

Dynamic representations of faces in the human ventral visual stream link visual features to behaviour

Diana C. Dima*¹ and Krish D. Singh¹

¹*Cardiff University Brain Research Imaging Centre (CUBRIC), School of Psychology, Cardiff University, Cardiff CF24 4HQ, United Kingdom*

Abstract

Humans can rapidly extract information from faces even in challenging viewing conditions, yet the neural representations supporting this ability are still not well understood. Here, we manipulated the presentation duration of backward-masked facial expressions and used magnetoencephalography (MEG) to investigate the computations underpinning rapid face processing. Multivariate analyses revealed two stages in face perception, with the ventral visual stream encoding facial features prior to facial configuration. When presentation time was reduced, the emergence of sustained featural and configural representations was delayed. Importantly, these representations explained behaviour during an expression recognition task. Together, these results describe the adaptable system linking visual features, brain and behaviour during face perception.

Keywords: face perception; magnetoencephalography (MEG); multivariate pattern analysis (MVPA); representational similarity analysis (RSA)

*Corresponding author: dimadc@cardiff.ac.uk (D.C. Dima)

1 Introduction

2 Our highly specialized face processing abilities are thought to be supported by feature-
3 based face detection followed by configural processing (Calder, Young, Keane, & Dean,
4 2000; Maurer, Grand, & Mondloch, 2002). However, it is still unclear how the brain
5 efficiently represents a high-dimensional array of relevant facial features, and how this
6 helps accomplish a wide range of behavioural goals.

7 Although there is disagreement on the exact sequence of processing stages, face percep-
8 tion is generally thought to progress from isolated features to first-order configuration (the
9 feature positioning common across all faces) and second-order configuration (the identity-
10 specific spacing between features), with holistic processing linking these into a gestalt
11 (Farah, Wilson, & Tanaka, 1998; Harris & Aguirre, 2008; Piepers & Robbins, 2012). On
12 the other hand, some behavioural goals, such as identity recognition, may rely on facial
13 features and not on holistic perception (Visconti Di Oleggio Castello, Wheeler, Cipolli, &
14 Gobbini, 2017).

15 Furthermore, although the neural correlates of face perception have been reliably
16 mapped in space and time, there is little agreement on how, where, and when specific
17 computations are implemented. Both modular and distributed neural codes are thought
18 to support face perception, with different computations being implemented within each of
19 the ventral face-responsive areas (Grill-Spector, Weiner, Gomez, Stigliani, & Natu, 2018;
20 Freiwald, Duchaine, & Yovel, 2016). For efficient information extraction, faces may be
21 represented along low-dimensional axes based on features or topology (Henriksson, Mur,
22 & Kriegeskorte, 2015; Leopold, O'Toole, Vetter, & Blanz, 2001); for example, a sparse
23 identity code has been shown to predict neural responses to faces in primates (Chang &
24 Tsao, 2017). However, it remains an open question how such codes adapt to task require-
25 ments and viewing conditions, and the dynamics of face feature representations are not
26 well understood.

27 Here, we focused on the temporal dynamics of face representations during a chal-
28 lenging expression discrimination task, by combining magnetoencephalography (MEG)
29 with multivariate pattern analyses and a rapid presentation paradigm. We manipulated
30 the presentation duration of backward-masked faces, some of which were shown outside
31 awareness, to disentangle face detection from expression processing. This allowed us to
32 evaluate the impact of limiting visual input on representational dynamics, while keeping
33 task demands constant. We used source-space representational similarity analysis (RSA)
34 and variance partitioning to evaluate the contribution of visual features to MEG responses

35 and behaviour.

36 We found that among the visual features tested, facial features and configuration
37 were most strongly represented in the ventral stream and contributed to behaviour. The
38 temporal dynamics of these representations changed in response to stimulus duration,
39 suggesting that it is important to study visual feature coding in dynamic contexts and
40 with high temporal resolution. Finally, despite a behavioural effect, a neural response to
41 faces outside of awareness did not encode any of the stimulus features tested, highlighting
42 the qualitative distinction between face detection and face categorization.

43 **Methods**

44 **Participants**

45 The participants were 25 healthy volunteers (16 female, age range 19-42, mean age 25.6
46 \pm 5.39). All volunteers gave written consent to participate in the study in accordance
47 with The Code of Ethics of the World Medical Association (Declaration of Helsinki). All
48 procedures were approved by the ethics committee of the School of Psychology, Cardiff
49 University.

50 **Stimuli**

51 The stimulus set consisted of 20 faces with angry, neutral and happy expressions (10
52 female faces; model numbers: 2, 6, 7, 8, 9, 11, 14, 16, 17, 18, 22, 23, 25, 31, 34, 35,
53 36, 38, 39, 40) from the NIMSTIM database (Tottenham et al., 2009). The eyes were
54 aligned using automated eye detection as implemented in the Matlab Computer Vision
55 System toolbox (Mathworks, Inc., Natick, Massachusetts). An oval mask was used to crop
56 the faces to a size of 378×252 pixels subtending 2.6×3.9 degrees of visual angle. All
57 images were converted to grayscale. Their spatial frequency was matched by specifying
58 the rotational average of the Fourier amplitude spectra as implemented in the SHINE
59 toolbox (Willenbockel et al., 2010), and Fourier amplitude spectra for all faces were set to
60 the average across the face set.

61 Masks and control stimuli were created by scrambling the phase of all face images in
62 the Fourier domain. This was achieved by replacing the phase information in each of the
63 images with phase information from a white noise image of equal size (Perry & Singh,
64 2014). To ensure matched low-level properties between face and control stimuli, pixel
65 intensities were normalized between each image and its scrambled counterpart, using the

66 minimum and maximum pixel intensity of the scrambled image.

67 **Experimental design**

68 At the start of each trial, a white fixation cross was centrally presented on an isoluminant
69 gray background. Its duration was pseudorandomly chosen from a uniform distribution
70 between 1.3 and 1.6 s. A face stimulus was then centrally presented with a duration of
71 either 10 ms, 30 ms or 150 ms; the stimulus was followed by a phase-scrambled mask with
72 a duration of 190 ms, 170 ms or 50 ms respectively (for a constant total stimulus duration
73 of 200 ms). In each block, 10 trials contained no face; instead, a phase-scrambled control
74 stimulus was flashed for 10 ms and followed by another mask.

75 After a 500 ms delay intended to dissociate face perception from response preparation,
76 participants had to correctly select the expression they had perceived out of three alter-
77 natives presented on screen (Figure 3A). They had 1.5 seconds to make a button press; if
78 they were sure that no face had been presented, they could refrain from responding. The
79 mapping of the response buttons to emotional expressions changed halfway through the
80 experiment so as to ensure that emotional expression processing would not be confounded
81 by specific motor preparation effects.

82 Next, participants had to rate how clearly they had seen the face using a 3-point scale
83 starting from 0. They were instructed to only select 0 if no face had been perceived, 1 if
84 they had perceived a face but not clearly, and 2 if they had clearly perceived the face. They
85 had 2 seconds to make this response. Note that since the expression discrimination task
86 was not forced-choice, references to awareness in this paper refer exclusively to subjective
87 awareness, as indicated by perceptual ratings.

88 In each of four blocks, each face was presented once with each of the three possible
89 stimulus durations. We thus collected 80 trials per condition, except for the control
90 condition (containing no face) which only had 40 trials.

91 **Data acquisition**

92 All participants with one exception acquired a whole-head structural MRI using a 1 mm
93 isotropic Fast Spoiled Gradient-Recalled-Echo pulse sequence.

