

1

1 **Title:** Testing the possibility of model-based Pavlovian control of attention to threat.

2 **Abbreviated title:** Instant attention capture by inferred threat

3 **Authors:** D Talmi^a, M Slapkova^a, MJ Wieser^b

4 **Affiliations:**

5 ^aDivision of Neuroscience and experimental psychology, University of Manchester, Oxford Road, M13
6 9PL, Manchester, UK

7 ^bErasmus School of Social and Behavioural Sciences, Department of Psychology, Education, and Child
8 Studies, Erasmus University Rotterdam , Burgemeester Oudlaan 50, 3062 PA Rotterdam, Netherlands

9

10

11 **Acknowledgements.** We thank N. Chater and C. Frith for the discussions that inspired this paradigm,
12 A. Wilkinson for help with data collection, J. Taylor and C. Charalambous for help with data analysis,
13 and M. Bauer and L. Hunter for helpful comments. DT acknowledges the support of ESRC First Grant
14 ES/I010424/1.

15

16 **Address for correspondence**

17 Deborah Talmi, Division of neuroscience and experimental psychology, School of Biological Sciences,
18 University of Manchester, Manchester, UK, M139PL.

19 Telephone: 0161 275 1968

20 Email: Deborah.Talmi@manchester.ac.uk

21

22

Abstract

23 Signals for reward or punishment attract attention preferentially, a principle termed ‘value-
24 modulated attention capture’ (VMAC). The mechanisms that govern the allocation of attention
25 resources can be described with a terminology that is more often applied to the control of overt
26 behaviours, namely, the distinction between instrumental and Pavlovian control, and between
27 model-free and model-based control. While instrumental control of VMAC can be either model-free
28 or model-based, it is not known whether Pavlovian control of VMAC can be model-based. To decide
29 whether this is possible we measured Steady-State Visual Evoked Potentials (SSVEPs) while 20
30 healthy adults took part in a novel task. During the learning stage participants underwent aversive
31 threat conditioning with two CSs, one that predicted pain (CS+) and one that predicted safety (CS-).
32 Instructions given prior to the test stage in the task allowed participants to infer whether novel,
33 ambiguous CSs (new_CS+/ new_CS-) were threatening or safe. Correct inference required combining
34 stored internal representations and new propositional information, the hallmark of model-based
35 control. SSVEP amplitudes quantified the amount of attention allocated to novel CSs on their very
36 first presentation, before they were ever reinforced. We found that SSVEPs were higher for new_CS+
37 than new_CS-. This result is potentially indicative of model-based Pavlovian control of VMAC, but
38 additional controls are necessary to verify this conclusively. This result underlines the potential
39 transformative role of information and inference in emotion regulation.

40

41 Introduction

42 Regulating established emotional responses upon receipt of novel information can be adaptive. For
43 example, it would be advantageous if, when a patient is told that their new bottle of medication has
44 fewer side effects than the one they used for some years prior, they immediately down-regulated
45 their feeling of anxiety about administering this medication, as well as concomitant cognitive,
46 physiological and behavioural responses, such as increased attention to the bottle. In this example
47 the bottle of medication has become, through the years of being paired with unpleasant side effects,
48 a signal for aversive outcomes. Signals for reward or punishment are known to attract attention
49 preferentially, a principle termed Value-Modulated Attention Capture (VMAC, Le Pelley et al., 2016).

50 The mechanisms that govern attention allocation can be productively described with a terminology
51 more commonly applied to the control of overt behaviours – the distinction between
52 instrumental/Pavlovian control, and model-free/model-based control (Dayan & Berridge, 2014). The
53 distinction between instrumental and Pavlovian control has to do with the dependencies between
54 behaviour and outcome (Mackintosh, 1983). In an instrumental learning task outcomes depend on
55 agents' behaviour, so agents act in order to increase utility – either increase the likelihood of reward,
56 or decrease the likelihood of punishment. An example for an instrumental control of VMAC is
57 increased attention to stimuli when told that doing so will be remunerated. By contrast, Pavlovian
58 control refers to behaviour that is triggered by stimuli that predict reward or punishment, even when
59 the outcomes are independent of the agent's behaviour, such as freezing in response to a tone that
60 predicts a pain. While Pavlovian control of VMAC was established in the case of rewarding stimuli,
61 there are also demonstrations of the same effect with aversive stimuli (Van Damme, Crombez,
62 Hermans, Koster, & Eccleston, 2006; L. Wang, Yu, & Zhou, 2013; Wentura, Müller, & Rothermund,
63 2014). One of the surest ways to be convinced that a particular behaviour is controlled by a

64 Pavlovian, rather than an instrumental, process is when it incurs a loss (Dayan, Niv, Seymour, & Daw,
65 2006). Pavlovian control of VMAC was elegantly demonstrated in an experiment that used an
66 omission schedule, where attending distractors that signalled reward magnitude resulted in the
67 omission of the reward (Le Pelley, Pearson, Griffiths, & Beesley, 2015). Because increased attention
68 to distractors that predicted high (compared to low) reward decreased the likelihood of reward and
69 was never itself rewarded, VMAC in that experiment could not be attributed to instrumental control.
70 Instead, the findings revealed a Pavlovian control of VMAC. Subsequent work showed that these
71 effects extended to separate tasks (Bucker & Theeuwes, 2017).

72 The distinction between instrumental and Pavlovian control is orthogonal to the distinction between
73 model-based and model-free control (Daw, Niv, & Dayan, 2005; Dayan & Berridge, 2014). A model,
74 according to Dayan and Berridge, is an internal representation of stimuli, situations and
75 environmental circumstances, which supports prospective cognition. Model-based responses are
76 executed when we infer, based on our model of the environment, that responding in a particular way
77 would maximise our expected utility. Model-based control can be contrasted to model-free control,
78 which depends on the accumulated average experience agents have with the outcomes associated
79 with a particular environmental state. Model-based control allows us to respond flexibly to a volatile,
80 changing environment; model-free control gives us the wisdom of the average experience. The
81 example above for instrumental control, where participants attend stimuli when told they will be
82 rewarded for doing so, is likely to be an example for *model-based* instrumental control. This is
83 because propositional information in instructions shape our model of the environment; we can take
84 up a suggestion or follow an instruction regardless of previous experience in a task (Olsson & Phelps,
85 2004). Model-based instrumental responses, such as those informed by instructions, can become
86 model-free if they are repeatedly executed and reinforced (Yin & Knowlton, 2006). For example, with

87 repeated pairing between attention to certain stimuli and reward attainment participants acquire a
88 habit to attend to those stimuli. The model-free nature of this behaviour is demonstrated when
89 participants continue to pay preferential attention to these stimuli even when further reinforcement
90 is unlikely (Luque et al., 2017).

91 Here we ask whether model-based Pavlovian control of VMAC is possible. Previous experiments in
92 animals have demonstrated model-based Pavlovian control of overt behaviour. For example, placing
93 animals in entirely new states, such as a salt-deprived state, instantly transforms the learned aversive
94 value of a lever that predicts a salty taste (Robinson & Berridge, 2013). The opening example
95 demonstrates what Model-based Pavlovian control of VMAC might look like in humans: the
96 information on the new medication revises the patient's model of the environment, yielding new
97 inferences that instantaneously transform the value they assign to the bottle and, consequently, the
98 attention she pays to it. Not much is currently known about Model-based Pavlovian control in
99 humans, although a recent study suggested that a model-based algorithm fitted conditioned threat
100 response in the amygdala better than model-free algorithms (Prévost, McNamee, Jessup, Bossaerts,
101 & O'Doherty, 2013), and the distinction between model-based and model-free Pavlovian control
102 resembles the one between cognitive and emotional control of Pavlovian responses in aversive
103 conditioning tasks (Sevenster, Beckers, & Kindt, 2012).

104 It is not, at present, known whether model-based Pavlovian control of VMAC is possible. Because
105 Pavlovian control of VMAC was evident even when participants had plenty of opportunity to learn
106 that the way they were allocating their attention was detrimental, and even when they were fully
107 informed about the nature of the omission schedule (Pearson, Donkin, Tran, Most, & Le Pelley,
108 2015), the Pavlovian control of VMAC may always be model-free (Le Pelley et al., 2016). The same
109 conclusion appears to be supported by findings that instructed extinction did not modulate the

110 classically-conditioned potentiated startle responses (Sevenster et al., 2012). Our aim was to test this
111 contention by using an optimised task. Evidence for model-based, Pavlovian control of VMAC will
112 confirm, in the domain of internal resource allocation, the distinction between model-free and
113 model-based Pavlovian control of incentive behaviour.

114 In the *conditioning stage* of the task participants passively viewed two Conditioned Stimuli (CSs),
115 which fully predicted a painful electric shock (CS+) or the absence of shock (CS-). We measured the
116 Steady-State Visual Evoked Potential (SSVEP), a validated neural signal of visual attention (Matthias
117 M. Müller, Teder-Sälejärvi, & Hillyard, 1998; Norcia, Appelbaum, Ales, Cottureau, & Rossion, 2015).
118 The SSVEP is known to be sensitive to VMAC, specifically value-modulated attentional engagement
119 and disengagement processes that occur from 500ms after stimulus presentation (Miskovic & Keil,
120 2013; Wieser, Miskovic, & Keil, 2016), so we expected greater SSVEP amplitudes for the CS+ than the
121 CS-. The *test stage* included a single presentation of two ambiguous, novel CSs, which we refer to as
122 “New_CSs” to contrast them to the “old_CSs” that participants experienced during the conditioning
123 stage. New_CSs were constructed such that their value could not be predicted by previous
124 experience. Before the test, participants received propositional information that, when combined
125 with their learned internal representation of the CSs, allowed them to infer the prospective value of
126 new_CSs. We only measured attention to the very first presentation of the new_CSs, before they
127 were ever reinforced (the entire task was repeated several times, but new stimuli were used in each
128 repetition).

129 While not a formal test of the model-based or the Pavlovian nature of this form of control over
130 VMAC, we argue that it would be difficult to account for differential attention to new_CSs in any
131 other way, an issue we return to in the discussion section. We hypothesised that participants would
132 be able to utilise stored information jointly with propositional information to control attention

133 allocation, and therefore expected greater SSVEP amplitudes for the new_CS+ compared to the
134 new_CS-.

