

1 Distinguishing neural correlates of prediction errors on perceptual content and
2 detection of content

3 Abbreviated title: Content and detection prediction errors

4 Nadine Dijkstra^{*1}, Oliver Warrington¹, Peter Kok¹ and Stephen M. Fleming^{1,2,3}

5

6 1. Wellcome Centre for Human Neuroimaging, UCL Queen Square Institute of Neurology, University
7 College London, United Kingdom, WC1N 3AR

8 2. Max Planck UCL Centre for Computational Psychiatry and Aging Research, University College
9 London, United Kingdom, WC1B 5EH

10 3. Department of Experimental Psychology, University College London, United Kingdom, WC1H 0AP

11

12 * Corresponding author

13

14 Contact: n.dijkstra@ucl.ac.uk

15

16

17 Conflict of interest statement: The authors declare to have no conflict of interest.

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37 Accounting for why discrimination between different perceptual contents is not always
38 accompanied conscious detection of that content remains a challenge for predictive
39 processing theories of perception. Here we test a hypothesis that detection is supported
40 by a distinct inference within generative models of perceptual content. We develop a
41 novel visual perception paradigm that probes such inferences by manipulating both
42 expectations about stimulus content (stimulus identity) and detection of content
43 (stimulus presence). In line with model simulations we show that both content and
44 detection expectations influence reaction times on a categorisation task. By combining a
45 no-report version of our task with functional neuroimaging we reveal that violations of
46 expectations (prediction errors; PEs) about perceptual content and detection are
47 supported by posterior and prefrontal cortex in qualitatively different ways: within
48 posterior sensory cortex, activity patterns diverge only on trials with a content PE, but
49 within these trials, further divergence is seen for detection PEs. In contrast, within
50 prefrontal cortex, activity patterns diverge only on trials with a detection PE, but within
51 these trials, further divergence is seen for content PEs. These results suggest rich
52 encoding of both content and detection prediction errors and highlight a distributed
53 neural basis for inference on content and detection of content in the human brain.

54
55

56 Our perceptual experience is characteristically limited: at any given moment in time we are aware of
57 only a subset of perceptual inputs. Such failures of awareness do not necessarily reflect failures of
58 sensory processing. For instance, in a widely cited example of a dissociation between perceptual
59 performance and awareness, a patient with blindsight is still able to respond above-chance to the
60 identity of a stimulus, despite not seeing that stimulus (Persaud et al., 2011; Weiskrantz et al., 1974).
61 Similar dissociations have been documented in otherwise healthy subjects using techniques such as
62 masking, where the content of stimuli which are rendered invisible nevertheless continues to exert an
63 impact on behaviour (Dehaene et al., 2001; Marcel, 1983; Peters & Lau, 2015, although see (Meyen
64 et al., 2022). Within the framework of perceptual decision-making, dissociations between
65 performance and awareness can be modelled as a distinction between discrimination – categorising
66 some aspect of stimulus identity – and detection – responding as to whether a stimulus is perceived
67 or not (Azzopardi & Cowey, 1997; Green & Swets, 1966; Peters & Lau, 2015).

68 Predictive processing offers a powerful and general computational framework for modelling
69 perception (Hohwy & Seth, 2020; Marvan & Havlík, 2021). Within this framework, the content of
70 perception is realized by combining prior knowledge (expectations) with incoming sensory evidence
71 (Bastos et al., 2012; Friston et al., 2006; Kersten et al., 2004; Kok et al., 2013). Mismatches between
72 expectations and evidence about a particular feature result in prediction errors – tell-tale signatures
73 of inference on that feature. Previous accounts have suggested that what we become aware of is
74 determined by specific aspects of perceptual inference, for example, the perceptual hypothesis with
75 the highest posterior probability (Hohwy, 2012) or the updating of perceptual hypotheses by
76 unexpected signals (Hobson & Friston, 2014). However, these accounts struggle to accommodate

77 dissociations between high-fidelity discrimination performance – presumably reflecting intact
78 perceptual inference – and detection judgments (Lau, 2022).

79 An alternative proposal is that detection arises from inferences that are distinct from
80 inferences about content. We recently proposed a computational architecture (the higher-order state
81 space (HOSS) model) in which a higher-order global inference about the presence or absence of first-
82 order perceptual content supports detection judgements (Lau, 2019; Lau, 2007; Morales, 2022). This
83 processing step is proposed to be distinct to bottom-up salience or attention (Fleming, 2020). This
84 model builds on a large body of prior work that associates awareness with changes in higher-order
85 cognitive processes, including global workspace and higher-order theories of consciousness (Brown,
86 2015; Dehaene & Changeux, 2011; Lau & Rosenthal, 2011; Mashour et al., 2020). A more general
87 question that goes beyond these theories of consciousness is whether inferences on content and
88 detection of content rely on distinct processes in the human brain. Addressing this question would
89 provide initial empirical constraints on the architecture of predictive processing theories of perception
90 and consciousness.

91 In this study, we investigated to what extent inferences about content and detection of
92 content correlate with distinct neural substrates measured with fMRI. In what follows, we refer to
93 predictions about perceptual content that are relevant for discrimination as “content expectations”,
94 and to predictions about the presence (vs. absence) of content that are relevant for detection as
95 “detection expectations”. We developed a novel experimental paradigm in which we independently
96 manipulated expectations about perceptual content (stimulus identity) and detection (whether
97 stimulus content will be present or absent). In line with a neural hierarchy supporting detection
98 inferences, we hypothesised that prediction errors about perceptual content would be localized to
99 sensory cortex (Bastos et al., 2012; Kok et al., 2013) whereas prediction errors on detection would be
100 localized to prefrontal cortex (Merten & Nieder, 2012; van Vugt et al., 2018).

101 To preface our results, we first show in a behavioural experiment that both content and
102 detection expectations influence reaction times. Using a no-report version of the same task in
103 conjunction with neuroimaging (Tsuchiya et al., 2015), we show that content prediction errors are
104 predominantly encoded in sensory (visual) cortical areas whereas detection prediction errors are
105 predominantly encoded in prefrontal cortical areas. However, a strict separation between content and
106 detection of content is nuanced by findings of mutual interactions between the two types of
107 prediction error signals in both visual and prefrontal cortices. Taken together, our findings suggest
108 that inferences on content and detection of content rely on distinct but interacting neural substrates
109 in the human brain.

110

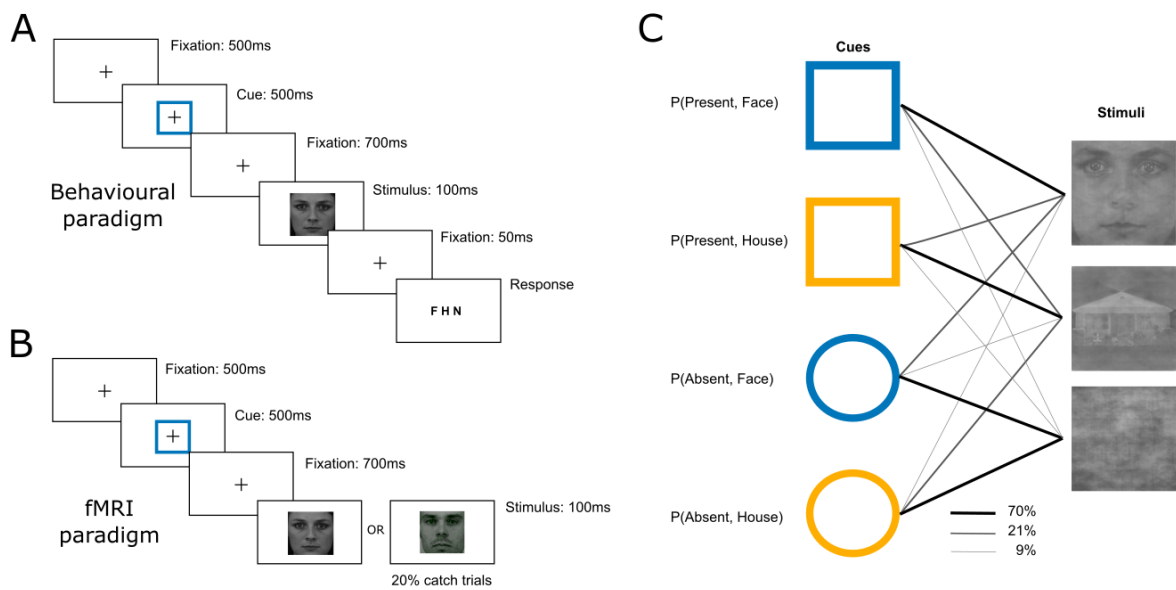
111 **Materials and Methods**

112 **Participants.** To determine our sample size, we assumed medium effect sizes (Cohen's $d = 0.05$) which
113 require 34 participants to achieve a power of 80% to detect an effect at an alpha level of 0.05. To allow
114 for drop-out, 36 participants gave written informed consent and participated in the study (mean age
115 26.4, SD 7.2). All participants were included in the behavioural analysis. 9 participants were excluded
116 from the MRI analyses: 2 participants took part in the behavioural session but could not be scanned,
117 6 were excluded due to low catch trial accuracy (below 70%), 1 was excluded due to a technical error
118 during response recording. 27 participants were included in the final fMRI analyses. Given that
119 dropout for the MRI session was higher than anticipated, this reduced our power to detect a medium
120 effect size to 71%. The study was approved by the University College London ethics committee
121 (approval number 8231_001). Participants were paid £8 for the behaviour session and £10 per hour
122 for the fMRI session.

123
124 **Stimuli.** The experiment was programmed in MATLAB R2019 (MathWorks) using Psychtoolbox
125 (version 3.0.16). In the behavioural session, stimuli were presented on a desktop monitor and in the
126 scanner, stimuli were presented via a projector at an approximate viewing distance of 58cm. Stimuli
127 consisted of 8 faces and 8 houses displayed in black and white and embedded in noise. The face stimuli
128 were selected from the Karolinska Directed Emotional Faces database (Lundqvist, D., Flykt, A., &
129 Öhman, 1998) with a 50-50 male-female ratio and a neutral expression. The house stimuli were
130 adapted from the Pasadena houses database collected by Helle and Perona (California Institute of
131 Technology, Pasadena, California). We controlled for differences in spatial frequency between stimuli
132 by calculating the frequency phase and magnitude for each stimulus in the set using Fourier
133 transformation, and then combining the mean spatial frequency magnitude over stimuli with each
134 individual image's frequency phase information to create images that are perceptually similar to the
135 original images but have the same spatial frequency profile. Noise was generated as matrices of
136 random numbers uniformly distributed between 0 and 1 and combined with each stimulus image in a
137 weighted sum with the weight on the stimulus image controlled by a visibility parameter. Face and
138 house stimuli were each presented at 0.9 visibility. Noise images consisted solely of noise. To ensure
139 participant engagement in the fMRI session, 20% of the trials were catch trials in which 10% of image
140 pixels were turned green. Participants were instructed to press a button when they saw a green tinge
141 to the image (more details below).