94 Whole-head MEG recordings were made using a 275-channel CTF radial gradiometer
95 system (CTF, Vancouver, Canada) at a sampling rate of 1200 Hz. Four of the sensors
96 were turned off due to excessive sensor noise. An additional 29 reference channels were
97 recorded for noise rejection purposes and the primary sensors were analyzed as synthetic

98 third-order gradiometers (Vrba & Robinson, 2001).

99 Stimuli were presented using a ProPixx projector system (VPixx Technologies, Saint-
100 Bruno, Canada) with a refresh rate set to 100 Hz. Images were projected to a screen with
101 a resolution of 1920 x 1080 pixels situated at a distance of 1.2 m from the participant.
102 Recordings were made in four blocks of approximately 15 minutes each, separated by short
103 breaks. The data were collected in 2.5 s epochs beginning 1 s prior to stimulus onset.

104 Participants were seated upright while viewing the stimuli and electromagnetic coils
105 were attached to the nasion and pre-auricular points on the scalp in order to continuously
106 monitor head position relative to a fixed coordinate system on the dewar. To help co-
107 register the MEG data with the participants' structural MRI scans, we defined the head
108 shape of each subject using an ANT Xensor digitizer (ANT Neuro, Enschede, Netherlands).
109 An Eyelink 1000 eye-tracker system (SR Research, Ottawa, Canada) with a sampling rate
110 of 1000 Hz was used to track the subjects' right pupil and corneal reflex.

111 Behavioural analysis

112 In order to assess the effects of stimulus duration and face expression on behaviour, we
113 calculated confusion matrices based on expression discrimination responses to each stim-
114 ulus category (Figure 3). Performance was quantified as proportion correct trials after
115 excluding trials with no response, and a rationalized arcsine transformation was applied
116 prior to statistical analysis (Studebaker, 1985). We then performed a 3×3 repeated-
117 measures ANOVA with factors *Duration* (levels: 10 ms, 30 ms, 150 ms) and *Expression*
118 (levels: angry, happy, neutral).

119 MEG multivariate pattern analysis (MVPA)

120 To test for differences between conditions present in multivariate patterns, we used a linear
121 Support Vector Machine (SVM) classifier with L2 regularization and a box constraint $c =$
122 1. The classifier was implemented in Matlab using LibLinear (Fan, Chang, Hsieh, Wang,
123 & Lin, 2008) and the Statistics and Machine Learning Toolbox (Mathworks, Inc.). We
124 performed binary classification on (1) responses to neutral faces versus scrambled stimuli
125 (face decoding); (2) all three pairs of emotional expressions (expression decoding).

126 For face decoding, time-resolved classification was performed separately for each stim-
127 ulus duration. To assess the presence of subjectively non-conscious responses, the classi-
128 fication of faces presented for 10 ms was performed after excluding any trials reported as
129 containing a face. To ensure that decoding results were not biased by stimulus repetitions

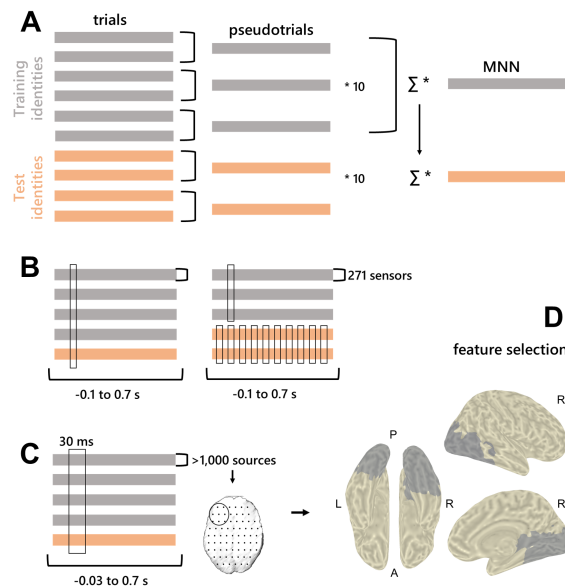


Figure 1: Overview of the MVPA analysis pipeline. **A.** Trial averaging and multivariate noise normalization (MNN) procedure. Σ is the error covariance matrix. **B.** Sensor-space time-resolved decoding (left) and temporal generalization (right). **C.** Source-space searchlight decoding procedure. **D.** Sources included in the representational similarity analysis based on face vs. scrambled classification results. P: posterior; A: anterior; L: left; R: right.

130 or recognition of face identities across the training and test sets, cross-exemplar five-fold
 131 cross-validation was used to assess classification performance: the classifier was trained on
 132 16 of the 20 face identities and 8 of the 10 scrambled images, and tested on the remaining
 133 4 faces and 2 scrambled exemplars.

134 To assess similarities between responses across stimulus duration conditions, face cross-
 135 decoding was also performed, whereby a decoder was trained on 150 ms faces and tested
 136 on 30 ms faces and vice-versa. The analysis was repeated for all pairs of conditions, using
 137 cross-exemplar cross-validation to ensure true generalization of responses; the resulting
 138 accuracies were averaged across the two training/testing directions, which led to similar
 139 results.

140 The temporal structure of face responses was assessed through temporal generaliza-
 141 tion decoding (King & Dehaene, 2014). Classifier models were trained on each sampled
 142 time point between -0.1 and 0.7 s and tested on all time points in order to evaluate the
 143 generalizability of neural patterns over time at each stimulus duration. For this analysis,
 144 a cross-exemplar hold-out procedure was used to speed up computation (the training and
 145 test sets each consisted of 10 face identities/5 scrambled exemplars).

146 For expression decoding, classification was separately applied to all pairs of emotional
 147 expression conditions for each stimulus duration. As low trial numbers were a limitation

148 of the study design, we increased the power of our analysis by also pooling together trials
149 containing faces shown for 30 ms and 150 ms (which were shown to share representations
150 in the cross-decoding analysis). Performance was evaluated using five-fold cross-exemplar
151 cross-validation. Note that splitting the datasets according to perceptual rating led to
152 largely similar results (Supplementary Figure 3).

153 To achieve equal class sizes in face decoding, face trials were randomly subsampled
154 (after cross-exemplar partitioning) to match the number of scrambled trials. For expression
155 classification, trial numbers did not significantly differ between conditions after artefact
156 rejection ($F(1.92, 46.18) = 0.15, P = 0.85, \eta^2 = 0.0062$).

157 MEG sensor-level analyses

158 MEG data were analyzed using Matlab (Mathworks, Inc.) and the Fieldtrip toolbox
159 (Oostenveld, Fries, Maris, & Schoffelen, 2011). Prior to analysis, trials containing excessive
160 eye or muscle artefacts were excluded based on visual inspection, as were trials exceeding 5
161 mm in head motion (quantified as the displacement of any head coil between two sampled
162 time points). Using eyetracker information, we also excluded trials containing saccades
163 and fixations away from stimulus or blinks during stimulus presentation. A mean of 8.71%
164 $\pm 9.4\%$ of trials were excluded based on this procedure.

165 For all analyses, MEG data were downsampled to 300 Hz and baseline corrected using
166 the 500 ms before stimulus onset. A low-pass filter was applied at 100 Hz and a 50 Hz
167 comb filter was used to remove the mains noise and its harmonics.

168 To improve SNR (Grootswagers, Wardle, & Carlson, 2017), each dataset was divided
169 into 20 equal partitions and pseudo-trials were created by averaging the trials in each
170 partition. This procedure was repeated 10 times with random assignment of trials to
171 pseudo-trials and was performed separately for the training and test sets.