135

136

Materials and methods

137 **Participants**

138 Twenty seven undergraduate students from the University of Manchester participated in the study in
139 exchange for course credits. None of the participants reported personal or family history of photic
140 epilepsy, none were taking centrally-acting medication, none had a history of psychiatric or
141 neurological disorders, and all had normal or corrected-to-normal vision. The experiment was
142 approved by the University of Manchester ethics committee. Three participants did not complete the
143 study and one participant did not exhibit an SSVEP signal. Three participants were excluded because
144 they failed the contingency awareness criterion (see below). This resulted in a total of 20 participants
145 (6 males, mean age 19.5, SD=1.15).

146 **Apparatus.**

147 The amplitude of the electric skin stimulation which served as a US (see below) was controlled by a
148 Digitimer DS5, an Isolated Bipolar Constant Current Stimulator (Digitimer DS5 2000, Digitimer Ltd.,
149 Welwyn Garden City, UK). The DS5 produces an isolated constant current stimulus proportional to a
150 voltage applied at its input. For reasons of participant safety this stimulator is limited to delivering a
151 maximum of 10V/10mA. Participants received the stimulation through a ring electrode built in-house
152 (Medical Physics, Salford Royal Hospital) attached to the DS5. To ensure adequate conductance
153 between the electrode and the skin, the back of each participant's hand was prepared with Nuprep
154 Skin Preparation Gel and Ten20 Conductive Paste prior to attaching the electrode. Transcutaneous

155 electrical stimulation activates myelinated A β somatosensory fibres as well as A δ nociceptive fibres
156 (Hird, Jones, Talmi, & El-Deredy, 2018) and can therefore cause both a sensation of touch and a
157 sensation of pain.

158 The experiment was implemented using the Psych toolbox on a Matlab (The Mathworks Inc., Natick,
159 MA, USA) platform. The voltage inputs to the DS5 were sent from Matlab through a data acquisition
160 interface (National Instruments, Austin, TX, USA). The behavioural ratings were taken on Microsoft
161 Excel.

162 **Materials**

163 **CSs.** Stimuli resembled Navon figures (Navon, 1977), in that they were composed of global and local
164 shapes where the outline of the large 'global' shape was created out of smaller 'local' shapes. To
165 create these stimuli we first created 48 unique shapes using Adobe Illustrator, each with a black
166 outline and white filling. These shapes were divided into 24 pairs so that the two shapes in each pair
167 were visually dissimilar (e.g. an arrow and a star). Each pair was used to create a subset of 4 Navon
168 figures, as shown in Figure 1. Two were congruent (e.g. a global arrow made of local arrows, a global
169 star made of local stars), and two incongruent (e.g. a global arrow made of local stars, a global star
170 made of local arrows). In total, the experiment used 24 such four-figure subsets (96 Navon figures).
171 All figures were created and presented in grayscale to minimise differences in colours and luminance.
172 Four-figure subsets were randomly allocated to experimental block. The two congruent figures were
173 randomly allocated to the "old_CS+" and "old_CS-" conditions. There were three types of task blocks,
174 as described below, termed global, local, and control blocks. The new_CS+ in 'Global' blocks used the
175 global attribute of the old_CS+ and the local attribute of the old_CS-. The new_CS+ in 'Local' blocks
176 used the global attribute of the old_CS- and the local attribute of the old_CS+.

177 figure subsets were used for the 2 practice blocks. The figures in practice blocks were created from 4
178 letters with one four-figure subset consisting of 'H' and 'O', and the other one 'Z' and 'I'.

179 **US.** The majority of studies of VMAC use rewarding USs, but there is also evidence for VMAC with
180 aversive outcomes, including pain (e.g. Wang et al., 2013). Here the US was a painful electric
181 stimulation delivered to the back of the right hand (see Apparatus, above).

182

183 **Procedure**

184 On arriving at the laboratory, participants were given an information sheet informing them of the
185 justification for the study and of the use of electrical stimulation. After they signed the consent form,
186 participants were fitted with the electroencephalogram (EEG) cap, and sat in a dimly-lit and sound-
187 attenuated room, 90 cm in front of the monitor screen, where an electrode was attached to the
188 dorsum of their right hand. Once the electrode was attached the participants undertook a series of
189 procedures, described below, in the following order: pre-experiment rating of materials, pain
190 calibration, habituation, experimental task, and post-experiment rating of materials.

191 ***Pre- and post-experiment rating of likability and contingency.*** Participants were presented with all
192 of the figures and rated how much they liked each one using a 5-point Likert scale (likability rating
193 task). They then saw all figures again, and guessed, by entering a percentage, how likely each figure
194 was to be followed by a painful stimulation (contingency rating task). The order of the figures in each
195 rating task was randomised for each participant. The likability and contingency rating tasks were
196 repeated at the end of the experiment.

197 ***Pain calibration.*** This procedure ensured that participant could tolerate the stimulation, and that the
198 stimulations were psychologically equivalent across participants. During this procedure participants

199 received a series of stimulations starting from 0.4mA, and incrementing by 0.4mA at each step.
200 Participants rated each stimulation on a scale from 0 – 10 where a score of 0 reflected not being able
201 to feel the stimulation, 3 reflected a stimulation level that was on the threshold of being painful, 7
202 related to a stimulation that was deemed painful, but still tolerable, and 10 related to ‘unbearable
203 pain’. The scaling procedure was terminated once the participant reported the level of stimulation as
204 being equivalent to ‘7’ on the scale. This calibration procedure was performed twice to allow for
205 initial habituation/sensitisation to the stimulation. The power levels that induced a rating of ‘7’ on
206 the second run of the calibration procedure were used as US. The US was therefore a painful but
207 tolerable sensation.

208 ***Habituation and method of CS presentation.*** Participants passively viewed a randomised list of all of
209 the CS figures. CS figures were displayed at the centre of the screen, a 17” monitor with a resolution
210 of 1024x768 pixels and a refresh rate of 60Hz. The duration of the presentation of each CS was
211 3,300ms. That time included 66 on-off cycles in which CS figures were displayed on a uniform white
212 background for 33.3ms (‘on’) and the screen turned black for 16.6ms (‘off’), resulting in a 20Hz
213 flickering display. The inter-trial interval between CSs was 2,500ms, during which the screen was
214 white.

215 ***Experimental task.*** We designed a novel task to reveal model-based Pavlovian control of VMAC. A
216 schematic of the task is shown in Figure 1. The task progressed through two stages - a conditioning
217 stage with 24 trials and a test stage with 4 trials, which are described in detail below. Each trial
218 included the presentation of a CS; when this was a CS+, the trial always terminated with US delivery.
219 Crucially, the logic of the task necessitated an extremely brief test stage that yielded only a single
220 trial for the contrast of interest. This was necessary in order to ensure that VMAC could not be
221 controlled through a model-free process; once new_CSs were reinforced, that reinforcement could

222 inform the value assigned to new_CSs in their second presentation. Therefore, we needed to
223 measure attention to new_CSs upon their first presentation, before they were reinforced, to prevent
224 any possibility that threat value could be informed by the experience of reinforcement. This
225 requirement led us to measure attention using SSVEPs (Matthias M. Müller et al., 1998). SSVEPs have
226 excellent SNR compared to ERP and time-frequency analysis of EEG measurements (Norcia et al.,
227 2015; Wieser et al., 2016), which can even allow a measure of effects at the single-trial level (Keil et
228 al., 2008). This is because the neurons that generate the SSVEP signal oscillate at the driving
229 frequency, which is precise and known a-priori, such that it is less affected by background EEG noise.
230 Even within the narrow band, the time-locking of oscillation phase to the stimulus (here, the CS)
231 decreases the impact of non-experimentally-driven oscillations. In addition, a-priori knowledge of the
232 frequency band and the electrodes sensitive to it prevents the need to correct for multiple
233 comparisons.

234 The task was repeated once in each of 24 task blocks. Each task block used novel stimuli, as described
235 in the material section, preventing the transfer of learning across blocks. Each task block lasted 2.5
236 minutes, with a 5-second break between blocks. Participants practiced the experimental task before
237 it commenced in two practice blocks, using the practice materials described above.

238 Before the experimental task began participants were given instructions for the experimental task.
239 They were asked to fixate on the fixation cross throughout each block, to observe the figures, and to
240 pay attention to the relationships between the figures and the pain. To encourage compliance,
241 participants were told that their memory for these associations will be tested. This instruction does
242 not privilege memory for the CS+ compared to the CS-, and therefore cannot be responsible for
243 observed threat responses.

244 *Conditioning stage.* During the conditioning stage, participants learned that one figure (old_CS+) 245 always predicted pain but another (old_CS-) was safe. CSs were fully predictive of their respective 246 outcomes to reduce any effects of stimulus predictability and of uncertainty, which are tightly 247 intertwined with the effect of value on attention control (Le Pelley et al., 2016). We used previous 248 trial-by-trial dissection of threat effects on the SSVEP signal (Wieser, Miskovic, Rausch, & Keil, 2014) 249 to decide how many conditioning trials were necessary in the conditioning stage. They observed a 250 significant modulation of the SSVEP by aversive reinforcement was observed after 5-10 conditioning 251 trials. Therefore, here we used 12 conditioning trials with each CS. The old_CS+ figure and the 252 old_CS- figure were presented 12 times each, at a random order. The details of how each CS was 253 presented was the same as during habituation, but here, when old_CS+ figures were presented, the 254 US was delivered during the very last cycle, at the offset of the last 'on' screen.

255 *Test stage.* After the conditioning stage was completed, participants viewed one of three possible 256 instructions for 10s. In the experimental condition the instruction was the word 'global' or the word 257 'local'. These words indicated the terms under which the US was to be delivered in the test stage, 258 namely, whether the global or local attribute of the old_CS+ would be reinforced. In the control 259 condition the instruction was a meaningless alphanumeric string, which gave participants no 260 information as to which attribute of the old_CS+ would be reinforced.

261 Four trials were presented after the instructions. The first two included the presentation of old_Css 262 (their order was randomised), and the last two the presentation of new_CSs (their order was also 263 randomised). New_CSs were the "other" two figures from the same four-figure subset from which 264 the old_CSs were drawn. As can be appreciated from examining Figure 1, each of the new_CSs 265 consisted of one previously-reinforced attribute and one previously-safe attribute. The old_CS+ and 266 the new_CS+ were reinforced; the old_CS- and the new_CS- were not.