142
143 **Task and procedure.** Participants attended the laboratory on two different days within a 5-day
144 window. On the first day they completed a behavioural training session in which the task and cue-

145 stimulus associations were learned (Fig. 1A). On the second day participants completed the fMRI
 146 session during which they performed the same task but now without requiring a response (Fig. 1B).
 147 To facilitate investigation of whether prediction errors on detection are encoded differently from
 148 prediction errors on content, we independently manipulated the probability of whether any
 149 perceptual content (face or house) would be presented or not (the detection prior), and whether such
 150 content would be a face or a house (the content prior). These probabilities were reflected in the shape
 151 (detection) and colour (content) of the cues (Fig. 1C). For example, an orange circle indicates that most
 152 likely no stimulus would be presented, but that if one were presented, it would likely be a house. To
 153 prevent neural correlates of predictions being contaminated by responses to the physical cues, cue-
 154 stimulus mappings within each level were swapped halfway through the experiment. For instance, the
 155 high-probability presence / high-probability face cue could be represented by a blue rectangle in the
 156 first half of the experiment, and an orange circle in the second half. Both mappings were also used
 157 during the behavioural training session to familiarize participants with possibility of such switches.
 158



159
 160 **Figure 1. Experimental paradigm.** (A) Behavioural paradigm. Trials consisted of 500ms fixation followed by
 161 500ms cue, 700ms fixation, 100ms stimulus. After another 50ms fixation participants had to indicate whether
 162 they saw a face (F), house (H) or noise (N) using the 'a', 's' or 'd' keyboard keys respectively. (B) The fMRI
 163 paradigm was the same as the behavioural paradigm except that no response was required, ensuring task-
 164 related activations would not reflect reporting requirements. To ensure participants remained engaged with the
 165 task, they were instructed to press a button when the stimulus contained green pixels, which was the case on
 166 20% of the trials. (C) The shape of the cue indicated the probability a stimulus would be present or not (detection
 167 expectation) whereas the colour indicated whether that stimulus was likely to be a face or a house (content
 168 expectation). For example, a blue circle indicated a high probability that no stimulus would be presented but
 169 that if a stimulus were presented, it would likely be a face.

170
 171 **Behavioural session.** At the start of the experiment, participants were trained on the task and the
 172 relationship between the cues and stimuli in 3 phases. Within each phase, task instructions were

173 presented to participants via a self-paced PowerPoint presentation interleaved with 20 practice trials,
174 resulting in 60 practice trials in total. In phase 1, participants were introduced to how the colour of
175 cues indicated stimulus content (e.g. blue: '70% chance of being a face', orange: '70% change of being
176 a house'). During each trial, a cue was presented for 500ms, followed by a fixation cross for 700ms,
177 and then the stimulus for 100ms (Fig. 1A). After another brief 50ms fixation, a response screen
178 appeared indicating the different response options (in this phase, only 'F' for face and 'H' for house
179 responses were available). Participants were instructed to use the 'a' and 's' keys to indicate whether
180 they saw a face or a house respectively. After pressing the button, their selected category, 'F' or 'H',
181 was highlighted for 200ms before continuing to the next trial.

182 In phase 2, the noise category (absence of content) was added and participants were
183 introduced to the shape dimension of the cues (e.g. square: '70% likely to be a picture (face or house)',
184 rectangle: '70% likely to be no picture/noise'). They then practiced 20 trials of face-house-noise
185 categorisation with just the presence/absence shape cues. 'N' was added to the response options and
186 participants used the 'd' key to select this 'no picture' answer. In the final phase 3, both cue
187 dimensions were combined and their meaning was explained (e.g. blue square: '70% likely to be a
188 picture AND that picture will likely be a *face*'). Then, participants completed two blocks of the main
189 task with these cues, each taking approximately 6 minutes. After this, the cue-stimulus mappings were
190 swapped, the new cues were explained and another similar training session was completed. Finally,
191 the participants completed two blocks with these new cue-stimulus mappings.

192

193 **fMRI session.** To ensure that any correlates of prediction errors were not confounded by response
194 requirements in the MRI scanner, participants performed a no-report version of the task (Fig. 1B).
195 During the initial setup scans, participants were presented with instructions reminding them of the
196 cue-stimulus associations and introducing the catch trial task. To ensure participants continued to pay
197 attention to the stimuli, we introduced a target detection task: participants were asked to detect
198 stimuli in which green pixels were intermixed within the image (face, house or noise). Participants
199 performed 73 practice trials with the first cue-stimulus mapping, lasting ~4 minutes. They then
200 completed 3 × 8-minute blocks of the main task under the first cue-stimulus mapping, each with 145
201 trials. After each block, the scanner was stopped and participants were asked whether they needed a
202 break. Halfway through the experiment, while their structural scan was obtained, participants were
203 reminded about the second cue-stimulus relationship via another set of instructions and 73 practice

204 trials. They then completed a further 3×8 -minute blocks of the main task under the second cue-
205 stimulus mappings. In total, participants performed 870 trials of the main task.

206 Scanning took place at the Wellcome Centre for Human Neuroimaging, University College
207 London, using a 3 Tesla Siemens Prisma MRI scanner with a 64-channel head coil. We acquired
208 structural images using an MPRAGE sequence (1x1x1 mm voxels, 176 slices, in plane FoV =
209 256x256mm²), followed by a double-echo FLASH (gradient echo) sequence with TE1 = 10ms and TE2
210 = 12.46ms (64 slices, slice thickness = 2mm, gap = 1mm, in plane FoV = 192 x 192mm², resolution = 3
211 x 3mm²) that was later used for field inhomogeneity correction. Functional scans were acquired using
212 a 2D EPI sequence, optimized for regions near the orbito-frontal cortex (3x3x3mm voxels, TR = 3.36s,
213 TE = 30ms, 48 slices tilted by 30 degrees with respect to the T > C axis, matrix size = 64x72, Z-shim =
214 1.4).

215
216 **Model simulations.** We used core functions of the Higher-Order State Space (HOSS) model
217 (<https://github.com/smfleming/HOSS>) to simulate the expected pattern of prediction errors in our
218 experiment (Fleming, 2020). The model is instantiated as a probabilistic graphical model, where nodes
219 correspond to unknown variables and the graph structure indicates dependencies between variables
220 (Fig. 2A). The model is generative, such that higher levels of the hierarchy generate expectations over
221 variables in the layers below. The highest level, the detection (A) state, is a simple scalar such that
222 higher probabilities lead to the activation of content (face or house) states in the W layer below. W
223 is a $1 \times N$ vector that encodes the relative probabilities of each of N discrete perceptual states. Here,
224 $N = 3$, reflecting the 3 possible stimulus categories of face, house or noise. To simulate multivariate
225 sensory data X , we drew samples from one of three multivariate normal distributions conditioned on
226 W : ‘noise’ with $\mu = [0.5 \ 0.5]$, ‘face’ with $\mu = [1.5 \ 0.5]$ or ‘house’ with $\mu = [0.5 \ 1.5]$. The covariance
227 matrix was specified as $\Sigma = [0.1 \ 0; 0 \ 0.1]$. The locations of samples in evidence space are arbitrary;
228 what is important is the mapping between stimulus categories and the detection state. The likelihood
229 of X given W is then:

230

$$231 \quad P(X = x|W) \sim N(\mu_W, \Sigma)$$

232

233 Upon receipt of a sample of X , the model can be inverted to compute the posteriors over A and W by
234 marginalising:

235

$$236 \quad P(A|X = x) \propto \sum_W P(A)P(W|A)P(X = x|W)$$

237
$$P(W|X = x) \propto \sum_A P(A)P(W|A)P(X = x|W)$$

238

239 We simulated four possible prior states, reflecting the four cues in the experiment: $p(A) =$
 240 $0.8 \& p(W_{face} = 0.8)$, $p(A) = 0.8 \& p(W_{face} = 0.2)$, $p(A) = 0.2 \& p(W_{face} = 0.8)$, $p(A) =$
 241 $0.2 \& p(W_{face} = 0.2)$. For each cue-target combination, prediction error was computed at both the
 242 detection (A) and content (W) layers as the Kullback-Leibler (KL) divergence between the prior and
 243 posterior distributions. The simulated categorisation response was determined by the W -state with
 244 the highest poster probability. Only correct trials were used to calculate the prediction errors. We
 245 simulated 300 trials per cue-target combination and 30 participants in total.

246

247 **Behavioural analysis.** We first tested whether there was a congruency effect of the content and
 248 detection cues by comparing valid versus invalid trials for both cue types using simple t-tests on both
 249 accuracy and reaction time (RT). Trials with RTs faster than 200ms or slower than 2s were removed
 250 prior to analysis. For the content congruency effects, noise trials were ignored. To investigate the
 251 effects of W -level (content) and A -level (detection) prediction errors in more detail, we ran a linear
 252 mixed-effects analysis using MATLAB's (R2021b) 'lme.m' function, with the following model:

253

254
$$\log(RT) \sim KL_W + KL_A + KL_W \times KL_A + (1 | participant)$$

255

256 where KL_W and KL_A are proxies for the qualitative patterns expected for W and A -level prediction
 257 errors per cue-target combination respectively, calculated using the HOSS model.

258

259 **fMRI pre-processing.** Data pre-processing followed the procedure described in (Mazor et al., 2020;
 260 Morales et al., 2018): Imaging analysis was performed using SPM12 (Statistical Parametric Mapping;
 261 www.fil.ion.ucl.ac.uk/spm). The first five volumes of each run were discarded to allow for T1
 262 stabilization. Functional images were realigned and unwarped using local field maps (Andersson et al.,
 263 2001) and then slice-time corrected (Sladky et al., 2011). Each participant's structural image was
 264 segmented into gray matter, white matter, CSF, bone, soft tissue, and air/background images using a
 265 nonlinear deformation field to map it onto template tissue probability maps (Ashburner & Friston,
 266 2005). This mapping was applied to both structural and functional images to create normalized images
 267 in Montreal Neurological Institute (MNI) space. Normalized images were spatially smoothed using a
 268 Gaussian kernel (6 mm FWHM). We set a within-run 4 mm affine motion cut-off criterion. Pre-
 269 processing and construction of first- and second-level models used standardized pipelines and scripts
 270 available at <https://github.com/metacoglab/MetaLabCore/>.

271

272 **Univariate analysis.** To test where in the brain activation correlated with W and A-level prediction
273 errors, we performed univariate analyses within SPM12 in MATLAB R2021b. Main effects of A
274 (detection) and W (content) level prediction errors were characterised using the model-predicted KL
275 divergence per trial type (Fig. 2). The general linear model (GLM) contained one regressor aligned to
276 the onset of the stimulus with two parametric modulators, one for each type of prediction error. The
277 onset of the cues and responses were included as nuisance regressors, as were movement
278 parameters, their first derivatives and the mean amplitudes of voxels containing white matter and
279 cerebral spinal fluid (CSF). Regressors were specified per run. Significance testing was implemented at
280 the group-level with a t-test of each KL regressor against 0. Correction for multiple comparisons was
281 applied at the cluster-level ($P < 0.05$, family-wise error corrected), using a cluster-forming threshold
282 of $P < 0.001$, uncorrected. Effects were small-volume corrected using either (a) a posterior mask that
283 was generated by combining the following regions from the AAL atlas (Destrieux et al., 2010): all
284 occipital regions, inferior temporal gyrus, calcarine, cuneus, and lingual gyrus or (b) a frontal mask that
285 included: all frontal regions, rectus, insula and anterior cingulate. For the representational similarity
286 analysis (see below), beta weights per cue and target combination were estimated by running a
287 separate GLM with the same nuisance regressors as before but now instead of the KL divergence,
288 including a condition regressor per cue-target combination (12 in total), centred on the target onset.

