning and resulting in less pronounces preference modification

828 effect. To examine this hypothesis, we design a novel CAT procedure, which was tested in a

829 preliminary experiment and an additional larger preregistered replication experiment. Cue

830 anticipation was manipulated by introducing two training conditions within the CAT task - in the

831 100%-contingency condition, the Go stimuli perfectly anticipated Go cue onset, while in the

832 50% contingency condition, the Go stimuli only anticipated the Go cue in half of the trials.

833 Manipulating cue-contingency during CAT induced the expected effects both in the

834 training and in the probe task. During training, participant relied less on early anticipatory

835 responses and transitioned slower to anticipatory responses in the 50% contingency condition.

836 This effect was captured by the $\theta_{slope_i}$ estimates, which were larger in the 100% contingency

837 condition compared with the 50% contingency condition. More importantly, in the subsequent

838 probe phase, a consistent pattern of stronger preference modification effect was observed for the

839 100% contingency Go stimuli. Thus, manipulating training difficulty has also affected future

840 preference modification measured in the probe phase. When individual differences in preference

841 modification were examined, $\theta_{slope_i}$ marker accounted for a large portion of variability in

842 choices, while no significant difference was found between the two contingency conditions

843 above and beyond the $\theta_{slope_i}$ effect. This suggest that the more robust preference modification

844 effect observed for the 100% contingency Go stimuli, was captured by the stark differences in

845 $\theta_{slope_i}$ learning marker. Furthermore, the fact that intentionally manipulating the training

846 procedure in Study 2 induced differential learning parameter estimates as well as differential

847 behavioral change effect, further supports the hypothesis of a causal nature of the association

848 with the behavioral change. Thus, it is plausible to deduce that shaping training at the individual

849 item level, which was captured by the individualized learning parameter, has driven the induced

850 differential behavior modification effect, as evaluated in the subsequent probe phase. It is

851 important to note that although the cue-contingency was different between the two conditions,

852 the actual exposure-time to the stimuli was identical. Thus, the enhanced preference for Go

853 stimuli in the 100% contingency condition could not be accounted for by mere exposure effect [6].

854

855 **What is the mechanism?**

856 These findings, showing a consistent link between motor-planning during training at the

857 individual items level and preference modification in the subsequent binary choice probe phase,

858    correspond with previous works. These works showed that a rapid response is a crucial feature

859    for preference modification with CAT[13], and neural finding showing that increased striatal and

860    premotor activity during training, were associated with more robust preference modification

861    effect in the probe phase[18]. Thus, the results suggest that learning efficacy in the training phase is

862    manifested as stimulus-specific motor planning. Importantly, all participants were alert and

863    attentive throughout the training phase, as demonstrated by the negligible rates of non-response.

864    However, it evident that mere attention is not sufficient to induce strong preference modification,

865    but rather an attention that can be translated to action, suggesting a unique valuation pathway in

866    the absence of external reinforcements, putatively based on parieto-frontal circuits involving

867    attention with motor planning [5].

868    To produce an anticipatory response for a Go stimulus, participants must also rely on

869    memory. Thus, it is possible that the anticipatory response pattern captured in our data identified

870    individualized difference in memory which in turn mediated a change in preferences[23]. This

871    hypothesis is in line with previous findings with non-reinforced training, which identified a

872    positive association between preferences modification and enhanced memory both at the

873    participant- level [15] and at the stimulus- level[19,21]. Since memory for the trained Go stimuli was

874    not examined in the current work, it would be interesting to examine in future experiments

875    whether $\theta_{slope_i}$ computational marker correlates with declarative memory and does it provide

876    independent predictive power of subsequent non-reinforced preference modification above and

877    beyond memory measurements. Future imaging studies could also examine the new task design

878    with fMRI to identify whether differential neural activation patterns characterize the two

879    contingency conditions.

880    One question that is often raised in non-reinforced preference modification research is

881    whether the behavioral effect of choosing Go stimuli reflects an internal change in the value of

882    the chosen Go stimuli, or rather habitual automated responses to choose the Go stimuli

883    originating from mechanisms similar to operant conditioning. On the one hand, CAT

884    experiments showed that not all cue-response associations result in enhanced preference for Go

885    stimuli – CAT was found to be less effective in enhancing preferences for stimuli of negative

886    affective valence [15] and low-value snack food items [12], suggesting that the stimulus pretraining

887    value interacts with the effect of CAT on preferences. In addition non-challenging CAT designs

888    such as CAT where the cue starts with the stimulus onset or when the response association was

33

889   made as a block of stimuli to which participants were required to respond also failed to induce

890   preference change[13], which provide evidence that the mechanism impacting preferences requires

891   more than a simple motor response association.  On the other hand, some findings elude to the

892   involvement of rapid impulsive decision making mechanisms - CAT experiments in which

893   participants were asked to make slow well considered decisions found that CAT effect on

894   preferences diminished when participants made slow non-impulsive choices. This result further

895   resonates the finding that during the CAT probe phase, participants were faster in their choices

896   when Go stimuli were chosen [17,41].

897   To test these two competing notions that CAT non-reinforced preference modification is

898   driven by impulsive motor response habitual mechanisms versus internal value representation

899   change, we examined in our data whether choice RT could account for the enhanced likelihood

900   of choosing Go over NoGo stimuli during probe. Our results showed most of the times that

901   indeed faster choice RTs were associated with increased likelihood of choosing Go stimuli, in

902   agreement with past findings. However, and more importantly, we found that the computational

903   marker was not correlated with choice RT, and maintained its predictive power of choice above

904   and beyond choice RT. Thus, our results indicate that while some aspects of non-reinforced

905   preference change could be attributed to rapid impulsive choices, our proposed computational

906   marker tracks an independent preference modification mechanism, which could putatively be

907   indicative of value representation change. Future work could attempt to find behavioral or neural

908   correlates which could corroborate this hypothesis.

909   An explanation offered by the model for the modification of preferences revolves around

910   the presence of an internal reinforcement process during the training phase. Following the initial

911   exposure to stimulus-cue contingencies, participants attempted to respond based on the stimuli

912   rather than the cue itself, displaying anticipatory responses. When participants correctly pressed

913   earlier, an internal feedback mechanism was activated, further enhancing the learning process.

914   This internal reinforcement can be seen as a broader mechanism for facilitating learning. It is

915   challenging to differentiate this internal reinforcer from memory or attention processes, and our

916   study does not aim to do so. Instead, the current model integrates these mechanisms and

917   incorporates the internal reinforcer as an integral component of the learning process. Drawing on

918   previous research and considering the correlation between brain regions associated with the

919   reward system [18], we propose that this inner mechanism plays a role in reinforcement. However,

920  since no external feedback was provided during the training task, our hypothesis suggests that
921  the reinforcement mechanism is likely to be internal in nature. The current work's unique
922  contribution establishing a new passive marker for non-reinforced learning could be tested in
923  future neuroimaging studies. These studies could shade light on whether the marker is associated
924  more with parietal attention mechanisms, temporal memory-related regions or striatal reward and
925  motor-learning neural mechanisms (ref).

926  In addition to a basic understanding of the cognitive mechanisms of non-reinforced
927  learning, identifying a learning marker based on early training runs could be used to design
928  powerful prediction tool for learning efficacy, with relevance to other paradigms and situations
929  that share the basic training/transfer structure of the current studies. A naïve approach to evaluate
930  learning-efficacy during training tasks such as CAT might suggest to introduce direct
931  measurements of behavior-change throughout the training process, e.g. by asking participants to
932  make active value-based choices [42,43]. However, it has been shown that the mere act of choice
933  could induce a longitudinal effect of enhanced preference for the chosen stimuli [8,9,44] and
934  introduced bias in the effect of preference modification following non-reinforced training [21].
935  Therefore, probing preference within a preference modification training, is likely to alter the
936  learning process and undermine its validity. Moreover, in more standard learning procedures,
937  such as conditioning-based interventions, to examine the efficacy of learning researchers observe
938  the subject's response to the conditioned stimulus when the associated unconditioned stimulus is
939  omitted [45–52]. Over repetition of probe tests in all of these procedures might initiate a new
940  learning procedure – where the conditioned stimulus is no longer associated with the
941  unconditioned stimulus, which would eventually lead to extinction- learning [53–56]. Thus,
942  evaluating learning by active probing of the behavior-change effect could interfere with the
943  learning process. However, a passive marker for learning which can predict behavioral change
944  efficacy based on earlier training data, does not evoke the drawbacks of introducing an active
945  choice task to reveal preferences. Using a passive learning marker which is evaluated based on
946  independent training data that preceded the probe phase, overcomes this obstacle. The temporal
947  primacy of $\theta_{slope_i}$ estimation could provide a prediction tool for future preference modification
948  in CAT, without tainting the results with a direct evaluating of preferences using a choice task.
949  An interesting hypothesis could anticipate that the 50% cue-contingency would resemble partial
950  reinforcement learning schedule, and thus will have potentially higher long term sustainability

35

951  and resistance to extinction [57]. While the current work did not examine long term maintenance of

952  CAT effect and its association cue-contingency, this could be an interesting subject to test in

953  future studies.

954        The benefits of real-time individualized markers for learning could also be of great value

955  for additional learning and behavior change procedures beyond CAT. In many learning tasks

956  (including CAT), the behavioral change effect on choices is measured in a separate probe phase,

957  proceeding the training phase. For example, in experiments testing pain perception, participants

958  are trained to associate neutral stimuli with differential level of painful heat stimulation before

959  probing the impact of interventions such as drug versus placebo administration [58,59]. Studies

960  examining habit formation may use lengthy free-operant learning protocol, in which participants

961  repeatedly perform an action (such as pressing a button) to gain food rewards, in order to test in a

962  later probe phase whether the participants demonstrate habitual behavior and continue to perform

963  the action even when the reward is devalued (e.g., pressing to get food when satiated) [60–62]. And

964  even outside the field of value-based decision making, in the clinical psychological wellbeing

965  domain, attention bias modification (ABM) procedures use computerized attention training

966  interventions similar in nature to CAT to treat depression and anxiety [63–65]. In such experimental

967  settings, where training is either unpleasant or exhausting, with plausible odds that learning

968  would not be well-established if under-trained, monitoring learning efficacy individually based

969  on the training data could assist in optimizing training efficacy and efficiency. While the current

970  work focuses on CAT, we assert that the basic mechanisms identified by the current work, could

971  be adapted to accommodate a need in a wider range of learning procedures. With an appropriate

972  adaptation to the desired experimental design, this unique feature of a passive computational

973  marker which is predictive of subsequent change, opens a new path for potential interventions

974  that will allow more efficient training via monitoring and real-time feedback [66,67].