267 Participants did not see the new_CSs before, so without the instructions they could not predict which
268 one would be reinforced. The only way for participants to predict whether the US will follow the
269 new_CS+ or the new_CS- was to infer this from the instructions by drawing on their memory of
270 old_CSs. For example, after the instruction 'local' participants who remembered the local attribute of
271 the old_CS+ could infer that (1) the global attribute of new_CSs did not determine whether the US
272 will be delivered or not (2) the US will follow any new_CS with the same local attribute as the
273 old_CS+. Taken together, such participants would predict pain after the new_CS+ but not after the
274 new_CS-. New_CSs were reinforced in accordance with the instructions, confirming participants'
275 expectations. Old_Css were reinforced in accordance both with the contingencies established during
276 the conditioning stage and the instructions.

277 While in previous research a newly-acquired conditioned response could be observed within 5-10
278 trials with each CS (Wieser et al., 2014), the test stage in each task block here only yielded just one
279 trial in each condition. To increase SNR the same structure described thus far – a conditioning stage
280 followed by the test stage - was repeated in each task block. 16 task blocks were allocated to the
281 experimental condition (8 with 'global' and 8 with 'local' instructions) and 8 were allocated to the
282 control condition.

283 **EEG recording and analysis**

284 **EEG recording.** Continuous EEG recordings were obtained from a 64-channel cap with in-build
285 electrodes (Biosemi Active Two) using the 10-20 configuration system. Data were digitised at a rate
286 of 2048Hz and filtered online between 0.1 and 100 Hz. The recording was referenced online to the
287 Common Mode Sense active electrode. The Driven Right Leg passive electrode was used as ground.
288 The impedance was kept below 40k Ω . Eye movement and blinks were recorded from horizontal and
289 vertical electro-oculogram.

290 **Preprocessing of EEG data.** Data were analysed using SPM12. They were converted from their native
291 format and then filtered with three 2nd order Butterworth IIR zero-phase forward and reverse filters:
292 a 1Hz highpass, a 80Hz lowpass, and a 49.50Hz-50.5Hz notch filter to remove mainline noise. Data
293 were downsampled to 200Hz and re-referenced to the average of all electrodes. Eye blinks and
294 saccades were marked on the VEOG and HEOG channels (or Fp1 in two participants) using an
295 automatic algorithm that was thresholded separately for each participant.

296 Further pre-processing was conducted for the purpose of complex demodulation (see below). Data
297 were segmented between -600ms before the onset of CSs to 3250ms after CS onset (Just before the
298 offset of the CS/US onset, 3300ms after CS onset). Segments where the following artefacts were
299 present on occipital channels (Oz, POz, O1, O2, O3, O4) were rejected: jumps greater than 150 μ V;
300 peak-to-peak differences greater than 250 μ V; flat segments. Channels where more than 20% of the
301 trials were rejected were excluded from analysis. This left, on average, 282.62 learning trials and 7.83
302 test trials with each CS in each condition. Artefacts associated with eye blinks and saccades were
303 corrected using the singular value decomposition (SVD) technique implemented in SPM12 which
304 captures eye artefacts and removes the associated component.

305 **Complex demodulation.** Threat effects were operationalised as an increased response to the CS+
306 compared to CS-. Previous work suggested that threat effect are more pronounced later in the
307 presentation of the CS, because the threat response is greater when threat is imminent, and the
308 perceptual processing of predictive sensory features of the CS is enhanced only when the US is
309 imminent (Miskovic & Keil, 2012; Paterson & Neufeld, 1987). Complex demodulation was therefore
310 carried out to determine exactly when threat effects were present. Complex demodulation was
311 conducted using SPM12 on the entire segment, with the multitaper method, a hanning window, and
312 a resolution of 1Hz. Data in each condition were averaged using robust averaging, a method that

313 down-weights outliers (Litvak et al., 2011). The signal from electrodes Oz and POz was extracted
314 around the driving frequency (19-21Hz). These data were averaged across all 12 trials in each
315 condition and the 24 blocks of the experimental task (288 trials for each CS). In agreement with
316 Miskovic and Keil (2012), threat effects were greatest during the second half of the presentation of
317 the CS (Figure 2). An examination of the topographies associated with threat supported our selection
318 of electrodes of analysis. We used these results to constrain our spectral analysis.

319 **Spectral analysis.** Based on the results of the complex demodulation step, spectral analysis was
320 conducted using the spatiotemporal window of 1500-3000ms from CS onset, at Oz and POz. We
321 followed the method in Keil (2008) to maximise sensitivity to differences of interest by minimising 20Hz signal
322 that does not keep in phase with the stimulation. Within the time and spatial windows defined above, 26
323 windows of 250ms (5 cycles of the SSVEP) were segmented for each trial. The first window started at 1500ms
324 from CS onset and each subsequent window started 50ms (1 SSVEP cycle) later; the last window started
325 2750ms from CS onset. For each trial, these 26 windows were then averaged in time, resulting in averages that
326 corresponded to single trials. We then averaged across all of the single trial averages from the same condition
327 (Keil et al. did not perform this last step because they were interested in single-trial data). These condition-
328 wise averages were Fourier-transformed using the FFT function in Matlab.

329 **Statistical analysis.** Across participants, some of the data significantly diverged from normality,
330 according to the Shapiro-Wilk test, and all data were therefore log-transformed. Data from the
331 conditioning stage in each experimental block were binned (4 bins; 3 trials per bin), averaged across
332 the 24 task blocks, and analysed with a 2 (threat: CS+, CS-) x 4 (bins in the conditioning stage). The
333 main hypothesis concerned the response to new_CSs during the test stage. The response to new_CS+
334 was compared to the response to new_CS- across the 16 blocks in the experimental condition
335 (collapsing across the 8 blocks in the 'global' and the 8 blocks in the 'local' conditions) using a one-
336 sample t-test. These responses were also compared to the response to responses to the two

337 ambiguous CSs that followed the 8 control blocks where no instructions were given using two
338 Bonferroni-corrected two-sample t-tests (the averages that went into these comparisons comprised
339 of 16 trials in each condition).

340

341 **Results**

342 **Behavioural results**

343 **Sample selection.** Behavioural and EEG data were first checked to verify the presence of a
344 conditioned threat response. We excluded participants who may have disengaged from the task
345 based on their 'learning score', computed as the increased contingency ratings given to old_CS+
346 stimuli after the experiment compared to the pre-experiment rating given to the same stimuli.
347 Increased ratings must be based on learning that occurred during the experiment, and therefore
348 reflects at least a minimal level of engagement. The learning score purposefully ignores the
349 magnitude of the conditioned response, which could be computed as differential ratings of the
350 old_CS+ and old_CS-, in order not to bias the selection of the sample. To be included participants had
351 to show a numerical (above 0) increase in ratings. Based on this criterion, three participants were
352 excluded from analysis, leaving a sample of N=20.

353 **Manipulation check.** Contingency and liking ratings for the stimuli were entered into two separate 3-
354 way repeated-measures ANOVAs with the factors time (pre-experiment, post-experiment), threat
355 (CS+, CS-), and status (old, new). The results evidenced a conditioning effect (Figure 3). In both the
356 analysis of contingency ratings and in the (separate) analysis of liking ratings the 3-way interaction
357 between time (pre-experiment, post-experiment), threat (CS+, CS-), and status (old, new) was
358 significant (contingency: $F(1,19)=34.95$, $P<.001$, partial $\eta^2=.65$; liking $F(1,19)=4.79$, $p=.04$, partial

359 eta=.20). We unpacked this interaction by examining the old and new CSs separately. The 2-way
360 interaction was significant for old_CSs (contingency: $F(1,19)=59.98$, $p<.001$, partial eta=.76; liking:
361 $F(1,19)=17.89$, $p<.001$, partial eta=.48) as well as for new_CSs (contingency: $F(1,19)=4.65$, $p=.04$,
362 partial eta =.20; liking: $F(1,19)=9.16$, $p=.007$, partial eta=.32). The results showed that participants
363 expected the US more frequently following the CS+ and liked the CS+ less than the CS-, and that
364 these effects were stronger for old_CSs than for new_CSs, probably because old_CSs were repeated
365 multiple times, but new_CSs were only presented once.

366 EEG results

367 **Manipulation check.** Threat effects on the SSVEP during the conditioning stage were analysed with a
368 2 (threat) x 4 (trial bins) repeated-measures ANOVA. The successful manipulation of threat in this
369 experiment was evident in differential responses to the CS+ and CS- during the conditioning stage
370 (Figure 4), $F(1,19)=5.50$, $p=0.02$, partial $\eta^2=0.22$. Although numerically the magnitude of threat
371 effects was only observable after 3 trials, the interaction of threat and trial bins was not significant,
372 $F(3,57)<1$, suggesting that they remained consistent across each of the experimental blocks. The
373 main effect of binned trials, $F(3,57)=6.98$, $p<.001$, partial $\eta^2=0.27$, was also significant, denoting an
374 overall decrease in attention as the block progressed. To verify that threat responses were also
375 obtained in the test stage, responses to old_CSs in the control condition, where participants had no
376 reason to make any model-based inferences, were analysed with a one-tailed paired t-test,
377 contrasting old_CS+ and old_CS-. We observed differential responses to the old_CS+ and old_CS-,
378 $t(19)=2.36$, $p=.015$ (one-tailed), Cohen's $d = 0.53$. Because only the data from the conditioning stage
379 were pre-selected, the data from the test stage provide a useful confirmation.

380 **Main hypothesis.** Our main hypothesis was that responses to the new_CS+ would be greater than
381 responses to the new_CS-. The hypothesis was evaluated with a one-tailed paired t-test comparing

382 SSVEPs to the new_CS+ and new_CS-. As predicted, SSVEP amplitudes were higher during the
383 presentation of the new_CS+ compared to the new_CS-, $t(19)=2.22$, $p=.02$, Cohen's $d = 0.50$.
384 Additional 2-tailed t-tests, controlled for multiple comparisons with a p-value of 0.025, compared
385 each of these SSVEPS to the SSVEP elicited by ambiguous new_CSs in the control condition. SSVEP
386 amplitudes were equivalent during the presentation of the new_CS+ and the new_CSs in the control
387 condition, $t<1$, suggesting that participants experienced ambiguous figures as threatening when they
388 could not use the instructions to disambiguate them. This interpretation is supported by a significant
389 difference between the new_CSs in the control condition and the new_CS-, $t(19)=2.58$, $p=0.018$,
390 Cohen's $d = 0.58$, suggesting that VMAC was attenuated when participants knew that pain was
391 unlikely.