289

290 **Representational similarity analysis.** To investigate the representational structure of the different
291 types of prediction errors, we performed a searchlight representational similarity analysis (RSA). RSA
292 was performed using MATLAB R2021b and Timo Flesch's RSA toolbox
293 (https://github.com/TimoFlesch/fmri_utils/tree/master/RSA) in combination with custom MATLAB
294 code. We defined 5 different model RDMs ([Representational Dissimilarity Matrices](#)) encoding
295 dissimilarity between (1) the presented stimuli, (2) content priors, (3) detection priors, (4) content
296 prediction errors and (5) detection prediction errors (Fig. 5A). Neural RDMs were generated per
297 participant by calculating the Euclidean distance between the activation patterns of different
298 conditions. All within-run comparisons were set to NaN, distances along the diagonals of all model and
299 neural RDMs were also set to NaN, and the lower-triangles were transformed into distance vectors.
300 Per searchlight, a GLM was run to predict the neural RDM from the model RDMs. Group-level
301 inference was performed by testing the inferred beta weights per RDM regressor over participants
302 against 0 using a one-sample t-test. Correction for multiple-comparisons was performed using family-
303 wise error correction ($p < 0.05$) at the whole-brain level, using a cluster-forming threshold of $p < 0.001$,
304 uncorrected (as before).

305

306 **Results**

307 To independently manipulate predictions about perceptual content and detection of content, we
308 developed a novel perceptual discrimination task with compound cues (Fig. 1). In the behavioural
309 version of the experiment, the task was to infer whether a briefly shown stimulus was a face, a house
310 or noise (Fig. 1A). Preceding the stimulus was a compound cue in which the shape indicated the
311 probability of seeing a stimulus (face or house) versus noise, regardless of its identity – an expectation
312 about the detection of content, rather than content itself. In contrast, the colour of the cue indicated
313 the probability of a stimulus being a face or a house, regardless of whether it was likely to be present
314 – an expectation about content, rather than detection of content. For example, a blue circle indicated
315 that there was a high likelihood no stimulus would be shown (detection expectation) but that if a
316 stimulus was shown, it would likely be a face (content expectation; Fig. 1C).

317

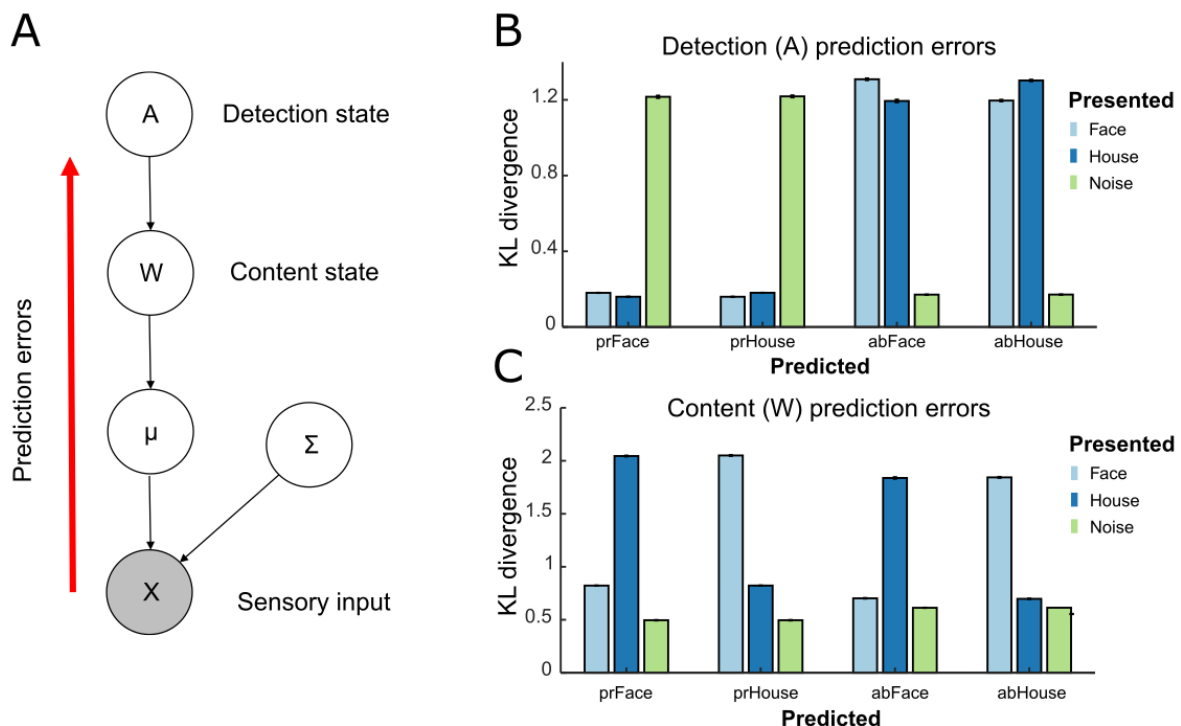
318 **Simulations predict diverging patterns of prediction error on content and detection.**

319 We used the higher-order state space (HOSS) model (Fleming, 2020) to simulate expected patterns of
320 prediction errors (PEs) for content and detection expectations in our experiment. HOSS specifies a
321 Bayesian network in which higher-order “detection states” (A) furnish expectations about the
322 presence of content, with content-specific “content states” (W) nested under the detection state layer
323 (Fig. 2A). In the current experiment, content states W denote [face, house, noise] and detection state
324 A encodes [present, absent] irrespective of content (face and house are mapped to “present”, and
325 noise to “absent”). We simulated the model with priors for the W and A layers set to the empirical
326 prior probabilities used to construct the compound cues. Upon receipt of a multivariate sensory input
327 X , the model is inverted and posterior probabilities over both content and detection of content can
328 be derived.

329 HOSS naturally nests high-dimensional perceptual content (the W layer) within a more
330 abstract state that tracks the magnitude or reliability of higher-dimensional perceptual signals.
331 However, HOSS is only one of several possible architectures that could support dissociable inferences
332 on content and detection of content – for instance, a “flat” architecture with no explicit representation
333 of global presence vs. absence may suffice (Whyte & Smith, 2021). The more general point is that
334 inferences on detection of content are proposed to be distinct from inferences on specific contents,
335 with the former being factorised with respect to the latter (Fleming, 2020).

336 This can be appreciated in the pattern of prediction errors simulated from the model (Figure
337 2). Simulated prediction errors (Kullback-Leibler divergences) within the detection and content layers
338 for each cue-target combination are shown in Figures 2B and 2C, respectively. In Figure 2B, inferences

339 about detection of content are predominantly sensitive to expectations about presence vs. absence,
 340 and not about specific stimulus contents. For example, a large detection PE is generated when the
 341 model expects to see something, but only noise is presented, irrespective of whether the content layer
 342 is expecting a face (prFace) or house (prHouse, Fig. 2B; two green bars on the left). In contrast, within
 343 the content layer, the largest PEs are observed when content expectations are violated, irrespective
 344 of detection expectations. For example, a large PE is generated when a face is expected but a house
 345 is presented, regardless of whether the detection layer expected to see a stimulus (prFace) or not
 346 (abFace, Fig. 2C; dark blue bars). Note that in the content layer, prediction errors for noise are low
 347 because a noise patch contains an absence of information about either feature (face or house).
 348 Strictly, in the model, noise is an absence of any input supporting either of the two features – whereas
 349 in our experimental paradigm, we use a noise patch to indicate the ‘absence’ of content. As noise does
 350 not contain meaningful content at the level of object categories such as faces and houses, this is a
 351 reasonable approximation to the simulations. Crucially, the shared variance between the A and W PEs
 352 is low (correlation of -0.14 between the two simulated condition vectors), showing that our
 353 experimental design allows us to independently investigate neural and behavioural correlates of
 354 content and detection PEs.
 355



356
 357 **Figure 2. Simulated prediction errors within content and detection layers of the higher-order state space**
 358 **(HOSS) model.** (A) Graphical representation of the HOSS model. Perceptual states W and detection state A are
 359 inferred based on sensory input X. Simulated prediction errors per cue-target combination within the A-
 360 detection layer (B) and W-content layer (C) are plotted. The x-axis reflects the cues, with the first two letters
 361 indicating whether the cue indicated high probability of the presence (pr) or absence (ab) of content (the

362 presence expectation), followed by whether it indicated a high probability of being face or house (the content
 363 expectation). The y-axis indicates the simulated KL-divergence between the prior and the posterior at the
 364 different levels of the model (a proxy for prediction error, or how much belief change is induced by a sensory
 365 sample). Error bars indicates standard errors of the mean (SEM) over simulation samples.

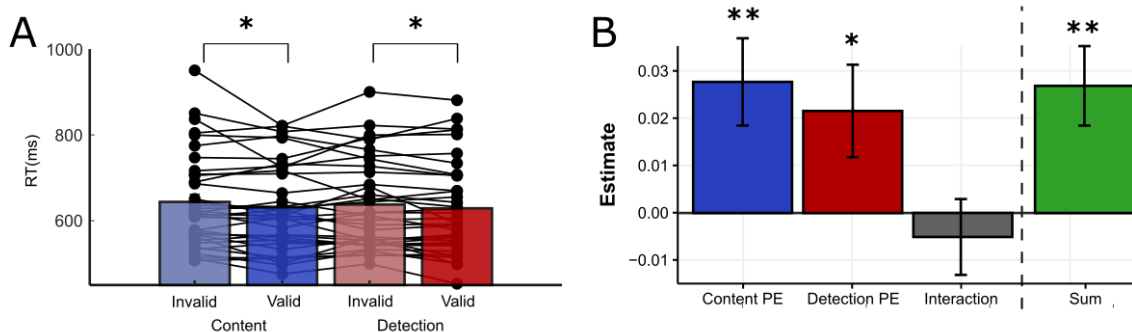
366

367 **Content and detection expectations both influence behaviour.**

368 We next investigated whether detection and content expectations separately modulated
 369 discrimination reaction times in a behavioural experiment. The logic of our approach is grounded in
 370 previous findings that expected stimuli lead to faster responses, whereas violations of expectations
 371 (prediction errors) lead to behavioural slowing (Brodersen et al., 2008; Carpenter & Williams, 1995;
 372 Mars et al., 2008). Thirty-six participants performed the behavioural version of the task (Fig. 1A). We
 373 first ran a model-free analysis to investigate whether there were congruency effects of content (only
 374 face and house trials) and detection cues on reaction time (RT). RTs were indeed faster for trials with
 375 a valid content cue ($M = 631.38\text{ms}$, $SD = 101.63\text{ms}$) compared to trials with an invalid content cue (M
 376 $= 644\text{ms}$, $SD = 109\text{ms}$; $t(35) = 2.12$, $p = 0.041$, $d = 0.12$) and also faster for trials with a valid detection
 377 cue ($M = 629\text{ms}$, $SD = 109\text{ms}$) compared to trials with an invalid presence cue ($M = 638\text{ms}$, $SD = 104\text{ms}$,
 378 $t(35) = 2.04$, $p = 0.049$, $d = 0.08$). The interaction between the two cue dimensions (content and
 379 detection) was not significant ($t(35) = -1.92$, $p = 0.063$).

