

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

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

*These authors contributed equally to this work

¹ Wellcome Centre for Human Neuroimaging, UCL

² Max Planck UCL Centre for Computational Psychiatry and Ageing Research

³ Department of Experimental Psychology, UCL

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

Abstract: A key goal of consciousness science is identifying neural signatures of being aware vs. unaware of simple stimuli. This is often investigated in the context of near-threshold detection, with reports of stimulus awareness being linked to heightened activation in a frontoparietal network. However, due to the fact that reports of stimulus presence are also associated with higher confidence than reports of stimulus absence, these results could be explained by frontoparietal regions encoding stimulus visibility, decision confidence or both. Consistent with this view, previously we showed that prefrontal regions encode confidence in decisions about target presence (Mazor, Friston & Fleming, 2020). Here, we further ask if prefrontal cortex also encodes information about stimulus visibility over and above confidence. We first show that, whereas stimulus identity was best decoded from the visual cortex, stimulus visibility (presence vs. absence) was best decoded from prefrontal regions. To control for effects of confidence, we then selectively sampled trials prior to decoding to equalize the confidence distributions between absence and presence responses. This analysis revealed that posterior medial frontal cortex encoded stimulus visibility over and above decision confidence. We interpret our findings as providing support for a representation of stimulus visibility in specific higher-order cortical circuits, one that is dissociable from representations of both decision confidence and stimulus identity.

Introduction

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

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

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

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

Methods

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

Participants

46 participants took part in the study (ages 18–36, mean = 24 ± 4). We applied the same subject- and block-wise exclusion criteria as in the original study. Specifically, participants were excluded for having low response accuracy, pronounced response bias, or insufficient variability in their confidence ratings. 35 participants met our pre-specified inclusion criteria (ages 18–36, mean = 24 ± 4 ; 20 females). All analyses are based on the included blocks from these 35 participants.

Design and procedure

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

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

Participants performed 5 experimental runs comprising one discrimination and one detection block, each of 40 trials, presented in random order. A bonus was awarded for accurate responses and confidence ratings.