975        The current work's approach in modeling learning using RT could also be applicable in

976  experimental designs which do not separate the probe from the training task, such as

977  reinforcement learning tasks, where learning can be evaluated as it progresses on a trial by trial

978  basis [68–71]. Since RT patterns could be indicative of the individual's confidence in her choice

979  during reinforcement learning tasks such as probabilistic selection [51,52] or multi-arm bandit task

980  [72,73], incorporating RT data could improve reinforcement learning models and capture learning

981  more accurately [74].

982    Some limitations were not directly tested in the current work and should be addressed in
983    future studies. Primarily, in Study 2 of the current work, we demonstrated that learning and
984    preference modification could be hindered using a more difficult/partial association procedure.
985    However, our model also predicts that easing the association procedure could enhance learning
986    efficacy and preference modification effect. Further work should empirically test this hypothesis,
987    for example by examining the effect of enhancing the Go cue saliency or attention during
988    predefined Go stimuli [27,75,76], or by manipulating the temporal features of training schedule [77].
989    Furthermore, by instructing participants to respond when they have enough confidence that a
990    Cue will follow, we might have introduced an undesired confound to the task in the subjectivity
991    of confidence each participant needed to respond; i.e., some participants might have learned well
992    the association between a stimulus and the Go cue but were reluctant to respond before they
993    validate that in a 50% contingency condition a cue will appear. Although such confound does not
994    undermine the validity of the findings, future studies should avoid it and aim to examine a
995    differential effect in which the difference between conditions reflect only task difficulty, for
996    example, by asking to respond to stimuli which at any point were associated with a cue, or by
997    clarifying in the instruction that a response based on guess is also valid.

998    Another limitation of the current work could raise a concern for a confound which might
999    undermine the ability to deduce causal relationship between learning and choice in our new
1000   experimental design. In two experiments, we found that participants demonstrated an enhanced
1001   preference modification effect for stimuli that were more consistently associated with Go cue. By
1002   maintaining a similar presentation time of both 50% and 100% contingency stimuli, we
1003   eliminated the confound of the standard mere exposure effect [6], as participant viewed the stimuli
1004   in both conditions for the same time duration and putatively had to maintain high alertness to
1005   both type of stimuli, which both required response with a high likelihood. However, one might
1006   argue that participants were more engaged with the 100% contingency stimuli, to which they
1007   pressed twice as much. Thus, an alternative explanation could claim that increased engagement
1008   (i.e., more press responses) is the true causal factor which induced a stronger preference
1009   modification effect, rather than the learning procedure identified by the computational marker.

1010   Our statistical analysis provides evidence that suggests that the learning marker had a
1011   stronger impact on choices than the level of exposure. When analyzing the factors predicting
1012   choice patterns, the training condition had no explanatory power above and beyond the

1013    computational marker of learning. Nonetheless, a dedicated experimental design could aim to

1014    directly discern between these two competing hypotheses. One such proposed training design

1015    could present the 50% contingency stimuli twice as often as 100% contingency stimuli. This

1016    would result in a training design with more consistent association for the 100% contingency

1017    condition (similarly to the current design), while maintaining equal engagement levels between

1018    the two conditions (manifested as identical number of press responses in both training

1019    conditions) and an increased exposure to the 50% contingency stimuli, which would be presented

1020    for twice the duration of time as the 100% contingency condition stimuli. Our theory

1021    hypothesizes that despite the increased exposure and similar engagement with the 50%

1022    contingency stimuli in this design, participants will still demonstrate faster learning in the

1023    consistent 100% contingency condition and would thus show stronger preference modification

1024    effect for the 100% contingency condition stimuli. Future work could try to run this dedicated

1025    design and provide empirical evidence which would settle whether our proposed model of

1026    learning overcome mere exposure effect combined with more equal engagement-level, which we

1027    did not fully address in the current work.

1028    We carefully documented the experimental choices in a preregistration before data was

1029    collected for an independent replication experiment. Utilizing a Bayesian computational

1030    framework, such as the one used here, provided great flexibility required to fit complicated

1031    theoretical model. We aimed to use a relatively straightforward model, with as similar features as

1032    possible in all studies. Future studies could take the liberty to use the openly accessible data from

1033    this work and try to improve the model by changing our current assumptions including

1034    distributions' shape, free parameters, and priors.

1035    In conclusion, the current work laid the foundations for understanding and quantifying

1036    non-externally reinforced learning at the individual item level via a novel Bayesian modeling

1037    approach of RT pattern during training. Using a large meta-analysis dataset of previous CAT

1038    studies and two new original experiments, we demonstrated a unique method to evaluate a

1039    learning marker at the individual item level with a robust predictive power of subsequent

1040    preference change. Thus, we propose that motor response patterns could provide a passive

1041    marker for value-change, which does not require direct measurement of preferences.

1042

1043    **Methods**

38

**Study 1 – CAT meta-analysis**

**Data collection and sample sizes.** For study 1 meta-analysis, the data from 29 previous CAT experiment was combined. The dataset included 21 experiments from published works by our research group and collaborators [12–16,18,19,22], as well as eight additional unpublished works, which were collected as part of the preparations for previously published manuscripts, or are planned to be published in the future. The experiments comprised of a median sample size of $n = 26$ ($M_{sample\ size} = 29.79$, range = [23, 70]; see Table 1 for information of the sample size in each experiment, and the supplementary data for more detailed description of the unpublished data).

Table 1. Experiment included in the meta-analysis.

| Exp. | Stimuli | n | Training runs | Go cue | Publication (exp. number in publication [a]) |
|------|---------|---|---------------|--------|----------------------------------------------|
| 1 | Fractal art | 25 | 12 | Auditory | Salomon et al., 2018 (2) |
| 2 | IAPS positive | 27 | 12 | Auditory | Salomon et al., 2018 (3) |
| 3 | IAPS negative | 28 | 12 | Auditory | Salomon et al., 2018 (4) |
| 4 | Snacks (IL) | 25 | 20 | Visual | Salomon et al., 2018 (5) |
| 5 | Snacks (IL) | 25 | 20 | Neg. auditory | Salomon et al., 2018 (6) |
| 6 | Faces | 25 | 20 | Auditory | Salomon et al., 2018 (7) |
| 7 | Fractals | 25 | 20 | Auditory | Salomon et al., 2018 (8) |
| 8 | IAPS positive | 29 | 20 | Visual | Salomon et al., 2018 (9) |
| 9 | IAPS negative | 32 | 20 | Visual | Salomon et al., 2018 (10) |
| 10 | Faces: politic. | 25 | 20 | Auditory | Unpublished |
| 11 | Faces: politic. | 39 | 20 | Auditory | Unpublished |
| 12 | Faces | 42 | 16 | Auditory | Salomon et al., 2019 |
| 13 | Faces: affective | 42 | 20 | Auditory | Unpublished |
| 14 | Faces: affective | 70 | 20 | Auditory | Unpublished |
| 15 | Fractals | 29 | 16 | Auditory | Aridan et al., 2019 |
| 16 | Snacks (US) | 26 | 12 | Visual | Unpublished |
| 17 | Snacks (IL) | 24 | 20 | Visual (Reward)[b] | Unpublished |
| 18 | Snacks (IL) | 30 | 16 | Auditory | Botvinik-Nezer, Bakkour, et al., 2021 (1) |
| 19 | Snacks (IL) | 25 | 16 | Auditory | Botvinik-Nezer, Bakkour, et al., 2021 (Pilot) |
| 20 | Snacks (IL) | 23 | 12 | Auditory | Unpublished |
| 21 | Snacks (IL) | 25 | 20 | Auditory | Unpublished |
| 22 | Snacks (US) | 29 | 12 | Auditory | Schonberg et al., 2014 (1) |
| 23 | Snacks (US) | 25 | 8 | Auditory | Schonberg et al., 2014 (2) |
| 24 | Snacks (US) | 25 | 12 | Auditory | Schonberg et al., 2014 (3) |
| 25 | Snacks (US) | 27 | 16 | Auditory | Schonberg et al., 2014 (4) |
| 26 | Snacks (US) | 26 | 16 | Auditory | Schonberg et al., 2014 (7) |
| 27 | Snacks (US) | 30 | 12 | Auditory | Bakkour et al., 2017 |
| 28 | Snacks (US) | 25 | 16 | Auditory | Bakkour et al., 2016 |
| 29 | Snacks (IL) | 36 | 16 | Auditory | Botvinik-Nezer, Salomon, et al., 2019 |

[a] In papers with multiple experiments, we note in parenthesis the experiment number or title as can be found in the related manuscript.
[b] In experiment 16 a visual cue was used indicating to participants they were awarded with a sum of money, which accumulated throughout the training phase.

All participants gave their informed consent to take part in the experiments. In most experiments, participants received monetary compensation for their time, and in few

experiments, some participants took part in the experiment in exchange for course credit. All experiments were approved by the ethical review board of the institutes where they were performed (Tel Aviv University, The Hebrew University of Jerusalem, University of Texas at Austin, and McGill University).