392

393

Discussion

394

395 During the learning stage of our Pavlovian conditioning task, we observed an increase in the
396 amplitude of the SSVEP signal towards stimuli with learned aversive value. These results are not
397 surprising, given much evidence that the valuation system can control attention allocation (Le Pelley
398 et al., 2016), but important because they confirm relatively limited evidence for Pavlovian control of
399 VMAC towards stimuli with aversive value (Van Damme et al., 2006; L. Wang et al., 2013; Wentura et
400 al., 2014; Wieser et al., 2014), which is less established than Pavlovian control of VMAC towards
401 reward, or the effect of value on instrumental control of VMAC. Uniquely, these results also suggest a
402 way to observe the neural evolution of VMAC across the learning process, and even in the first trial,

403 something that has not been possible using other neuroimaging investigation (Olsson & Phelps, 2007;
404 Phelps et al., 2001) or in animal models (Balcarras, Ardid, Kaping, Everling, & Womelsdorf, 2016).

405 Our key result was that SSVEP amplitudes were larger when participants were presented with a new
406 shape that they inferred predicted physical pain (the new_CS+), compared to a new shape that they
407 inferred predicted safety (the new_CS-). Because the amplitude of SSVEPs is higher for attended
408 stimuli compared to unattended ones (Matthias M. Müller et al., 1998; Muller et al., 1998), our
409 findings suggest that more attention was allocated to the new_CS+ compared to the new_CS-. We
410 argue that the differential attentional response to new_CS+ and new_CS- suggests a Pavlovian
411 model-based control of VMAC, an argument that we dissect in the 'theoretical considerations'
412 section, below. It has to be noted that while the SSVEP is known to be sensitive to VMAC such as in
413 aversive learning (Kastner-Dorn, Andreatta, Pauli, & Wieser, 2018; Miskovic & Keil, 2013; Wieser et
414 al., 2016) and to emotional stimuli in general (Keil et al., 2003; Keil, Moratti, Sabatinelli, Bradley, &
415 Lang, 2005; McTeague, Shumen, Wieser, Lang, & Keil, 2011; Wieser et al., 2016), heightened ssVEP
416 amplitudes are also found in response to increased working-memory load (Silberstein, Nunez,
417 Pipingas, Harris, & Danieli, 2001), and for attended relative to unattended stimuli (Morgan, Hansen,
418 & Hillyard, 1996) (Hillyard et al., 1997)(M. M. Müller, Malinowski, Gruber, & Hillyard, 2003; Matthias
419 M. Müller et al., 1998). Thus, the enhanced SSVEP amplitudes for may reflect any of these respective
420 processes.

421 By contrast to the experimental blocks, where the value of new_CSs could be inferred through a
422 combination of the instructions and stored memories of the learning stage, in control blocks no
423 instructions were given. Therefore, the threat value of control new_CSs was ambiguous, and
424 participants could not predict which one would be followed by pain. These ambiguous new_CSs in
425 the control condition attracted increased attention compared to the new_CS-. This result, which

426 suggests orienting towards ambiguous stimuli, accords with previous findings, where instructed,
427 ambiguous, novel CSs gave rise to increased physiological arousal and increased activation in the
428 amygdala and the insula (Phelps et al., 2001).

429 Our study confirms other demonstrations where instructions about Pavlovian contingencies
430 encourage responses that mimic the effect of associative learning through experience (reviewed in
431 Mitchell, De Houwer, & Lovibond, 2009). This is the first demonstration that propositional
432 information – a form of instruction – influences Pavlovian control of attention allocation. This
433 demonstration is particularly important because a previous experiment suggested that in Pavlovian
434 tasks attention allocation obeys associative learning principles, and is immune to propositional
435 knowledge. In Moratti and Keil's study (Moratti & Keil, 2009) the SSVEF during CS presentation
436 (steady-state visual evoked field, measured with MEG) increased with increased number of
437 sequentially reinforced CSs, not with increased US expectancy, which, in turn, only increased when
438 previous CSs have consistently not been reinforced. Indeed, other studies have also observed that
439 similar 'gambler's fallacy'-like paradigms give rise to conditioned responses that are based on model-
440 free, not model-based value (Clark, Manns, & Squire, 2001; Perruchet, 1985). The results of Moratti
441 and Keil, indicative of attention allocation, were particularly intriguing because they appeared to
442 contradict evidence that expectancy influences visual attention (Downing, 1988). Taken together, it
443 appears that associative mechanisms dominate attention allocation in the gambler's fallacy
444 paradigm, at least when using a delayed conditioning procedure (Clark et al., 2001), while in other
445 paradigms – including those using delayed conditioning, as we did here – propositional information
446 holds more sway. It is possible that the balance between these mechanisms is affected by the
447 certainty of each system in its threat evaluation (Daw et al., 2005).

448 While attentional responses were affected by propositional information here, it is possible that other
449 classically-conditioned responses were not. In particular, because attention allocation was affected in
450 the first trial it bears stronger resemblance to US-expectancy ratings and to classically-conditioned
451 skin conductance responses, which were immediately influenced by instructed extinction, than to
452 potentiated startle, which was not (Sevenster et al., 2012). Further research is required to examine
453 this potential dissociation using our paradigm, and to verify whether instructed extinction, like
454 instructed threat, also alters VMAC instantaneously.

455 By using neural measures to index control of VMAC we move a little closer to understanding how
456 propositional information is implemented at the level of the neurobiological mechanism. Increased
457 SSVEPs during the presentation of threatening stimuli is thought to be driven by re-entrant
458 connections from the amygdala, ACC and OFC, which amplify the processing of adaptive information
459 (Miskovic & Keil, 2012). Repeated pairing between a stimulus and pain can change the neural
460 representation of the pain-predicting stimulus. For example, repeated pairing between a tone and a
461 painful shock change the tuning frequency of neurons that encode these tones, and stimulation of
462 the amygdala is sufficient to produce this effect (Chavez, McGaugh, & Weinberger, 2013). Here,
463 however, such a process could not occur because attention modulation was manifested before the
464 reinforcement itself. A meta-analysis of studies of instructed fear found that the dorsomedial
465 prefrontal cortex is uniquely associated with a conscious appraisal process (Mechias, Etkin, & Kalisch,
466 2010). Similarly, the same region has been shown to dynamically modulate model-free valuation in
467 the OFC, striatum, and hippocampus (Li, Delgado, & Phelps, 2011). It is therefore likely that increased
468 response to the new_CS+ was due to projections from the dorsomedial prefrontal cortex to the OFC
469 and ACC, regions that are strongly connected to the amygdala and able to modulate its activity (Lee,

470 Heller, van Reekum, Nelson, & Davidson, 2012; Schiller & Delgado, 2010), with downstream re-
471 entrant effects in the visual cortex.

472 The constrained data yield of the paradigm should be acknowledged as a limitation of this study.
473 While the effect sizes in all of the statistical tests were all of a 'medium' size, according to Cohen's
474 classification (Cohen, 1988), the study should be replicated in order to increase confidence in this
475 novel result. For the same reason, we could not explore the influence of 'dimension' (global or local)
476 in the results we obtained, because this would have halved the number of trials that we could
477 analyse.

478

479 **Theoretical considerations.**

480 We argue that increased SSVEPs to new_CS+ compared to new_CS- in this experiment was likely due
481 to a Pavlovian, not an instrumental process. The paradigm was entirely passive; pain outcomes were
482 independent of participants' behaviour or how they allocated their attention. Participants could not
483 benefit from allocating differential attention to specific CSs. Indeed, participants were told explicitly,
484 and also knew through experience across the 24 blocks of the task, that the stimulation levels were
485 pre-determined and that they could therefore not influence it. Participants had no reason to allocate
486 more attention to the new_CS+ in order to increase success in the post-experiment rating task,
487 because they have already completed it once before they started the experiment, and knew,
488 therefore, that performance would benefit equally from attending all of the stimuli.

489 Although the control of VMAC here could not influence objective outcomes, it is possible, in
490 principle, that it incurred some *internal* benefit. Specifically, paying extra attention to threatening CSs
491 here may have decreased subjective pain. It is important to consider this possibility because

492 expected and experienced pain are not true reflections of objective tissue damage. Instead, pain
493 experience is strongly modulated by pain expectations (Atlas & Wager, 2012; Berns et al., 2006;
494 Morley, Vlaeyen, & Schrooten, 2012; Paterson & Neufeld, 1987; Tabor, Thacker, Moseley, & Körding,
495 2017; Vlaev, Seymour, Dolan, & Chater, 2009), which modulate endogenous analgesic mechanisms
496 (Anchisi & Zanon, 2015; Tracey, 2010; Wager et al., 2004), and experimental pain expectations are
497 themselves influenced by pre-existing individual biases (Hoskin et al., n.d.).

498 Yet closer scrutiny suggests that it is unlikely that attending the new_CS+ triggered endogenous
499 analgesia. While participants needed to attend experimental new_CSs to decipher exactly which one
500 predicted pain and which one did not, this need not result in *differential* attention to the two
501 new_CSs. The US was completely predictable, always of the same intensity, and presented at the
502 same time, so attending the new_CS+ is unlikely to have altered pain expectations meaningfully. It is
503 possible that attending the new_CS+ increased the precision of expectations (Kok, Rahnev, Jehee,
504 Lau, & de Lange, 2012). However, expecting high pain with greater certainty would increase
505 subjective pain, not decrease it (Hird et al., 2018). In fact, much evidence suggests that distraction,
506 not attention, is an effective pain-coping strategy (Buhle, Stevens, Friedman, & Wager, 2012;
507 Eccleston, 1995; Sharar et al., 2016; Weiss, Dahlquist, & Wohlheiter, 2011).

508 To fully establish that a response is controlled by a Pavlovian process researchers have utilised
509 omission schedules, where the response incurs a cost (Le Pelley et al., 2016; Mackintosh, 1983). This
510 method is achievable for model-free Pavlovian responses, but in the case of model-based Pavlovian
511 control known eventual costs should, by definition, alter the world-model that inspired the responses
512 in the first place. It is therefore potentially tricky to utilise this method to test that control was
513 Pavlovian and model-based. In summary, although the feeling of pain is malleable, and we have not
514 conclusively demonstrated that attention was controlled through a Pavlovian rather than an

515 instrumental process, we can presently think of no a-priori reason that attending the new_CS+ in this
516 experiment would be advantageous. The intuitive sense that we would all *want* to look – that we
517 might not be able to attend anything else – is perhaps simply the reflection of Pavlovian
518 misbehaviour (Dayan et al., 2006).