172 To improve data quality, we performed multivariate noise normalization (MNN; Guggen-
173 mos, Sterzer, and Cichy, 2018). The time-resolved error covariance between sensors was
174 calculated based on the covariance matrix (Σ) of the training set (X) and used to nor-
175 malize both the training and test sets, in order to downweight MEG channels with higher
176 noise levels (Equation 1).

$$X^* = \Sigma^{-\frac{1}{2}} X \quad (1)$$

177 In sensor-level MVPA analyses, all 271 MEG sensors were included as features and de-
178 coding was performed for each sampled time point between -0.1 and 0.7 s around stimulus

179 onset.

180 MEG source-space analyses

181 For source analyses, each participant's MRI ($N=24$) was coregistered to the MEG data
182 by marking the fiducial coil locations on the MRI and aligning the digitized head shape
183 to the MRI with Fieldtrip. MEG data were projected into source space using a vectorial
184 Linearly Constrained Minimum Variance (LCMV) beamformer (Van Veen, van Dronge-
185 len, Yuchtman, & Suzuki, 1997). To reconstruct activity at locations equivalent across
186 participants, a template grid with a 10 mm isotropic resolution was defined using the
187 MNI template brain and was warped to each individual MRI. The covariance matrix was
188 calculated based on the average of all trials across conditions bandpass-filtered between
189 0.1 and 100 Hz; this was then combined with a single-shell forward model to create an
190 adaptive spatial filter, reconstructing each source as a weighted sum of all MEG sensor
191 signals (Hillebrand, Singh, Holliday, Furlong, & Barnes, 2005). To alleviate the depth bias
192 in MEG source reconstruction, beamformer weights were normalized by their vector norm
193 (Hillebrand, Barnes, Bosboom, Berendse, & Stam, 2012). To improve data quality, MNN
194 was included in the source localization procedure, by multiplying the normalized beam-
195 former filters by the error covariance matrix to ensure that sensors with higher noise levels
196 were downweighted. Next, the sensor-level data were multiplied by the corresponding
197 weighted filters in order to reconstruct the time-courses of virtual sensors at all loca-
198 tions in the brain. This resulted in three time-courses for each source, containing each
199 of the three dipole orientations, which were concatenated for use in the MVPA analysis
200 in order to maximize classification performance (Gohel, Lim, Kim, Kwon, & Kim, 2018).
201 Preprocessing (baseline correction and downsampling) was performed as for sensor-level
202 analyses.

203 A searchlight approach was used in source-space classification, whereby clusters with
204 a 10 mm radius were entered separately into the decoding analysis. To exclude sources
205 outside the brain and in the cerebellum, we restricted our searchlight analysis to sources
206 included in the 90-region Automated Anatomical Labelling (AAL) atlas (Tzourio-Mazoyer
207 et al., 2002). Given the 10 mm resolution of our sourcemodel, this amounted to a maximum
208 of 27 neighbouring sources being included as features (mean 26.9, median 27, SD 0.31).
209 Source-space subliminal face decoding was performed on 30 ms time windows with a 3
210 ms overlap using the time windows identified in sensor-space decoding in order to reduce
211 computational cost. We also performed supraliminal face decoding (150 ms faces vs.

212 scrambled stimuli) in order to identify a face-responsive ROI for use in the RSA analysis.
213 This was accomplished by identifying searchlights achieving a cross-subject accuracy above
214 the 99.5th percentile ($P < 0.005$, 66 searchlights; Figure 1).

215 **Significance testing**

216 We evaluated decoding performance using the averaged accuracy across subjects (propor-
217 tion correctly classified trials) and assessed its significance through randomization testing
218 (Nichols & Holmes, 2001).

219 For sensor-level decoding, 1,000 label shuffling iterations across the training and test
220 sets were used to estimate the null distribution using the time point achieving maximum
221 average accuracy in the MVPA analysis (Dima, Perry, & Singh, 2018). Omnibus correction
222 for multiple comparisons was applied across tests, time points and sources where applicable
223 (Nichols & Holmes, 2001; Singh, Barnes, & Hillebrand, 2003), with a supplementary false
224 discovery rate correction applied for tests where the null distribution was not separately
225 estimated. To avoid spurious effects, a threshold of 5 consecutive significant time points
226 (5^2 in 2D temporal generalization maps) was imposed. For source-space decoding, 100
227 randomization iterations were performed for each source cluster and subject in order to
228 reduce computational cost, which were randomly combined into 1000 whole-brain group
229 maps (Stelzer, Chen, & Turner, 2013). A minimal extent of three consecutive time windows
230 with a FDR-corrected $P < 0.005$ was applied.

231 **Representational Similarity Analysis (RSA)**

232 **Neural patterns and analysis framework**

233 To interrogate the content of neural representations in space and time, we performed
234 representational similarity analysis (RSA). For this analysis, MEG data were source re-
235 constructed as described above and trials were sorted according to expression and face
236 identity. RSA was performed separately for each stimulus duration and only trials con-
237 taining faces were included in the analysis. We tracked representational dynamics using a
238 searchlight analysis restricted to the occipitotemporal sources identified in face decoding,
239 with a temporal resolution of 30 ms. The exclusion of responses to scrambled stimuli from
240 the RSA ensured that feature selection was based on an orthogonal contrast (Figure 1).

241 To create MEG representational dissimilarity matrices (RDMs), we calculated the
242 squared cross-validated Euclidean distance between all pairs of face stimuli (Guggenmos
243 et al., 2018). Note that as the data were multivariately noise-normalized, this is equivalent

244 to the squared cross-validated Mahalanobis distance (Walther et al., 2016). For each
245 participant, the data were split into a training set (the first 2 sessions) and a test set (the
246 last 2 sessions). The two stimulus repetitions contained in each set were averaged, and
247 these were averaged across subjects to create training and test sets. To compute the cross-
248 validated Euclidean distance between two stimulus patterns (X^* , Y^*), we calculated the
249 dot products of pattern differences based on the training set and the test set (Equation 2).
250 This procedure has the advantage of increasing the reliability of distance estimates in the
251 presence of noise.

$$d^2(X^*, Y^*) = \sum_{i=1}^n (X_i^* - Y_i^*)_{train} (X_i^* - Y_i^*)_{test} \quad (2)$$

252 The spatiotemporally resolved MEG RDMs were then correlated with several model
253 RDMs to assess the contribution of different features to neural representations. In an initial
254 analysis, we calculated Spearman's rank correlation coefficients between each model RDM
255 and the MEG RDM (Nili et al., 2014). To further investigate the unique contribution
256 of each model, we entered the significantly correlated models based on visual features of
257 the images into a partial correlation analysis, where each model's correlation to the MEG
258 data was recalculated after partialling out the contribution of the other models.

259 Note that a model based on behaviour, which was also represented in the MEG data
260 for all stimulus duration conditions, was not included in the partial correlation analysis;
261 the rationale is that we were interested in the contribution of each visual property in-
262 dependently of the others, but we did not expect a unique contribution of behaviour in
263 the absence of expression-related visual properties, and partialling out the behavioural
264 model from the visual models would not be easily interpretable. Instead, we preferred
265 to independently describe the correlations between behaviour and visual features, brain
266 and behaviour, and brain and visual features, as the three main factors of interest in our
267 analysis.

268 **Model RDMs**

269 We investigated the temporal dynamics of face perception by assessing the similarity
270 between MEG patterns and 9 models quantifying behaviour and facial/visual properties
271 (Figure 2).

