

Research Articles: Behavioral/Cognitive

Dissociating the neural correlates of subjective visibility from those of decision confidence

<https://doi.org/10.1523/JNEUROSCI.1220-21.2022>

Cite as: J. Neurosci 2022; 10.1523/JNEUROSCI.1220-21.2022

Received: 14 June 2021

Revised: 14 December 2021

Accepted: 17 January 2022

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.jneurosci.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

1 Dissociating the neural correlates of
2 subjective visibility from those of
3 decision confidence

4 Matan Mazor*¹, Nadine Dijkstra*¹ & Stephen M. Fleming^{1,2,3}

5

6 Abbreviated title: Dissociating visibility from confidence

7

8 *These authors contributed equally to this work

9 ¹ Wellcome Centre for Human Neuroimaging, UCL

10 ² Max Planck UCL Centre for Computational Psychiatry and Ageing Research

11 ³ Department of Experimental Psychology, UCL

12

13

14

15

16

17

18

19

20 Correspondence should be addressed to Nadine Dijkstra n.dijkstra@ucl.ac.uk and Matan Mazor
21 mtnmzor@gmail.com

22

23

24

25 Abstract: A key goal of consciousness science is identifying neural signatures of being aware
26 vs. unaware of simple stimuli. This is often investigated in the context of near-threshold
27 detection, with reports of stimulus awareness being linked to heightened activation in a
28 frontoparietal network. However, due to reports of stimulus presence typically being associated
29 with higher confidence than reports of stimulus absence, these results could be explained by
30 frontoparietal regions encoding stimulus visibility, decision confidence or both. In an exploratory
31 analysis, we leverage fMRI data from 35 human participants (20 females) to disentangle these
32 possibilities. We first show that, whereas stimulus identity was best decoded from the visual
33 cortex, stimulus visibility (presence vs. absence) was best decoded from prefrontal regions. To
34 control for effects of confidence, we then selectively sampled trials prior to decoding to equalize
35 confidence distributions between absence and presence responses. This analysis revealed
36 striking differences in the neural correlates of subjective visibility in prefrontal cortex regions of
37 interest, depending on whether or not differences in confidence were controlled for. We interpret
38 our findings as highlighting the importance of controlling for metacognitive aspects of the
39 decision process in the search for neural correlates of visual awareness.

40
41 Significance statement: While much has been learned over the past two decades about the
42 neural basis of visual awareness, the role of the prefrontal cortex remains a topic of debate. By
43 applying decoding analyses to functional brain imaging data, we show that prefrontal
44 representations of subjective visibility are contaminated by neural correlates of decision
45 confidence. We propose a new analysis method to control for these metacognitive aspects of
46 awareness reports, and use it to reveal confidence-independent correlates of perceptual
47 judgments in a subset of prefrontal areas.

48
49
50

51 Introduction

52 In neuroimaging studies of visual perception, frontal and parietal cortices typically show
53 stronger activation when participants report being aware rather than unaware of a visual
54 stimulus (Sahraie et al., 1997; Dehaene et al., 2001; Fisch et al., 2009; Koivisto & Revonsuo,
55 2010). This finding is a cornerstone of several influential theories of awareness (e.g., *Global*
56 *Neuronal Workspace*: Dehaene, Sergent & Changeux, 2003; Dehaene., Changeux, &
57 Naccache, 2011; *Higher Order Thought*: Lau & Rosenthal, 2011; Brown, Lau, & LeDoux, 2019),
58 and is central to recent debates about the specific role of these regions in the generation of
59 subjective experience (Boly et al., 2017; Odegaard, Knight & Lau, 2017; Michel & Morales,
60 2020; Raccach, Block & Fox, 2021).

61 However, reports of awareness and unawareness of a visual stimulus differ not only in
62 terms of whether a stimulus was visible or not, but also in other cognitive factors (Bayne &
63 Hohwy, 2013). Specifically, when asked to rate their subjective confidence in near-threshold
64 detection, participants' confidence in decisions about stimulus presence is reliably higher than in
65 decisions about stimulus absence (Mazor, Friston & Fleming, 2020; Mazor, Moran & Fleming,
66 2021). This confidence asymmetry between judgments of presence and absence makes
67 interpreting frontoparietal activations in reports of visual awareness difficult: they may reflect
68 stimulus visibility, subjective confidence in the percept (which is higher when a stimulus is
69 detected), or both.

70 Consistent with the idea that frontoparietal activations found to correlate with awareness
71 might reflect confidence, the same regions associated with awareness reports are also found to
72 be implicated in reports of subjective confidence. For example, a coordinate-based meta-
73 analysis revealed that dorsolateral prefrontal cortex, lateral parietal cortex, and posterior medial
74 frontal cortex show a reliable parametric modulation of confidence (Vacarro & Fleming, 2018) -
75 all regions that have been associated with subjective visibility in previous studies (Sahraie et al.,
76 1997; Dehaene et al., 2001; Lau & Passingham, 2008; Fisch et al., 2009; Koivisto & Revonsuo,
77 2010). Importantly, these regions encode subjective confidence not only in perceptual decisions,
78 but also in memory-based (Morales, Lau & Fleming, 2018) and value-based (De-Martino et al.,
79 2013) decisions, suggesting that their link to subjective confidence is not solely in virtue of their
80 role in tracking subjective visibility.

81 Here, we set out to systematically dissociate the neural correlates of visibility and
82 confidence, and ask to what extent neural representations within a frontoparietal network track
83 one or both of these variables. To address this question, we performed a series of exploratory
84 analyses on neuroimaging data collected during performance-matched visual detection and
85 discrimination tasks with subjective confidence ratings (originally reported in Mazor et al., 2020).
86 We first asked where in the brain can we decode the presence or absence of a visual target
87 stimulus (a sinusoidal grating) from multivariate spatial activity patterns during the detection
88 task. By comparing these results against similar decoding of stimulus identity (grating
89 orientation) in a performance-matched discrimination task, we could control for non-specific
90 neural contributions to perceptual decision-making and report. Critically, by leveraging trial-wise
91 confidence ratings we were able to equate differences in subjective confidence between
92 conditions, allowing us to isolate neural representations associated with stimulus visibility. To
93 anticipate our results, we find that a number of prefrontal representations of stimulus visibility

94 are confounded with representations of confidence, but that a confidence-independent
95 representation of perceptual content is present in posterior medial frontal cortex (pMFC). Our
96 approach provides a novel method for controlling for such confidence effects in future studies of
97 visual awareness.