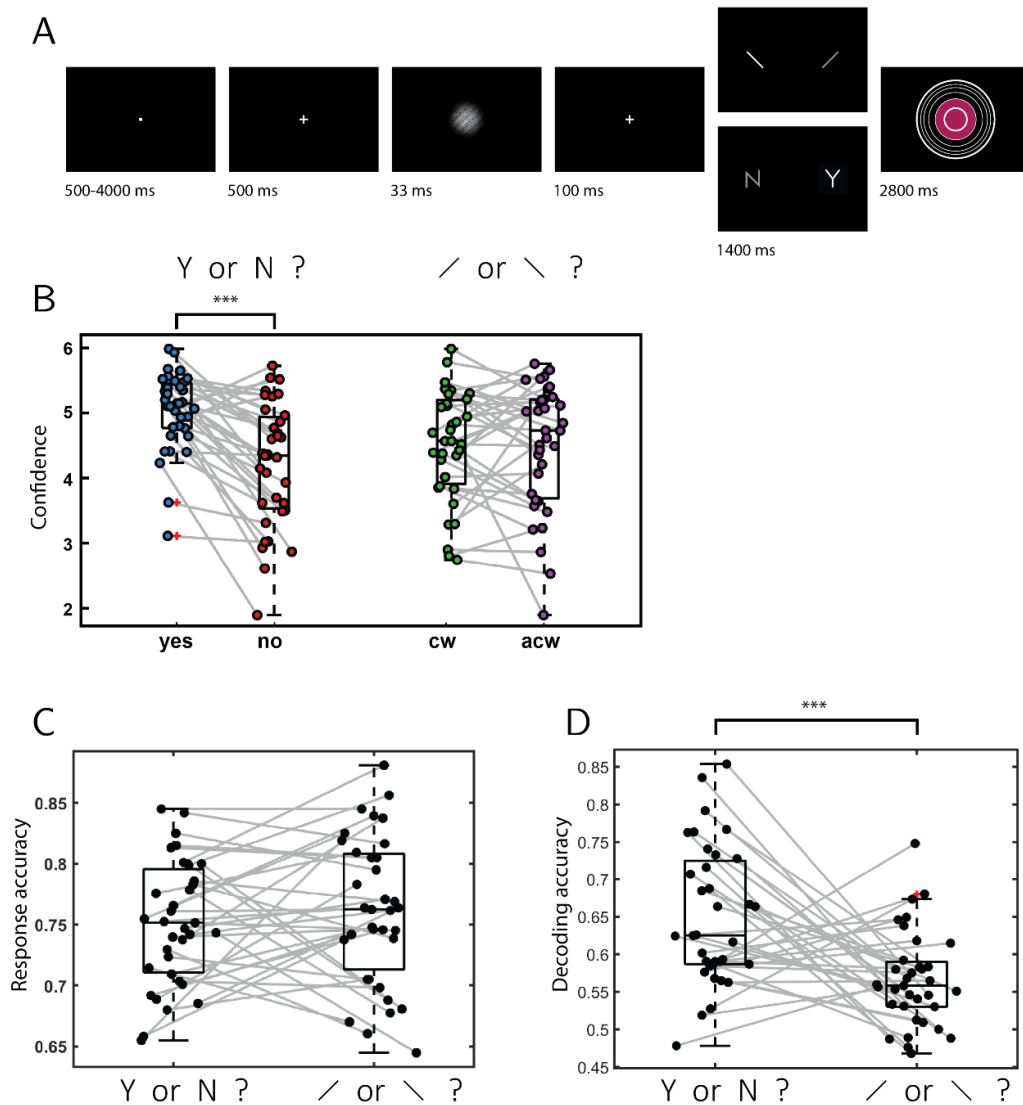


Figure 1: Experimental design and behavioural results. A: In discrimination trials, participants made discrimination judgments about clockwise and anticlockwise tilted noisy gratings, and then rated their subjective confidence by controlling the size of a colored circle. In detection judgments, decisions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

Scanning parameters

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

Analysis

Preprocessing

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

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

Multivariate analysis

Only correct trials were used for decoding. Stimulus presence (present vs. absent) was decoded during detection blocks, and stimulus identity (clockwise vs. anticlockwise orientation) during discrimination blocks. Both decoding analyses used an LDA (Linear Discriminant Analysis) classifier with leave-one-run-out cross-validation and a searchlight radius of 4 voxels (~257 voxels per searchlight). Significance testing was done using permutation testing to generate the empirical null-distribution. We followed the approach suggested by (Stelzer, Chen, & Turner, 2013) for

searchlight MVPA measurements which uses a combination of permutation testing and bootstrapping to generate chance distributions for group studies. Per participant, 25 permutation maps were generated by permuting the class labels within each run. Group-level permutation distributions were subsequently generated by bootstrapping over these 25 maps, i.e. randomly selecting one out of 25 maps per participant. 10000 bootstrapping samples were used to generate the group null-distribution per voxel and per comparison. *P*-values were calculated per searchlight or ROI as the right-tailed area of the histogram of permuted accuracies from the mean over participants. We corrected for multiple comparisons in the searchlight analyses using whole-brain FDR-correction. Cluster correction was performed, ensuring that voxels were only identified as significant if they belonged to a cluster of at least 50 significant voxels (Dijkstra, Bosch, & van Gerven, 2017).

Results

Decoding of stimulus presence and orientation

We first searched for multivariate activation patterns that encoded information about stimulus orientation (in discrimination) and stimulus presence/visibility (in detection). Stimulus orientation could be reliably decoded only from the visual cortex (Fig. 2). In contrast, information about stimulus presence was identified in parietal and prefrontal brain regions, including the dorsolateral prefrontal cortex, the middle frontal gyrus, and the precuneus (see Fig. 2, left panel; for unthresholded classification maps, see neurovault.org/collections/9912).