380 We next ran a model-based linear mixed-effects regression analysis predicting reaction times
 381 from the simulated content and detection PEs obtained from the HOSS model (entered as random
 382 effects) with random intercepts for each participant (Fig. 3B). Both content PEs ($\beta = 0.028$,
 383 $t(16508.07) = 3.01$, $p = 0.0026$; Fig. 3B, blue bar) as well as detection PEs ($\beta = 0.022$, $t(16508.04) =$
 384 2.21 , $p = 0.027$; Fig. 3B, red bar) led to significant increases in reaction time. The interaction between
 385 content and detection PEs was again not significant ($t(16508.04) = -0.63$, $p = 0.53$; Fig. 3B, grey bar).

386



387

388 **Figure 3. Behavioural results.** (A) Reaction times (RTs) separated by valid and invalid expectations on content
 389 (blue) and detection (red). Points represent individual participants (B) Results of a linear mixed-effects analysis
 390 using simulated prediction errors (Fig. 2) as predictors of response times. The green bar reflects the beta
 391 estimate for a model including a predictor consisting of the sum of the simulated content and detection

392 prediction errors on each trial. * $p < 0.05$, ** $p < 0.005$, *** $p < 0.0005$. Error bars reflect standard errors of the
393 mean (SEM).

394

395 To explore whether the effects of content and detection PEs on reaction times were additive, we also
396 tested another model in which the PEs of the two components were summed (i.e. the sum of the K-L
397 divergences displayed in Fig. 2A and B) to create one combined PE regressor per trial. This summed
398 predictor was associated with a significant increase in reaction time, as expected ($\beta = 0.027$,
399 $t(16541) = 3.19$, $p = 0.0014$; Fig. 3B, green bar). Model comparison indicated that a model containing
400 the summed PE (BIC = 417.81) was a more parsimonious explanation of the reaction time data than a
401 model with independent PE terms (BIC = 426.23).

402 Taken together, our analysis of reaction time effects reveals that, at the level of behaviour,
403 detection and content PEs both lead to significant slowing, and that their influence is best modelled
404 as a linear sum of the two PE terms. On the basis of behavioural data alone, we are unable to conclude
405 whether inferences on content and detection of content are supported by distinct (neural)
406 computations. Therefore, in order to investigate this, we next turned to neuroimaging to ask how
407 these two types of PEs are encoded in the brain.

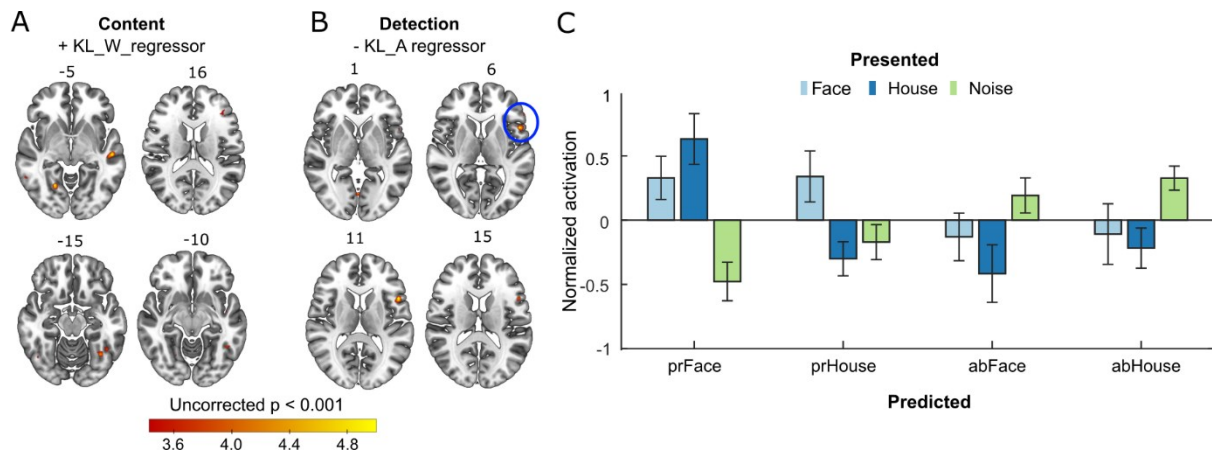
408

409 **Neural correlates of content and detection prediction errors.**

410 Twenty-seven of the participants who had completed the behavioural experiment went on to perform
411 a no-report version of the same task while undergoing whole-brain functional neuroimaging (Fig. 1B).
412 Participants passively viewed the compound cues and stimuli without being required to explicitly
413 categorise each stimulus. Instead, to ensure attention, a button press was required on catch trials
414 (indicated by green pixels within the stimulus, 20% trials). Given that the catch manipulation was
415 independent of our main effects of interest, these trials were included in subsequent analyses.

416 We pursued two complementary analysis approaches that aimed to identify a) univariate
417 signals and b) multivariate patterns covarying with either content or detection PEs. First, to identify
418 univariate brain activity modulated by detection and content PEs, we ran a whole-brain GLM entering
419 the simulated K-L divergences from the detection and content layers of the HOSS model as regressors.
420 Content and detection PEs correlated with activation in different brain areas (Fig. 4). To test our a
421 priori hypotheses that content PEs would be observed in visual sensory areas and detection PEs in
422 prefrontal areas, we applied small-volume corrections based on posterior and frontal masks.

423



424

425 **Figure 4. Univariate correlates of prediction errors on detection and content.** (A,B) Brain areas that
 426 significantly correlated positively with the KL W regressor (A) and negatively with the KL A regressor (B)
 427 thresholded at $p < 0.001$ uncorrected. (C) Activation profile of the region that showed a significant ($p < 0.05$, FWE-
 428 small-volume corrected) effect of the KL A regressors, circled in blue in (B). Activation is z-scored per participant
 429 to account for large variations in mean amplitude between participants. Error bars reflect SEM.

430

431 A positive effect of the content PE regressor in the fusiform gyrus (Fig. 4A) did not survive correction
 432 for multiple comparisons and we therefore refrain from interpreting it further ($t(26) = 4.77$, *cluster-*
 433 *level* $p_{\text{FWE-corrected}} = 0.052$). The detection PE regressor showed a *negative* correlation with activation in
 434 the left inferior frontal cortex when applying small-volume correction within a frontal mask (IFC; $t(26)$
 435 $= 5.55$, *cluster-level* $p_{\text{FWE-corrected}} = 0.026$; Fig. 4C), close to voxels showing (uncorrected) effects of the
 436 content PE regressor. There was no significant detection PE effect when applying small-volume
 437 correction within a posterior mask (all p-values > 0.338). Within this region, when presence was
 438 expected (prFace and prHouse) activation tended to be higher when this expectation was confirmed
 439 and content was presented (irrespective of whether this stimulus was a face or a house) compared to
 440 when noise was presented (Fig. 4D). Conversely, when no content was expected (abFace or abHouse),
 441 activation was higher when this expectation was confirmed, and noise was presented, compared to
 442 when a face or house was presented (Fig. 4D). These results show that activity in left IFC decreases
 443 when detection predictions are violated.

444

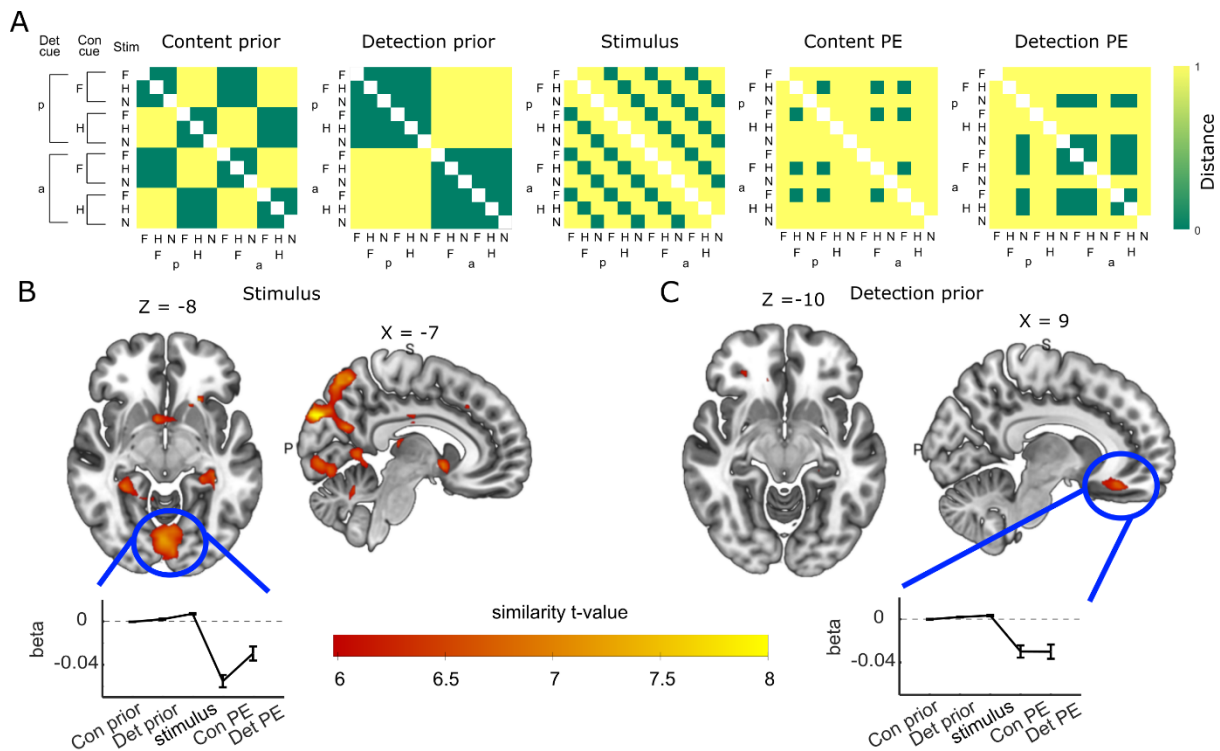
445 **Rich representations of content and detection prediction errors in sensory and prefrontal brain**
 446 **networks.**

447 While univariate analyses can reveal low-dimensional neural signatures of content and detection
 448 prediction errors, they are blind to changes in distributed neural codes that might support prior and
 449 posterior beliefs in richer representational spaces. Therefore, we next investigated whether prediction
 450 errors were also encoded in multivariate patterns rather than univariate amplitude differences, using
 451 representational similarity analysis (RSA). We formulated representational dissimilarity matrices

452 (RDMs) that identified dissimilarities between conditions in either the prior (cue), stimulus or
453 prediction error, with prior and prediction error RDMs being specified separately for the content and
454 detection layers of the model (Fig. 5A). These RDMs were designed to test whether the activity
455 patterns in one set of conditions sharing a specific feature (e.g. all conditions in which a face was
456 presented) were more similar to each other (lower distance between them) than to conditions with a
457 different feature.