272 To create behavioural model RDMs, we calculated the number of error responses made
273 by each participant to each stimulus and summed these up to create a cross-subject be-
274 havioural RDM. For each stimulus duration, we created separate behavioural RDMs by

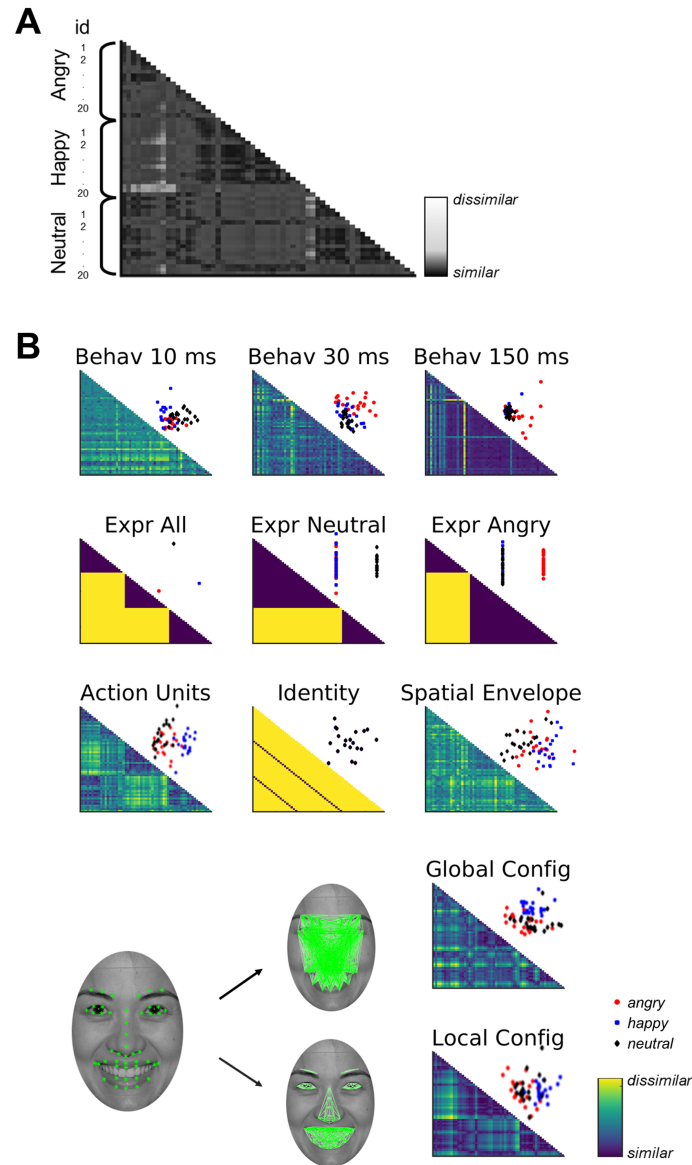


Figure 2: Models used in RSA analysis. **A**. Example model RDM: each model maps pairwise dissimilarities between faces, which are sorted according to expression and identity. **B**. Model RDMs showing predicted distances between all pairs of stimuli. 2D multidimensional scaling (MDS) plots are shown above each model to visualize how the three expression categories are organized according to each model. For the local and global configuration models, we also show the facial landmarks and the within-feature/between-feature distances used to create each model. Behav: behavioural models; Expr: high-level expression models (all-vs-all, neutral-vs-others, and angry-vs-others); Config: face configuration models.

275 calculating pairwise cross-validated Euclidean distances between error response patterns,
276 using a cross-session training/test split as described above.

277 To create face configuration RDMs, we first used OpenFace (Baltrusaitis, Robinson, &
278 Morency, 2016) to automatically detect and label face landmarks. The software created
279 68 2D landmarks for each face. We removed landmarks corresponding to the face outline
280 and the 2 outermost eyebrow landmarks, to account for cases in which these landmarks
281 were cropped out by the oval mask used in the MEG stimulus set. The final landmark
282 set consisted of 47 coordinates for 6 facial features (eyes, eyebrows, nose, and mouth),
283 which were visually inspected to ensure that they were correctly marked. To capture
284 feature-based (local) facial configuration, we calculated within-feature pairwise Euclidean
285 distances between landmarks (Figure 2B). To quantify global face configuration, we cal-
286 culated between-feature Euclidean distances (the distances between each landmark and
287 all landmarks belonging to different facial features). Distances were then concatenated to
288 create feature vectors describing each face in terms of its local/global configuration, and
289 Euclidean distances between them gave the final configural model RDMs. These mod-
290 els correspond to the featural and configural stages in classic models of face perception
291 (Diamond & Carey, 1986; Piepers & Robbins, 2012).

292 To create a high-level identity model, we assigned distances of 0 to pairs of face identi-
293 ties repeated across emotional expression conditions, and distances of 1 to pairs of different
294 face identities. We used a similar strategy to create high-level emotional expression mod-
295 els. An all-versus-all model was created by assigning distances of 0 to all faces belonging
296 to the same emotional expression condition, and distances of 1 to pairs of faces differing
297 in emotion. We also tested a neutral-versus-others model by assigning distances of 0 to all
298 emotional faces (happy + angry), and an angry-versus-others model by assigning distances
299 of 0 to all benign faces (happy + neutral).

300 To account for variability in expression that is not captured by such high-level binary
301 representations, we also tested a model based on Action Units. Action Units quantify
302 changes in expression by categorizing facial movements (Ekman & Friesen, 1977). We used
303 OpenFace (Baltrusaitis et al., 2016) to automatically extract the intensity of 12 Action
304 Units in our image set (Supplementary Table 4), and we calculated pairwise Euclidean
305 distances between these intensities for all pairs of faces in our stimulus set to obtain an
306 Action Unit RDM.

307 Finally, a spatial envelope model was created in order to capture image characteristics
308 using the GIST descriptor (Oliva & Torralba, 2001). This procedure extracts 512 values

309 per image by applying a series of Gabor filters at different orientations and positions,
310 and thus quantifies the average orientation energy at each spatial frequency. To obtain
311 the spatial envelope RDM, we calculated pairwise Euclidean distances between all images
312 using the GIST values.

313 Finally, models were subject to multidimensional scaling (MDS) to visualize how each
314 model represents the similarity between facial expressions in a 2D space (Figure 2).

315 **Significance testing**

316 To assess the significance of spatiotemporally resolved correlation maps, we used a ran-
317 domization approach. Model RDMs were shuffled 1,000 times and correlations were re-
318 computed for each of the 66 searchlights using the time window achieving the maximal
319 correlation coefficient across models for each of the stimulus duration conditions. Since
320 negative correlations were not expected and would not be easily interpretable, P-values
321 were calculated using a one-sided test (Furl, Lohse, & Pizzorni-Ferrarese, 2017). To correct
322 for multiple comparisons, P-values were omnibus-corrected by creating a maximal distri-
323 bution of randomized correlation coefficients across searchlights, models and conditions,
324 and FDR and cluster-corrected across timepoints ($\alpha = 0.05$, thresholded at 3 consecutive
325 time windows).

326 **Variance partitioning**

327 To gain more insight into the relationship between behavioural responses, expression cat-
328 egories and face configuration models, we used a variance partitioning approach (Greene,
329 Baldassano, Esteva, Beck, & Fei-fei, 2016; Groen et al., 2018). For each stimulus duration
330 condition, the corresponding behavioural RDM was entered into a hierarchical multiple
331 linear regression analysis, with three model RDMs as predictors: the two facial configura-
332 tion models and the most correlated high-level expression model (10 ms: neutral-vs-others;
333 30 and 150 ms: angry-vs-others). These models were selected to reduce the predictor space
334 before performing variance partitioning. To quantify the unique and shared variance con-
335 tributed by each model, we calculated the R^2 value for every combination of predictors (i.e.
336 all three models together, each pair of models separately, and each model separately). The
337 EulerAPE software was used for visualization (Micallef and Rodgers, 2014; Figure 6B).