Based on these maps, we decided to focus our subsequent analyses on four regions of interest (ROIs; see Fig. S1): an occipital ROI, defined using the AICHA atlas as 'occipital mid' regions (Joliot et al., 2015), and three prefrontal ROIs which were also used in Mazor et al (2020): the posterior medial frontal cortex (pmFC; an 8 mm sphere around MNI coordinates [0, 17, 46]), Brodmann area 46, and the lateral frontopolar cortex (BA46 and FPI; both defined based on a connectivity-based parcellation; Neubert et al., 2014). Within these four ROIs, stimulus orientation could be decoded significantly from the occipital ($M = 0.54$, $SD = 0.09$, $p < 0.0001$) and FPI ROIs ($M = 0.51$, $SD = 0.06$, $p = 0.04$). In contrast, stimulus presence could be decoded from pmFC ($M = 0.53$, $SD = 0.08$, $p = 0.0009$), area 46 ($M = 0.54$, $SD = 0.06$, $p < 0.0001$) and FPI ROIs ($M = 0.52$, $SD = 0.07$, $p = 0.015$), but not from the occipital ROI ($M = 0.51$, $SD = 0.07$, $p = 0.11$). Classification accuracy showed a significant ROI x task interaction ($F(3,32) = 5.31$, $p = 0.004$; see Fig. 2, right panel), suggesting that stimulus identity and stimulus presence are encoded differentially across ROIs.

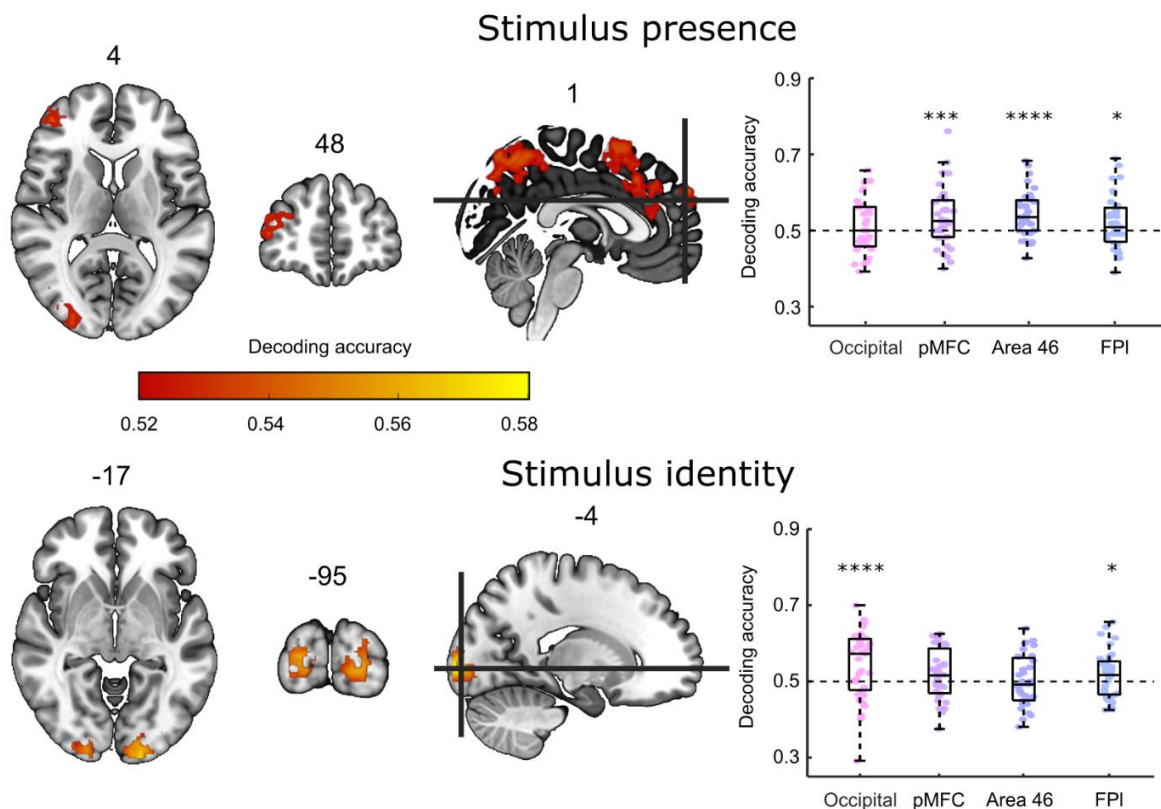


Figure 2: **Decoding of stimulus presence and stimulus identity.** Left: whole brain searchlight decoding, corrected for multiple comparisons at the cluster level using a nonparametric permutation test. Right: classification accuracy in the four regions of interest. *: $p < 0.5$, **: $p < 0.01$, ***: $p < 0.001$, ****: $p < 0.0001$.