458 For the prediction error RDMs, we hypothesized that if an area coded for PEs at a given level,
459 the presence of a PE would lead to convergence towards a specific activity pattern. In contrast, in the
460 absence of a PE at that level, activity patterns would reflect noise, and be random (uncorrelated).
461 Therefore, for both content and detection PEs, we assumed that conditions with a PE would be similar
462 to other conditions with a PE, whereas conditions without a PE would be dissimilar both to each other
463 and to conditions with a PE (Fig. 5A). Note that for all RSA analyses, the diagonal elements were
464 removed before computing similarity with neural data.

465 We found strong positive correlations between the stimulus RDMs and posterior brain regions
466 (Fig. 5B, Appendix A - Table 1) indicating that within these regions, stimuli belonging to the same
467 category were encoded in similar activation patterns, irrespective of prior expectations. Furthermore,
468 we found that the detection prior RDM showed significant positive correlations with activity patterns
469 in the ventromedial prefrontal cortex (vmPFC; Fig. 5C, Appendix A - Table 2), indicating that in this
470 region, the activity patterns of conditions in which presence was expected were more similar to each
471 other than to conditions in which stimulus absence was expected, and vice-versa, irrespective of
472 expectations about content (face or house). We did not find any significant correlation with the
473 content-level prior RDM.



474

475 **Figure 5. RDM hypotheses and similarity coding of stimulus- and prior-related information.** (A)
 476 Representational Dissimilarity Matrices (RDMs) reflecting different hypotheses about the similarity of neural
 477 patterns of different conditions based on encoding of content prior, detection prior, stimulus, content prediction
 478 error and detection prediction error. Darker colours indicate higher similarity, i.e. lower distance, between
 479 conditions. Only the Stimulus and Detection Prior RDMs showed significant positive correlations with brain
 480 activity. (B) Similarity in stimulus category encoding (FWE-corrected at $p < 0.05$): conditions in which stimuli of
 481 the same category were presented were more similar than those in which different stimuli were presented. Beta
 482 values for region of interest (ROI) based on this contrast encircled in blue are shown below. All RDMs significantly
 483 predicted activation patterns in this ROI, except the content prior, with all p -values < 0.006 (uncorrected). (C)
 484 Similarity in detection prior encoding (FWE-corrected at $p < 0.05$): conditions in which presence was expected
 485 were more similar to each other than to conditions in which absence was expected and vice-versa. Beta values
 486 for the blue encircled ROI are shown below. All RDMs significantly predicted activation patterns in this ROI,
 487 except the content prior, with all p -values < 0.0006 (uncorrected). Statistical maps in (B) and (C) are thresholded
 488 at $p < 0.001$ uncorrected; see Tables 2 and 3 for details of clusters surviving whole-brain correction including the
 489 EVC and vmPFC ROIs.

490

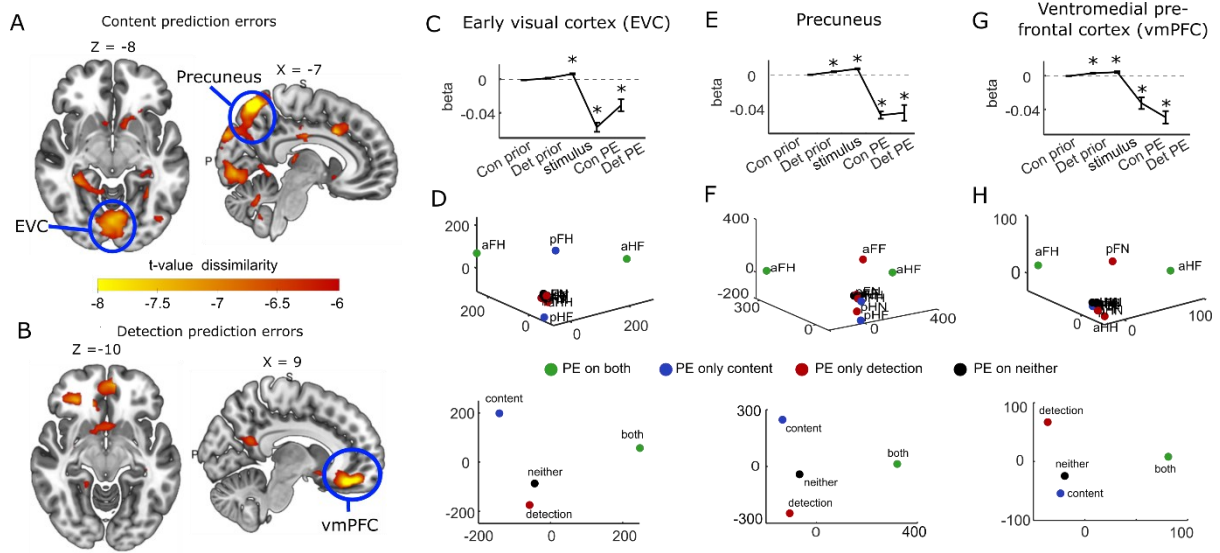
491 When inspecting the beta weights for each of the five RDMs in selected brain regions, we observed
 492 that activity patterns in both early visual cortex (EVC) as well as vmPFC showed strong *negative*
 493 relationships with both our hypothesized content and detection PE RDMs (see Fig. 5 beta plots below
 494 panels B and C). Negative relationships indicate that the conditions that are hypothesized to show
 495 similar patterns of activity in fact show dissimilar patterns of activity in the brain, and vice versa.
 496 Within the RSA literature, both similarity- and dissimilarity-based neural coding schemes have been
 497 hypothesized (Kriegeskorte et al., 2008). For example, in the fusiform face area, faces are encoded as
 498 more similar to other faces than to houses (positive correlation with a category membership RDM)
 499 whereas in face-selective parts of the inferotemporal cortex (IT), faces are encoded as more *dissimilar*
 500 to other faces, indicating exemplar encoding, revealing that this region is sensitive to the identity of

501 the face (negative correlation with a category membership RDM (Kriegeskorte et al., 2008). In the
502 current context, negative correlations with the hypothesized PE RDMs would therefore be in line with
503 the corresponding brain region being sensitive to the identity of the PE (e.g. whether it tracks a
504 violation of face or house predictions).

505 Therefore, to further explore these negative correlations, we expanded our search to examine
506 negative correlations with content and detection PE RDMs in a whole-brain searchlight analysis. This
507 approach indeed revealed significant dissimilarity PE encoding in several brain areas (Fig. 6). Content
508 PE-related patterns were predominant in posterior sensory regions (Fig. 6A, Appendix A – Table 3) and
509 detection PE-related patterns were predominant in prefrontal regions (Fig. 6B, Appendix A -Table 4).
510 However, directly comparing detection PE RDMs against content PE RDMs in a whole-brain contrast
511 revealed no significant differences between the two PE maps. Together these results indicate that
512 while neural patterns associated with content and detection PEs are predominantly expressed in
513 different parts of the brain, this distinction is graded rather than discrete.

514 To further characterize the relationship between content and detection PEs within different
515 brain areas, we examined activity patterns within functional ROIs selected for their dominance of
516 content- or detection-PE effects: the early visual cortex (EVC) and the Precuneus for content PEs and
517 the ventromedial prefrontal cortex (vmPFC) for detection PEs (encircled in blue in Fig. 6A&B). Note
518 that effects of the other RDMs in these ROIs cannot be explained by collinearity between the RDMs,
519 as the maximum correlation between regressors was low (0.11). Besides the expected content PE
520 effect, both the EVC and the Precuneus also showed a significant negative correlation with the
521 detection PE RDM (EVC: $t(26) = -3.68, p = 0.0011$; Precuneus: EVC: $t(26) = -4.93, p = 0.0008$; Fig 6C&E).
522 Furthermore, the vmPFC ROI defined based on the detection PE effect also showed a significant
523 negative correlation with the content PE RDM ($t(26) = -3.99, p = 0.0005$). This suggests that the brain
524 regions that are modulated by content PEs are also sensitive to detection PEs and vice versa.

525



526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

Figure 6. Dissimilarity coding of prediction errors. (A) Dissimilarity in the encoding of content prediction errors: conditions in which there was a prediction error on content (e.g. a face was expected and a house was presented) were more dissimilar/less similar to each other than to conditions without a prediction error on content. (B) Dissimilarity in the encoding of detection prediction errors: conditions in which there was a prediction error on detection (e.g. noise was expected but a face or house was presented) were more dissimilar/less similar to each other than to conditions without a prediction error on detection. (C, E, G) Beta values for the three blue encircled regions of interest (ROI) for each of the RDMs. * indicates the RDM significantly predicted activation pattern in this ROI (Bonferroni corrected). (D, F, H) Multidimensional scaling (MDS) of activation patterns for each cue-target combination in the three ROIs. Top: MDS in 3 dimensions with each cue-target combination plotted separately. Bottom: MDS in 2 dimensions in which conditions are grouped together based on whether they contain prediction errors on both content and detection (green), only on content (blue), only on detection (red) or on neither (black). Statistical maps in (A) and (B) are thresholded at $p < 0.001$ uncorrected; see Tables 4 and 5 for details of clusters surviving whole-brain correction including the EVC, precuneus and vmPFC ROIs.

542

543

544

545

546

547

548

549

To further characterize the representational structure of activity profiles in these regions, we performed classical multi-dimensional scaling (MDS) on the similarity of the activation patterns. MDS visualizes the similarity in neural patterns between conditions by projecting the data into a lower dimensional space in which similar conditions are plotted nearby to each other. We note that such visualisations were expected to recapitulate the RDM used to identify the ROI – for instance, we expect that within vmPFC we should see detection PEs as being encoded distinct from other trial types. However, it is possible that the nature of voxel patterns within these candidate ROIs diverge in other interesting ways.

550

551

552

553

554

555

Within the EVC, all conditions that did not contain a content PE were represented as similar to each other (Fig. 6D, red and black dots) whereas conditions containing a content PE were dissimilar to each other as well as to other conditions, as expected from the content PE RDM (Fig. 6C, green and blue dots). Interestingly, however, and in line with the observation that the detection PE RDM also showed a significant effect in this region, content PEs that also violated detection predictions (green dots) were dissimilar both to each other and to prediction errors on content only (blue dots). Together,

556 this suggests that in the EVC, once a content prediction is violated, activation diverges according to (a)
557 the exact type of content violation (face > house or house > face) and (b) whether the detection
558 prediction is also violated or not. Note that in the EVC, prediction errors on detection only (red dots)
559 were represented as similar to conditions in which no prediction error occurred (black dots),
560 suggesting that the EVC is only modulated by detection PEs when content predictions are violated.