**Stimuli.** The different experiments included in Study 1 examined the effect of CAT on preferences for various stimuli (see table 1 for a summary of all experiments). In most of the experiments, the stimuli set comprised of images of familiar local snack food items, popular in the US (eight experiments) or in Israel (eight experiments), which participants received for actual consumption as part of the experiment. Other experiments used face stimuli of unfamiliar figures posing neutral expression from the Siblings dataset [40] (two experiments), unfamiliar faces with neutral and happy expression from the Karolinska directed emotional faces dataset [78] (two experiments) or familiar faces of famous Israeli politicians (two experiments). Unfamiliar abstract stimuli of fractal art [79] were used in three experiments. Two experiments included positive affective stimuli from the international affective picture system [80,81] dataset and two experiments included negative affective stimuli from the IAPS dataset.

The face and snack stimuli were modified using Photoshop to remove any background features and create visual standardization within each experiment. In each experiment, all stimuli were colored images of identical dimensions, where the subject of the image (the snack or the face) was centered in the middle of the image frame and the background was replaced with homogenous background (either black or gray). In experiments with fractal art and IAPS stimuli, the images were only cropped to identical pixel dimensions.

*Go stimuli used in the CAT task.* In most experiments (23 out of 29) a neutral auditory cue was used as Go cue during the cue-approach training task. In the remaining experiments a neutral visual cue (Experiments 4, 8, 9, and 16), an aversive auditory cue (Experiment 5), or a visual cue indicating a reward (Experiment 17), was used. Both types of visual cues (neutral and reward related) comprised of a semi-transparent Gabor shape, presented on top of the associated Go stimuli. In experiment 17, the appearance of the visual cue indicated to the participant that she or he had won an additional amount of money, which accumulated across the training task.

**Procedure.** While the different experiments diverged in several aspects of their procedure, all experiments maintained three key phases: an initial-preferences evaluation task,

1090    followed by a training task, and a probe task (see Figure 1 for illustration of the procedural

1091    design).

1092           ***Initial preferences evaluation task.*** In the first task of the experimental procedure,

1093    participants were exposed to the complete set of stimuli for the first time and were required to

1094    indicate their subjective preferences using one of two tasks. To evaluate participants' preferences

1095    in experiments using consumable snack food stimuli, participants performed a Becker-DeGroot-

1096    Marschak (BDM) auction procedure [39]. Participants were allocated with either 3 USD (in

1097    experiments conducted in the US) or 10 ILS (in experiments conducted in Israel; approximately

1098    equivalent to 3.1 USD) which were used to bid on the snack-food stimuli, presented one by one.

1099    Participants were informed that at the end of the experiment, one of the trials will be randomly

1100    selected, for which the computer will generate a counter bid. If the participants' bid was higher

1101    than that of the computer, they were required to purchase the snack for the lower price bided by

1102    the computer. Participants purchased the snack for actual consumption at the end of the

1103    experiment. However, if the computer's bid was higher, participants got to keep the allocated

1104    sum of money. Participant were explicitly instructed that the best strategy for the task is to

1105    indicate their true subjective preference. Prior to their participation, participants were asked to

1106    fast for at least 3 hours, to make sure they were hungry and incentivized to purchase the food

1107    items, according to their subjective preferences.

1108           In experiments with of non-consumable stimuli (such as fractals and faces), which are

1109    less appropriate to be evaluated using a monetary scale, preferences were evaluated using a

1110    binary ranking procedure. Participants were presented in each trial a random pair of stimuli and

1111    were required to choose the stimuli they prefer better. Based on the idea of choice transitivity,

1112    binary choices were quantified to produce subjective value ranks using the Colley Matrix

1113    ranking procedure [82].

1114           Following the BDM or binary choice task, stimuli were ranked-ordered according to

1115    subjective preferences. The ranks of the stimuli were used as a basis to form two value-groups -

1116    one of high-value stimuli (above-median rank) and a group of low-value stimuli (below median

1117    rank). The size of the value groups differed between experiments (having 8, 12 or 16 stimuli per

1118    value group). In each of the value-groups, half of the stimuli were allocated to be associate with

1119    the Go stimuli and response in the subsequent CAT task (Go stimuli), and the other half of

1120    stimuli was allocated to be presented in CAT without Go cue (NoGo stimuli).

1121    Allocation for Go and NoGo stimuli within each value category maintained an equal

1122 mean rank for Go and NoGo stimuli - e.g., for a high-value group consisting of eight stimuli, the

1123 stimuli ranked 8, 11, 12, and 15 were allocated to be Go stimuli, while stimuli ranked 9, 10, 13,

1124 14 were allocated to be NoGo stimuli, such that both allocations were characterized with a mean

1125 rank = 11.5. The Go / NoGo allocation was counterbalanced across participants.

1126    ***Cue Approach Training (CAT).*** In the CAT task, stimuli were presented individually on

1127 the screen center for a fixed duration of one second (except for Experiment 27, where the

1128 duration was extended to 1.2 seconds for compatibility with fMRI scanning protocol). In each

1129 experiment, all training stimuli were presented once in each training run, thus the number of

1130 training runs indicate the number of stimulus repetitions. A fixed proportion of stimuli per

1131 experiment (commonly 30% of stimuli; range 25%-40%) were Go stimuli. When a Go stimulus

1132 appeared, a delayed Go cue appeared after a Go signal delay. Participants were asked to respond

1133 to the Go cue with a button press as rapidly as possible, before the stimulus offset (one second

1134 after the stimulus onset). The Go signal delay was adjusted according to the participant's

1135 performance - a failure to respond before stimulus offset resulted in 50ms shortening of the Go

1136 signal delay (thus reducing the task difficulty for the following trial), while a successful response

1137 on time resulted in 16.667ms increase of the next Go signal delay (making the task more

1138 difficult; 1:3 ratio of signal delay increase to decrease). The Go signal delay commonly started at

1139 750ms, and ranged around 700ms (across experiments, $M = 691.78$, $SD = 102.77$).

1140    The training phase in each experiment, consisted of 8 to 20 training runs (see Table 1).

1141 Participants were not informed in advance of the contingency between Go signal and Go stimuli.

1142 However, as they were repeatedly exposed to the same stimuli, in the later runs of the task, some

1143 of the participants were able to identify the Go stimulus and cue association, thus producing

1144 accurate faster responses, sometimes even preceding the Go cue onset. This measurement of RT

1145 to Go stimuli trials was used in this current work as our main behavioral measurement for the

1146 CAT learning model. The details of the model are further discussed below.

1147    ***Probe.*** In the final probe phase, preference modification following CAT was evaluated.

1148 The probe phase was usually preceded by a 'filler' task such as filling up questionnaires or

1149 ranking liking for fractal art images. This task generally provided some time for consolidation

1150 and a dissociation between the CAT and probe phases. In the probe task, participants were

1151 presented with pairs of stimuli of the similar initial subjective value (both high-value stimuli or

42

1152  both low-value). In each pair, one of the stimuli was a Go stimulus and the other was a NoGo

1153  stimulus. Preference modification was evaluated as the proportion of trials in which participants

1154  chose the Go stimulus over the NoGo stimulus, above and beyond the expected 50% chance

1155  level. This measurement was used as the main outcome variable which indicated preference

1156  modification following CAT.

1157  **Analysis.** In Study 1, we aimed to identify a marker for learning in the CAT task. To

1158  achieve this goal, we examined RT patterns in the task. In an exploratory analysis of the meta-

1159  analysis data, we identified that as training progressed, mean RT in the task was reduce. We

1160  hypothesized that as the training task progresses, participants that were able to identify and learn

1161  the stimulus-cue contingency pattern, and thus could generate faster Go responses which do not

1162  rely on the delayed Go cue onset, but rather on the Go stimulus onset. We formally modeled this

1163  process in a single experiment before testing it on the entire meta-analysis dataset.

1164  ***Computational model of CAT.*** To model our proposed cognitive mechanism we utilized

1165  a Bayesian modeling approach of the RT in the CAT task, using the R implementation of Stan

1166  programing language [83]. RTs were modeled as a mixture of two gaussian distributions (Equation

1167  1) – one Gaussian of shorter mean RT, representing the early anticipatory responses generated

1168  when the stimulus-cue association is predicted by the participants; and a second Gaussian with a

1169  later mean RT following the cue-onset (i.e., a distribution representing standard cue-dependent

1170  responses).

$$RT_{t,i} - Cue_{t,i} = N(\mu_1, \sigma_{\varepsilon_1})(\theta_{t,i}) + N(\mu_2, \sigma_{\varepsilon_2})(1 - \theta_{t,i}) \qquad 1$$
$$\theta_{t,i} = \Phi(-3.1 + \theta_{slope_i} Run_{t,i})$$
$$\theta_{slope_i} \sim N(\theta_{slope}, \sigma_{\theta_{slope}})$$
$$\mu_1 < 0, \ \mu_2 > 0; \ t - trial \ index; \ i - participant \ index$$

1171  The $\theta_{t,i}$ mixture probability of the two Gaussian distribution was used to determine the

1172  proportion of trials participants were expected to produce early anticipatory response (1). It was

1173  defined using a linear function of time-dependent $\theta_{slope_i}$ parameter, which was multiplied by the

1174  scaled training run index (training repetition), scaled to 0-1, where 0 and 1 indicating the first

1175  and last (20th) training repetition, respectively. This $\theta_{slope_i}$ parameter was fitted individually for

1176  every participant and was therefore defined conceptually as the individualized learning

1177  parameter of interest. To scale the linear function to a range suitable for proportions (0 – 1), we

1178  used the normal CDF as a link function.

1179    The baseline proportion of early anticipatory responses (at the first training run) was

1180    fixed at the value $\theta_0 = -3.1$, corresponding with probability of 0.1% to generate an anticipatory

1181    response at the very first run. Thus, the individualized learning parameter could be interpreted as

1182    the sole parameter effecting the mixture proportion for the two Gaussian distributions, e.g.,

1183    $\theta_{slope_i} = 3.1$ would be interpreted as a $\Phi(0) = 0.5$, 50% proportion of anticipatory responses at

1184    the final (20$^{\text{th}}$) training run. The individually fitted $\theta_{slope_i}$ parameter was determined in the

1185    current work as the computational marker of learning and was our main parameter of interest.