Confidence-matching via downsampling

Prefrontal decoding of stimulus presence, but not orientation, is consistent with the idea that subjective visibility is represented in a frontoparietal network. However, given the fact that activity in prefrontal cortex is also sensitive to variation in confidence (Vacarro & Fleming, 2018), and the fact that we found a significant difference in confidence between classes during detection but not discrimination (Fig. 1B), it is plausible that prefrontal decoding of detection also reflects representations of confidence, instead of visibility. Indeed, mean confidence in detection 'yes' responses (correct responses only) was 5.03 on a 1-6 scale and significantly higher than 4.21 for detection 'no' responses ($t(34)=5.86$, $p < 0.001$). In contrast, confidence was similar for clockwise (4.28) and anticlockwise (4.25) discrimination responses ($t(34)=0.31$, $p=0.76$). Consistent with detection-specific confidence differences, a linear classifier reliably separated detection 'yes' and 'no' responses based on decision confidence alone (mean cross-validated classification accuracy = 0.65), but performed worse when trained to classify discrimination responses based on confidence (mean cross-validated classification accuracy = 0.57; $t(34)=3.88$, $p < 0.001$ for a paired t-test testing the difference in classification accuracy between detection and discrimination; see Fig. 1D).

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

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

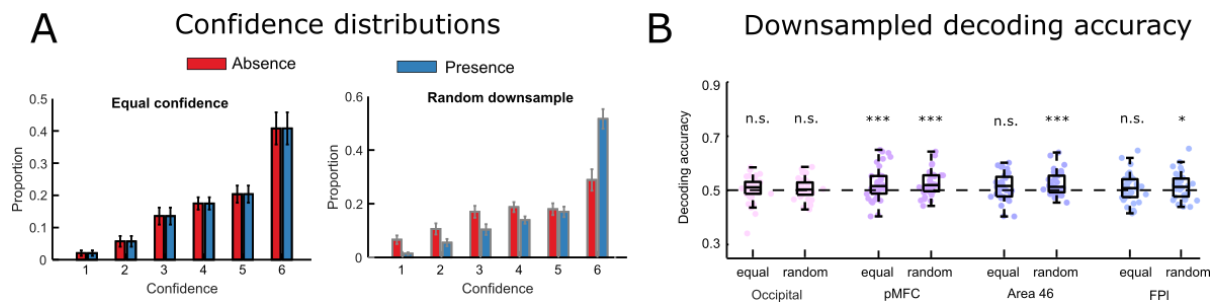


Figure 3: Stimulus presence downsampling analysis. A: for each participant, trials were deleted until confidence distributions were matched for target present and target absent responses. As a control analysis, we repeated this procedure with random downsampling, deleting the same number of trials irrespective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

When equalizing confidence, classification accuracy for decoding stimulus presence remained significant in pMFC ($M = 0.52$, $SD = 0.06$, $p = 0.002$). However, decoding was no longer significant after equalizing the confidence distributions in FPI ($M = 0.51$, $SD = 0.05$, $p =$

0.11), and only marginally significant in area 46 ($M=0.51$, $SD=0.05$, $p=0.07$). In both regions, decoding was still significant after random downsampling (FPI: $M = 0.52$, $SD = 0.05$, $p=0.02$; area 46: $M=0.53$, $SD=0.04$, $p=0.0017$). A decrease in classification accuracy after equalizing confidence relative to random downsampling was marginally significant in area 46 ($t(34)=-1.733$, $p=0.09$), but not in the FPI ROI ($t(34)=-1.615$, $p=0.11$). In the pMFC ROI, classification accuracy for confidence-matched and random downsampling was highly similar (0.524 and 0.525, $t(34)=-0.20$, $p=0.84$). Taken together, these results show that in pMFC, but not area 46 and FPI, stimulus presence/visibility can be reliably decoded independent of differences in decision confidence.