561 Within the vmPFC, in contrast, all conditions that did not contain a detection PE were
562 represented as similar to each other (Fig. 6F, blue and black dots) whereas conditions containing a
563 detection PE were dissimilar to each other as well as to other conditions, as expected from the
564 detection PE RDM (Fig. 6F, green and red dots). However, and going beyond the hypothesis RDM, MDS
565 showed that trials that additionally violated content predictions (green dots) were also dissimilar to
566 each other and to prediction errors on detection only (red dots), in line with the observation that
567 content PEs also modulate activity patterns in this region. Note that in the vmPFC, in contrast to the
568 EVC, conditions with prediction errors on content only (blue dots) were represented as similar to
569 conditions in which no prediction error occurred (black dots), suggesting the vmPFC is only modulated
570 by content PEs when detection predictions are violated as well. Finally, the (pre-)cuneus showed a
571 significant effect of both the content and detection PE RDMs and closer inspection of the activity
572 patterns in this region (Fig. 6 H) showed that, in contrast to the EVC and vmPFC, in the (pre-)cuneus
573 all conditions with any kind of PE were associated with diverging activity patterns.

574

575 **Discussion**

576 In this study we set out to test whether inferences about perceptual content (*what* is perceived) and
577 inferences about detection of that content (*whether* something is perceived) are underpinned by
578 distinct neural substrates. To this end, we developed a novel experimental paradigm that used
579 compound cues to separately induce expectations about perceptual content and detection of that
580 content. We found that both content and detection expectations influenced reaction times, with
581 higher probability stimuli being identified more quickly. Using a no-report version of this paradigm in
582 conjunction with fMRI, we found that prediction errors on content correlated most strongly with
583 posterior visual brain areas, whereas prediction errors on detection correlated most strongly with
584 prefrontal brain areas. However, contrary to our hypothesis, these representations were not
585 orthogonal. Instead, prediction errors on one level gated the expression of prediction errors at the
586 other level. Taken together, our results suggest that inferences on content and detection of content
587 rely on distinct but interacting neural computations.

588 We observed a negative univariate effect of detection prediction errors in IFC, indicating that
589 in this region, activation was higher when a detection prediction was confirmed. One possibility is that

590 this region encodes a detection prior, with activity being strengthened in conditions in which these
591 priors are reinforced by matching input. However, within most neuronal models of predictive
592 processing, priors and prediction errors are assumed to be encoded within the same brain region
593 (Bastos et al., 2012). An alternative possibility is that a confirmation effect in IFC reflects a signature
594 of perceptual confidence (Cortese et al., 2016; Hilgenstock et al., 2014; Shekhar & Rahnev, 2018),
595 which is likely to be higher when predictions are confirmed compared to when they are violated.
596 Further work is needed to identify neural substrates supporting putative detection-specific confidence
597 signals, and distinguish these from other aspects of metacognition (Mazor et al., 2020, 2022).

598 Furthermore, our representational similarity analyses (RSA) revealed that prediction errors at
599 both content and detection levels were encoded as being dissimilar to each other. The pattern of
600 dissimilarity that we found indicates a sharp distinction in the initial trigger for PE coding in EVC and
601 prefrontal cortex – in EVC, the trigger for coding PEs is a violation of content expectations, whereas in
602 vmPFC, the trigger for coding PEs is a violation of detection expectations. Specifically, a divergence in
603 activity patterns in EVC is triggered by a content PE – here, a violation in the expectation of face or
604 house. Once a content PE is triggered, then EVC activity patterns go on to represent the type of PE
605 within the full compound cue space, tracking violations of both content and detection expectations.
606 In contrast, a divergence in activity patterns in vmPFC is triggered by a detection PE – whether an
607 absence expectation has been violated by stimulus presence, or whether a presence expectation has
608 been violated by stimulus absence. Once a detection PE is triggered, then vmPFC activity patterns also
609 go on to represent the type of PE within the full compound cue space, inheriting information about
610 content (face vs. house) violations. The Precuneus shows an intermediate effect, representing
611 diverging patterns for every type of PE. This pattern dissimilarity suggests that prediction errors are
612 encoded in an exemplar specific way, similar to individual faces in IT (Kriegeskorte et al., 2008). Further
613 work is needed to fully understand what this dissimilarity-based coding implies for the computational
614 underpinnings of inferences about perceptual content and detection.

615 As our neuroimaging results were obtained in the absence of reports, they provide evidence
616 in favour of an architecture in which detection prediction errors are automatically elicited even under
617 passive viewing conditions. Our results are consistent with other studies observing prefrontal
618 correlates of subjective detection in the absence of overt report (Hatamimajoumerd et al., 2022). We
619 note that the focus of the current study is on distinguishing between neural signatures of inference
620 on content, and detection of content, and did not set out to measure variation in subjective perception
621 or awareness. However, our results bear on the possible neural architectures supporting predictive
622 processing accounts consciousness, as inferences on detection – i.e. whether subjects are “aware” or
623 “unaware” of particular stimulus features – are the cornerstone of conscious reportability. We also

624 note that our findings cannot be explained by mere stimulus effects, as our focus here is on how the
625 stimulus interacts with an experimentally-manipulated expectation, thereby generating a prediction
626 error signal. Future studies are necessary to investigate whether the inferences on content and
627 detection we identify here relate to changes in conscious experience while keeping stimulus input
628 near threshold (Frith et al., 1999; Leopold & Logothetis, 1996).

629 One example of such an approach is a recent MEG study which revealed a neural signature of
630 the content of false percepts in the occipital lobe, whereas confidence in stimulus detection was
631 reflected in a parieto-frontal network (Haarsma, Hetenyi, et al., 2024). Combining such a false percept
632 paradigm with the type of compound content-detection cue employed here (Haarsma, Kaltenmaier,
633 et al., 2024) is a promising avenue to investigate how the neural correlates identified here relate to
634 fluctuations in conscious experience. To ensure that effects are not due to reports in that case, future
635 research could develop no-report read-outs of perceptual content such as eye movements (Frassle et
636 al., 2014; Frässle et al., 2013) to capture aspects of both discrimination and detection. Finally, It would
637 also be interesting to seek to causally intervene on regions (such as vmPFC) exhibiting signatures of
638 detection prediction error (for instance, using multivariate neurofeedback (Taschereau-Dumouchel et
639 al., 2021)), and ask whether and how such interventions alter conscious experience.

640 Several ideas have been advanced to accommodate conscious awareness (in the form of
641 subjective detection) within a predictive processing framework (Clark et al., 2019; Doerig et al., 2020;
642 Fleming, 2020; Hobson & Friston, 2014, 2012; Hohwy et al., 2008; Hohwy & Seth, 2020). One
643 instantiation of such an architecture proposes that conscious detection arises from inferences deep
644 within a perceptual hierarchy. Interestingly, in line with this hierarchical view, detection PEs were
645 preferentially localised to a vmPFC region overlapping with the default model network (DMN) (Raichle,
646 2015) which is proposed to occupy a deep position within a cognitive hierarchy (Margulies et al., 2016).
647 Other work has linked the vmPFC to carrying information about latent (unobservable) perceptual
648 spaces, such as hidden states governing task structure, or links between arbitrary stimuli on a graph
649 (Park et al., 2020; Schuck et al., 2016). Another key node of the DMN, the Precuneus, was also evident
650 in the PE RSA analysis, and showed clear detection- as well as content-related prediction error effects
651 (Raichle, 2015). Other research has shown that the Precuneus is modulated by both the level of
652 awareness (Bisenius et al., 2015; Cavanna, 2007; Kjaer et al., 2001) and represents stimulus content
653 (Doesburg et al., 2009) – consistent with it inhabiting an intermediate position in a perceptual
654 hierarchy.

655 In conclusion, using a novel experimental paradigm we show that prediction errors on
656 perceptual content and detection of content are encoded in distinct but interacting activity patterns
657 in the human brain. These results are consistent with a proposal that detection may require distinct

658 neural computations that go beyond those required for inferences on content itself. More generally,
659 our findings provide a framework for future empirical and theoretical studies that incorporate and
660 model detection and discrimination as distinct dimensions within powerful predictive processing
661 accounts of perception and cognition.

662

663 **Data Availability**

664 Second-level maps are uploaded on NeuroVault at: <https://neurovault.org/collections/14778/>.

665

666 **Code Availability**

667 All analysis code is available on: <https://github.com/NadineDijkstra/HSFPA>

668

669 **Author contribution statement**

670 Conceptualization: N.D., P.K. & S.M.F.; Data Curation: N.D. & O.W.; Formal Analysis: N.D. & O.W.;
671 Funding Acquisition: S.M.F.; Investigation: N.D. & O.W.; Methodology: N.D., O.W. & S.M.F.; Project
672 Administration: N.D. & O.W.; Resources: S.M.F.; Software: N.D. & O.W.; Supervision: N.D. & S.M.F.;
673 Validation: N.D.; Visualization: N.D. & O.W.; Writing – original draft: N.D.; Writing – review & editing:
674 N.D., O.W., P.K. & S.M.F.

675

676 **Gender citation bias**

677 Number of DOIs categorized: 54, Number of DOIs due to missing author data: 7. Proportion per
678 category: MM 0.87; WM: 0.074; MW: 0.037; WW: 0.019. GCBI per category: MM 1.139; WM -0.769;
679 MW -0.678; WW: -0.884.

680

681 **References**

682 Andersson, J. L. R., Hutton, C., Ashburner, J., Turner, R., & Friston, K. (2001). Modeling Geometric

683 Deformations in EPI Time Series. *NeuroImage*, 13(5), 903–919.

684 <https://doi.org/10.1006/NIMG.2001.0746>

685 Ashburner, J., & Friston, K. J. (2005). Unified segmentation. *NeuroImage*, 26(3), 839–851.

686 <https://doi.org/10.1016/J.NEUROIMAGE.2005.02.018>

687 Azzopardi, P., & Cowey, A. (1997). Is blindsight like normal, near-threshold vision? *Proceedings of the*

688 *National Academy of Sciences of the United States of America*, 94(25).