519 During the test stage participants could combine the propositional information provided to them in
520 the instructions (about the dimension that will be reinforced - global or local) with their stored
521 representation of the reinforced attribute (the particular shape that predicted the US during the
522 conditioning stage) to form a prediction of the value of new_CSs. That they have, indeed, done so is
523 evident in the differential attention they allocated to new_CSs in experimental blocks. While we did
524 not test the model-based nature of control of VMAC formally, e.g. by using a two-step task (Otto,
525 Skatova, Madlon-Kay, & Daw, 2015), increased attention to the new_CS+ compared to the new_CS-
526 here involved “prospective cognition, formulating and pursuing explicit possible future scenarios
527 based on internal representations of stimuli, situations and environmental circumstances” – the
528 hallmark of model-based Pavlovian control according to Dayan and Berridge (2014, p. 5).

529 While new_CSs were constructed such that previous experience would render them equally
530 ‘threatening’ and ‘safe’, and render it difficult for a model-free algorithm to implement differential
531 responses to these CSs, it is possible that the instructions created a model which then trained a
532 model-free controller, as in schemes such as Dyna (Sutton, Szepesvári, Geramifard, & Bowling, 2008)
533 or “biased” learning (Doll, Jacobs, Sanfey, & Frank, 2009), or, alternatively, trigger an existing model-
534 free controller. Regarding the first alternative, it is clear that the limits of model-free reinforcement
535 learning are constantly expanding with the introduction of meta-reinforcement learning (J. X. Wang
536 et al., 2016). Future computational work could explore whether indirect reinforcement learning can
537 link stored knowledge and propositional information within a few seconds and a single trial to

538 influence the value predicted in a novel state. In the “biased” learning scheme, for example, the
539 training of the model-free controller relied on a modulation of the value of reinforcers (Doll et al.,
540 2009), so it may not be realistic to expect such training in measurements that take place prior to any
541 reinforcement, as in the present paradigm.

542 The second alternative is that the instructions created a model that triggered model-free habits to
543 control VMAC. Organisms may habitually attend more highly valued (compared to neutral) stimuli
544 preferentially because this incurs reward over the long-term, even though it does not do so in a
545 particular scenario. Attention to valued stimuli could improve the encoding of CSs, strengthen their
546 memory traces, and thus facilitate optimal decisions when the opportunity arises to act on the same
547 stimuli (Lieder, Griffiths, & Hsu, 2018). Additionally, prediction error minimisation – something that is
548 considered globally adaptive (Pezzulo, Rigoli, & Friston, 2015) - may be facilitated if the excellent
549 encoding of previously-valued stimuli increases the precision of the model we have of the world
550 around us. The model constructed by the instructions could therefore simply indicate to the system
551 which stimuli are likely to have large absolute value, as well as which have value which is highly
552 uncertain; but the actual attentional control may be carried out through the habitual mechanism.

553 Finally, it is possible that VMAC was, indeed, controlled by a model-based, Pavlovian process. In that
554 case, following the experimental instructions, participants may have constructed a model of the
555 new_CSs and their predictive value by combining the new propositional information and stored
556 internal representations. It is possible that while they viewed the instructions, participants recalled
557 old_CSs and generalised their aversive value to imagined new_CSs. It is also possible that participants
558 volitionally inhibited the representation of the global or local dimensions of memorised learned CSs
559 as well as actual new_CSs, to support generalisation from the learning to the test stage. Dayan and

560 Berridge (2014) discuss such recall and revaluation processes as mechanisms that allow model-based
561 Pavlovian control.

562

563 At the experiential level, increased attention to the new_CS+ suggests that the information given to
564 participants worked as an emotion regulation technique – it rendered ambiguous stimuli instantly
565 threatening. Drawing the connection between model-based and model-free control, on the one
566 hand, and cognitive and emotional control, on the other (Sevenster et al., 2012) can help the quest to
567 ground emotion regulation and behaviour change techniques more tightly in computational theories
568 (Etkin, Büchel, & Gross, 2016).

569

570

571 **References**

- 572 Anchisi, D., & Zanon, M. (2015). A Bayesian perspective on sensory and cognitive integration in pain
573 perception and placebo analgesia. *PLoS ONE*. <http://doi.org/10.1371/journal.pone.0117270>
- 574 Atlas, L. Y., & Wager, T. D. (2012). How expectations shape pain. *Neuroscience Letters*.
575 <http://doi.org/10.1016/j.neulet.2012.03.039>
- 576 Balcarras, M., Ardid, S., Kaping, D., Everling, S., & Womelsdorf, T. (2016). Attentional Selection Can Be
577 Predicted by Reinforcement Learning of Task-relevant Stimulus Features Weighted by Value-
578 independent Stickiness. *Journal of Cognitive Neuroscience*, *28*(2), 333–349.
579 http://doi.org/10.1162/jocn_a_00894
- 580 Berns, G. S., Chappelow, J., Cekic, M., Zink, C. F., Pagnoni, G., & Martin-Skurski, M. E. (2006).
581 Neurobiological substrates of dread. *Science (New York, N.Y.)*, *312*(5774), 754–8.
582 <http://doi.org/10.1126/science.1123721>
- 583 Bucker, B., & Theeuwes, J. (2017). Pavlovian reward learning underlies value driven attentional
584 capture. *Attention, Perception, & Psychophysics*, *79*(2), 415–428.
585 <http://doi.org/10.3758/s13414-016-1241-1>
- 586 Buhle, J. T., Stevens, B. L., Friedman, J. J., & Wager, T. D. (2012). Distraction and Placebo.
587 *Psychological Science*, *23*(3), 246–253. <http://doi.org/10.1177/0956797611427919>
- 588 Chavez, C. M., McGaugh, J. L., & Weinberger, N. M. (2013). Activation of the basolateral amygdala
589 induces long-term enhancement of specific memory representations in the cerebral cortex.
590 *Neurobiology of Learning and Memory*, *101*, 8–18. <http://doi.org/10.1016/j.nlm.2012.12.013>
- 591 Clark, R. E., Manns, J. R., & Squire, L. R. (2001). Trace and delay eyeblink conditioning: contrasting
592 phenomena of declarative and nondeclarative memory. *Psychological Science*: *A Journal of the*
593 *American Psychological Society / APS*, *12*(4), 304–308. <http://doi.org/10.1111/1467-9280.00356>
- 594 Cohen, J. (1988). *Statistical power analysis for the behavioral sciences*. *Statistical Power Analysis for*
595 *the Behavioral Sciences* (2nd ed.). Lawrence Erlbaum Associates, Hillsdale, NJ.
596 <http://doi.org/10.1234/12345678>
- 597 Daw, N. D., Niv, Y., & Dayan, P. (2005). Uncertainty-based competition between prefrontal and
598 dorsolateral striatal systems for behavioral control. *Nature Neuroscience*, *8*(12), 1704–1711.
599 <http://doi.org/10.1038/nn1560>
- 600 Dayan, P., & Berridge, K. C. (2014). Model-based and model-free Pavlovian reward learning:
601 revaluation, revision, and revelation. *Cognitive, Affective & Behavioral Neuroscience*, *14*(2), 473–
602 92. <http://doi.org/10.3758/s13415-014-0277-8>
- 603 Dayan, P., Niv, Y., Seymour, B., & Daw, N. D. (2006). The misbehavior of value and the discipline of
604 the will. *Neural Networks*, *19*(8), 1153–1160. <http://doi.org/10.1016/j.neunet.2006.03.002>
- 605 Doll, B. B., Jacobs, W. J., Sanfey, A. G., & Frank, M. J. (2009). Instructional control of reinforcement
606 learning: A behavioral and neurocomputational investigation. *Brain Research*, *1299*, 74–94.
607 <http://doi.org/10.1016/j.brainres.2009.07.007>

- 608 Downing, C. J. (1988). Expectancy and visual-spatial attention: Effects on perceptual quality. *Journal*
609 *of Experimental Psychology: Human Perception and Performance*, 14(2), 188–202.
610 <http://doi.org/10.1037/0096-1523.14.2.188>
- 611 Eccleston, C. (1995). The attentional control of pain: methodological and theoretical concerns. *Pain*,
612 63(1), 3–10. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/8577487>
- 613 Etkin, A., Büchel, C., & Gross, J. J. (2016). Emotion regulation involves both model-based and model-
614 free processes. *Nature Reviews Neuroscience*, 17(8), 532–532.
615 <http://doi.org/10.1038/nrn.2016.79>
- 616 Hillyard, S. a, Hinrichs, H., Tempelmann, C., Morgan, S. T., Hansen, J. C., Scheich, H., & Heinze, H. J.
617 (1997). Combining steady-state visual evoked potentials and f MRI to localize brain activity
618 during selective attention. *Human Brain Mapping*, 5(4), 287–292.
619 [http://doi.org/10.1002/\(SICI\)1097-0193\(1997\)5:4<287::AID-HBM14>3.0.CO;2-B](http://doi.org/10.1002/(SICI)1097-0193(1997)5:4<287::AID-HBM14>3.0.CO;2-B)
- 620 Hird, E. J., Jones, A. K. P., Talmi, D., & El-Deredy, W. (2018). A comparison between the neural
621 correlates of laser and electric pain stimulation and their modulation by expectation. *Journal of*
622 *Neuroscience Methods*, 293, 117–127. <http://doi.org/10.1016/j.jneumeth.2017.09.011>
- 623 Hoskin, R., Berzuini, C. C., Acosta-Kane, D., El-Deredy, W., Guo, H., & Talmi, D. (n.d.). Sensitivity to
624 Pain Expectations: A Bayesian Model of Individual Differences. *Cognition*.
- 625 Kastner-Dorn, A. K., Andreatta, M., Pauli, P., & Wieser, M. J. (2018). Hypervigilance during anxiety
626 and selective attention during fear: Using steady-state visual evoked potentials (ssVEPs) to
627 disentangle attention mechanisms during predictable and unpredictable threat. *Cortex, Online*.
628 <http://doi.org/10.1016/j.cortex.2018.05.008>
- 629 Keil, A., Gruber, T., Müller, M. M., Moratti, S., Stolarova, M., Bradley, M. M., & Lang, P. J. (2003). Early
630 modulation of visual perception by emotional arousal: evidence from steady-state visual evoked
631 brain potentials. *Cognitive, Affective & Behavioral Neuroscience*, 3(3), 195–206.
632 <http://doi.org/10.3758/CABN.3.3.195>
- 633 Keil, A., Moratti, S., Sabatinelli, D., Bradley, M. M., & Lang, P. J. (2005). Additive effects of emotional
634 content and spatial selective attention on electrocortical facilitation. *Cerebral Cortex*, 15(8),
635 1187–1197. <http://doi.org/10.1093/cercor/bhi001>
- 636 Keil, A., Smith, J. C., Wangelin, B. C., Sabatinelli, D., Bradley, M. M., & Lang, P. J. (2008).
637 Electrocortical and electrodermal responses covary as a function of emotional arousal: A single-
638 trial analysis. *Psychophysiology*, 45(4), 516–523. <http://doi.org/10.1111/j.1469-8986.2008.00667.x>
- 640 Kok, P., Rahnev, D., Jehee, J. F. M., Lau, H. C., & de Lange, F. P. (2012). Attention Reverses the Effect
641 of Prediction in Silencing Sensory Signals. *Cerebral Cortex*, 22(9), 2197–2206.
642 <http://doi.org/10.1093/cercor/bhr310>
- 643 Le Pelley, M. E., Mitchell, C. J., Beesley, T., George, D. N., Wills, A. J., & Le Pelley, M. (2016). Attention
644 and associative learning in humans: An integrative review. *Psychological Bulletin*, 142(10), 1111–
645 1140. <http://doi.org/10.1037/bul0000064>
- 646 Le Pelley, M. E., Pearson, D., Griffiths, O., & Beesley, T. (2015). When Goals Conflict With Values:

- 647 Counterproductive Attentional and Oculomotor Capture by Reward-Related Stimuli
648 Predictiveness-Driven Attentional Capture. *Experimental Psychology*, 144(1), 158–171.
649 <http://doi.org/10.1037/xge0000037>
- 650 Lee, H., Heller, A. S., van Reekum, C. M., Nelson, B., & Davidson, R. J. (2012). Amygdala-prefrontal
651 coupling underlies individual differences in emotion regulation. *NeuroImage*, 62(3), 1575–1581.
652 <http://doi.org/10.1016/j.neuroimage.2012.05.044>
- 653 Li, J., Delgado, M. R., & Phelps, E. A. (2011). How instructed knowledge modulates the neural systems
654 of reward learning. *Proceedings of the National Academy of Sciences*, 108(1), 55–60.
655 <http://doi.org/10.1073/pnas.1014938108>
- 656 Lieder, F., Griffiths, T. L., & Hsu, M. (2018). Overrepresentation of extreme events in decision making
657 reflects rational use of cognitive resources. *Psychological Review*, 125(1), 1–32.
658 <http://doi.org/10.1037/rev0000074>
- 659 Litvak, V., Mattout, J., Kiebel, S., Phillips, C., Henson, R., Kilner, J., ... Friston, K. (2011). EEG and MEG
660 data analysis in SPM8. *Computational Intelligence and Neuroscience*, 2011, 852961.
661 <http://doi.org/10.1155/2011/852961>
- 662 Luque, D., Beesley, T., Morris, R. W., Jack, B. N., Griffiths, O., Whitford, T. J., & Le Pelley, M. E. (2017).
663 Goal-Directed and Habit-Like Modulations of Stimulus Processing during Reinforcement
664 Learning. *The Journal of Neuroscience*: The Official Journal of the Society for Neuroscience,
665 37(11), 3009–3017. <http://doi.org/10.1523/JNEUROSCI.3205-16.2017>
- 666 Mackintosh, N. J. (1983). *Conditioning and associative learning*. Clarendon Press. Retrieved from
667 https://books.google.co.uk/books/about/Conditioning_and_Associative_Learning.html?id=a8x9AAAAMAAJ&redir_esc=y&hl=en
668
- 669 McTeague, L. M., Shumen, J. R., Wieser, M. J., Lang, P. J., & Keil, A. (2011). Social vision: Sustained
670 perceptual enhancement of affective facial cues in social anxiety. *NeuroImage*, 54(2), 1615–
671 1624. <http://doi.org/10.1016/j.neuroimage.2010.08.080>
- 672 Mechias, M.-L., Etkin, A., & Kalisch, R. (2010). A meta-analysis of instructed fear studies: Implications
673 for conscious appraisal of threat. *NeuroImage*, 49(2), 1760–1768.
674 <http://doi.org/10.1016/J.NEUROIMAGE.2009.09.040>
- 675 Miskovic, V., & Keil, A. (2012). Acquired fears reflected in cortical sensory processing: a review of
676 electrophysiological studies of human classical conditioning. *Psychophysiology*, 49(9), 1230–41.
677 <http://doi.org/10.1111/j.1469-8986.2012.01398.x>
- 678 Miskovic, V., & Keil, A. (2013). Perceiving threat in the face of safety: excitation and inhibition of
679 conditioned fear in human visual cortex. *The Journal of Neuroscience*: The Official Journal of
680 the Society for Neuroscience, 33(1), 72–8. <http://doi.org/10.1523/JNEUROSCI.3692-12.2013>
- 681 Mitchell, C. J., De Houwer, J., & Lovibond, P. F. (2009). The propositional nature of human associative
682 learning. *Behavioral and Brain Sciences*. <http://doi.org/10.1017/S0140525X09000855>
- 683 Moratti, S., & Keil, A. (2009). Not What You Expect: Experience but not Expectancy Predicts
684 Conditioned Responses in Human Visual and Supplementary Cortex. *Cerebral Cortex December*,
685 19, 2803–2809. <http://doi.org/10.1093/cercor/bhp052>

- 686 Morgan, S. T., Hansen, J. C., & Hillyard, S. A. (1996). Selective attention to stimulus location
687 modulates the steady-state visual evoked potential. *Proceedings of the National Academy of*
688 *Sciences*, 93(10), 4770–4774. <http://doi.org/10.1073/pnas.93.10.4770>
- 689 Morley, S., Vlaeyen, J. W., & Schrooten, M. G. S. (2012). Psychological interventions for chronic pain:
690 reviewed within the context of goal pursuit. *Pain Management*, 2(March), 1–10.
691 <http://doi.org/http://dx.doi.org/10.2217/pmt.12.2>
- 692 Müller, M. M., Malinowski, P., Gruber, T., & Hillyard, S. A. (2003). Sustained division of the
693 attentional spotlight. *Nature*, 424(6946), 309–312. <http://doi.org/10.1038/nature01812>
- 694 Muller, M. M., Picton, T. W., Valdes-Sosa, P., Riera, P., Teder-Salejarvi, W., & Hillyard, S. A. (1998).
695 Effects of spatial selective attention on the steady-state visual evoked potential in the 20-28 Hz
696 range. *Cognitive Brain Research*, 6, 249–261.
- 697 Müller, M. M., Teder-Sälejärvi, W., & Hillyard, S. A. (1998). The time course of cortical facilitation
698 during cued shifts of spatial attention. *Nature Neuroscience*, 1(7), 631–634.
699 <http://doi.org/10.1038/2865>
- 700 Navon, D. (1977). Forest before trees: the precedence of global features in visual perception.
701 *COGNITIVE PSYCHOLOGY*, 9, 353–383.
- 702 Norcia, A. M., Appelbaum, L. G., Ales, J. M., Cottureau, B. R., & Rossion, B. (2015). The steady-state
703 visual evoked potential in vision research: A review. *Journal of Vision*, 15(6), 4.
704 <http://doi.org/10.1167/15.6.4>
- 705 Olsson, A., & Phelps, E. A. (2004). Learned fear of “unseen” faces after pavlovian, observational, and
706 instructed fear. *Psychological Science*, 15(12), 822–828. [http://doi.org/10.1111/j.0956-](http://doi.org/10.1111/j.0956-7976.2004.00762.x)
707 [7976.2004.00762.x](http://doi.org/10.1111/j.0956-7976.2004.00762.x)
- 708 Olsson, A., & Phelps, E. A. (2007). Social learning of fear. *Nature Neuroscience*, 10(9), 1095–1102.
709 <http://doi.org/10.1038/nn1968>
- 710 Otto, A. R., Skatova, A., Madlon-Kay, S., & Daw, N. D. (2015). Cognitive Control Predicts Use of Model-
711 based Reinforcement Learning. *Journal of Cognitive Neuroscience*, 27(2), 319–333.
712 http://doi.org/10.1162/jocn_a_00709
- 713 Paterson, R. J., & Neufeld, R. W. J. (1987). Clear Danger: Situational Determinants of the Appraisal of
714 Threat. *Psychological Bulletin*. <http://doi.org/10.1037/0033-2909.101.3.404>
- 715 Pearson, D., Donkin, C., Tran, S. C., Most, S. B., & Le Pelley, M. E. (2015). Cognitive control and
716 counterproductive oculomotor capture by reward-related stimuli. *Visual Cognition*, 6285(May
717 2015), 1–26. <http://doi.org/10.1080/13506285.2014.994252>
- 718 Perruchet, P. (1985). A pitfall for the expectancy theory of human eyelid conditioning. *The Pavlovian*
719 *Journal of Biological Science*: *Official Journal of the Pavlovian*, 20(4), 163–170.
720 <http://doi.org/10.1007/BF03003653>
- 721 Pezzulo, G., Rigoli, F., & Friston, K. (2015). Active Inference, homeostatic regulation and adaptive
722 behavioural control. *Progress in Neurobiology*, 134, 17–35.
723 <http://doi.org/10.1016/J.PNEUROBIO.2015.09.001>