98 Methods

99 This is an exploratory analysis of neuroimaging data, originally reported in Mazor et al.
100 (2020). For a more elaborate description of the experimental design and behavioural findings,
101 see Mazor et al. (2020).

102 Participants

103 46 participants took part in the study (ages 18–36, mean = 24 ± 4). We applied the same
104 subject- and block-wise exclusion criteria as in the original study. Specifically, participants were
105 excluded for having low response accuracy, pronounced response bias, or insufficient variability
106 in their confidence ratings. 35 participants met our pre-specified inclusion criteria (ages 18–36,
107 mean = 24 ± 4 ; 20 females). We pre-registered a sample size of 35 to maximize statistical
108 power given resource limitations. This allowed us to detect a medium effect in a paired-samples
109 t-test (cohen's $d = 0.49$) with a power of 80%. All analyses are based on the included blocks
110 from these 35 participants.

111 Pre-registration was time-locked by initializing the pseudorandom number generator with
112 a hash of our pre-registered protocol folder (link:
113 github.com/matanmazor/detectionVsDiscrimination_fmri/tree/master/protocol_folder) prior to
114 determining the order and timing of experimental events (Mazor, Mazor & Mukamel, 2019).
115 Importantly, this pre-registration was motivated by a different set of hypotheses (tested in
116 Mazor, Friston & Fleming, 2020). The results we present here are derived from a data-driven,
117 exploratory set of analyses.

118 Experimental Design and Statistical Analysis

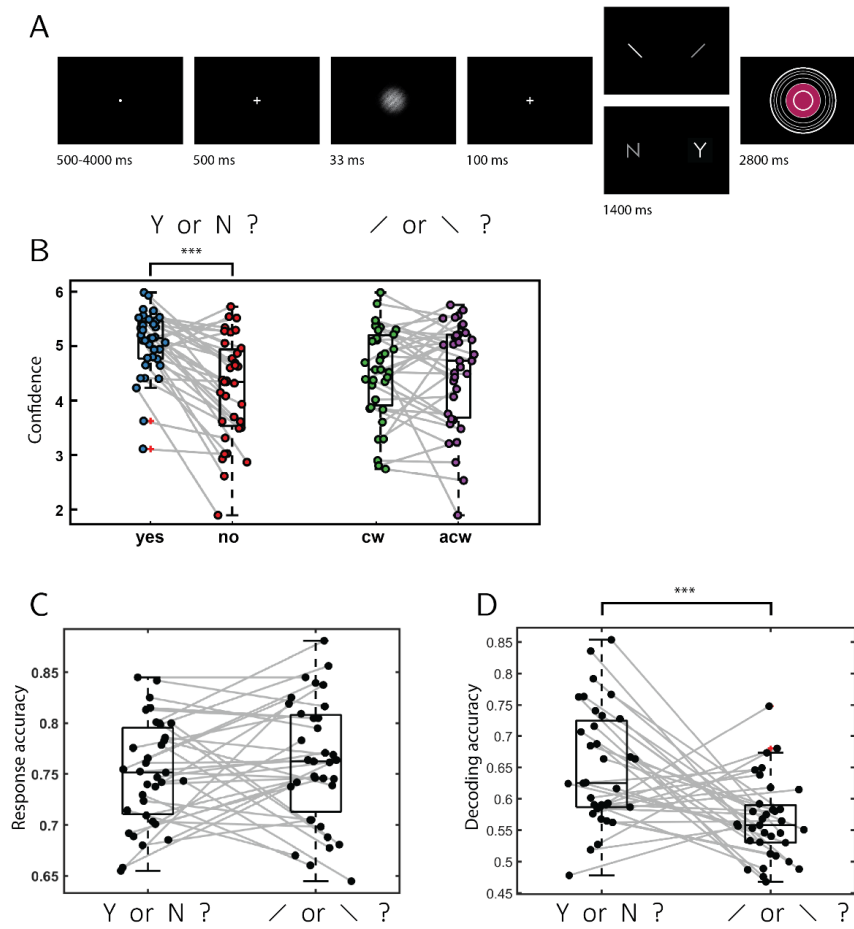
119 Design and procedure

120 Trials started with a fixation cross (500 milliseconds), followed by a presentation of a
121 stimulus for 33 milliseconds. In discrimination trials, the stimulus was a circle of diameter 3°
122 containing randomly generated white noise, merged with a sinusoidal grating (2 cycles per
123 degree; oriented 45° or -45°). In half of the detection trials, stimuli did not contain a sinusoidal
124 grating and consisted of random noise only. After stimulus offset, participants used their right-
125 hand index and middle fingers to make a perceptual decision about the orientation of the grating
126 (discrimination blocks), or about the presence or absence of a grating (detection blocks; see
127 Fig. 1, top panel). Response mapping was counterbalanced between blocks which means that
128 significant decoding of decisions cannot reflect motor representations.

129 Immediately after making a decision, participants rated their confidence on a 6-point
130 scale by using two keys to increase or decrease their reported confidence level with their left-
131 hand thumb. Confidence levels were indicated by the size and color of a circle presented at the
132 center of the screen. The initial size and color of the circle was determined randomly at the
133 beginning of the confidence rating phase. The mapping between color and size to confidence
134 was counterbalanced between participants: for half of the participants high confidence was
135 mapped to small, red circles, and for the other half high confidence was mapped to large, blue
136 circles. The perceptual decision and the confidence rating phases were restricted to 1500 and
137 2500 milliseconds, respectively. No feedback was delivered to subjects about their
138 performance. Trials were separated by a temporally jittered rest period of 500-4000
139 milliseconds.

140 Prior to the scanning day, participants underwent a behavioral session in which task
141 difficulty was adjusted independently for the detection and discrimination tasks, targeting around
142 70% accuracy. We achieved this by adaptively adjusting the stimulus signal-to-noise ratio (SNR)
143 every 10 trials (increasing the SNR if accuracy fell below 60%, and decreasing it if accuracy
144 exceeded 80%). Task difficulty was further calibrated within the scanner environment at the
145 beginning of the scanning session, during the acquisition of anatomical (MP-RAGE and
146 fieldmap) images, using a similar procedure. Upon completion of the calibration phase,
147 participants performed 5 experimental runs comprising one discrimination and one detection
148 block, each of 40 trials, presented in random order. A bonus was awarded for accurate
149 responses and confidence ratings according to the following formula: $\sum_{i=1}^N accuracy_i \times$
150 $confidence_i$, where *accuracy* equals 1 for correct responses and -1 for incorrect responses, and
151 *confidence* is the reported confidence level on a scale of 1-6.