When decoding stimulus identity in the discrimination task, confidence-matching had no effect on classification accuracy relative to random downsampling (see Fig. S2). This is consistent with there already being little difference in the (behavioural) confidence distributions between the two response types in discrimination blocks. Importantly, in pMFC, we observed no significant classification of stimulus identity, regardless of whether the analysis used confidence-matched data or not. In other words, in this prefrontal ROI, we were able to decode visibility (independently of confidence) but not identity.

Discussion

What role the prefrontal cortex plays in visual awareness is much debated (e.g. Aru, Bachmann, Singer & Melloni, 2012; Boly et al., 2017). Here, we investigated whether prefrontal areas encode the visibility of a faint stimulus independently of stimulus identity and decision confidence. We first showed that prefrontal ROIs tracked stimulus presence during a detection task but not stimulus identity during a discrimination task, consistent with prefrontal involvement in encoding of stimulus visibility. However, because seeing a stimulus is associated with higher confidence than not seeing a stimulus, this asymmetry could also reflect (potentially domain-general) confidence coding in frontal areas. To investigate this possibility, we tested whether decoding of stimulus presence remained significant after controlling for differences in confidence. We found that such decoding was indeed still possible in pMFC, but not in area 46 and FPI. Taken together, these results suggest that pMFC encodes stimulus visibility over and above both stimulus identity and decision confidence.

Conceptually, visibility and decision confidence appear similar. They can both be defined in terms of precision: the precision of a visual percept in the first case, and the precision with which a decision is made in the second (Denison et al., 2017). Empirically, neural correlates of visibility and decision confidence overlap, specifically in the dorsolateral prefrontal cortex (dlPFC) but also in medial prefrontal, parietal, and insular cortices (Vacarro & Fleming, 2018). Notwithstanding this conceptual and empirical overlap, visibility and confidence are not one and the same thing. Critically, within a Bayesian framework, decision confidence is defined as the probability correct of a particular response, and should therefore be sensitive not only to the precision of sensory representations, but also response requirements (Pouget et al., 2016; Bang & Fleming, 2018). Accordingly, visibility judgments scale with stimulus contrast even in trials in which participants make erroneous decisions, but confidence judgments show a different profile,

and are sensitive to stimulus contrast only for correct responses (Rausch and Zehelsteiner, 2016).

Despite a theoretical distinction between confidence and visibility, neuroimaging findings of visual awareness have often not been able to separate their respective contributions to differential brain activation. For example, it has not been possible to determine whether the dorsolateral prefrontal cortex is more active on aware versus unaware trials because it is sensitive to subjective visibility, or because participants are generally more confident in their decisions when they are aware of a stimulus. In an exploratory analysis of existing imaging data, we found that visibility was encoded independently of confidence in pMFC, but not in the more anterior Brodmann area 46 and lateral frontopolar cortex.

As reported in Mazor et al. (2020), univariate analysis of this data indicated a similar parametric modulation of confidence for detection and discrimination responses in pMFC. Specifically, a similar modulation of confidence in decisions about target presence and absence indicate that univariate signal in this region also scales with decision confidence. Univariate analysis did not reveal a pMFC modulation of visibility, which would manifest as an interaction of confidence and class in detection (because visibility is negatively correlated with confidence in 'no' responses, but positively correlated with confidence in 'yes' responses). However, a pre-registered cross-classification analysis revealed shared multivariate representations for discrimination confidence and detection responses indicating whether a stimulus is seen or not in pMFC and area 46 (Mazor et al., 2020; Appendix 8). We previously interpreted these findings as indicating that multivariate spatial activation patterns in area 46 and pMFC hold information about stimulus visibility, because like detection responses, confidence during discrimination might also track stimulus visibility (it is easier to determine what it is when you see it clearer). Our current results corroborate this finding with respect to pMFC, and further show that above chance cross-classification in this region is not merely driven by differences in subjective confidence between 'yes' and 'no' responses during detection. Taken together, these results suggest that both confidence and visibility are represented in different components of the pMFC signal.