689 <https://doi.org/10.1073/pnas.94.25.14190>

690 Bastos, A. M., Usrey, W. M., Adams, R. A., Mangun, G. R., Fries, P., & Friston, K. J. (2012). Canonical
691 Microcircuits for Predictive Coding. *Neuron*, *76*(4), 695–711.
692 <https://doi.org/10.1016/j.neuron.2012.10.038>

693 Bisenius, S., Trapp, S., Neumann, J., & Schroeter, M. L. (2015). Identifying neural correlates of visual
694 consciousness with ALE meta-analyses. *NeuroImage*, *122*, 177–187.
695 <https://doi.org/10.1016/j.neuroimage.2015.07.070>

696 Brodersen, K. H., Penny, W. D., Harrison, L. M., Daunizeau, J., Ruff, C. C., Duzel, E., Friston, K. J., &
697 Stephan, K. E. (2008). Integrated Bayesian models of learning and decision making for
698 saccadic eye movements. *Neural Networks*, *21*(9).
699 <https://doi.org/10.1016/j.neunet.2008.08.007>

700 Brown, R. (2015). The HOROR theory of phenomenal consciousness. *Philosophical Studies*, *172*(7),
701 1783–1794. <https://doi.org/10.1007/s11098-014-0388-7>

702 Carpenter, R. H. S., & Williams, M. L. L. (1995). Neural computation of log likelihood in control of
703 saccadic eye movements. *Nature*, *377*(6544). <https://doi.org/10.1038/377059a0>

704 Cavanna, A. E. (2007). The Precuneus and Consciousness. *CNS Spectrums*, *12*(7), 545–552.
705 <https://doi.org/10.1017/S1092852900021295>

706 Clark, A., Friston, K., & Wilkinson, S. (2019). Bayesing qualia: Consciousness as inference, not raw
707 datum. *Journal of Consciousness Studies*, *26*(9–10), 19–33.

708 Cortese, A., Amano, K., Koizumi, A., Kawato, M., & Lau, H. (2016). Multivoxel neurofeedback
709 selectively modulates confidence without changing perceptual performance. *Nature*
710 *Communications*, *7*(1), 1–18. <https://doi.org/10.1038/ncomms13669>

711 Dehaene, S., & Changeux, J. P. (2011). Experimental and Theoretical Approaches to Conscious
712 Processing. *Neuron*, *70*(2), 200–227. <https://doi.org/10.1016/j.neuron.2011.03.018>

713 Dehaene, S., Naccache, L., Cohen, L., Bihan, D. L., Mangin, J. F., Poline, J. B., & Rivière, D. (2001).
714 Cerebral mechanisms of word masking and unconscious repetition priming. *Nature*
715 *Neuroscience*, *4*(7), 752–758. <https://doi.org/10.1038/89551>

716 Destrieux, C., Fischl, B., Dale, A., & Halgren, E. (2010). Automatic parcellation of human cortical gyri
717 and sulci using standard anatomical nomenclature. *NeuroImage*, *53*(1), 1–15.
718 <https://doi.org/10.1016/j.neuroimage.2010.06.010>

719 Doerig, A., Schurger, A., & Herzog, M. H. (2020). Hard criteria for empirical theories of consciousness.
720 *Cognitive Neuroscience*, 1–22. <https://doi.org/10.1080/17588928.2020.1772214>

721 Doesburg, S. M., Green, J. J., McDonald, J. J., & Ward, L. M. (2009). Rhythms of consciousness:
722 Binocular rivalry reveals large-scale oscillatory network dynamics mediating visual
723 perception. *PloS One*, *4*(7), e6142. <https://doi.org/10.1371/journal.pone.0006142>

724 Fleming, S. M. (2020). Awareness as inference in a higher-order state space. *Neuroscience of*
725 *Consciousness*, *2020*(1). <https://doi.org/10.1093/nc/niz020>

726 Frassle, S., Sommer, J., Jansen, A., Naber, M., & Einhauser, W. (2014). Binocular Rivalry: Frontal
727 Activity Relates to Introspection and Action But Not to Perception. *Journal of Neuroscience*,
728 *34*(5), 1738–1747. <https://doi.org/10.1523/JNEUROSCI.4403-13.2014>

729 Frässle, S., Sommer, J., Naber, M., Jansen, A., & Einhäuser, W. (2013). Neural Correlates of Binocular
730 Rivalry as measured in fMRI are partially confounded by observers' active report. *Journal of*
731 *Vision*, *13*(9), 937. <https://doi.org/10.1167/13.9.937>

732 Friston, K., Kilner, J., & Harrison, L. (2006). A free energy principle for the brain. *Journal of Physiology*
733 *Paris*, *100*(1–3), 70–87. <https://doi.org/10.1016/j.jphysparis.2006.10.001>

734 Frith, C., Perry, R., & Lumer, E. (1999). The neural correlates of conscious experience: An
735 experimental framework. *Trends in Cognitive Sciences*, *3*(3), 105–114.
736 [https://doi.org/10.1016/s1364-6613\(99\)01281-4](https://doi.org/10.1016/s1364-6613(99)01281-4)

737 Green, D. M., & Swets, J. A. (1966). Signal detection theory and psychophysics. In *John Wiley* (Vol. 5).
738 [https://doi.org/10.1016/0022-460x\(67\)90197-6](https://doi.org/10.1016/0022-460x(67)90197-6)

739 Haarsma, J., Hetenyi, D., & Kok, P. (2024). *Shared and diverging neural dynamics underlying false and*
740 *veridical perception* (p. 2023.11.16.567367). bioRxiv.
741 <https://doi.org/10.1101/2023.11.16.567367>

742 Haarsma, J., Kaltenmaier, A., Fleming, S. M., & Kok, P. (2024). *Expectations about presence enhance*
743 *the influence of content-specific expectations on low-level orientation judgements* (p.
744 2024.02.22.581334). bioRxiv. <https://doi.org/10.1101/2024.02.22.581334>

745 Hatamimajoumerd, E., Ratan Murty, N. A., Pitts, M., & Cohen, M. A. (2022). Decoding perceptual
746 awareness across the brain with a no-report fMRI masking paradigm. *Current Biology*,
747 32(19), 4139-4149.e4. <https://doi.org/10.1016/J.CUB.2022.07.068>

748 Hilgenstock, R., Weiss, T., & Witte, O. W. (2014). You'd Better Think Twice: Post-Decision Perceptual
749 Confidence. *NeuroImage*, 99. <https://doi.org/10.1016/j.neuroimage.2014.05.049>

750 Hobson, J. A., & Friston, K. (2014). Consciousness, dreams, and inference: The cartesian theatre
751 revisited. *Journal of Consciousness Studies*, 21(1–2), 6–32.

752 Hobson, J. A., & Friston, K. J. (2012). Waking and dreaming consciousness: Neurobiological and
753 functional considerations. *Progress in Neurobiology*, 98(1), 82–98.
754 <https://doi.org/10.1016/J.PNEUROBIO.2012.05.003>

755 Hohwy, J. (2012). Attention and conscious perception in the hypothesis testing brain. *Frontiers in*
756 *Psychology*, 3(April), 96. <https://doi.org/10.3389/fpsyg.2012.00096>

757 Hohwy, J., Roepstorff, A., & Friston, K. (2008). Predictive coding explains binocular rivalry: An
758 epistemological review. *Cognition*, 108(3), 687–701.
759 <https://doi.org/10.1016/j.cognition.2008.05.010>

760 Hohwy, J., & Seth, A. (2020). Predictive processing as a systematic basis for identifying the neural
761 correlates of consciousness. *Philosophy and the Mind Sciences*, 1(II).
762 <https://doi.org/10.33735/phimisci.2020.ii.64>

763 Kersten, D., Mamassian, P., & Yuille, A. (2004). Object perception as Bayesian inference. *Annual*
764 *Review of Psychology*, 55, 271–304.
765 <https://doi.org/10.1146/annurev.psych.55.090902.142005>

766 Kjaer, T. W., Nowak, M., Kjaer, K. W., Lou, A. R., & Lou, H. C. (2001). Precuneus-prefrontal activity
767 during awareness of visual verbal stimuli. *Consciousness and Cognition*, *10*(3), 356–365.
768 <https://doi.org/10.1006/ccog.2001.0509>

769 Kok, P., Brouwer, G. J., van Gerven, M. A. J., & de Lange, F. P. (2013). Prior expectations bias sensory
770 representations in visual cortex. *Journal of Neuroscience*, *33*(41).
771 <https://doi.org/10.1523/JNEUROSCI.0742-13.2013>

772 Kriegeskorte, N., Mur, M., & Bandettini, P. (2008). Representational similarity analysis—Connecting
773 the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, *2*(November), 4.
774 <https://doi.org/10.3389/neuro.06.004.2008>

775 Lau, H. (2019). Consciousness, Metacognition, & Perceptual Reality Monitoring. *PsychArxiv*, 1–17.

776 Lau, H. (2022). In Consciousness we Trust. In *In Consciousness we Trust*.
777 <https://doi.org/10.1093/oso/9780198856771.001.0001>

778 Lau, H. C. (2007). A higher order Bayesian decision theory of consciousness. *Progress in Brain*
779 *Research*, *168*, 35–48. [https://doi.org/10.1016/S0079-6123\(07\)68004-2](https://doi.org/10.1016/S0079-6123(07)68004-2)

780 Lau, H., & Rosenthal, D. (2011). Empirical support for higher-order theories of conscious awareness.
781 *Trends in Cognitive Sciences*, *15*(8), 365–373. <https://doi.org/10.1016/j.tics.2011.05.009>

782 Leopold, D. A., & Logothetis, N. K. (1996). Activity changes in early visual cortex reflect monkeys'
783 percepts during binocular rivalry. *Nature*, *379*, 549–553. <https://doi.org/10.1038/379549a0>

784 Lundqvist, D., Flykt, A., & Öhman, A. (1998). *The Karolinska Directed Emotional Faces—KDEF, CD*
785 *ROM from Department of Clinical Neuroscienc.*

786 Marcel, A. J. (1983). Conscious and unconscious perception: Experiments on visual masking and
787 word recognition. *Cognitive Psychology*, *15*(2). [https://doi.org/10.1016/0010-](https://doi.org/10.1016/0010-0285(83)90009-9)
788 [0285\(83\)90009-9](https://doi.org/10.1016/0010-0285(83)90009-9)

789 Margulies, D. S., Ghosh, S. S., Goulas, A., Falkiewicz, M., Huntenburg, J. M., Langs, G., Bezgin, G.,
790 Eickhoff, S. B., Castellanos, F. X., Petrides, M., Jefferies, E., & Smallwood, J. (2016). Situating
791 the default-mode network along a principal gradient of macroscale cortical organization.