152
153



154
 155 **Figure 1: Experimental design and behavioural results.** A: In discrimination trials, participants
 156 made discrimination judgments about clockwise and anticlockwise tilted noisy gratings, and then rated
 157 their subjective confidence by controlling the size of a colored circle. In detection judgments, decisions
 158 were made about the presence (Y) or absence (N) of a grating in noise. B: mean confidence as a function
 159 of response for the 35 participants. Confidence in detection 'yes' responses was significantly higher than
 160 in 'no' responses. No significant difference was observed between confidence in discrimination
 161 responses (cw: clockwise, acw: anticlockwise). C: Response accuracy was not different between the two
 162 tasks. D: Decoding accuracy for a classifier trained to classify response (yes or no in detection, clockwise
 163 or anticlockwise in discrimination) based on confidence ratings alone. Decoding accuracy was
 164 significantly higher for detection than for discrimination. ***: $p < 0.001$.

165 Scanning parameters

166 Scanning took place at the Wellcome Centre for Human Neuroimaging, London, using a
167 3 Tesla Siemens Prisma MRI scanner with a 64-channel head coil. We acquired structural
168 images using an MPRAGE sequence (1×1×1 mm voxels, 176 slices, in plane FoV = 256×256
169 mm²), followed by a double-echo FLASH (gradient echo) sequence with TE1 = 10 ms and TE2
170 = 12.46 ms (64 slices, slice thickness = 2 mm, gap = 1 mm, in plane FoV = 192 × 192 mm²,
171 resolution = 3 × 3 mm²) that was later used for field inhomogeneity correction. Functional scans
172 were acquired using a 2D EPI sequence, optimized for regions near the orbito-frontal cortex
173 (3×3×3 mm voxels, TR = 3.36 s, TE = 30 ms, 48 slices tilted by -30 degrees with respect to the
174 T > C axis, matrix size = 64×72, Z-shim = -1.4).

175 Preprocessing

176 Data preprocessing followed the procedure described in Morales et al. (2018): Imaging
177 analysis was performed using SPM12 (Statistical Parametric Mapping;
178 www.fil.ion.ucl.ac.uk/spm). The first five volumes of each run were discarded to allow for T1
179 stabilization. Functional images were realigned and unwarped using local field maps
180 (Andersson et al., 2001) and then slice-time corrected (Sladky et al., 2011). Each participant's
181 structural image was segmented into gray matter, white matter, CSF, bone, soft tissue, and
182 air/background images using a nonlinear deformation field to map it onto template tissue
183 probability maps (Ashburner and Friston, 2005). This mapping was applied to both structural
184 and functional images to create normalized images in Montreal Neurological Institute (MNI)
185 space. Normalized images were spatially smoothed using a Gaussian kernel (6 mm FWHM).
186 We set a within-run 4 mm affine motion cutoff criterion.

187
188 To extract trial-wise activation estimates, we used SPM to fit a design matrix to the
189 preprocessed images. The design matrix included a regressor for each experimental trial, as
190 well as nuisance regressors for instruction screens and physiological parameters. Trials were
191 modeled as 33 millisecond boxcar functions, locked to the presentation of the stimulus, and
192 convolved with a canonical hemodynamic response function. Trial-wise beta estimates were
193 then used in multivariate analysis.

195 Multivariate analysis

196 Only correct trials were used for decoding (75% and 76% of trials from included blocks in
197 the detection and discrimination tasks, respectively). We chose to limit our decoding analysis to
198 correct trials in order not to conflate effects of subjective confidence with those of objective
199 accuracy, or stimulus type. However, we found that qualitatively similar results are obtained
200 when analyzing all trials (unthresholded whole brain maps are available in this study's
201 NeuroVault collection: neurovault.org/collections/9912/).

202 Stimulus presence (present vs. absent) was decoded during detection blocks, and
203 stimulus identity (clockwise vs. anticlockwise orientation) during discrimination blocks. Both
204 decoding analyses used an LDA (Linear Discriminant Analysis) classifier with leave-one-run-out
205 cross-validation and a searchlight radius of 4 voxels (~257 voxels per searchlight). Significance

206 testing was done using permutation testing to generate the empirical null-distribution. We
207 followed the approach suggested by Stelzer, Chen, & Turner (2013) for searchlight MVPA
208 measurements which uses a combination of permutation testing and bootstrapping to generate
209 chance distributions for group studies. Per participant, 25 permutation maps were generated by
210 permuting the class labels within each run. Group-level permutation distributions were
211 subsequently generated by bootstrapping over these 25 maps, i.e. randomly selecting one out
212 of 25 maps per participant. 10000 bootstrapping samples were used to generate the group null-
213 distribution per voxel and per comparison. *P*-values were calculated per searchlight or ROI as
214 the right-tailed area of the histogram of permuted accuracies from the mean over participants.
215 We corrected for multiple comparisons in the searchlight analyses using whole-brain FDR-
216 correction. A cluster-extent threshold was applied, ensuring that voxels were only identified as
217 significant if they belonged to a cluster of at least 50 significant voxels (Dijkstra, Bosch, & van
218 Gerven, 2017).