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

Our results with respect to the lateral frontopolar cortex (FPI) are more difficult to interpret. We found that this area did not represent stimulus presence over and above confidence, but that it did represent stimulus identity, even after controlling for confidence differences between the different stimulus classes. Several factors may have contributed to these results. First, our observation that the FPI does not encode visibility irrespective of

confidence does not mean that this region cannot play a role in visual awareness. In target absence trials, participants can sometimes be fully aware of the absence of a target – a case where visibility is low, but awareness (of absence) is high (Mazor & Fleming, 2020). Therefore, if FPI tracks content-invariant aspects of visual awareness, its activation may not differentiate between target presence and target absence. However, a representation of stimulus identity in FPI suggests that this area might also encode stimulus content. We are not aware of previous reports of decoding of visual content from the frontopolar cortex. Moreover, a recent meta-analysis reported no known effects of intracranial electrical stimulation of the frontopolar cortex on spontaneous reports of visual experience (Raccah, Block & Fox, 2021). Given the relatively modest effect sizes in FPI decoding of stimulus identity ($M=0.51$) in comparison to the more robust encoding of stimulus identity in occipital cortex ($M=0.55$), we are cautious in over-interpreting this surprising result. Future studies are necessary to explore to what extent FPI truly represents stimulus identity, and/or contributes to visual awareness.

Finally, when considering the implications of these findings for the study of visual awareness and its neural correlates, it is important to note the difference between subjective reports of stimulus awareness, and decisions about the presence or absence of a target stimulus in a perceptual detection task. While the first is a subjective decision about the contents of one's perception, the second is a report of one's beliefs about the state of the external world. Consequently, these two types of decisions draw on different sets of prior beliefs and expectations. For example, in detection, but not in subjective visibility reports, participants may adjust their decision criterion when noticing that they haven't detected a stimulus in a long time. Furthermore, participants may base their detection responses not on the visibility of a stimulus, but on other visual and non-visual cues (adopting a different *criterion content*; Kahneman, 1968). Our findings are based on the analysis of detection decisions, and their generalizability to subjective awareness reports is an open empirical question.

To conclude, an exploratory data analysis revealed that pMFC encodes stimulus visibility independent of task response and confidence. Our results support a functional dissociation between the neural correlates of visibility, confidence, and stimulus identity, thus serving to disentangle key contributions to the neural correlates of visual awareness.

Funding

N.D. is supported by a Rubicon grant from the Netherlands Organization for Scientific Research (NWO) [019.192SG.003] and SMF is funded by a Wellcome/Royal Society Sir Henry Dale Fellowship (206648/Z/17/Z) and a Philip Leverhulme Prize from the Leverhulme Trust. The Wellcome Centre for Human Neuroimaging is supported by core funding from the Wellcome Trust (206648/Z/17/Z).

The authors declare that there are no competing interests.

Aru, J., Bachmann, T., Singer, W., & Melloni, L. (2012). Distilling the neural correlates of consciousness. *Neuroscience & Biobehavioral Reviews*, 36(2), 737-746.

Bang, D., & Fleming, S. M. (2018). Distinct encoding of decision confidence in human medial prefrontal cortex. *Proceedings of the National Academy of Sciences*, 115(23), 6082-6087.

Bayne, T., & Hohwy, J. (2013). Consciousness: theoretical approaches. *Neuroimaging of consciousness*, 23-35.

Boly, M., Massimini, M., Tsuchiya, N., Postle, B. R., Koch, C., & Tononi, G. (2017). Are the neural correlates of consciousness in the front or in the back of the cerebral cortex? Clinical and neuroimaging evidence. *Journal of Neuroscience*, 37(40), 9603-9613.

Brown, R., Lau, H., & LeDoux, J. E. (2019). Understanding the higher-order approach to consciousness. *Trends in cognitive sciences*, 23(9), 754-768.

Dehaene, S., Changeux, J. P., & Naccache, L. (2011). The global neuronal workspace model of conscious access: from neuronal architectures to clinical applications. *Characterizing consciousness: From cognition to the clinic?*, 55-84.