1186    See Supplementary Code S1 for a complete specification of the model parameters and priors.

1187    The model's parameters were evaluated using Markov-chain Monte Carlo (MCMC)

1188    gradient algorithm, implemented with RStan [83]. Each model was evaluated four independent

1189    times, with chains length of 2000 (1000 chain links burn-out time). The mean values of each

1190    parameter of the converged model were used as the parameter estimates and are reported along

1191    with 95% credible interval (CI). All reported results converged onto stable solution with $\hat{R} = 1$.

1192    To maintain stable interpretation of $\mu_1$ and $\mu_2$ as the centers of anticipatory and cue-

1193    dependent responses, respectively, an upper threshold of 150ms from cue onset was imposed on

1194    the $\mu_1$ parameter, and a lower threshold of 200ms was for $\mu_2$ parameter. Thus, in every chain $\mu_1$

1195    was imposed the role of the earlier anticipatory responses. Alternative formulation forms of the

1196    computational model were also considered and tested, including using 0ms as both upper and

1197    lower threshold for the distribution means, using a log-normal distribution instead of a normal

1198    distribution, and fitting a unique $\mu_{1_i}$ parameter for each participant. However, these models

1199    either did not converge when applying to the entire dataset or converged with some issues.

1200    Initially we used 0ms as an upper threshold for $\mu_1$, which caused the parameter to converge at

1201    this maximal threshold. Using a unique $\mu_{1_i}$ parameter for each participant was theoretically

1202    favorable (and indeed used in Study 2), but could not be applied in Study 1, potentially either

1203    due to computational challenges (fitting many more inter-dependent parameters) or due to the

1204    complex structure of the data (e.g., where anticipatory responses were intertwined with late

1205    responses due to changing go signal delay)

1206    ***Computational model of CAT with stimulus-level parameters.*** In an exploratory model,

1207    designed following the analysis of Study 2, we aimed to fit a stimulus-level learning parameter.

1208    For each stimulus$_s$ which was presented to participant$_i$, we fitted an individualized $\theta_{slope_{i,s}}$

1209   parameter, which was used instead of the participants-level individualized learning parameter

1210   ($\theta_{slope_i}$), see Supplementary Code S2. Initially, we attempted to model a participant-level

1211   dependence, by modeling $\theta_{slope_{i,s}}$ as a parameter derived from higher-level $\theta_{slope_i}$ participant-

1212   level parameter, similarly to the design in Study 2 (see below). However, such model with tens

1213   of thousands of $\theta_{slope_{i,s}}$ parameter estimates that covaried with hundreds of $\theta_{slope_i}$ parameters

1214   was computationally too demanding. To simplify the model structure, $\theta_{slope_{i,s}}$ parameter was

1215   modeled independently of participants' identity (i.e., not considering within-participant possible

1216   effect; see Supplementary Code S2).

1217        *Deviation from previous version and pre-registration.* In our pre-registration, we used an

1218   identical Stan model with one key difference – to enforce an order in which $\mu_1, \sigma_1$ would relate

1219   to early anticipatory RTs and $\mu_2, \sigma_2$ would relate to late cue-dependent RTs, a different

1220   restriction was set in which the upper and lower limit of $\mu_1$ and $\mu_2$, respectively, were set to 0

1221   (Cue onset), instead of the final limits set to 150ms and 200ms. Running the model with these

1222   different restrictions converged to a stable solution. Early anticipatory responses were modeled

1223   as having an earlier mean ($\mu_1$ = -0.04ms, 95%CI [-0.16, 0.0], $\sigma_{\varepsilon_1}$ = 283.02, 95%CI [280.56,

1224   285.62]) and late cue-dependent responses were very similar to the new model ($\mu_2$ = 287.71ms,

1225   95%CI [287.26, 288.16], $\sigma_{\varepsilon_1}$ = 76.97, 95%CI [76.62, 77.32]). Other parameters were also very

1226   similar ($\theta_{slope}$ = 3.49, 95%CI [3.20, 3.77]), with variation between participants ($\sigma_{\theta_{slope}}$ = 4.16,

1227   95% CI [3.93, 4.39]). All reported association with probe phase remained very similar.

1228        While the model showed no formal convergence errors (see trace plots of parameters in

1229   Supplementary Fig. S9), the $\mu_1$ parameter estimate seemed to have been affected by the artificial

1230   limitation which was imposed on its maximal value for technical reasons (to avoid inversion with

1231   $\mu_2$ across the different chains). Throughout the different chains, the final parameter converged

1232   around the upper limit set to 0 (the Cue onset). After several experimental attempts to change

1233   this limit, we found the increasing $\mu_1$ lower limit to 150ms resulted in a better solution, in which

1234   $\mu_1$ parameter estimate did not converge around the uppermost limit. This model is the final

1235   model we chose to use and report here.

1236        ***Probe analysis.*** To evaluate the effect of CAT on preferences we analyzed the proportion

1237   of trials in which participants chose the Go stimulus over the NoGo stimulus in the probe phase,

1238   using mixed-model logistic regression. As Go and NoGo stimuli were matched based on initial

1239      value, under the null hypothesis, participants were expected to choose Go stimuli at 50% of trials

1240      (log-odds = 0; odds = 1). The results of this analysis were of main interest in previous

1241      publications and is reported in the current work for unpublished data (see supplementary

1242      materials).

1243      While early work with CAT task showed a differential effect of value on CAT effect on

1244      preferences - i.e. CAT usually had induced more prominent preference modification for stimuli

1245      of initial high value, compared to stimuli of initial low value [12,13], more recent work with CAT

1246      found that this effect was not a dominant feature of CAT [15,18,19,22]. Thus, in the current work, we

1247      pooled together data of stimuli with both high- and low-initial value.

1248      ***Probe and individualized-learning parameter association.*** Our main analysis of interest

1249      aimed to evaluate the association of the CAT preference modification effect with the

1250      individualized learning parameter $\theta_{slope_i}$, derived from the computational modeling of CAT. In a

1251      mixed model logistic regression analysis implemented with lme4 R package [84], the number of

1252      trial each participant chose the Go versus NoGo stimuli was explained using the independent

1253      variable of $\theta_{slope_i}$ of each participant, as calculated in the preceding CAT task. The data was

1254      aggregated over all participants in the 29 experiments. To account for the aggregation of

1255      different experiments in the dataset, we also included random intercept and random slope terms

1256      for each of the 29 experiments, see logistic regression formula in Equation 2.

$$Choice \sim 1 + \theta_{slope_i} + \left(1 + \theta_{slope_i}|Exp.\right) \qquad 2$$

1257

1258      In an additional exploratory analysis (which was conceived after the analysis of Study 2),

1259      we modeled choices using a stimulus-specific $\theta_{slope_{i,s}}$ parameter. We aimed to examine if

1260      modeling a learning parameter at a more precise stimulus-level (contrary to the pre-registered

1261      participant-level $\theta_{slope_i}$ parameter), would serve as a more accurate predictor of future

1262      preference modification. For this aim we merged the $\theta_{slope_{i,s}}$ which were estimated for each

1263      contingency with the probe data. Choices (number of trials each Go stimulus$_s$ was chosen or not

1264      by each participant$_i$) were modeled with two additional models using $\theta_{slope_{i,s}}$ as an independent

1265      variable. In the first model, $\theta_{slope_{i,s}}$ parameter estimates were used instead of $\theta_{slope_i}$ estimates.

1266      Since $\theta_{slope_{i,s}}$ varied within participants, which were nested within experiment, a nested random

1267      slope effect for $\theta_{slope_{i,s}}$ was also added (Equation 3a). In two additional models, probe choices

1268    were explained using a full model with both $\theta_{slope_{i,s}}$ and $\theta_{slope_i}$ (Equation 3b), and a nested

1269    model with the same IVs except for the fixed and random effects of $\theta_{slope_{i,s}}$ (Equation 3c). Using

1270    a likelihood ratio test for nested models (comparing model 3b with the nested model of 3c), we

1271    examined whether $\theta_{slope_{i,s}}$ provided additional explanatory above and beyond $\theta_{slope_i}$.

$$3$$

$$Choice \sim 1 + \theta_{slope_{i,s}} + \left(1 + \theta_{slope_{i,s}} \middle| Exp./Participant\right) \quad (a)$$

$$Choice \sim 1 + \theta_{slope_{i,s}} + \theta_{slope_i} + \left(1 + \theta_{slope_{i,s}} \middle| Exp./Part.\right) + \left(\theta_{slope_i} \middle| Exp.\right) \quad (b)$$

$$Choice \sim 1 + \theta_{slope_i} + \left(1 + \theta_{slope_i} \middle| Exp.\right) \quad (c)$$

1272    ***Effect size estimation for mixed models with $R^2_{GLMM}$.*** For the reported logistic regression

1273    mixed model, we included an additional effect size score using a generalized linear mixed effect

1274    model (GLMM) $R^2$ estimate, based on the work by Nakagawa, and Schielzeth (2013),

1275    implemented with R's Multi-Model Inference (MuMIn) package [85–88]. Like $R^2$ in linear mixed

1276    model, $R^2_{GLMM}$ was used to quantify the relative proportion of variance accounted by the three

1277    generalized mixed model's variance components: variance explained by the fixed effects ($\sigma_f^2$),

1278    variance explained by the random effects ($\sigma_\alpha^2$), and unexplained residual variance ($\sigma_\varepsilon^2$). In

1279    accordance with the developers' suggestion [85,86], we report here two $R^2$ scores – the marginal $R^2$

1280    (equation 4a), which signifies the proportion of total variance accounted by the fixed effects

1281    ($\sigma_f^2$); and the conditional $R^2$ (equation 4b), which signifies the relative proportion of variance

1282    explained by both the fixed and random effects of the model.

$$4$$

$$R^2_{GLMM(m)} = \frac{\sigma_f^2}{\sigma_f^2 + \sigma_\alpha^2 + \sigma_\varepsilon^2} \quad (a)$$

$$R^2_{GLMM(c)} = \frac{\sigma_f^2 + \sigma_\alpha^2}{\sigma_f^2 + \sigma_\alpha^2 + \sigma_\varepsilon^2} \quad (b)$$

1283

1284      In mixed models examining the contribution of stimulus-level computational marker

1285    ($\theta_{slope_{i,s}}$; Equation 3) $R^2_{GLMM}$ values were similarly evaluated. However, it is important to note

1286    that the basic unit of analysis in these model (choices per Go stimulus within participant) differ

1287    from the unit of analysis in Equation 2 (choices per participant). Thus the $R^2_{GLMM}$ values are not

1288    comparable between the two models The $R^2_{GLMM(m)}$ used the participant-level model is the most

1289    similar in structure and interpretation to "traditional" $R^2$ in ordinary linear model, representing

1290    the proportion of variance explained by the (fixed) effects of interest, using the participant as the

1291    basic unit of analysis.