219

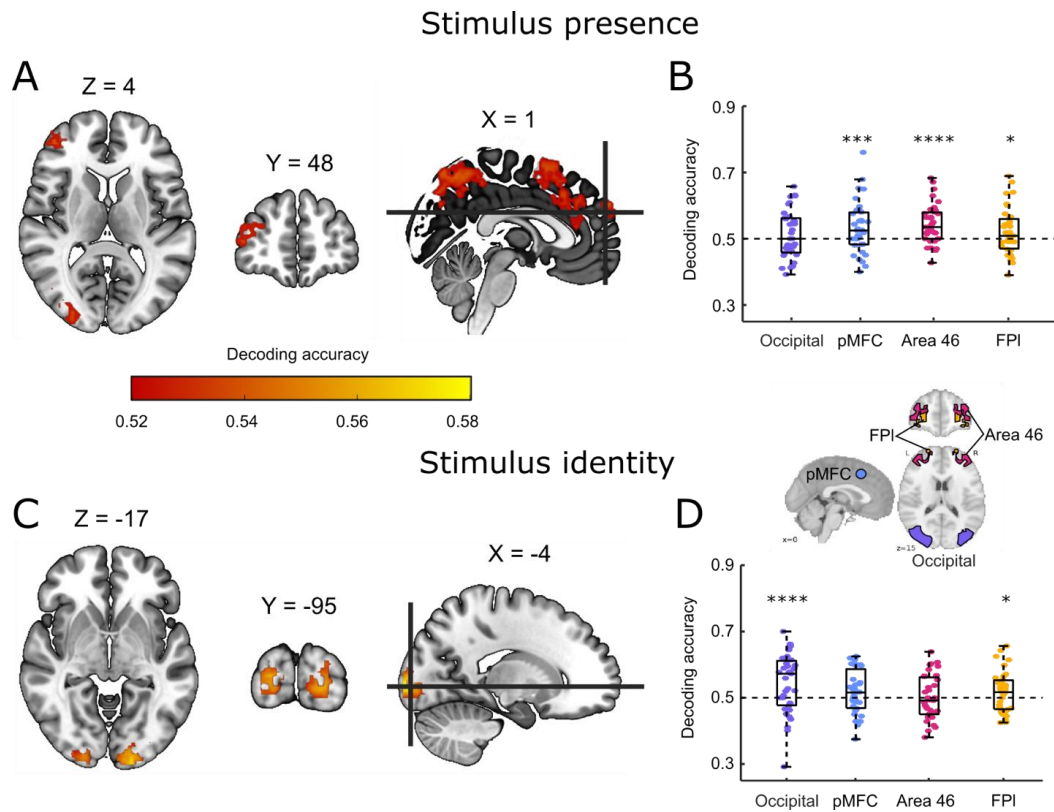
220 Results

221 Decoding of stimulus presence and orientation

222 We first searched for multivariate activation patterns that encoded information about
223 stimulus orientation (in discrimination) and stimulus presence/visibility (in detection). In a whole-
224 brain searchlight analysis, stimulus orientation could be reliably decoded only from the visual
225 cortex (Fig. 2C). In contrast, information about stimulus presence was identified in parietal and
226 prefrontal brain regions, including the dorsolateral prefrontal cortex, the middle frontal gyrus,
227 and the precuneus (see Fig. 2A, for unthresholded classification maps, see
228 neurovault.org/collections/9912).

229 Based on these maps, we decided to focus our subsequent analyses on four regions of
230 interest (ROIs): an occipital ROI, defined using the AICHA atlas as 'occipital mid' regions (Joliot
231 et al., 2015) and three prefrontal ROIs which were also used in Mazor et al (2020): posterior
232 medial frontal cortex (pMFC; an 8 mm sphere around MNI coordinates [0, 17, 46]), Brodmann
233 area 46, and lateral frontopolar cortex (BA46 and FPI; both defined based on a connectivity-
234 based parcellation; Neubert et al., 2014). Bilateral ROIs were defined as the union of the right
235 and left hemispheres.

236 Within these four ROIs, stimulus orientation could be decoded significantly from occipital
237 ($M = 0.54$, $SD = 0.09$, $p < 0.0001$) and FPI ROIs ($M = 0.51$, $SD = 0.06$, $p = 0.04$). In contrast,
238 stimulus presence could be decoded from pMFC ($M = 0.53$, $SD = 0.08$, $p = 0.0009$), area 46 (M
239 $= 0.54$, $SD = 0.06$, $p < 0.0001$) and FPI ROIs ($M = 0.52$, $SD = 0.07$, $p = 0.015$), but not from the
240 occipital ROI ($M = 0.51$, $SD = 0.07$, $p = 0.11$). Classification accuracy showed a significant ROI
241 x task interaction ($F(3,32) = 5.31$, $p = 0.004$; see Fig. 2, right panel), suggesting that stimulus
242 presence (Fig. 2B) and stimulus identity (Fig. 2D) are encoded differentially across ROIs. Post-
243 hoc contrasts revealed a significantly higher classification accuracy for detection compared to
244 discrimination in area 46 ($t(34) = 3.06$, $p < 0.005$), with no significant difference between detection
245 and discrimination decoding in the FPI, pMFC, or occipital ROIs.



246
 247 **Figure 2: Decoding of stimulus presence and stimulus identity.** A: whole brain searchlight
 248 decoding of stimulus presence versus absence in the detection task, correct responses only. B:
 249 decoding of stimulus presence versus absence in the occipital, pMFC, BA 46 and FPI ROIs. C:
 250 Whole brain searchlight decoding of stimulus identity in the discrimination task, correct
 251 responses only. D: decoding of stimulus identity in the four ROIs. Whole-brain maps are
 252 corrected for multiple comparisons at the voxel level with a cluster-size cutoff of 50 voxels. *:
 253 $p < 0.5$, **: $p < 0.01$, ***: $p < 0.001$, ****: $p < 0.0001$.
 254

255 **Behavioural analysis and confidence-based decoding**

256 As previously reported in Mazor et al. (2020), task performance was similar for detection
 257 (75% accuracy, $d' = 1.48$) and discrimination (76% accuracy, $d' = 1.50$). Repeated measures t-
 258 tests failed to detect a difference between tasks both in mean accuracy ($t(34) = -0.90$, $p = 0.37$,
 259 $d = 0.15$, $BF_{01} = 5.15$), and d' ($t(34) = -0.30$, $p = 0.76$, $d = 0.05$, $BF_{01} = 7.29$), indicating that
 260 performance was well matched. Within detection, participants were significantly more confident
 261 in 'yes' responses (mean confidence = 5.03 on a 1-6 scale) compared to 'no' responses (mean
 262 confidence = 4.21; $t(34) = 5.83$, $p < 0.001$, $d = 1.00$). In contrast, confidence in discrimination

263 'clockwise' responses (mean confidence =4.28) was not significantly different from confidence in
264 discrimination 'anticlockwise' responses (mean confidence =4.25; $t(34)=0.31$, $p=0.76$, $d=0.05$).

265 This absence of a significant difference between confidence in discrimination responses
266 may indicate that a typical participant rated confidence similarly for discrimination 'clockwise'
267 and 'anticlockwise' responses. Alternatively, it may be that some participants showed a bias
268 towards higher confidence in 'clockwise' responses and others showed a bias towards higher
269 confidence in 'anticlockwise' responses. Deciding between these two alternatives is important
270 for interpreting our multi-voxel pattern analysis of discrimination responses: if single participants
271 were consistently more confident in one of the two discrimination responses, above chance
272 classification accuracy for discrimination may still be driven by differences in decision
273 confidence, even if such differences average out at the group level (Gilron et al., 2017).