Dehaene, S., Naccache, L., Cohen, L., Le Bihan, D., Mangin, J. F., Poline, J. B., & Rivière, D. (2001). Cerebral mechanisms of word masking and unconscious repetition priming. *Nature neuroscience*, 4(7), 752-758.

Dehaene, S., Sergent, C., & Changeux, J. P. (2003). A neuronal network model linking subjective reports and objective physiological data during conscious perception. *Proceedings of the National Academy of Sciences*, 100(14), 8520-8525.

De Martino, B., Fleming, S. M., Garrett, N., & Dolan, R. J. (2013). Confidence in value-based choice. *Nature neuroscience*, 16(1), 105-110.

Denison, R. N. (2017). Precision, Not Confidence, Describes the Uncertainty of Perceptual Experience: Comment on John Morrison's "Perceptual Confidence". *Analytic Philosophy*, 58(1), 58-70.

Fisch, L., Privman, E., Ramot, M., Harel, M., Nir, Y., Kipervasser, S., ... & Malach, R. (2009). Neural "ignition": enhanced activation linked to perceptual awareness in human ventral stream visual cortex. *Neuron*, 64(4), 562-574.

Fleming, S. M., Huijgen, J., & Dolan, R. J. (2012). Prefrontal contributions to metacognition in perceptual decision making. *Journal of Neuroscience*, 32(18), 6117-6125.

Fleming, S. M., Van Der Putten, E. J., & Daw, N. D. (2018). Neural mediators of changes of mind about perceptual decisions. *Nature neuroscience*, 21(4), 617-624.

Joliot, M., Jobard, G., Naveau, M., Delcroix, N., Petit, L., Zago, L., ... & Tzourio-Mazoyer, N. (2015). AICHA: An atlas of intrinsic connectivity of homotopic areas. *Journal of neuroscience methods*, 254, 46-59.

- Kahneman, D. (1968). Method, findings, and theory in studies of visual masking. *Psychological Bulletin*, 70(6p1), 404.
- Koivisto, M., & Revonsuo, A. (2010). Event-related brain potential correlates of visual awareness. *Neuroscience & Biobehavioral Reviews*, 34(6), 922-934.
- Lau, H. C., & Passingham, R. E. (2006). Relative blindsight in normal observers and the neural correlate of visual consciousness. *Proceedings of the National Academy of Sciences*, 103(49), 18763-18768.
- Lau, H., & Rosenthal, D. (2011). Empirical support for higher-order theories of conscious awareness. *Trends in cognitive sciences*, 15(8), 365-373.
- Mazor, M., Friston, K. J., & Fleming, S. M. (2020). Distinct neural contributions to metacognition for detecting, but not discriminating visual stimuli. *ELife*, 9, e53900.
- Mazor, M., & Fleming, S. M. (2020). Distinguishing absence of awareness from awareness of absence. *Philosophy and the Mind Sciences*, 1(II).
- Michel, M., & Morales, J. (2020). Minority reports: Consciousness and the prefrontal cortex. *Mind & Language*, 35(4), 493-513.
- Molenberghs, P., Trautwein, F. M., Böckler, A., Singer, T., & Kanske, P. (2016). Neural correlates of metacognitive ability and of feeling confident: a large-scale fMRI study. *Social cognitive and affective neuroscience*, 11(12), 1942-1951.
- Morales, J., Lau, H., & Fleming, S. M. (2018). Domain-general and domain-specific patterns of activity supporting metacognition in human prefrontal cortex. *Journal of Neuroscience*, 38(14), 3534-3546.
- Odegaard, B., Knight, R. T., & Lau, H. (2017). Should a few null findings falsify prefrontal theories of conscious perception?. *Journal of Neuroscience*, 37(40), 9593-9602.
- O'Reilly, J. X., Schüffegen, U., Cuell, S. F., Behrens, T. E., Mars, R. B., & Rushworth, M. F. (2013). Dissociable effects of surprise and model update in parietal and anterior cingulate cortex. *Proceedings of the National Academy of Sciences*, 110(38), E3660-E3669.
- Pouget, A., Drugowitsch, J., & Kepecs, A. (2016). Confidence and certainty: distinct probabilistic quantities for different goals. *Nature neuroscience*, 19(3), 366.
- Racch, O., Block, N., & Fox, K. C. (2021). Does the prefrontal cortex play a necessary role in consciousness? Insights from intracranial electrical stimulation of the human brain. *Journal of Neuroscience*, 1(41).
- Rausch, M., & Zehetleitner, M. (2016). Visibility is not equivalent to confidence in a low contrast orientation discrimination task. *Frontiers in psychology*, 7, 591.
- Sahraie, A., Weiskrantz, L., Barbur, J. L., Simmons, A., Williams, S. C. R., & Brammer, M. J. (1997). Pattern of neuronal activity associated with conscious and unconscious processing of visual signals. *Proceedings of the National Academy of Sciences*, 94(17), 9406-9411.