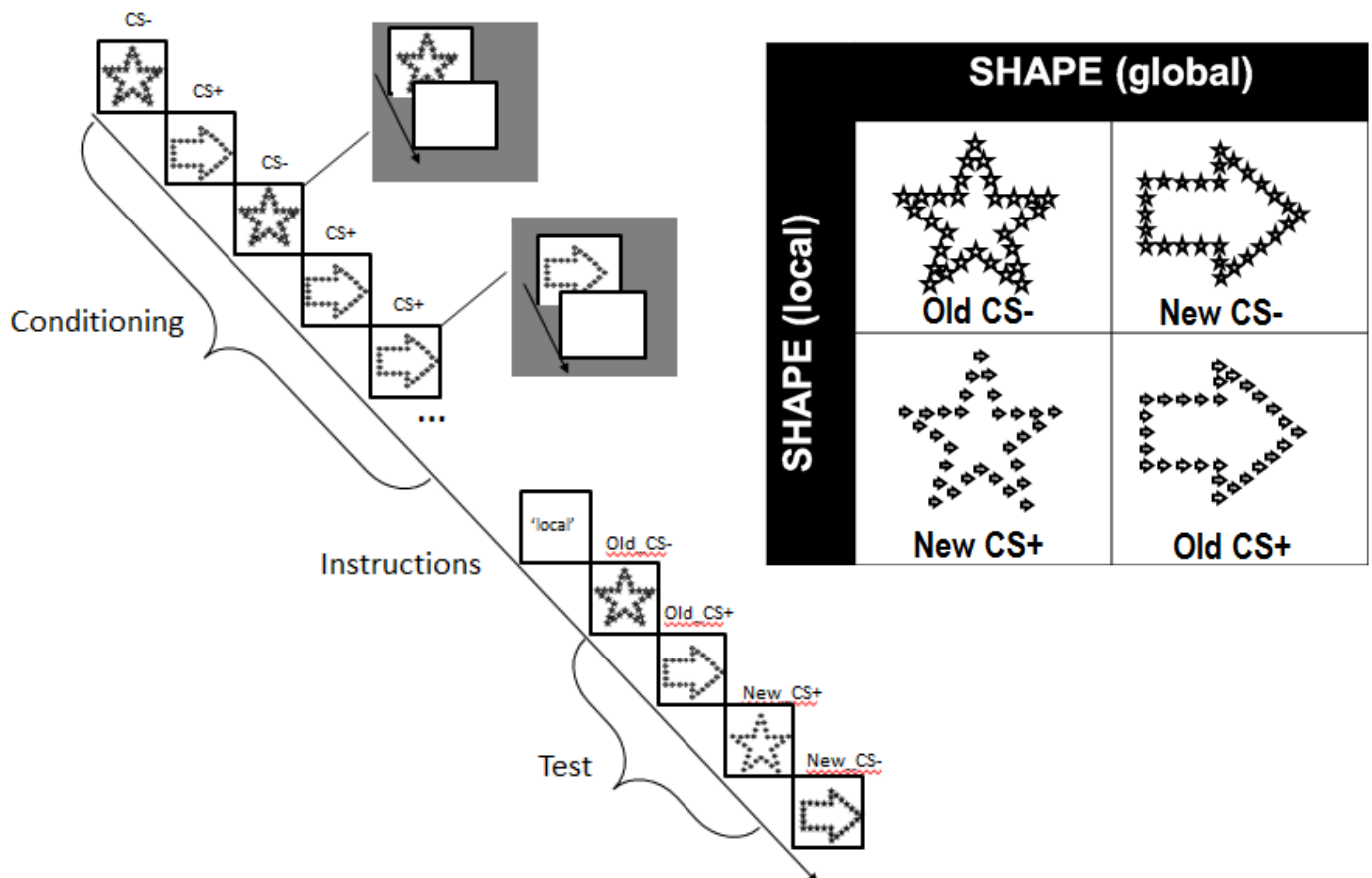
- 724 Phelps, E. A., O'Connor, K. J., Gatenby, J. C., Gore, J. C., Grillon, C., & Davis, M. (2001). Activation of
725 the left amygdala to a cognitive representation of fear. *Nature Neuroscience*, 4(4), 437–441.
726 <http://doi.org/10.1038/86110>
- 727 Prévost, C., McNamee, D., Jessup, R. K., Bossaerts, P., & O'Doherty, J. P. (2013). Evidence for model-
728 based computations in the human amygdala during Pavlovian conditioning. *PLoS Computational*
729 *Biology*, 9(2), e1002918. <http://doi.org/10.1371/journal.pcbi.1002918>
- 730 Robinson, M. J. F., & Berridge, K. C. (2013). Instant transformation of learned repulsion into
731 motivational “wanting.” *Current Biology*, 23(4), 282–289.
732 <http://doi.org/10.1016/j.cub.2013.01.016>
- 733 Schiller, D., & Delgado, M. R. (2010). Overlapping neural systems mediating extinction, reversal and
734 regulation of fear. *Trends in Cognitive Sciences*. <http://doi.org/10.1016/j.tics.2010.04.002>
- 735 Sevenster, D., Beckers, T., & Kindt, M. (2012). Instructed extinction differentially affects the
736 emotional and cognitive expression of associative fear memory. *Psychophysiology*, 49(10),
737 1426–1435. <http://doi.org/10.1111/j.1469-8986.2012.01450.x>
- 738 Sharar, S. R., Alamdari, A., Hoffer, C., Hoffman, H. G., Jensen, M. P., & Patterson, D. R. (2016).
739 Circumplex Model of Affect: A Measure of Pleasure and Arousal During Virtual Reality
740 Distraction Analgesia. *Games for Health Journal*, 5(3), 197–202.
741 <http://doi.org/10.1089/g4h.2015.0046>
- 742 Silberstein, R. B., Nunez, P. L., Pipingas, A., Harris, P., & Danieli, F. (2001). Steady state visually evoked
743 potential (SSVEP) topography in a graded working memory task. In *International Journal of*
744 *Psychophysiology* (Vol. 42, pp. 219–232). [http://doi.org/10.1016/S0167-8760\(01\)00167-2](http://doi.org/10.1016/S0167-8760(01)00167-2)
- 745 Sutton, R. S., Szepesvári, C., Geramifard, A., & Bowling, M. (2008). Dyna-Style Planning with Linear
746 Function Approximation and Prioritized Sweeping. *Proceedings of the Twenty-Fourth Conference*
747 *on Uncertainty in Artificial Intelligence*, 528–536. Retrieved from
748 <https://arxiv.org/ftp/arxiv/papers/1206/1206.3285.pdf>
- 749 Tabor, A., Thacker, M. A., Moseley, G. L., & Körding, K. P. (2017). Pain: A Statistical Account. *PLoS*
750 *Computational Biology*. <http://doi.org/10.1371/journal.pcbi.1005142>
- 751 Tracey, I. (2010). Getting the pain you expect: Mechanisms of placebo, nocebo and reappraisal
752 effects in humans. *Nature Medicine*. <http://doi.org/10.1038/nm.2229>
- 753 Van Damme, S., Crombez, G., Hermans, D., Koster, E. H. W., & Eccleston, C. (2006). The role of
754 extinction and reinstatement in attentional bias to threat: A conditioning approach. *Behaviour*
755 *Research and Therapy*, 44(11), 1555–1563. <http://doi.org/10.1016/j.brat.2005.11.008>
- 756 Vlaev, I., Seymour, B., Dolan, R. J., & Chater, N. (2009). The price of pain and the value of suffering.
757 *Psychological Science*, 20(3), 309–317. <http://doi.org/10.1111/j.1467-9280.2009.02304.x>
- 758 Wager, T. D., Rilling, J. K., Smith, E. E., Sokolik, A., Casey, K. L., Davidson, R. J., ... Cohen, J. D. (2004).
759 Placebo-induced changes in fMRI in the anticipation and experience of pain. *Science (New York,*
760 *N.Y.)*, 303(5661), 1162–1167. <http://doi.org/10.1126/science.1093065>
- 761 Wang, J. X., Kurth-Nelson, Z., Tirumala, D., Soyer, H., Leibo, J. Z., Munos, R., ... Botvinick, M. (2016).

- 762 Learning to reinforcement learn. *arXiv:1611.05763*, 1–17. Retrieved from
763 <http://arxiv.org/abs/1611.05763>
- 764 Wang, L., Yu, H., & Zhou, X. (2013). Interaction between value and perceptual salience in value-driven
765 attentional capture. *Journal of Vision*, *13*(2013), 1–13. <http://doi.org/10.1167/13.3.5>.doi
- 766 Weiss, K. E., Dahlquist, L. M., & Wohlheiter, K. (2011). The Effects of Interactive and Passive
767 Distraction on Cold Pressor Pain in Preschool-aged Children. *Journal of Pediatric Psychology*,
768 *36*(7), 816–826. <http://doi.org/10.1093/jpepsy/jsq125>
- 769 Wentura, D., Müller, P., & Rothermund, K. (2014). Attentional capture by evaluative stimuli: Gain-
770 and loss-connoting colors boost the additional-singleton effect. *Psychonomic Bulletin and*
771 *Review*, *21*(3), 701–707. <http://doi.org/10.3758/s13423-013-0531-z>
- 772 Wieser, M. J., Miskovic, V., & Keil, A. (2016). Steady-state visual evoked potentials as a research tool
773 in social affective neuroscience. *Psychophysiology*. <http://doi.org/10.1111/psyp.12768>
- 774 Wieser, M. J., Miskovic, V., Rausch, S., & Keil, A. (2014). Different time course of visuocortical signal
775 changes to fear-conditioned faces with direct or averted gaze: A ssVEP study with single-trial
776 analysis. *Neuropsychologia*, *62*(1), 101–110.
777 <http://doi.org/10.1016/j.neuropsychologia.2014.07.009>
- 778 Yin, H. H., & Knowlton, B. J. (2006). The role of the basal ganglia in habit formation. *Nature Reviews.*
779 *Neuroscience*, *7*(6), 464–76. <http://doi.org/10.1038/nrn1919>
- 780

781 **FIGURE 1. EXPERIMENTAL TASK.**

782 **Left:** timeline of a single block in the experimental task, including the conditioning, instructions, and
783 test stages. **Right:** CSs in this block were drawn from the 4-figure subset crossing the local and global
784 dimensions of the star and arrow shapes.