274 To decide between these two alternatives, we trained and tested an LDA classifier to
275 predict participants' decisions from their confidence ratings only. We used the same leave-one-
276 run-out cross-validation procedure as in our MVPA analysis. This was done separately for the
277 two tasks and for each participant. Confidence ratings successfully predicted detection
278 responses, in line with a difference in mean confidence between detection 'yes' and 'no'
279 responses ($M=0.65$, $t(34)=9.70$, $p<0.001$, $d=1.64$). Importantly, an LDA classifier also separated
280 discrimination responses based on decision confidence ($M=0.57$, $t=6.25$, $p<0.001$, $d=1.06$), but
281 to a lesser extent than in detection ($t(34)=3.88$, $p<0.001$, $d=0.67$ for a paired t-test testing the
282 difference in classification accuracy between detection and discrimination). These analyses
283 further emphasise the need to control for confidence when interpreting above-chance
284 classification of detection and discrimination responses in higher-order brain regions in our data,
285 as these may reflect person-specific differences in mean confidence between the two
286 responses. Our next set of analyses was designed to control for this potential confound.

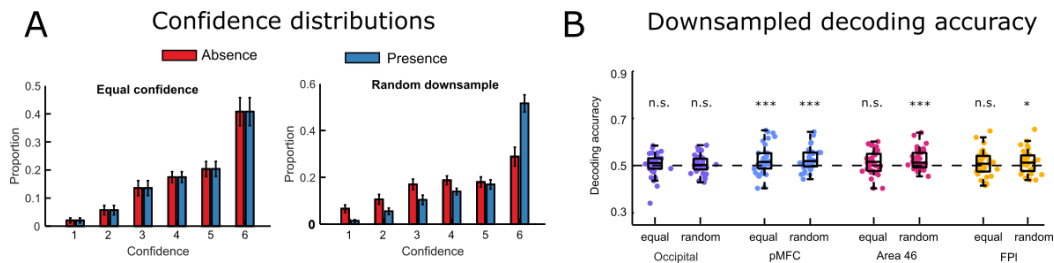
287 Confidence-matching via downsampling

288 Prefrontal decoding of stimulus presence is consistent with the proposal that subjective
289 visibility is represented in a frontoparietal network. However, it is also plausible that prefrontal
290 decoding of detection reflects representations of confidence, instead of visibility. This alternative
291 interpretation is in line with the finding that activity in prefrontal cortex is sensitive to variation in
292 confidence (Vacarro & Fleming, 2018), and with our observation that confidence varied between
293 detection decisions more than between discrimination decisions.

294 In our next analysis we therefore set out to determine whether our prefrontal ROIs would
295 continue to represent stimulus presence *after controlling for decision confidence*. Having trial-
296 wise confidence ratings allowed us to perfectly match not only mean confidence, but the entire
297 distribution of confidence ratings for target present and target absent responses, and quantify
298 the effect this had on classification accuracy. This was achieved by downsampling: for each
299 participant and for each task, we selectively deleted trials until the two response categories had
300 an equal number of trials for each confidence level (see Fig. 3A, left histogram). For example, if
301 a participant had 15 trials in which they gave a confidence rating of 6, out of which only 3 were
302 target absent trials, we randomly deleted 9 target-present trials in which the participant gave a
303 confidence rating of 6, resulting in an equal number of confidence-6 trials for each response
304 category. By then applying our presence/absence decoding analysis to these downsampled

305 data, we were able to obtain a “downsampled” decoding accuracy, which reflected the ability of
 306 a classifier to determine stimulus presence vs. absence from activation patterns, after removing
 307 differences in confidence.

308 To make sure any change in decoding accuracy was not simply due to a reduction in
 309 trial number, we also repeated this procedure with random instead of confidence-based
 310 downsampling, resulting in a second ‘random downsampled’ decoding accuracy value for each
 311 ROI. Importantly, this procedure of random downsampling ensures that the trial numbers in the
 312 two classes are the same as in the equalized confidence analysis, while keeping any confidence
 313 differences intact (see Fig. 3A, right histogram). Because there are multiple ways in which a
 314 dataset could be downsampled, for both types of analyses we repeated the procedure 25 times
 315 to take into account the variance created by selective sampling and then averaged decoding
 316 accuracy over these different downsampled sets. Finally, for statistical testing we created null
 317 distributions by following the same downsampling procedure on label-shuffled datasets.
 318



319
 320
 321 **Figure 3: Stimulus presence downsampling analysis.** A: for each participant, trials were deleted until
 322 confidence distributions were matched for target present and target absent responses. As a control
 323 analysis, we repeated this procedure with random downsampling, deleting the same number of trials
 324 irrespective of confidence ratings. B: presence/absence classification accuracy in the four ROIs for the
 325 equal confidence and random downsampling datasets. *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$, ****:
 326 $p < 0.0001$
 327

328 When equalizing confidence, classification accuracy for decoding stimulus presence
 329 remained significant in pMFC ($M = 0.52$, $SD = 0.06$, $p = 0.002$). However, decoding was no
 330 longer significant after equalizing the confidence distributions in FPI ($M = 0.51$, $SD = 0.05$, $p =$
 331 0.11), and only marginally significant in area 46 ($M = 0.51$, $SD = 0.05$, $p = 0.07$). In both regions,
 332 decoding was still significant after random downsampling (FPI: $M = 0.52$, $SD = 0.05$, $p = 0.02$;
 333 area 46: $M = 0.53$, $SD = 0.04$, $p = 0.0017$). A decrease in classification accuracy after equalizing
 334 confidence relative to random downsampling was marginally significant in area 46 ($t(34) = -1.733$,
 335 $p = 0.09$, $d = 0.29$), but not in the FPI ROI ($t(34) = -1.615$, $p = 0.11$, $d = 0.27$). In the pMFC ROI,
 336 classification accuracies for the confidence-matched and random downsamples were highly
 337 similar (0.524 and 0.525 , $t(34) = -0.20$, $p = 0.84$). Taken together, these results show that in
 338 pMFC, but not area 46 and FPI, stimulus presence/visibility can be reliably decoded
 339 independent of differences in decision confidence.

340 When decoding stimulus identity in the discrimination task, confidence-matching had no
 341 effect on classification accuracy relative to random downsampling (downsampled classification

342 accuracy in the occipital ROI: $M = 0.55$, $SD = 0.07$; FPI: $M = 0.52$, $SD = 0.06$; pMFC: $M = 0.51$,
343 $SD = 0.06$; area 46: $M = 0.51$, $SD = 0.05$, all pairwise comparisons with non-downsampled
344 accuracy $p > 0.28$). This is consistent with there already being little difference in the
345 (behavioural) confidence distributions between the two response types in discrimination blocks.
346 Importantly, in pMFC, we observed no significant classification of stimulus identity, regardless of
347 whether the analysis used confidence-matched data or not (downsampled classification
348 accuracy: $M = 0.52$, $SD = 0.06$, $p = 0.2$). In other words, in this prefrontal ROI, we were able to
349 decode visibility (independently of confidence) but not identity.