1292

1293    **Study 2: Novel CAT design**

1294    **Stimuli.** In the two novel experiments, the stimuli set included images of 80 unfamiliar

1295    faces from the Siblings dataset [40], 40 male and 40 female characters. As in previous CAT studies

1296    with this stimuli set [15,18], the images were processed in Photoshop, so that all stimuli were

1297    cropped to identical size (400×500 pixels), with the character's pupils positioned at the same

1298    spatial coordinates and a homogeneous gray background. The characters poised similar neutral

1299    expression and had minimal salient artificial characteristics such as jewelry, make-up, or distinct

1300    facial hair. Previous CAT studies with these stimuli found that a similar preferences

1301    enhancement for both high- and low-value faces which were associated with a Go cue during

1302    CAT [15,18]. Meaning, CAT had similar effect on preferences both for stimuli of initial low-value

1303    and of initial high-value. Using face stimuli allowed us to pool together a larger sample of

1304    stimuli of different initial values, under the reasonable assumption that the effect would show

1305    similar pattern across the different initial value categories.

1306    During the CAT task, a visual Go cue was associated with some of the face stimuli. The

1307    Go cue was identical to the one used in a previous study [15], it comprised of a semi-transparent

1308    Gabor image (38×38 pixels, alpha = 0.7), which appeared at the center of the screen, on top of

1309    the face stimuli.

1310    **Procedure.** Like previous CAT experiments, the two novel preliminary and replication

1311    experiments included three phases: a baseline preference evaluation task, a modified cue-

1312    approach training task, and a post training probe phase which examined preference modification.

1313    ***Initial preference evaluation task.*** To evaluate initial preferences, participants underwent

1314    a binary forced-choice task between random pairs of stimuli, as in previous CAT experiments

1315    with non-consumable stimuli [15,18]. The task included 400 unique choice trials (meaning, no

1316    choice between the same two faces repeated more than once), during which each of the 80 face

1317    stimuli was presented exactly 10 times. Each choice trial lasted 3000ms, of which participants

1318    were given a 2000ms time window to make their choice. Choice trials were followed by a 500ms

1319    confirmation screen, showing a green frame around the chosen stimulus, and a fixation cross,

1320    which was presented for the remaining trial duration as inter-stimulus interval (ISI), for at least

1321  500ms. In case no choice was made during the 2000ms time-window, a screen saying "You must

1322  respond faster" appeared for 500ms, followed by a 500ms ISI.

1323  Binary choices were transformed into individual preference scores using Colley ranking

1324  algorithm [82]. Stimuli were ranked based on each participant's individual initial preferences from

1325  the highest value (1) to lowest value (80). Ranks were used to categorize the stimuli into 10

1326  equal-size value groups (each containing eight stimuli; ranks 1-8, ranks 9-16, etc.). The value

1327  categories and internal ranks within each category were used to allocate conditions in the

1328  subsequent training and probe task, ensuring initial values were balanced across 100% and 50%

1329  contingency conditions, as well as across Go and NoGo stimuli within each value-category,

1330  which would be pitted against each other in the subsequent probe phase.

1331  ***Cue-approach training.*** The CAT task consisted of 20 training runs, in each run all 80

1332  stimuli were presented in a random order for 1000ms each, with a 500ms ISI. In 30% of trials a

1333  visual semi-transparent Go cue appeared 850ms following the stimulus onset for 100ms at the

1334  center screen position, on top of the face stimulus. Unlike previous CAT experiments, the Go cue

1335  onset did not change throughout the training task. Thus, actual RT was consistent with effective

1336  RT (RT from cue onset).

1337  Three Go association conditions were included in this modified version of the CAT – 16

1338  stimuli were always presented with the Go cue (Go stimuli, 100% contingency condition), 16

1339  stimuli were associated with the Go cue during half of the presentations (Go stimuli, 50%

1340  contingency condition), and the rest of the stimuli were never followed with a Go cue (NoGo

1341  stimuli). Participants were instructed to respond to Go stimuli by pressing a keyboard button as

1342  fast as they can, before the face-stimulus disappear. Participant were told in the instruction of the

1343  three Go cue contingencies, and that they may respond when they see the or even when they

1344  anticipate the cue will appear shortly (see the task instruction as presented to the participants in

1345  Supplementary Fig. S5). Unlike previous CAT studies, participants were explicitly told of the Go

1346  contingency and were thus encouraged to initiate anticipatory responses preceding the actual cue

1347  onset, while maintaining high accuracy by only responding when they have sufficient confidence

1348  that a cue will follow.

1349  Stimuli were allocated with a go contingency condition based on initial preferences.

1350  Stimuli were categorized into 10 equal sized initial value categories, each containing eight

1351  stimuli. The stimuli in the highest and lowest initial value categories were allocated to be all

1352 NoGo stimuli. Of the middle eight value categories, four were allocated to be 100% contingency

1353 and 50% contingency condition category. Within each of the eight middle categories, half of the

1354 stimuli were Go stimuli and the other half were NoGo stimuli. The role allocation between

1355 categories and within them was designed so that initial values were balanced across the groups

1356 E.g., within the second value categories, stimuli ranked 9, 12, 13, and 16 were allocated to be Go

1357 stimuli, while stimuli ranked 10, 11, 14, and 15 were allocated to be NoGo stimuli (thus, each

1358 group had the same mean rank of 12.5; see illustration in Figure 9). The conditions allocation by

1359 initial value was counter-balanced across participants.

1360

| Value category | Ranks | Contingency for Go stimuli | | Value rank | Go allocation |
|---|---|---|---|---|---|
| 1 (highest) | 1 – 8 | All NoGo | | 9 | 100% Go |
| 2 | 9 – 16 | 100% | | 10 | NoGo |
| 3 | 17 – 24 | 50% | | 11 | NoGo |
| 4 | 25 – 32 | 50% | | 12 | 100% Go |
| 5 | 33 – 40 | 100% | | 13 | 100% Go |
| 6 | 41 – 48 | 100% | | 14 | NoGo |
| 7 | 49 – 56 | 50% | | 15 | NoGo |
| 8 | 57 – 64 | 50% | | 16 | 100% Go |
| 9 | 65 – 72 | 100% | | | |
| 10 (lowest) | 73 – 80 | All NoGo | | | |

1361 Figure 9. Go stimuli allocation illustration. Both allocations to Go contingency condition (100%
1362 versus 50%) and Go/NoGo stimuli were balanced based on initial subjective value. Go stimuli of
1363 50% contingency were associated in 10 out of the 20 training runs with the Go cue. Different
1364 designs were counter-balanced across participants by switching between 100% and 50%
1365 contingency categories positions and within category Go/NoGo positions.
1366

1367 All Go stimuli of the 50% contingency condition were associated with the Go cue in the

1368 last two runs as well as in eight additional runs (in total - 10 out of 20 run). This was done for

1369 potential future implementation of the task in an MRI scanner, which have not been done at the

1370 time of this manuscript write-up.

1371 *Probe.* In the probe phase, preference modification effect was examined using a binary

1372 forced choice task, in which participants chose their preferred stimulus of pairs Go versus NoGo

1373 stimuli, of similar initial value. Each Go stimulus was pitted against the four NoGo stimuli

1374 within the same initial value category (e.g., Go stimuli ranked 9, 12, 13, and 16 were pitted

1375 against NoGo ranked 10, 11, 14, and 15). Participants had 1500ms to make their choice, which

1376 was followed by a 500ms confirmation feedback (or a message prompting faster response, in

1377 case no choice was made), and an ISI of varying duration (1000ms – 9500ms, $M$ = 3000ms),

1378 drawn from a truncated exponential distribution with 100ms precision. Like in previous CAT

1379 experiments with non-consumable stimuli (and in contrast to CAT experiments with consumable

1380 stimuli), choices in the probe phase were not incentive compatible. Previous studies showed that

1381 CAT effect is consistent both across incentive compatible choices of consumable stimuli, as well

1382 as in non-incentive compatible choices [15].

1383 Choices of the 100% contingency condition and 50% contingency condition were merged

1384 across the different value categories and analyzed with a mixed-logistic regression model. We

1385 hypothesized that participants would choose the Go stimuli over NoGo stimuli, and that this

1386 preference modification effect would be more robust for the 100% contingency condition. We

1387 also hypothesized that individual $\theta_{slope_i}$ parameter from computational modeling of the training

1388 task, would predict which participants would demonstrate stronger preferences modification

1389 effect.

1390 **Participants.** In the preliminary experiment and replication experiment $n$ = 20 and $n$ = 59

1391 valid participants completed the experiment and were included in the analysis, respectively. The

1392 sample size of the preliminary experiments was based on the minimal sample size used in

1393 previous CAT studies. The sample size required for the replication sample was based on a power

1394 analysis using the results of the preliminary experiment. A sample size of $n$ = 59 was expected to

1395 be sufficient to achieve 95% power to detect a significant correlation effect ($\alpha$ = 0.05) between

1396 probe choices and the $\theta_{slope_i}$ parameter estimates, in both 100% and 50% contingency

1397 conditions. The power analysis and the resulting sample sized was documented in the pre-

1398 registration and its supplementary materials (https://osf.io/nwr4v).