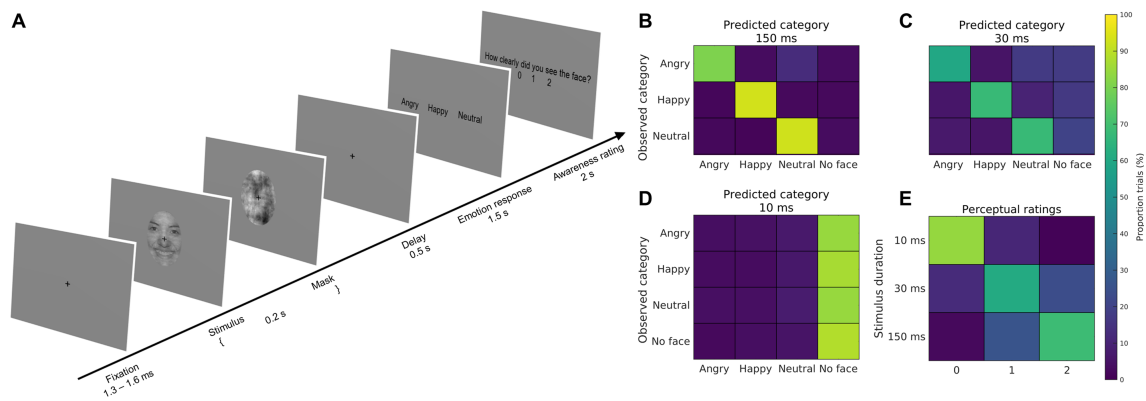


Figure 3: Overview of the experimental paradigm and behavioural results. **A**. Stimuli were presented on screen for 150 ms, 30 ms, or 10 ms, and were followed by a 50 ms, 170 ms, or 190 ms scrambled mask. **B-D**. Confusion matrices mapping the average proportion of trials receiving each of the possible responses (X-axis) out of the trials belonging to each category (Y-axis). "No response" trials were excluded for statistical analysis, but are shown here as representing a "no face" response. **E**. Perceptual ratings for each stimulus duration summarized as average proportion of trials.

338 Results

339 Perception and behaviour

340 We assessed the effects of stimulus duration and face expression on behaviour using a 3×3
 341 repeated-measures ANOVA with factors *Duration* (levels: 10 ms, 30 ms, 150 ms) and *Ex-*
 342 *pression* (levels: angry, happy, neutral) on rationalized arcsine-transformed accuracies
 343 (see Methods). Stimulus duration had a strong effect on expression discrimination perfor-
 344 mance, with average performance not exceeding chance level at 10 ms ($33.45\% \pm 2.99$)
 345 and rising well above chance at 30 and 150 ms ($78.62\% \pm 2.11$ and $91.83\% \pm 1$ re-
 346 spectively). This was reflected in a significant main effect of duration in the ANOVA
 347 ($P < 0.0001$, $F(1.21, 29.06) = 221.05$, $\eta^2 = 0.9$). Face expression had a weak ef-
 348 fect, with angry faces categorized less accurately than both happy and neutral faces
 349 ($P = 0.046$, $F(1.95, 46.71) = 3.33$, $\eta^2 = 0.12$), and with no significant interaction effect
 350 ($P = 0.23$, $F(1.74, 41.83) = 1.53$, $\eta^2 = 0.06$).

351 Participants found the task challenging, as reflected in the perceptual awareness rat-
 352 ings: 84.5% of the 10 ms trials were rated as not containing a face (Figure 3E). This
 353 suggests that participants were complying with the task with respect to both expression
 354 discrimination and perceptual rating. Importantly, for faces presented for 10 ms, there was
 355 no difference in accuracy between expressions ($P = 0.43$, $F(1.65, 39.5) = 0.8$) or between
 356 any pair of cells in the confusion matrix ($P = 0.6$, $F(3.42, 82.07) = 0.64$), suggesting that
 357 faces presented at this duration were equally likely to be categorized as any expression.

358 **Spatiotemporal dynamics of face perception**

359 To investigate face processing as a function of stimulus duration, we performed within-
360 subject cross-identity decoding of responses to faces vs. scrambled stimuli. The analysis
361 included three components: sensor-level time-resolved classification to evaluate the pro-
362 gression of condition-related information; sensor-level temporal generalization to assess
363 the temporal structure of this information; and source-space decoding to obtain spatial
364 information about subliminal responses to faces (Figure 1).

365 We first decoded responses to neutral faces vs. scrambled stimuli using data from all
366 MEG sensors, separately for each stimulus duration. In the case of faces presented for
367 10 ms, any trials reported as containing a face were excluded, to ensure that we assessed
368 responses outside of subjective awareness. Scrambled stimuli could be discriminated from
369 faces presented for 150 and 30 ms starting as early as 100 ms (Figure 4A). After the initial
370 peak in performance, decoding accuracy decreased, but remained well above chance for
371 the remainder of the decoding time window. For faces presented for 10 ms and reported
372 as not perceived, there was only a weak increase in decoding performance, which reached
373 significance at 147 ms and dropped back to chance level after ~350 ms (Supplementary
374 Table 1). To assess how well face representations generalized across stimulus durations, we
375 repeated this analysis by training and testing on stimulus exemplars presented for different
376 amounts of time (Figure 4B). Decoding accuracy was high when cross-decoding between
377 30 ms and 150 ms faces, with two increases in performance at M170 latencies (100-200
378 ms) and after 300 ms. On the other hand, representations only generalized to 10 ms faces
379 for a limited time window corresponding to the M170 component.

380 Using temporal generalization decoding (King & Dehaene, 2014), we investigated the
381 temporal structure of face responses, and we found that this changed with stimulus dura-
382 tion. For faces presented for 150 ms, successful temporal generalization started at ~93 ms
383 in a diagonal pattern suggestive of transient representations, with more sustained repre-
384 sentations (square patterns) arising at M170 latencies and after 300 ms (Figure 4D-E). For
385 30 ms stimuli, a transient representation pattern started at ~110 ms after stimulus onset
386 and sustained representations only arose later (~400 ms). Early processing thus appears
387 to be heavily biased by stimulus presentation duration, with 30 ms faces failing to elicit a
388 stable representation at M170 latencies. For faces presented for 10 ms, only few transient
389 clusters survived correction for multiple comparisons, with the largest one occurring after
390 200 ms.