350 Discussion

351 What role the prefrontal cortex plays in visual awareness is much debated (e.g. Aru,
352 Bachmann, Singer & Melloni, 2012; Boly et al., 2017). Here, we investigated whether prefrontal
353 areas encode the visibility of a faint stimulus independently of stimulus identity and decision
354 confidence. We first showed that a subset of prefrontal ROIs (pMFC and area 46) tracked
355 stimulus presence during a detection task but not stimulus identity during a discrimination task,
356 consistent with prefrontal involvement in encoding of stimulus visibility. Furthermore,
357 classification accuracy was significantly higher for stimulus presence than for stimulus identity in
358 area 46. However, because seeing a stimulus is associated with higher confidence than not
359 seeing a stimulus, this asymmetry could also reflect confidence coding in frontal areas. To
360 investigate this possibility, we tested whether decoding of stimulus presence remained
361 significant after controlling for differences in confidence. We found that such decoding was
362 indeed still possible in pMFC, but not in area 46. Taken together, these results suggest that
363 pMFC, in contrast to area 46, encodes stimulus visibility over and above decision confidence.
364 Furthermore, pMFC, unlike occipital regions, did not significantly encode stimulus identity, either
365 when allowing confidence to freely vary, or when controlling for confidence in a downsampling
366 analysis.

367 However, it is important to note that the interpretation of a “pure visibility” signal in pMFC
368 is nuanced by a lack of significant difference between classification accuracies for stimulus
369 presence and identity in this region. In other words, while we can decode stimulus visibility but
370 not identity in pMFC, we cannot conclude that the decoding of these two quantities are
371 themselves significantly different. Therefore, one viable alternative interpretation of our results
372 might be that pMFC encodes a low-dimensional projection of rich perceptual input onto a
373 decision axis: one that separates clockwise from anticlockwise gratings in discrimination blocks,
374 and noise patches with and without a grating in detection blocks. Nevertheless, regardless of
375 the nuance required when interpreting results in individual prefrontal ROIs, our results make
376 clear that what may appear to be neural signatures of visibility in prefrontal cortex (e.g. in whole-
377 brain searchlight decoding, such as in Figure 2) may on closer inspection be more closely
378 related to differences in decision confidence.

379 Conceptually, visibility and decision confidence appear similar. They can both be defined
380 in terms of precision: the precision of a visual percept in the first case, and the precision with
381 which a decision is made in the second (Denison et al., 2017). Empirically, neural correlates of
382 visibility and decision confidence overlap, specifically in the dorsolateral prefrontal cortex
383 (dlPFC) but also in medial prefrontal, parietal, and insular cortices (Vacarro & Fleming, 2018).

384 Notwithstanding this conceptual and empirical overlap, visibility and confidence are not one and
385 the same thing. Critically, within a Bayesian framework, decision confidence is defined as the
386 probability correct of a particular response, and should therefore be sensitive not only to the
387 precision of sensory representations, but also response requirements (Pouget et al., 2016; Bang
388 & Fleming, 2018). Accordingly, visibility judgments scale with stimulus contrast even in trials in
389 which participants make erroneous decisions, but confidence judgments show a different profile,
390 and are sensitive to stimulus contrast only for correct responses (Rausch and Zehelsteiner,
391 2016).

392 Despite a theoretical distinction between confidence and visibility, neuroimaging findings
393 of visual awareness have often not been able to separate their respective contributions to
394 differential brain activation. For example, it has not been possible to determine whether the
395 dorsolateral prefrontal cortex is more active on aware versus unaware trials because it is
396 sensitive to subjective visibility, or because participants are generally more confident in their
397 decisions when they are aware of a stimulus. In an exploratory analysis of existing imaging
398 data, we found that an apparent encoding of stimulus visibility in area 46 and lateral frontopolar
399 cortex disappeared when controlling for subjective confidence. In contrast, pMFC encoding of
400 visibility remained significant.

401 As reported in Mazor et al. (2020), univariate analysis of this data indicated a similar
402 parametric modulation of confidence for detection and discrimination responses in pMFC.
403 Specifically, a similar modulation of confidence in decisions about target presence and absence
404 indicate that univariate signal in this region also scales with decision confidence. Univariate
405 analysis did not reveal a pMFC modulation of visibility, which would manifest as an interaction of
406 confidence and class in detection (because visibility is negatively correlated with confidence in
407 'no' responses, but positively correlated with confidence in 'yes' responses). However, a pre-
408 registered cross-classification analysis revealed shared multivariate representations for
409 discrimination confidence and detection responses indicating whether a stimulus is seen or not
410 in pMFC and area 46 (Mazor et al., 2020; Appendix 8). We previously interpreted these findings
411 as indicating that multivariate spatial activation patterns in area 46 and pMFC hold information
412 about stimulus visibility, because like detection responses, confidence during discrimination
413 might also track stimulus visibility (it is easier to determine what something is when you see it
414 more clearly). Our current results corroborate this finding with respect to pMFC, and further
415 show that above chance cross-classification in this region is not merely driven by differences in
416 subjective confidence between 'yes' and 'no' responses during detection. Taken together, these
417 results suggest that pMFC signal carries information not only about subjective confidence, but
418 also about perceptual content, be it stimulus visibility, stimulus identity, or both.

419 Activation in pMFC is commonly found to correlate negatively with subjective confidence,
420 or positively with uncertainty (Fleming, Huijgen & Dolan, 2012; Molenberghs et al., 2016;
421 Vacarro & Fleming, 2018; Mazor, Friston & Fleming, 2020). In a recent study we found that
422 univariate pMFC activation tracked the effect of decision difficulty, although it was insensitive to
423 the precision of perceptual information in a motion perception task, which was instead tracked in
424 posterior parietal regions (Bang & Fleming, 2018). Other work has shown that the pMFC is
425 important for signaling when decisions or beliefs should be updated on the basis of new
426 information (Fleming et al., 2018; O'Reilly et al., 2013). Novel paradigms may be necessary to
427 further disentangle pMFC contributions to encoding stimulus visibility, and to relate this putative

428 computational role to the encoding of other types of (perceptual and non-perceptual)
429 uncertainty.