792 *Proceedings of the National Academy of Sciences of the United States of America*, 113(44).
793 <https://doi.org/10.1073/pnas.1608282113>

794 Mars, R. B., Debener, S., Gladwin, T. E., Harrison, L. M., Haggard, P., Rothwell, J. C., & Bestmann, S.
795 (2008). Trial-by-trial fluctuations in the event-related electroencephalogram reflect dynamic
796 changes in the degree of surprise. *Journal of Neuroscience*, 28(47).
797 <https://doi.org/10.1523/JNEUROSCI.2925-08.2008>

798 Marvan, T., & Havlík, M. (2021). Is predictive processing a theory of perceptual consciousness? *New*
799 *Ideas in Psychology*, 61, 100837. <https://doi.org/10.1016/J.NEWIDEAPSYCH.2020.100837>

800 Mashour, G. A., Roelfsema, P., Changeux, J. P., & Dehaene, S. (2020). Conscious Processing and the
801 Global Neuronal Workspace Hypothesis. *Neuron*, 105(5), 776–798.
802 <https://doi.org/10.1016/j.neuron.2020.01.026>

803 Mazor, M., Dijkstra, N., & Fleming, S. M. (2022). Dissociating the Neural Correlates of Subjective
804 Visibility from Those of Decision Confidence. *Journal of Neuroscience*, 42(12), 2562–2569.
805 <https://doi.org/10.1523/JNEUROSCI.1220-21.2022>

806 Mazor, M., Friston, K., & Fleming, S. (2020). Distinct neural contributions to metacognition for
807 detecting, but not discriminating visual stimuli. *eLife*, 9, 1–34.
808 <https://doi.org/10.1101/853366>

809 Merten, K., & Nieder, A. (2012). Active encoding of decisions about stimulus absence in primate
810 prefrontal cortex neurons. *Proceedings of the National Academy of Sciences*, 109(16), 6289–
811 6294. <https://doi.org/10.1073/pnas.1121084109>

812 Meyen, S., Zerweck, I. A., Amado, C., von Luxburg, U., & Franz, V. H. (2022). Advancing research on
813 unconscious priming: When can scientists claim an indirect task advantage? *Journal of*
814 *Experimental Psychology. General*, 151(1), 65–81. <https://doi.org/10.1037/xge0001065>

815 Morales, J. (2022). Introspection Is Signal Detection. *British Journal for the Philosophy of Science*.
816 <https://doi.org/10.1086/715184>

817 Morales, J., Lau, H., & Fleming, S. M. (2018). Domain-general and domain-specific patterns of activity
818 supporting metacognition in human prefrontal cortex. *Journal of Neuroscience*, *38*(14),
819 3534–3546. <https://doi.org/10.1523/JNEUROSCI.2360-17.2018>

820 Park, S. A., Miller, D. S., Nili, H., Ranganath, C., & Boorman, E. D. (2020). Map Making: Constructing,
821 Combining, and Inferring on Abstract Cognitive Maps. *Neuron*, *107*(6).
822 <https://doi.org/10.1016/j.neuron.2020.06.030>

823 Persaud, N., Davidson, M., Maniscalco, B., Mobbs, D., Passingham, R. E., Cowey, A., & Lau, H. (2011).
824 Awareness-related activity in prefrontal and parietal cortices in blindsight reflects more than
825 superior visual performance. *NeuroImage*, *58*(2).
826 <https://doi.org/10.1016/j.neuroimage.2011.06.081>

827 Peters, M. A. K., & Lau, H. (2015). Human observers have optimal introspective access to perceptual
828 processes even for visually masked stimuli. *eLife*, *4*(OCTOBER2015).
829 <https://doi.org/10.7554/eLife.09651>

830 Raichle, M. E. (2015). The Brain's Default Mode Network. *Annual Review of Neuroscience*, *38*.
831 <https://doi.org/10.1146/annurev-neuro-071013-014030>

832 Schuck, N. W., Cai, M. B., Wilson, R. C., & Niv, Y. (2016). Human Orbitofrontal Cortex Represents a
833 Cognitive Map of State Space. *Neuron*, *91*(6). <https://doi.org/10.1016/j.neuron.2016.08.019>

834 Shekhar, M., & Rahnev, D. (2018). Distinguishing the roles of dorsolateral and anterior PFC in visual
835 metacognition. *Journal of Neuroscience*, *38*(22). [https://doi.org/10.1523/JNEUROSCI.3484-](https://doi.org/10.1523/JNEUROSCI.3484-17.2018)
836 [17.2018](https://doi.org/10.1523/JNEUROSCI.3484-17.2018)

837 Sladky, R., Friston, K. J., Tröstl, J., Cunnington, R., Moser, E., & Windischberger, C. (2011). Slice-timing
838 effects and their correction in functional MRI. *NeuroImage*, *58*(2), 588–594.
839 <https://doi.org/10.1016/J.NEUROIMAGE.2011.06.078>

840 Taschereau-Dumouchel, V., Cortese, A., Lau, H., & Kawato, M. (2021). Conducting decoded
841 neurofeedback studies. *Social Cognitive and Affective Neuroscience*, *16*(8).
842 <https://doi.org/10.1093/scan/nsaa063>

- 843 Tsuchiya, N., Wilke, M., Frässle, S., & Lamme, V. A. F. (2015). No-Report Paradigms: Extracting the
844 True Neural Correlates of Consciousness. *Trends in Cognitive Sciences*, 19(12), 757–770.
845 <https://doi.org/10.1016/j.tics.2015.10.002>
- 846 van Vugt, B., Dagnino, B., Vartak, D., Safaai, H., Panzeri, S., Dehaene, S., & Roelfsema, P. R. (2018).
847 The threshold for conscious report: Signal loss and response bias in visual and frontal cortex.
848 *Science*, eaar7186. <https://doi.org/10.1126/science.aar7186>
- 849 Weiskrantz, L., Warrington, E. K., Sanders, M. D., & Marshall, J. (1974). Visual capacity in the
850 hemianopic field following a restricted occipital ablation. *Brain*, 97(1), 709–728.
851 <https://doi.org/10.1093/brain/97.1.709>
- 852 Whyte, C. J., & Smith, R. (2021). The predictive global neuronal workspace: A formal active inference
853 model of visual consciousness. *Progress in Neurobiology*, 199, 101918.
854 <https://doi.org/10.1016/j.pneurobio.2020.101918>

855

856 **Conflict of interest**

857 The authors declare no conflict of interest.

858

859 **Acknowledgments**

860 ND is supported by the European Union’s Horizon 2020 research and innovation programme under
861 the Marie Skłodowska-Curie grant agreement No. 882832. PK is supported by a Wellcome/Royal
862 Society Sir Henry Dale Fellowship (218535/Z/19/Z) and a European Research Council (ERC) Starting
863 Grant (948548). SMF is a CIFAR Fellow in the Brain, Mind and Consciousness Program, and supported
864 by a Wellcome/Royal Society Sir Henry Dale Fellowship (206648/Z/17/Z) and UK Research and
865 Innovation (UKRI) under the UK government’s Horizon Europe funding guarantee [selected as ERC
866 Consolidator, grant number 101043666]. The Wellcome Centre for Human Neuroimaging is supported
867 by core funding from the Wellcome Trust (203147/Z/16/Z). The Max Planck UCL Centre is a joint
868 initiative supported by UCL and the Max Planck Society. For the purposes of Open Access, the author
869 has applied a CC-BY public copyright license to any author accepted manuscript version arising from
870 this submission.

871

872

873

874 **Appendix A**

875 **Table 1. Stimulus-encoding RSA clusters.** All clusters are significant at $p < 0.05$ FWE corrected for multiple
 876 comparisons within the whole brain volume, with a cluster-forming threshold of $p < 0.001$ uncorrected. Clusters
 877 are included in the table if they surpass a threshold of 50 voxels. T-values, labels and coordinates are given for
 878 the peak within each cluster.

N voxels	t-value	AAL label	X Y Z
7541	9	Cuneus R	8 -74 40
6226	7.8	Hippocampus R	20 -39 1
5986	7.8	Cerebellum L	-1 -79 -8
5921	8.2	Cerebellum Crust1 R	38 -74 -22
4289	8.6	Hippocampus L	-23 -38 4
3068	6.9	Cerebellum L	-32 -64 -24
1129	7.4	-	-4 19 37
1121	7	Vermis 8	-5 -63 x -31
877	7	-	2 9 -5
812	7.2	Cingulum Mid R	3 -16 34
612	6.8	Frontal Mid L	-43 34 32
560	6.7	Cerebellum L	-32 -43 -27
498	7.1	Precuneus L	-10 -43 -27
372	7.8	Parietal Sub L	-22 -64 51
347	7.8	-	-23 22 -5
126	6.4	-	-15 -2 25
107	7	Cerebellum R	25 -53 -41
104	6.4	Cerebellum R	-14 -64 -36
93	6.7	Frontal Inf Oper L	-43 9 23
79	6.3	Cerebellum L	-12 -56 -48
55	6.6	Frontal Inf L	-43 23 22

879

880 **Table 2. Detection prior RSA clusters.** All clusters are significant at $p < 0.05$ FWE corrected for multiple
 881 comparisons within the whole brain volume, with a cluster-forming threshold of $p < 0.001$ uncorrected. T-values,
 882 labels and coordinates are given for the peak within each cluster.

N voxels	t-value	AAL label	X Y Z
318	7.2	Rectus R	10 32 -18
119	6.3	Frontal Mid Orb R	28 39 -12

883

884

885 **Table 3. Content prediction error RSA clusters.** All clusters are significant at $p < 0.05$ FWE corrected for multiple
886 comparisons within the whole brain volume, with a cluster-forming threshold of $p < 0.001$ uncorrected. Clusters
887 are included in the table if they surpass a threshold of 50 voxels. T-values, labels and coordinates are given for
888 the peak within each cluster.

N voxels	t-value	AAL label	X Y Z
10443	10.1	Precuneus	8 -72 46
8147	8.2	Putamen	22 -40 -2
4441	7.9	Cerebellum Crus R	39 -74 -22
3506	8.7	Cingulum Mid L	1 18 68
1951	9.4	Parietal Sup L	-20 -69 50
1899	8.0	Cingulum Mid R	2 -17 34
1370	7.4	Parietal Sup R	31 -55 57
902	6.8	-	-17 18 -5
368	6.7	Occipital Inf L	-43 -75 -10
358	6.4	-	-14 -2 24
299	6.5	Frontal Sup R	22 54 32
206	6.3	Frontal Mid L	-41 33 33
144	6.3	Frontal Sup Medial L	-1 60 15
140	6.7	Precuneus L	-10 44 72
122	6.5	Frontal Mid L	-26 40 30
115	6.7	Frontal Inf Oper L	-42 8 24
65	6.2	Caudate R	-7 13 8

889

890 **Table 4. Detection prediction error RSA clusters.** All clusters are significant at $p < 0.05$ FWE corrected for
891 multiple comparisons within the whole brain volume, with a cluster-forming threshold of $p < 0.001$ uncorrected.
892 Clusters are included in the table if they surpass a threshold of 50 voxels. T-values, labels and coordinates are
893 given for the peak within each cluster.

N voxels	t-value	AAL label	X Y Z
6599	8.9	Rectus R	2 40 15
4616	8.7	Cerebellum	-31 -51 24
2488	7.7	Angular L	-39 -69 45
1757	7.0	Cerebellum R	21 -44 -19
1340	6.7	Precuneus R	4 -56 21
1252	7.4	Frontal Mid Orb R	33 40 -11
1077	7.2	Frontal Mid L	-39 13 34
609	7.1	Vermis	-3 -64 -31

302	6.6	Amygdala R	28 1 -14
269	6.7	Frontal Sub Orb L	-33 59 -2
223	7.9	-	-23 -18 20
221	6.8	-	16 -2 25
145	6.3	Caudate L	13 15 10
139	6.9	Precuneus L	-12 -70 39
120	6.5	Frontal Mid L	-26 46 39
99	6.5	Frontal Mid R	32 50 12
97	6.5	Cerebellum L	-14 -53 -14
89	6.6	Cerebellum R	26 -50 -42
86	6.7	-	-16 -36 42
67	6.9	Hippocampus L	-35 -35 -6

894

895

896