1399 Three quality-assurance measurements were used as exclusion criteria, based on previous

1400 studies with CAT [15,18,22] – (1) low variability of Colley scores in the initial preference evaluation

1401 task (which indicate intransitive choice pattern), (2) proportion of false alarm during training,

1402 and (3) proportion of missed Go trials during training. In each experiment, we excluded

1403 participants with extremely low transitivity, high false alarm, or high miss rate (defined as 3SD

1404 from the group mean). The exclusion criteria were pre-registered for the replication experiment,

1405 along with the planned sample size.

1406        In the preliminary experiment, no participant was excluded due to the mentioned above

1407    exclusion criteria (Transitivity score: $M = 0.205$, $3SD$ cutoff $= 0.122$, min valid score $= 0.137$;

1408    False alarm rate: $M = 3.62\%$, $3SD$ cutoff $= 15.80\%$, max valid score $= 15.27\%$; Miss rate: $M =$

1409    $6.33\%$, $3SD$ cutoff $= 24.13\%$, max valid score $= 21.04\%$). In the replication experiment, five

1410    participants were excluded due to these pre-registered exclusion criteria: Two participants had

1411    low transitivity score, one participant had high rate of false-alarm, and two participants had high

1412    rate of missed Go trials (Transitivity score: $M = 0.213$, $3SD$ cutoff $= 0.141$, $M_{valid} = 0.216$, min

1413    valid score $= 0.143$; False alarm rate: $M = 4.74\%$, $3SD$ cutoff $= 39.00\%$, $M_{valid} = 3.49\%$, max

1414    valid score $= 37.59\%$; Miss rate: $M = 7.59\%$, $3SD$ cutoff $= 40.22\%$, $M_{valid} = 5.97\%$, max valid

1415    score $= 38.75\%$).

1416        **Analysis.** Like the procedural design, the analyses of the two experiments were also

1417    identical, and generally resembled that of the meta-analysis study, introducing some analysis

1418    improvements which could not be applied in the meta-analysis.

1419        *CAT computational model.* The general goal and design of the Bayesian computational

1420    framework in the two new experiments resembled that of Study 1. Participants' RTs were

1421    modeled as a mixture of two Gaussian distributions - one distribution of late, cue-dependent

1422    responses, and another distribution of earlier anticipatory responses. A $\theta_{t,i}$ mixture proportion

1423    determined the probability of making an anticipatory response by the participant$_i$ at trial$_t$. The $\theta_{t,i}$

1424    probability was modeled as a monotonic function of time, using as an individual rate parameter.

1425    Unlike the meta-analysis model, in Study 2, we introduced two distinct contingency conditions,

1426    thus for each participant two $\theta_{slope_i}$ parameters were modeled, namely one for the 100%

1427    contingency condition, and one for the 50% contingency condition. Another discrepancy

1428    between the meta-analysis model and Study 2 model, was an introduction of individual

1429    parameter for anticipatory responses center ($\mu_{1_i}$), which allowed flexibility in modeling

1430    individual differences in anticipatory responses onsets (Equation 5). This model improvement

1431    was possible thanks to the more homogenous nature of the data. Each experiment (preliminary

1432    and replication) was modeled separately. See the complete model and priors in Supplementary

1433    Code S3.

$$RT_{t,i} = N(\mu_{1_i}, \sigma_{\varepsilon_1})(\theta_{t,i}) + N(\mu_2, \sigma_{\varepsilon_2})(1 - \theta_{t,i}) \qquad 5$$

$$\theta_{t,i} = \Phi\left[-3.1 + I_{cond.t,i}\theta_{100\%slope_i}Run_{t,i} + (1 - I_{cond.t,i})\theta_{50\%slope_i}Run_{t,i}\right]$$

$$\theta_{100\%slope_i} \sim N(\theta_{100\%}, \sigma_{\theta_{100\%}}); \; \theta_{50\%slope_i} \sim N(\theta_{50\%}, \sigma_{\theta_{50\%}}); \; \mu_{1_i} \sim N(\mu_1, \sigma_{\mu_1})$$

Computational marker of non-reinforced learning

$$\mu_{1_i} < CueOnset + 100, \ \mu_2 > CueOnset;$$
$$t - trial\ index;\ i - participant\ index$$
$$I_{cond.t,i} - condition\ index\ (100\%\ versus\ 50\%)$$

1434      In an additional exploratory analysis (not pre-registered), we aimed to create an even

1435    more accurate learning computational marker by modeling a $\theta_{slope_{i,s}}$ for every stimulus$_s$ that

1436    participant$_i$ was trained with. Assuming that some variability in training speed could be

1437    measured not only between participants, but also withing each participant. Since the actual

1438    stimuli which were allocated to be Go stimuli varied between participants (i.e., for each

1439    participant the condition and role of a certain image stimulus was different), we assumed no

1440    mutual information is shared between the same stimulus index of different participants. Stimuli

1441    and participants were treated as random effects (Equation 6).

$$RT_{t,i} = N\big(\mu_{1_i}, \sigma_1\big)\big(\theta_{t,i}\big) + N(\mu_2, \sigma_2)(1 - \theta_{t,i}) \qquad 6$$
$$\theta_{t,i} = \Phi\big(-3.1 + I_{stim.t,i}\,\theta_{slope_{i,s}} Run_{t,i}\big)$$
$$\mu_{1_i} < CueOnset, \ \mu_2 > CueOnset;$$
$$t - trial\ index;\ i - participant\ index;\ s - stimulus\ index$$
$$I_{stim.t} - stimulus\ index\ (for\ each\ contingency\ condition\ separately)$$

1442      In the process of testing this novel approach, the model did not converge when presented

1443    with the full data which included both 100% and 50% contingency conditions. However, when

1444    trained twice, using each contingency-condition as an independent dataset, the model did

1445    converge. See the complete model and priors in Supplementary Code S4.

1446      ***Probe.*** As in previous CAT studies, preference modification was evaluated following

1447    CAT. Participants' preferences for Go over NoGo stimuli were categorized into two conditions

1448    corresponding with the CAT conditions – of 100% contingency and 50% contingency

1449    conditions. Trials of the different value groups were pooled together and analyzed with a mixed

1450    model logistic regression. In a post-hoc analysis we validated that the initial value had no

1451    significant impact on the conclusions of the analysis (see supplementary materials).

1452      Based on the results of the meta-analysis, we hypothesized that better learning

1453    (manifested as larger $\theta_{slope_i}$ parameter estimate) would induce stronger preference modification

1454    effect. Thus, we also hypothesized that a stronger preference modification effect would be

1455    observed for the 100% contingency condition, compared with the 50% contingency condition.

1456    To examine this hypothesis, we run a one-sided mixed model logistic regression with an

Computational marker of non-reinforced learning

1457   additional fixed and participant-based random explanatory variables (EVs) of the contingency

1458   condition (Equation 7).

$$Choice \sim 1 + Contingency + (1 + Contingency | Participant) \qquad 7$$

1459   ***Probe choice prediction based on individual learning parameter.*** Most importantly, we

1460   hypothesized that using the $\theta_{slope_i}$ parameter estimate for learning would predict future

1461   preference modification effect observed during the subsequent probe phase. To test this

1462   hypothesis, we introduced the $\theta_{slope_i}$ as an additional independent variable in the mixed model

1463   logistic regression. Choices were modeled using contingency condition (both as fixed effect and

1464   as a random slope between participants), $\theta_{slope_i}$, and an interaction term. This model is

1465   equivalent to modeling choices using an intercept and $\theta_{slope_i}$ slope separately for each

1466   contingency condition with a random intercept modeled within participants (Equation 8a).

1467       In an additional (not preregistered) exploratory analysis, we modeled a stimulus-specific

1468   $\theta_{slope_{i,s}}$ parameter. We aimed to examine if modeling a learning parameter at a more detailed

1469   stimulus-level (extending the pre-registered participant-level $\theta_{slope_i}$ parameter), would serve as a

1470   more accurate predictor of future preference modification. For this aim we merged the $\theta_{slope_{i,s}}$

1471   which were estimated for each contingency with the probe data. Choices were modeled with two

1472   additional models using a $\theta_{slope_{i,s}}$ as an independent variable. In the first model, $\theta_{slope_{i,s}}$

1473   parameter estimates were used instead of $\theta_{slope_i}$ estimates. Since $\theta_{slope_{i,s}}$ varied within

1474   participant, a random slope effect for $\theta_{slope_{i,s}}$ was also added (Equation 8b). In another model,

1475   probe choices were explained using a full model with both $\theta_{slope_{i,s}}$ and $\theta_{slope_i}$ (Equation 8c).

1476   This model was used to examine whether $\theta_{slope_{i,s}}$ provided additional explanatory power in a

1477   likelihood ratio test for nested models (comparing the model of 8c with the nested model of 8a).

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad 8$$
$$Choice \sim 1 + I_{cond.} * \theta_{slope_i} + (1 + I_{cond.} | Participant) \qquad (a)$$
$$Choice \sim 1 + I_{cond.} * \theta_{slope_{i,s}} + \left(1 + I_{cond.} * \theta_{slope_{i,s}} | Participant\right) \qquad (b)$$
$$Choice \sim 1 + I_{cond.} * \theta_{slope_i} + I_{cond.} * \theta_{slope_{i,s}} + \left(1 + I_{cond.} * \theta_{slope_{i,s}} | Part.\right) \qquad (c)$$
$$i - participant\ index;\ s - stimulus\ index$$
$$I_{cond.} - contingency\ condition\ indicator\ IV$$

1478       When results with $\theta_{slope_{i,s}}$ were analyzed, we found that some participants were fitted a

1479   $\theta_{slope_{i,s}}$ estimates of extremely small variability (SD<0.1, meaning their $\theta_{slope_{i,s}}$ estimates were

1480    effectively fixed; see Supplementary Fig. S6). In such cases a mixed model with $\theta_{slope_{i,s}}$ random

1481    slope (Equation 8b and 8c), resulted in convergence warnings (this phenomenon was not

1482    observed in Study 1). Thus, for these analyses, we excluded all participants which demonstrated

1483    such low variability of $\theta_{slope_{i,s}}$ estimates in either one of the two contingency conditions. This

1484    resulted in exclusion of nine participants in the preliminary experiment, and 14 participants in

1485    the replication experiment (out of 20 and 59 participants, respectively). For the nested model

1486    likelihood ratio analysis (comparing a full model to model without $\theta_{slope_{i,s}}$ independent

1487    variables), both models were examined using the same subset of 11 and 45 valid participants.