391 Finally, we spatially localized the subliminal response to faces in source space. All par-

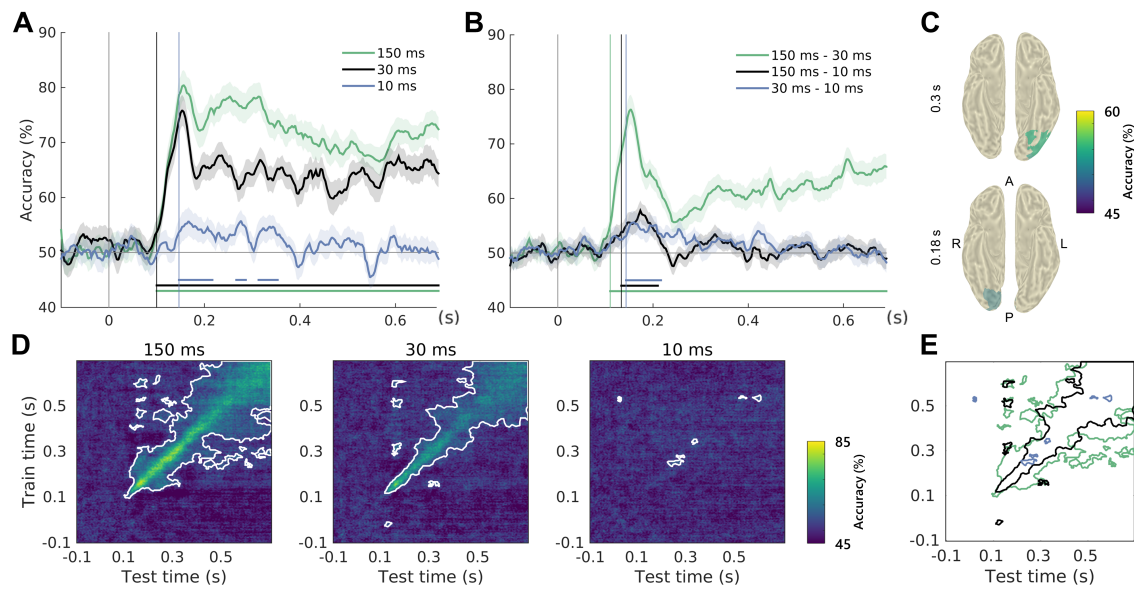


Figure 4: Face vs. scrambled decoding results. **A.** Sensor-space time-resolved decoding accuracy for all stimulus durations. Colour-coded vertical bars mark above-chance decoding onset and horizontal lines show significant time windows ($P < 0.05$, corrected). **B.** Sensor-space time-resolved cross-decoding for all pairs of stimulus durations. Cross-validation was performed across exemplars and accuracies were averaged over the two training/test directions. **C.** Sources achieving above-chance decoding of 10 ms faces outside awareness at M170 latencies in source space ($P < 0.005$, corrected). **D.** Sensor-space temporal generalization accuracy and significant clusters (white contours; $P < 0.05$, corrected) for all stimulus durations. **E.** Significant temporal generalization clusters for all three stimulus durations, showing more sustained representations of faces presented for 150 ms (legend as in A).

392 ticipants with one exception acquired a structural MRI, which was used to source localize
393 the MEG data using a Linearly Constrained Minimum Variance (LCMV) beamformer
394 (Van Veen et al., 1997). We performed whole-brain searchlight classification of 10 ms
395 faces vs. scrambled stimuli ($N=24$), using source clusters with a radius of 10 mm and
396 time windows of 30 ms. Faces were successfully decoded in a right occipital area at M170
397 latencies (Figure 4C), with a later stage associated with ventral patterns.

398 **Temporal dynamics of expression perception**

399 Next, we performed sensor-level cross-identity decoding of all pairs of emotional expres-
400 sions separately for each stimulus duration. The analysis was performed similarly to the
401 time-resolved face decoding analysis described above.

402 The highest decoding performance was achieved on late responses to expressions pre-
403 sented for 150 ms (Figure 5A). Expressions presented for 30 ms also achieved above-chance
404 decoding, although these effects were more transient. We also performed this analysis on
405 pooled datasets (faces presented for 30 and 150 ms), as the face cross-decoding analysis
406 showed that responses generalized between these two categories (Figure 4B). This revealed
407 a multi-stage progression for all expressions, with transient early decoding at M100 laten-
408 cies and an increasing accuracy at later stages (Figure 5B). We found no above-chance
409 performance when decoding 10 ms expressions. This finding adds to emerging evidence
410 against the automatic processing of expression outside awareness (Koster, Verschuere,
411 Burssens, Custers, & Crombez, 2007; Pessoa, Japee, & Sturman, 2006; Hedger, Gray,
412 Garner, & Adams, 2016; Schlossmacher, Junghöfer, Straube, & Bruchmann, 2017), and
413 we explore potential reasons for this result below.

414 **Face representations in occipitotemporal cortex**

415 To interrogate the neural representations underpinning these pattern differences, we per-
416 formed representational similarity analysis (RSA) using a searchlight approach at the
417 source level (Su, Fonteneau, Marslen-wilson, & Kriegeskorte, 2012) in a face-responsive
418 area of interest determined using an orthogonal contrast. We investigated the representa-
419 tional dynamics of face perception by assessing the similarity between MEG patterns and
420 models quantifying behaviour, face features, face configuration, expression, identity and
421 visual properties, using both a Spearman's rank correlation (Nili et al., 2014) and partial
422 correlation (see Methods).

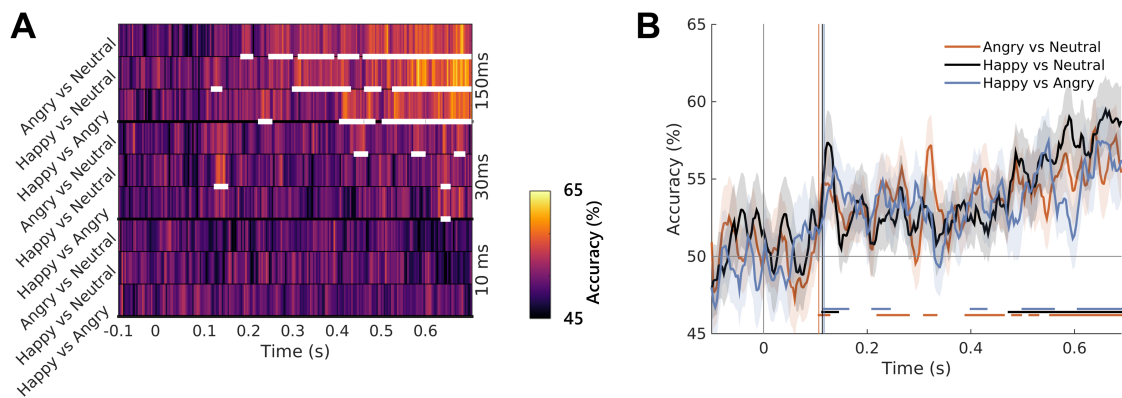


Figure 5: Expression decoding results. **A.** Time-resolved decoding accuracy for the three expression decoding problems and the three stimulus durations. White horizontal lines show significant time windows ($P < 0.05$, corrected). **B.** Time-resolved accuracy for the three expression decoding problems using the pooled datasets (30 + 150 ms).

423 Occipitotemporal cortex encodes behavioural responses

424 Among the other model RDMs tested, behavioural RDMs correlated most with the high-
425 level expression models (particularly the angry-vs-others model at 30 ms and 150 ms,
426 Spearman's $\rho = 0.29$ and $\rho = 0.34$). At 150 ms, the behavioural RDM also correlated with
427 the configural face models ($\rho = 0.22$ and $\rho = 0.18$). As expected based on performance,
428 behavioural RDMs at 10 ms did not correlate with the other two ($\rho = -0.05$ and $\rho =$
429 -0.09 respectively), while behavioural RDMs at 30 and 150 ms were positively correlated
430 ($\rho = 0.38$; Figure 6A).

431 Based on these links, face configuration, together with facial expression, appears to
432 partially explain behavioural responses. To test this, we performed a variance partitioning
433 analysis, using hierarchical multiple regression to quantify the unique and shared variance
434 in behaviour explained by facial configuration and high-level expression models. In the
435 10 ms condition, the neutral-vs-others model and the two configural models explained
436 25.1% of the variance; in the 30 ms and 150 ms conditions, the angry-vs-others model and
437 the configural models explained up to 45.7% of the variance in behaviour. Furthermore,
438 while the expression model contributed most of the variance, over 75% of this variance
439 was shared with the configural models. The unique contribution of the configural models
440 increased with stimulus duration (from ~2% at 10 ms, to ~20% at 150 ms). Together,
441 these results point to the role of face configuration in driving high-level representations
442 and behaviour. Note that we were unable to decode 10 ms expressions from the MEG
443 data; however, the variance partitioning analysis of behavioural responses in this condition
444 showed a contribution of both facial expression and configuration to behaviour.