Vaccaro, A. G., & Fleming, S. M. (2018). Thinking about thinking: A coordinate-based meta-analysis of neuroimaging studies of metacognitive judgements. *Brain and neuroscience advances*, 2, 2398212818810591.

Supplementary materials

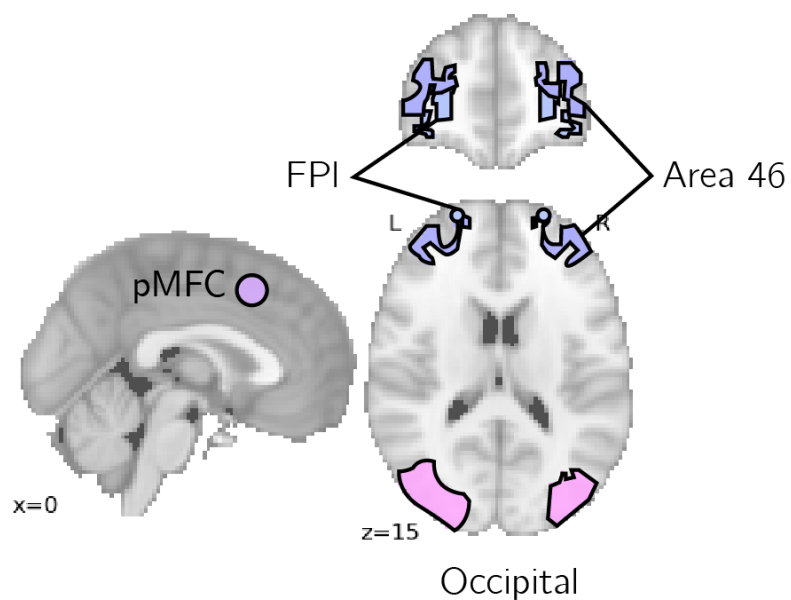


Figure S1: **Regions of interest.** The four regions of interest comprised an occipital ROI, pMFC, FPI and Brodmann area 46.

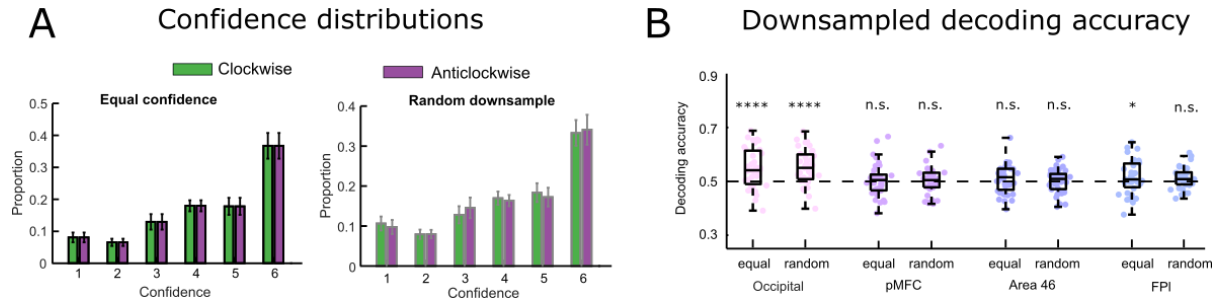


Figure S2: **Stimulus identity downsampling analysis.** Same conventions as in Figure 3