430 Our initial analysis specifically highlighted area 46 in the decoding of stimulus presence
431 *without* controlling for confidence differences. This pattern of results is consistent with area 46
432 contributing to detection confidence, whereas more posterior prefrontal cortex (pmFC) may
433 support visual detection responses, irrespective of differences in confidence. This result is in
434 keeping with previous observations TMS to area 46 leads to lower overall perceptual confidence
435 (a change in metacognitive bias), without affecting metacognitive sensitivity (Shekhar & Rahnev,
436 2018). In contrast, TMS applied to frontopolar cortex in Shekhar and Rahnev's study led to
437 increases in metacognitive sensitivity, without affecting confidence bias. We note that the
438 contribution of prefrontal cortical subregions to visual metacognitive sensitivity (the coupling
439 between confidence and accuracy) is difficult to assess using within-subject neuroimaging
440 methods applied here as it requires modeling confidence noise across many trials. It remains to
441 be determined whether the visual confidence signal in area 46 we observe here is specific to
442 perceptual judgments (Lau & Passingham, 2006), or generalises to different task domains
443 (Morales et al., 2018; Fleck et al., 2006).

444 Our results with respect to the lateral frontopolar cortex (FPI) are more difficult to
445 interpret. We found that this area did not represent stimulus presence over and above
446 confidence, but that it did represent stimulus identity, even after controlling for confidence
447 differences between the different stimulus classes. Several factors may have contributed to
448 these results. First, our observation that the FPI does not encode visibility irrespective of
449 confidence does not mean that this region cannot play a role in visual awareness. In target
450 absence trials, participants can sometimes be fully aware of the absence of a target – a case
451 where visibility is low, but awareness (of absence) is high (Mazor & Fleming, 2020). Therefore, if
452 FPI tracks content-invariant aspects of visual awareness, its activation may not differentiate
453 between target presence and target absence. However, a representation of stimulus identity in
454 FPI suggests that this area might also encode stimulus content. We are not aware of previous
455 reports of decoding of visual content from the frontopolar cortex. Moreover, a recent meta-
456 analysis reported no known effects of intracranial electrical stimulation of the frontopolar cortex
457 on spontaneous reports of visual experience (Racchah, Block & Fox, 2021). Given the relatively
458 modest effect sizes in FPI decoding of stimulus identity ($M=0.51$) in comparison to the more
459 robust encoding of stimulus identity in occipital cortex ($M=0.55$), we are cautious in over-
460 interpreting this surprising result. Future studies are necessary to explore to what extent FPI
461 truly represents stimulus identity, and/or contributes to visual awareness.

462 Finally, when considering the implications of these findings for the study of visual
463 awareness and its neural correlates, it is important to note the difference between subjective
464 reports of stimulus awareness, and decisions about the presence or absence of a target
465 stimulus in a perceptual detection task. While the first is a subjective decision about the
466 contents of one's perception, the second is a report of one's beliefs about the state of the
467 external world. Consequently, these two types of decisions draw on different sets of prior beliefs
468 and expectations. For example, in detection, but not in subjective visibility reports, participants
469 may adjust their decision criterion when noticing that they haven't detected a stimulus in a long
470 time. Furthermore, participants may base their detection responses not on the visibility of a
471 stimulus, but on other visual and non-visual cues (adopting different *criterion content*;

472 Kahneman, 1968). Our findings are based on the analysis of detection decisions, and their
473 generalizability to reports of subjective awareness is an open empirical question.

474 To conclude, an exploratory data analysis revealed that stimulus presence could be
475 decoded from prefrontal regions but that only the pMFC encoded stimulus presence after
476 controlling for decision confidence. Future hypothesis-driven investigation is needed to replicate
477 these exploratory results. We demonstrate the importance of controlling for confidence when
478 investigating reports of awareness versus unawareness and propose a novel analysis approach
479 to do so.

480

481

482 Funding

483 N.D. is supported by a Rubicon grant from the Netherlands Organization for Scientific
484 Research (NWO) [019.192SG.003]. SMF is funded by a Wellcome/Royal Society Sir Henry Dale
485 Fellowship (206648/Z/17/Z) and a Philip Leverhulme Prize from the Leverhulme Trust. The
486 Wellcome Centre for Human Neuroimaging is supported by core funding from the Wellcome
487 Trust (206648/Z/17/Z). This research was funded in whole or in part by the Wellcome Trust
488 (206648/Z/17/Z). For the purpose of Open Access, the author has applied a CC-BY public
489 copyright license to any author accepted manuscript version arising from this submission.

490

491 The authors declare that there are no competing interests.

- 492 Aru, J., Bachmann, T., Singer, W., & Melloni, L. (2012). Distilling the neural correlates of consciousness.
 493 *Neuroscience & Biobehavioral Reviews*, 36(2), 737-746.
 494
- 495 Bang, D., & Fleming, S. M. (2018). Distinct encoding of decision confidence in human medial prefrontal
 496 cortex. *Proceedings of the National Academy of Sciences*, 115(23), 6082-6087.
 497
- 498 Bayne, T., & Hohwy, J. (2013). Consciousness: theoretical approaches. *Neuroimaging of consciousness*,
 499 23-35.
 500
- 501 Boly, M., Massimini, M., Tsuchiya, N., Postle, B. R., Koch, C., & Tononi, G. (2017). Are the neural
 502 correlates of consciousness in the front or in the back of the cerebral cortex? Clinical and neuroimaging
 503 evidence. *Journal of Neuroscience*, 37(40), 9603-9613.
 504
- 505 Brown, R., Lau, H., & LeDoux, J. E. (2019). Understanding the higher-order approach to consciousness.
 506 *Trends in cognitive sciences*, 23(9), 754-768.
 507
- 508 Dehaene, S., Changeux, J. P., & Naccache, L. (2011). The global neuronal workspace model of
 509 conscious access: from neuronal architectures to clinical applications. *Characterizing consciousness:
 510 From cognition to the clinic?*, 55-84.
 511
- 512 Dehaene, S., Naccache, L., Cohen, L., Le Bihan, D., Mangin, J. F., Poline, J. B., & Rivière, D. (2001).
 513 Cerebral mechanisms of word masking and unconscious repetition priming. *Nature neuroscience*, 4(7),
 514 752-758.
 515
- 516 Dehaene, S., Sergent, C., & Changeux, J. P. (2003). A neuronal network model linking subjective reports
 517 and objective physiological data during conscious perception. *Proceedings of the National Academy of
 518 Sciences*, 100(14), 8520-8525.
 519
- 520 De Martino, B., Fleming, S. M., Garrett, N., & Dolan, R. J. (2013). Confidence in value-based choice.
 521 *Nature neuroscience*, 16(1), 105-110.
 522
- 523 Denison, R. N. (2017). Precision, Not Confidence, Describes the Uncertainty of Perceptual Experience:
 524 Comment on John Morrison's "Perceptual Confidence". *Analytic Philosophy*, 58(1), 58-70.
 525
- 526 Fisch, L., Privman, E., Ramot, M., Harel, M., Nir, Y., Kipervasser, S., ... & Malach, R. (2009). Neural
 527 "ignition": enhanced activation linked to perceptual awareness in human ventral stream visual cortex.
 528 *Neuron*, 64(4), 562-574.
 529
- 530 Fleck, M. S., Daselaar, S. M., Dobbins, I. G., & Cabeza, R. (2006). Role of prefrontal and anterior
 531 cingulate regions in decision-making processes shared by memory and nonmemory tasks. *Cerebral
 532 Cortex*, 16(11), 1623-1630.
 533
- 534 Fleming, S. M., Huijgen, J., & Dolan, R. J. (2012). Prefrontal contributions to metacognition in perceptual
 535 decision making. *Journal of Neuroscience*, 32(18), 6117-6125.
 536
- 537 Fleming, S. M., Van Der Putten, E. J., & Daw, N. D. (2018). Neural mediators of changes of mind about
 538 perceptual decisions. *Nature neuroscience*, 21(4), 617-624.
 539