1488    ***Effect size estimation for mixed models with $R^2_{GLMM}$.*** As in the meta-analysis study, we

1489    report for the logistic regression mixed models two $R^2$ estimates – $R^2_{GLMM(m)}$ and $R^2_{GLMM(c)}$,

1490    denoting respectively the marginal and conditional $R^2$ effect of the generalized linear model [85–

1491    88].

1492

1493    **References**

1494    1.    Staddon, J. E. R. & Cerutti, D. T. Operant Conditioning. *Annual Review of Psychology* **54**,

1495          115–144 (2003).

1496    2.    Rangel, A., Camerer, C. & Montague, P. R. A framework for studying the neurobiology of

1497          value-based decision making. *Nat. Rev. Neurosci.* **9**, 545–556 (2008).

1498    3.    Sutton, R. S. & Barto, A. G. *Reinforcement learning: An introduction.* (MIT press, 2018).

1499    4.    Niv, Y. Reinforcement learning in the brain. *J. Math. Psychol.* **53**, 139–154 (2009).

1500    5.    Schonberg, T. & Katz, L. N. A Neural Pathway for Nonreinforced Preference Change.

1501          *Trends in Cognitive Sciences* (2020). doi:10.1016/j.tics.2020.04.002

1502    6.    Zajonc, R. B. Attitudinal Effects of Mere Exposure. *J. Pers. Soc. Psychol.* **9**, 1–27 (1968).

1503    7.    Voigt, K., Murawski, C. & Bode, S. Endogenous formation of preferences: Choices

1504          systematically change willingness-to-pay for goods. *J. Exp. Psychol. Learn. Mem. Cogn.*

1505          **43**, 1872–1882 (2017).

1506    8.    Sharot, T., Velasquez, C. M. & Dolan, R. J. Do decisions shape preference? Evidence

1507          from blind choice. *Psychol. Sci.* **21**, 1231–1235 (2010).

1508    9.    Sharot, T., Fleming, S. M., Yu, X., Koster, R. & Dolan, R. J. Is Choice-Induced

1509          Preference Change Long Lasting? *Psychol. Sci.* **23**, 1123–1129 (2012).

1510    10.    Tom, G., Nelson, C., Srzentic, T. & King, R. Mere exposure and the endowment effect on
1511          consumer decision making. *J. Psychol. Interdiscip. Appl.* **141**, 117–125 (2007).

1512    11.    Rindfleisch, A. & Inman, J. J. Explaining the familiarity-liking relationship: Mere
1513          exposure, information availability, or social desirability? *Mark. Lett.* **9**, 5–19 (1998).

1514    12.    Schonberg, T. *et al.* Changing value through cued approach: an automatic mechanism of
1515          behavior change. *Nat. Neurosci.* **17**, 625–630 (2014).

1516    13.    Bakkour, A. *et al.* Mechanisms of choice behavior shift using cue-approach training.
1517          *Front. Psychol.* (2016). doi:10.3389/fpsyg.2016.00421

1518    14.    Bakkour, A., Lewis-Peacock, J. A., Poldrack, R. A. & Schonberg, T. Neural mechanisms
1519          of cue-approach training. *Neuroimage* **151**, 92–104 (2017).

1520    15.    Salomon, T. *et al.* The cue-approach task as a general mechanism for long term non-
1521          reinforced behavioral change. *Sci. Rep.* **8**, article number: 3614, 1-13 (2018).

1522    16.    Aridan, N., Pelletier, G., Fellows, L. K. & Schonberg, T. Is ventromedial prefrontal cortex
1523          critical for behavior change without external reinforcement? *Neuropsychologia* (2019).
1524          doi:10.1016/j.neuropsychologia.2018.12.008

1525    17.    Veling, H. *et al.* Training impulsive choices for healthy and sustainable food. *J. Exp.*
1526          *Psychol. Appl.* (2017). doi:10.1037/xap0000112

1527    18.    Salomon, T., Botvinik□Nezer, R., Oren, S. & Schonberg, T. Enhanced striatal and
1528          prefrontal activity is associated with individual differences in nonreinforced preference
1529          change for faces. *Hum. Brain Mapp.* **41**, 1043–1060 (2020).

1530    19.    Botvinik-Nezer, R., Bakkour, A., Salomon, T., Shohamy, D. & Schonberg, T. Memory for
1531          individual items is related to nonreinforced preference change. *Learn. Mem.* (2021).
1532          doi:10.1101/lm.053411.121

1533    20.    Zoltak, M. J., Veling, H., Chen, Z. & Holland, R. W. Attention! Can choices for low value
1534          food over high value food be trained? *Appetite* **124**, 124–132 (2018).

1535    21.    Chen, Z., Holland, R. W., Quandt, J., Dijksterhuis, A. & Veling, H. How preference
1536          change induced by mere action versus inaction persists over time. *Judgm. Decis. Mak.* **16**,
1537          201–237 (2021).

1538    22.    Botvinik-Nezer, R., Salomon, T. & Schonberg, T. Enhanced bottom-up and reduced top-
1539          down neural mechanisms drive long-lasting non-reinforced behavioral change. *Cereb.*
1540          *Cortex* **30**, 858–874 (2020).

1541    23.    Weber, E. U. & Johnson, E. J. Constructing Preferences From Memory. in *The*

1542            *Construction of Preference* (2009). doi:10.1017/cbo9780511618031.022

1543    24.    Weilbächer, R. A., Krajbich, I., Rieskamp, J. & Gluth, S. The influence of visual attention

1544            on memory-based preferential choice. *Cognition* **215**, (2021).

1545    25.    Kraemer, P. M., Weilbächer, R. A., Mechera-Ostrovsky, T. & Gluth, S. Cognitive and

1546            neural principles of a memory bias on preferential choices. *Curr. Res. Neurobiol.* **3**,

1547            100029 (2022).

1548    26.    Zoltak, M. J., Holland, R. W., Kukken, N. & Veling, H. Training choices toward low

1549            value options. *Judgm. Decis. Mak.* **15**, 254–265 (2020).

1550    27.    Krajbich, I., Armel, C. & Rangel, A. Visual fixations and the computation and comparison

1551            of value in simple choice. *Nat. Neurosci.* **13**, 1292–1298 (2010).

1552    28.    Krajbich, I. & Rangel, A. Multialternative drift-diffusion model predicts the relationship

1553            between visual fixations and choice in value-based decisions. *Proc. Natl. Acad. Sci.* **108**,

1554            13852–13857 (2011).

1555    29.    Krajbich, I. Accounting for attention in sequential sampling models of decision making.

1556            *Curr. Opin. Psychol.* **29**, 6–11 (2019).

1557    30.    Ridderinkhof, K. R., Ullsperger, M., Crone, E. A. & Nieuwenhuis, S. The role of the

1558            medial frontal cortex in cognitive control. *Science (80-. ).* **306**, 443–447 (2004).

1559    31.    Tanji, J. New concepts of the supplementary motor area. *Curr. Opin. Neurobiol.* **6**, 782–

1560            787 (1996).

1561    32.    O&apos;Doherty, J. *et al.* Dissociable roles of ventral and dorsal striatum in instrumental

1562            conditioning. *Sci. (New York, NY)* **304**, 452–454 (2004).

1563    33.    Pessiglione, M., Seymour, B., Flandin, G., Dolan, R. J. & Frith, C. D. Dopamine-

1564            dependent prediction errors underpin reward-seeking behaviour in humans. *Nature* **442**,

1565            1042–1045 (2006).

1566    34.    Clithero, J. A. & Rangel, A. Informatic parcellation of the network involved in the

1567            computation of subjective value. *Soc. Cogn. Affect. Neurosci.* (2013).

1568            doi:10.1093/scan/nst106

1569    35.    Kable, J. W. & Glimcher, P. W. The neurobiology of decision: consensus and controversy.

1570            *Neuron* **63**, 733–745 (2009).

1571    36.    Collins, A. G. E. & Frank, M. J. Opponent actor learning (OpAL): Modeling interactive

effects of striatal dopamine on reinforcement learning and choice incentive. *Psychol. Rev.* **121**, 337–366 (2014).

37. Tricomi, E. M., Delgado, M. R. & Fiez, J. A. Modulation of caudate activity by action contingency. *Neuron* **41**, 281–292 (2004).

38. Chew, B., Blain, B., Dolan, R. J. & Rutledge, R. B. A neurocomputational model for intrinsic reward. *J. Neurosci.* **41**, 8963–8971 (2021).

39. Becker, G. M., DeGroot, M. H. & Marschak, J. Measuring utility by a single□response sequential method. *Behav. Sci.* **9**, 226–232 (1964).

40. Vieira, T. F., Bottino, A., Laurentini, A. & De Simone, M. Detecting siblings in image pairs. *Vis. Comput.* **30**, 1333–1345 (2014).

41. Veling, H., Lawrence, N. S., Chen, Z., van Koningsbruggen, G. M. & Holland, R. W. What Is Trained During Food Go/No-Go Training? A Review Focusing on Mechanisms and a Research Agenda. *Curr. Addict. Reports* **4**, 35–41 (2017).