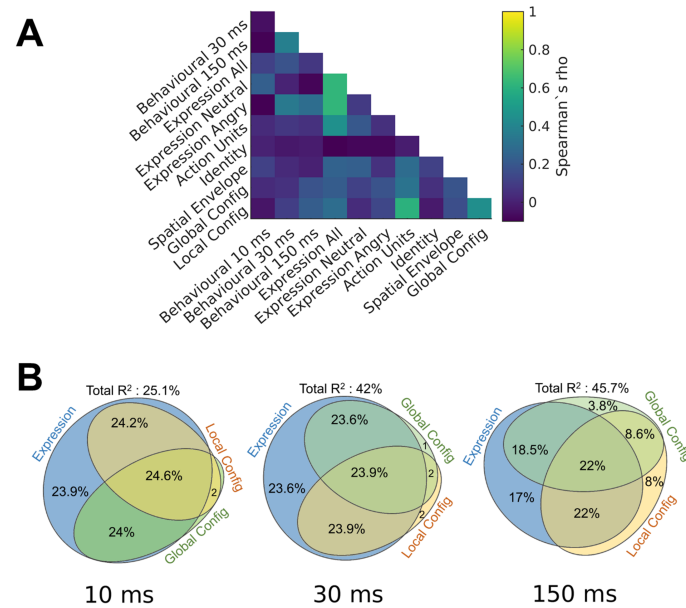


Figure 6: Relating behaviour to representational models. **A.** Model inter-correlations (Spearman's ρ). **B.** Variance partitioning results, showing the contributions of expression and face configuration models to behavioural responses at each stimulus duration. Values represent % of the total R^2 .

445 Behavioural RDMs showed the strongest and most sustained correlations with MEG
 446 patterns in ventral stream areas, including sources corresponding to the location of the
 447 fusiform face area (FFA) and occipital face area (OFA; Figure 7). Behavioural repre-
 448 sentations evolved differently in time for the three stimulus durations. For 10 ms faces,
 449 behaviour explained the data starting at 120 ms until the end of the analysis time window.
 450 Representations emerged similarly early for 150 ms faces and reached the noise ceiling be-
 451 fore decreasing again at 400 ms. For 30 ms faces, correlations were significant starting at
 452 210 ms in a relatively focal right temporal area. Patterns were more posterior for 10 ms
 453 faces and more extensive, including sources corresponding to the OFA and FFA, for 150
 454 ms faces.

455 The correlation time-courses suggest interesting differences in processing as a function
 456 of the information available: for clearly perceived faces, features relevant in behaviour
 457 are extracted between 120-400 ms, while behavioural responses for briefly presented faces
 458 appear to require sustained processing, as reflected by behaviour-related correlations not
 459 dropping back to zero. These results are in line with previous evidence of behavioural
 460 representations in ventral stream areas in scene and object perception (Walther, Caddigan,
 461 Fei-Fei, & Beck, 2009), and suggest that visual feature processing, even at early stages, is
 462 closely linked to behavioural goals.

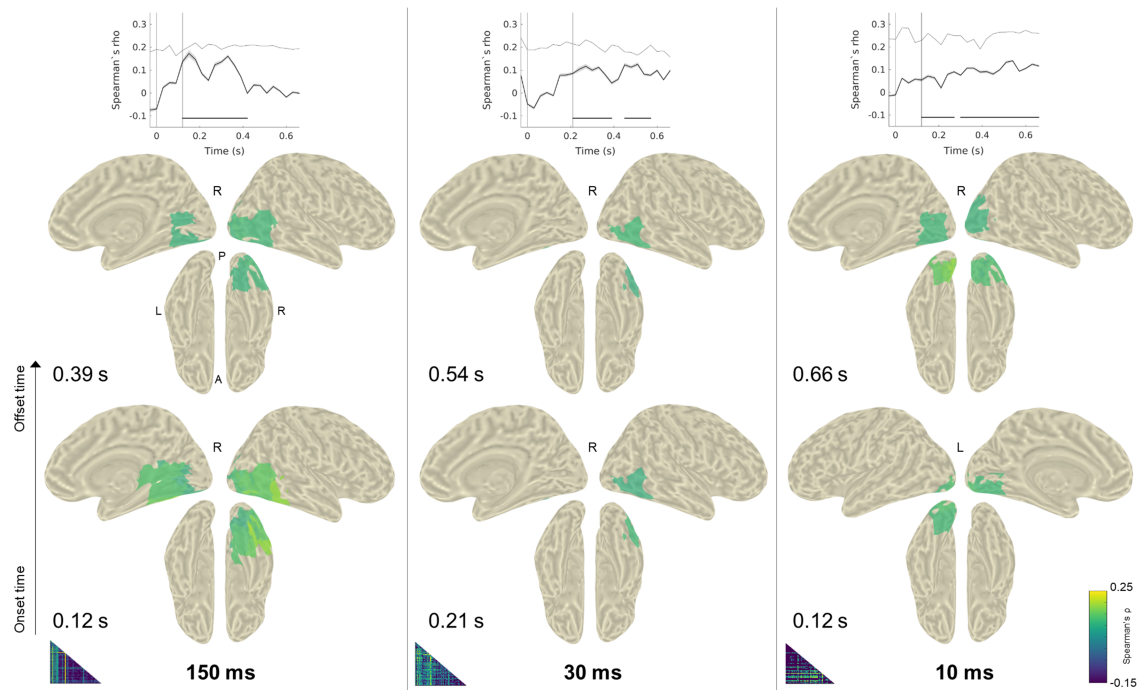


Figure 7: Correlations between MEG patterns and behavioural model RDMs for each stimulus condition duration (vertical columns). The top panels show correlation time-courses averaged across all significant searchlights; the noise ceiling is shown as a dotted horizontal line and is only approached in the 150 ms condition. The cortical maps show significant correlation coefficients for the first and last significant time windows (onset and offset times) on the inflated template MNI brain. The hemisphere shown is indicated with the letter R/L. Model RDMs are shown in the lower left corner of each column. See SourceMovies1 for movies showing the evolution of behavioural representations in time.

463 **Configural face processing from featural to relational**

464 The two face configuration models were also represented in the MEG patterns. In the
465 correlation analysis, the local and global configuration models explained representations
466 in partially overlapping areas of the ventral stream (corresponding to the right FFA lo-
467 cation), with local configuration representations arising earlier (at 120 ms for 150 ms
468 faces, and 360 ms for 30 ms faces). The RSA method used here favoured sustained cor-
469 relations over transient peaks; note that the global configuration model approached the
470 noise ceiling during a transient time window at M170 latencies for both 150 ms and 30
471 ms faces, suggesting a contribution of second-order characteristics, although this occurred
472 later than feature representations (Supplementary Figure 4). The partial correlation anal-
473 ysis revealed further differences between conditions: for 150 ms faces, the local and global
474 models made unique, successive contributions in explaining the data; conversely, for 30
475 ms faces we detected no unique contributions, suggesting that the extraction of configural
476 information from faces occurs differently in the absence of sufficient information. None of
477 the models significantly correlated with MEG patterns elicited by 10 ms faces.