- 540 Gilron, R., Rosenblatt, J., Koyejo, O., Poldrack, R. A., & Mukamel, R. (2017). What's in a pattern?
 541 Examining the type of signal multivariate analysis uncovers at the group level. *NeuroImage*, *146*, 113-
 542 120.
 543
- 544 Joliot, M., Jobard, G., Naveau, M., Delcroix, N., Petit, L., Zago, L., ... & Tzourio-Mazoyer, N. (2015).
 545 AICHA: An atlas of intrinsic connectivity of homotopic areas. *Journal of neuroscience methods*, *254*, 46-
 546 59.
 547
- 548 Kahneman, D. (1968). Method, findings, and theory in studies of visual masking. *Psychological*
 549 *Bulletin*, *70*(6p1), 404.
 550
- 551 Koivisto, M., & Revonsuo, A. (2010). Event-related brain potential correlates of visual awareness.
 552 *Neuroscience & Biobehavioral Reviews*, *34*(6), 922-934.
 553
- 554 Lau, H. C., & Passingham, R. E. (2006). Relative blindsight in normal observers and the neural correlate
 555 of visual consciousness. *Proceedings of the National Academy of Sciences*, *103*(49), 18763-18768.
 556
- 557 Lau, H., & Rosenthal, D. (2011). Empirical support for higher-order theories of conscious awareness.
 558 *Trends in cognitive sciences*, *15*(8), 365-373.
 559
- 560 Mazor, M., Friston, K. J., & Fleming, S. M. (2020). Distinct neural contributions to metacognition for
 561 detecting, but not discriminating visual stimuli. *ELife*, *9*, e53900.
 562
- 563 Mazor, M., & Fleming, S. M. (2020). Distinguishing absence of awareness from awareness of absence.
 564 *Philosophy and the Mind Sciences*, *1*(II).
 565
- 566 Mazor, M., Mazor, N., & Mukamel, R. (2019). A novel tool for time-locking study plans to
 567 results. *European Journal of Neuroscience*, *49*(9), 1149-1156.
 568
- 569 Mazor, M., Moran, R., & Fleming, S. M. (2021). Metacognitive asymmetries in visual
 570 perception. *Neuroscience of Consciousness*, *2021*(1), niab005.
 571
- 572 Michel, M., & Morales, J. (2020). Minority reports: Consciousness and the prefrontal cortex. *Mind &*
 573 *Language*, *35*(4), 493-513.
 574
- 575 Molenberghs, P., Trautwein, F. M., Böckler, A., Singer, T., & Kanske, P. (2016). Neural correlates of
 576 metacognitive ability and of feeling confident: a large-scale fMRI study. *Social cognitive and affective*
 577 *neuroscience*, *11*(12), 1942-1951.
 578
- 579 Morales, J., Lau, H., & Fleming, S. M. (2018). Domain-general and domain-specific patterns of activity
 580 supporting metacognition in human prefrontal cortex. *Journal of Neuroscience*, *38*(14), 3534-3546.
 581
- 582 Odegaard, B., Knight, R. T., & Lau, H. (2017). Should a few null findings falsify prefrontal theories of
 583 conscious perception?. *Journal of Neuroscience*, *37*(40), 9593-9602.
 584
- 585 O'Reilly, J. X., Schüffelgen, U., Cuell, S. F., Behrens, T. E., Mars, R. B., & Rushworth, M. F. (2013).
 586 Dissociable effects of surprise and model update in parietal and anterior cingulate cortex. *Proceedings of*
 587 *the National Academy of Sciences*, *110*(38), E3660-E3669.
 588

- 589 Pouget, A., Drugowitsch, J., & Kepecs, A. (2016). Confidence and certainty: distinct probabilistic
590 quantities for different goals. *Nature neuroscience*, 19(3), 366.
- 591
- 592 Raccach, O., Block, N., & Fox, K. C. (2021). Does the prefrontal cortex play a necessary role in
593 consciousness? Insights from intracranial electrical stimulation of the human brain. *Journal of*
594 *Neuroscience*, 1(41).
- 595
- 596 Rausch, M., & Zehetleitner, M. (2016). Visibility is not equivalent to confidence in a low contrast
597 orientation discrimination task. *Frontiers in psychology*, 7, 591.
- 598
- 599 Sahraie, A., Weiskrantz, L., Barbur, J. L., Simmons, A., Williams, S. C. R., & Brammer, M. J. (1997).
600 Pattern of neuronal activity associated with conscious and unconscious processing of visual signals.
601 *Proceedings of the National Academy of Sciences*, 94(17), 9406-9411.
- 602
- 603 Shekhar, M., & Rahnev, D. (2018). Distinguishing the roles of dorsolateral and anterior PFC in visual
604 metacognition. *Journal of Neuroscience*, 38(22), 5078-5087.
- 605
- 606 Vaccaro, A. G., & Fleming, S. M. (2018). Thinking about thinking: A coordinate-based meta-analysis of
607 neuroimaging studies of metacognitive judgements. *Brain and neuroscience advances*, 2,
608 2398212818810591.
- 609
- 610