42. Hands, D. W. Foundations of Contemporary Revealed Preference Theory. *Erkenntnis 2012 785* **78**, 1081–1108 (2012).

43. Vlaev, I., Chater, N., Stewart, N. & Brown, G. D. A. Does the brain calculate value? *Trends Cogn. Sci.* **15**, 546–554 (2011).

44. Nakamura, K. & Kawabata, H. I Choose, Therefore I Like: Preference for Faces Induced by Arbitrary Choice. *PLoS One* **8**, 1–8 (2013).

45. Barak, S. *et al.* Disruption of alcohol-related memories by mTORC1 inhibition prevents relapse. *Nat. Neurosci.* **16**, 1111–1117 (2013).

46. Joel, D. The signal attenuation rat model of obsessive–compulsive disorder: a review. *Psychopharmacol. 2006 1864* **186**, 487–503 (2006).

47. Tzschentke, T. M. Measuring reward with the conditioned place preference paradigm: a comprehensive review of drug effects, recent progress and new issues. *Prog. Neurobiol.* **56**, 613–672 (1998).

48. LaBar, K. S., Gatenby, J. C., Gore, J. C., LeDoux, J. E. & Phelps, E. A. Human Amygdala Activation during Conditioned Fear Acquisition and Extinction: a Mixed-Trial fMRI Study. *Neuron* **20**, 937–945 (1998).

49. De Houwer, J. Association learning of likes and dislikes: A review of 25 years of research on human evaluative conditioning. *Psychol. Bull.* **127**, 853 (2001).

1603    50.    Kawa, A. B., Bentzley, B. S. & Robinson, T. E. Less is more: prolonged intermittent

1604            access cocaine self-administration produces incentive-sensitization and addiction-like

1605            behavior. *Psychopharmacology (Berl).* **233**, 3587–3602 (2016).

1606    51.    Frank, M. J., Seeberger, L. C. & O'Reilly, R. C. By carrot or by stick: Cognitive

1607            reinforcement learning in Parkinsonism. *Science (80-. ).* **306**, 1940–1943 (2004).

1608    52.    Frank, M. J., Moustafa, A. A., Haughey, H. M., Curran, T. & Hutchison, K. E. Genetic

1609            triple dissociation reveals multiple roles for dopamine in reinforcement learning. *Proc.*

1610            *Natl. Acad. Sci. U. S. A.* **104**, 16311–16316 (2007).

1611    53.    Calcagnetti, D. J. & Schechter, M. D. Extinction of cocaine-induced place approach in

1612            rats: A validation of the "biased" conditioning procedure. *Brain Res. Bull.* **30**, 695–700

1613            (1993).

1614    54.    Bouton, M. E. Context, time, and memory retrieval in the interference paradigms of

1615            Pavlovian learning. *Psychol. Bull.* **114**, 80–99 (1993).

1616    55.    Bouton, M. E. M. E. Context and behavioral processes in extinction. *Learn. &amp; Mem.*

1617            *(Cold Spring Harb. NY)* **11**, 485–494 (2004).

1618    56.    Myers, K. M. & Davis, M. Behavioral and Neural Analysis of Extinction. *Neuron* **36**,

1619            567–584 (2002).

1620    57.    Nevin, J. A. Behavioral Momentum and the Partial Reinforcement Effect. *Psychol. Bull.*

1621            **103**, 44–56 (1988).

1622    58.    Atlas, L. Y., Bolger, N., Lindquist, M. A. & Wager, T. D. Brain mediators of predictive

1623            cue effects on perceived pain. *J. Neurosci.* **30**, 12964–12977 (2010).

1624    59.    Koban, L., Kusko, D. & Wager, T. D. Generalization of learned pain modulation depends

1625            on explicit learning. *Acta Psychol. (Amst).* **184**, 75–84 (2018).

1626    60.    Pool, E. R. *et al.* Determining the effects of training duration on the behavioral expression

1627            of habitual control in humans: a multi-laboratory investigation. *Learn. Mem.* **29**, 16–28

1628            (2022).

1629    61.    Tricomi, E., Balleine, B. W. & O&apos;Doherty, J. P. A specific role for posterior

1630            dorsolateral striatum in human habit learning. *Eur. J. Neurosci.* **29**, 2225–2232 (2009).

1631    62.    Dolan, R. J. & Dayan, P. Goals and habits in the brain. *Neuron* **80**, 312–325 (2013).

1632    63.    Browning, M., Holmes, E. A., Charles, M., Cowen, P. J. & Harmer, C. J. Using

1633            Attentional Bias Modification as a Cognitive Vaccine Against Depression. *Biol.*

1634      *Psychiatry* **72**, 572–579 (2012).

1635   64.   Hakamata, Y. *et al.* Attention Bias Modification Treatment: A Meta-Analysis Toward the

1636      Establishment of Novel Treatment for Anxiety. *Biol. Psychiatry* **68**, 982–990 (2010).

1637   65.   Bar-Haim, Y., Lamy, D., Pergamin, L., Bakermans-Kranenburg, M. J. & van IJzendoorn,

1638      M. H. Threat-related attentional bias in anxious and nonanxious individuals: a meta-

1639      analytic study. *Psychol. Bull.* **133**, 1–24 (2007).

1640   66.   Tempelaar, D. T., Rienties, B. & Giesbers, B. In search for the most informative data for

1641      feedback generation: Learning analytics in a data-rich context. *Comput. Human Behav.*

1642      **47**, 157–167 (2015).

1643   67.   Pardo, A. A feedback model for data-rich learning experiences. *Assess. Eval. High. Educ.*

1644      **43**, 428–438 (2018).

1645   68.   Leong, Y. C., Radulescu, A., Daniel, R., DeWoskin, V. & Niv, Y. Dynamic Interaction

1646      between Reinforcement Learning and Attention in Multidimensional Environments.

1647      *Neuron* **93**, 451–463 (2017).

1648   69.   Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and reward.

1649      *Science (80-. ).* **275**, 1593–1599 (1997).

1650   70.   Ratcliff, R. A theory of memory retrieval. *Psychol. Rev.* (1978). doi:10.1037/0033-

1651      295X.85.2.59

1652   71.   Smith, P. L. & Ratcliff, R. Psychology and neurobiology of simple decisions. *Trends*

1653      *Neurosci.* **27**, 161–168 (2004).

1654   72.   Niv, Y. *et al.* Reinforcement Learning in Multidimensional Environments Relies on

1655      Attention Mechanisms. *J. Neurosci.* **35**, 8145–8157 (2015).

1656   73.   Niv, Y., Edlund, J. A., Dayan, P. & O'Doherty, J. P. Neural Prediction Errors Reveal a

1657      Risk-Sensitive Reinforcement-Learning Process in the Human Brain. *J. Neurosci.* **32**,

1658      551–562 (2012).

1659   74.   Ballard, I. C. & McClure, S. M. Joint modeling of reaction times and choice improves

1660      parameter identifiability in reinforcement learning models. *J. Neurosci. Methods* **317**, 37–

1661      44 (2019).

1662   75.   Shimojo, S., Simion, C., Shimojo, E. & Scheier, C. Gaze bias both reflects and influences

1663      preference. *Nat. Neurosci.* **6**, 1317–1322 (2003).

1664   76.   Armel, K. C., Beaumel, A. & Rangel, A. Biasing simple choices by manipulating relative

1665  visual attention. *Judgm. Decis. Mak.* (2008).

1666  77.  Cepeda, N. J., Vul, E., Rohrer, D., Wixted, J. T. & Pashler, H. Spacing effects in learning:
1667  a temporal ridgeline of optimal retention. *Psychol. Sci.* **19**, 1095–1102 (2008).

1668  78.  Lundqvist, D., Flykt, A. & Öhman, A. Karolinska directed emotional faces. *Cogn. Emot.*
1669  (1998).

1670  79.  Fantastic Fractals. (2013).

1671  80.  Lang, P. J., Bradley, M. M. & Cuthbert, B. N. *International affective picture system*
1672  *(IAPS): Affective ratings of pictures and instruction manual. Technical Report A-8* (2008).

1673  81.  Lang, P. J., Bradley, M. M. & Cuthbert, B. N. International Affective Picture System
1674  (IAPS): Technical Manual and Affective Ratings. *NIMH Cent. Study Emot. Atten.* 39–58
1675  (1997). doi:10.1027/0269-8803/a000147

1676  82.  Colley, W. Colley's bias free college football ranking method: The Colley matrix
1677  explained. *Princet. Univ.* 1–23 (2002).

1678  83.  Stan Development Team. RStan: the R interface to Stan. (2020).

1679  84.  Bates, D., Mächler, M., Bolker, B. M. & Walker, S. C. Fitting linear mixed-effects models
1680  using lme4. *J. Stat. Softw.* **67**, 1–48 (2015).

1681  85.  Nakagawa, S. & Schielzeth, H. A general and simple method for obtaining R2 from
1682  generalized linear mixed-effects models. *Methods Ecol. Evol.* **4**, 133–142 (2013).

1683  86.  Nakagawa, S., Johnson, P. C. D. & Schielzeth, H. The coefficient of determination R2 and
1684  intra-class correlation coefficient from generalized linear mixed-effects models revisited
1685  and expanded. *J. R. Soc. Interface* **14**, (2017).

1686  87.  Johnson, P. C. D. Extension of Nakagawa & Schielzeth's R2GLMM to random slopes
1687  models. *Methods Ecol. Evol.* **5**, 944–946 (2014).

1688  88.  Bartoń, K. MuMIn: Multi-Model Inference. (2020).

1689

1690

1696    and the Fields-Rayant Minducate Learning Innovation Research Center. **Competing interests:**

1697    The authors declare they have no competing interests.

1698    **Data and materials availability:** All materials and analysis codes are available online at:

1699    https://github.com/tomsalomon/CAT_Individualized_Learning.    Preregistration    hypotheses,

1700    experimental design data and codes used for power analysis are available in the Open Science

1701    Foundation (OSF) depository (https://osf.io/nwr4v).