785

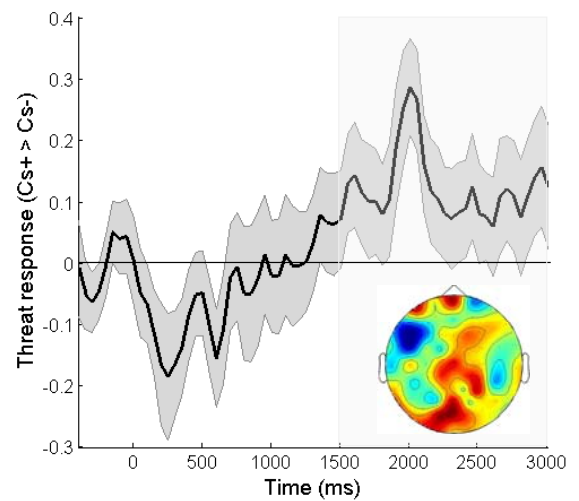
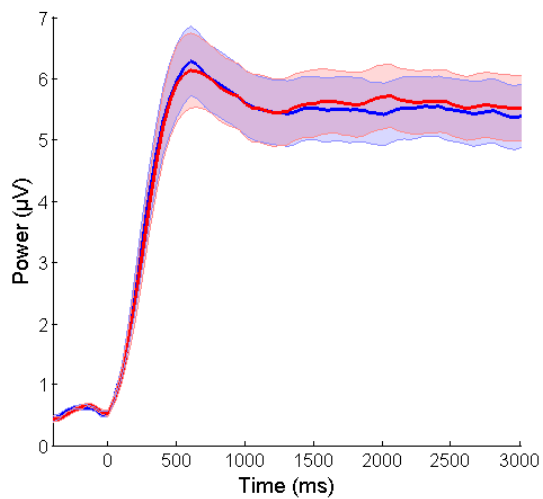


790 **FIGURE 2. THREAT EFFECTS IN THE CONDITIONING STAGE AS A FUNCTION OF TIME FROM CS**
791 **ONSET**

792 **Left.** SSVEP signal amplitudes in the conditioning stage for CS+ and CS-, extracted from occipital
793 electrodes Oz and POz, and averaged across 19-21Hz. Shaded areas plot the standard error.

794 **Right.** Threat effects in occipital electrodes Oz and POz, operationalised as the difference between
795 CS+ > CS-, are plotted as a function of time from CS onset, averaged across 19-21Hz. The time
796 window 1500-300ms (shaded grey) was used in all analysis of threat effects. The topology inset
797 shows that the threat effects across that time window.

798



799

800

801

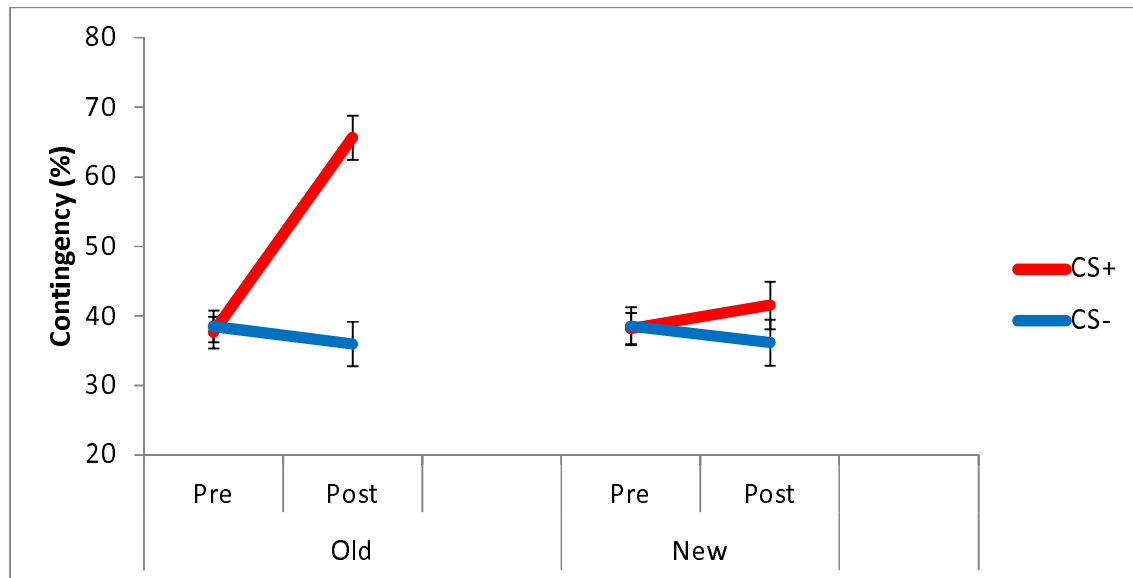
802

803 **FIGURE 3. CONTINGENCY AND LIKABILITY OF CONDITIONED STIMULI.**

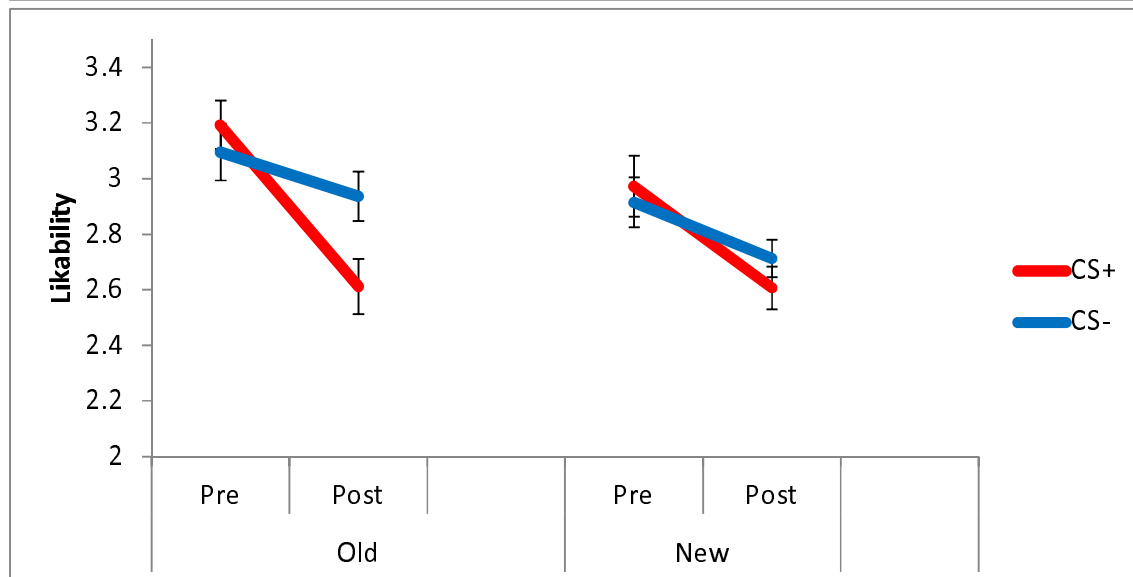
804 The contingency and likability ratings of stimuli used as CSs before and after the experimental task.

805 Error bars indicate the standard error of the mean.

806



807

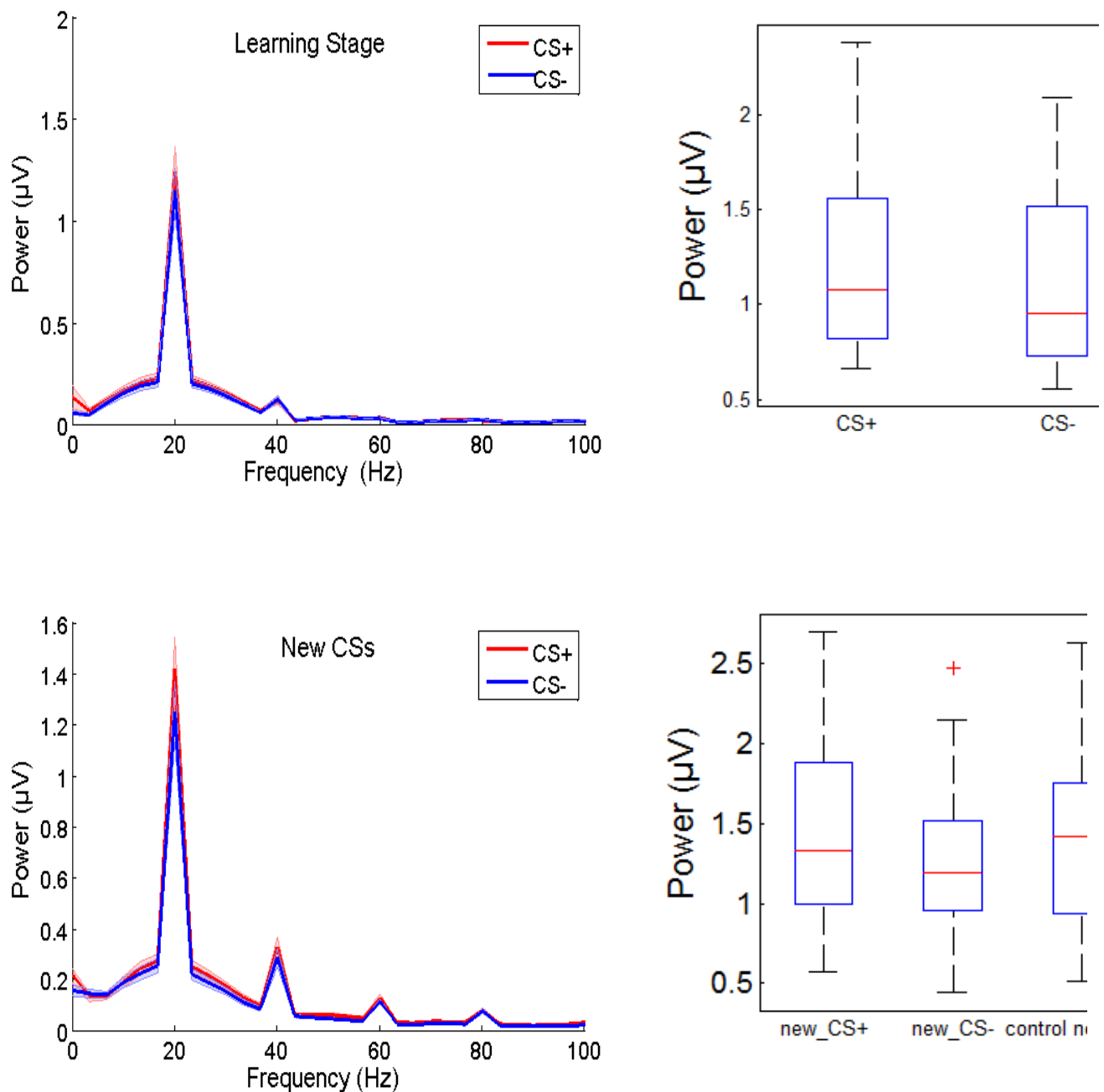


808

809 **FIGURE 4. SPECTRAL ANALYSIS OF THREAT EFFECTS**

810 **Left.** Spectral analysis of signal in the conditioning stage for CS+ and CS-, extracted from occipital
811 electrodes Oz and POz at the 1500-3000ms time window, showing that threat modulated the 20Hz
812 SSVEP signal and some of its harmonics.

813 **Right.** The magnitude of the 20Hz threat effect, showing the variability of this effect across
814 participants. The red line indicated the mean; the box indicates the inter-quartile range.



37

815