478 Note that although both internal (eyes, nose, mouth) and external (face shape, hair)
479 face features have been shown to contribute to neural responses to faces (Axelrod, 2010),
480 we focus here on internal features; for the purposes of this paper, external features were
481 excluded from the stimuli and we refer to the second-order configuration of distances
482 between internal features as "global configuration". Internal features are relevant to the
483 context of expression discrimination and have been shown to be more reliable even in facial
484 recognition contexts (e.g. Longmore, Liu, and Young, 2015).

485 **Transient representations of visual and high-level models**

486 Two other models elicited brief representations in the MEG data. For 150 ms faces, the
487 spatial envelope model explained left hemisphere occipital representations starting at ~400
488 ms, suggesting sustained processing of visual features, potentially based on feedback mech-
489 anisms. For 30 ms faces, a high-level expression model (neutral-vs.-others) was represented
490 in the MEG data starting at 300 ms (Figure 9). This can be speculatively explained by
491 the formation of task-related representations in the absence of sufficient information. On
492 the contrary, when faces are clearly presented, only models encoding face characteristics
493 are represented, while categorical models show no contribution to occipitotemporal repre-
494 sentations. Note that despite the role of facial features in explaining neural responses, the
495 Action Unit model RDM did not significantly correlate with the MEG patterns, probably

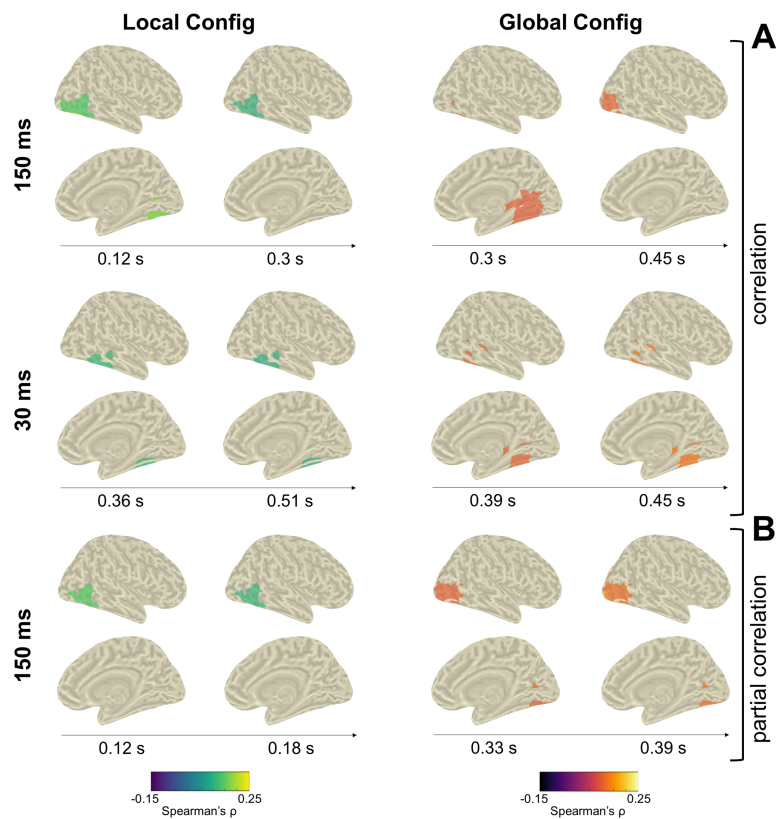


Figure 8: Significant correlations between MEG patterns and configural model RDMs. **A:** Correlation analysis results are significant for the 150 ms and 30 ms conditions. **B:** Partial correlation results are significant for the 150 ms condition. Only right hemisphere searchlights correlate with the configural models. Maps are shown for the onset and offset times of significant correlation. See SourceMovies2 for movies showing the evolution of behavioural representations in time.

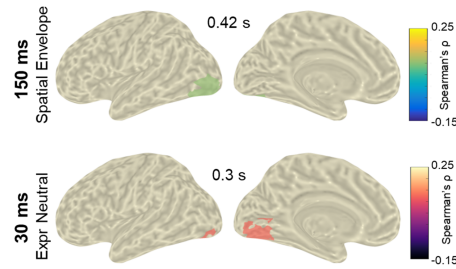


Figure 9: Significant correlations between: (1) MEG patterns for the 150 ms condition and the spatial envelope model RDM (**top**); (2) MEG patterns for the 30 ms condition and the high-level neutral-vs-others model (**bottom**). Only left hemisphere searchlights correlate with the two models. Maps are shown for the onset time of significant correlation, as clusters are sustained until offset (top: 0.54 s, bottom: 0.36 s).

496 due to the static and brief nature of our stimuli.

497 Although correlation coefficients between the models and neural data are generally low
498 (Supplementary Table 3), the noise ceiling shows that the maximal correlation possible
499 with our data is also low (mean $\rho=0.21$); this is not surprising, considering the low ρ -values
500 usually found in MEG RSA studies, and the fact that our paradigm involved complex, high-
501 level visual stimuli and a demanding task. In this case, the noise ceiling serves as a useful
502 benchmark for the explanatory power of our models. For example, the behavioural RDM
503 reaches the noise ceiling in the 150 ms condition, but not for briefer stimuli, suggesting
504 that behavioural representations fully explain the data when stimuli are clearly perceived.
505 The local configuration model also shows good explanatory power at its earliest stage, and
506 the same is true for the global model for a brief time window. Other significant models do
507 not reach the noise ceiling (Supplementary Figure 4); given the complex face processing
508 and task-related activity reflected by the MEG patterns, this is not surprising. In fact, the
509 explanatory power of the configural models at early stages (100-200 ms) is striking, as is
510 the strength of behavioural representations in ventral stream within 400 ms. Furthermore,
511 the initial peak in performance of the behavioural model overlaps with the peak of the local
512 configuration model. Together with the shared variance between configuration, expression
513 and behaviour shown in the variance partitioning analysis (Figure 6D), this points to the
514 role played by facial configuration in the extraction of emotional cues essential in the
515 expression discrimination task.

516 Discussion

517 The cross-identity decoding and representational similarity analyses described here con-
518 verge to highlight the dynamic nature of face representations in the ventral visual stream.

519 Face feature and face configuration representations link occipitotemporal neural patterns
520 and behavioural responses during an expression discrimination task, while their temporal
521 dynamics change to accommodate challenging viewing conditions.

522 In the time-resolved decoding analysis, a response to faces (150 ms and 30 ms) emerged
523 at ~100 ms, while faces shown outside of subjective awareness were decodable for a brief
524 time window (147 - 350 ms), in line with previous studies showing evidence of face per-
525 ception outside of awareness (Axelrod, Bar, & Rees, 2015). Temporal representations also
526 varied with stimulus duration: for 150 ms faces, a sustained representation emerged at
527 M170 latencies which was absent for 30 ms faces. This suggests that clearly presented faces
528 are perceived through a multi-stage process, while disrupted recurrent processing leads to
529 delayed stable representations. Although the M170 component decreases in amplitude
530 with face duration (Supplementary Figure 1), its duration does not predict such a marked
531 change in temporal structure, especially given the high decoding accuracy at this latency
532 obtained in both conditions in the time-resolved face decoding analysis. Trial-to-trial
533 variability, cited as another potential explanation for diagonal patterns (Vidaurre, Myers,
534 Stokes, Nobre, & Woolrich, 2018), is also not expected to systematically vary between
535 our conditions. On the other hand, sustained representations in temporal generalization
536 analyses are thought to be reflective of conscious perception and recurrent processes (De-
537 haene, 2016). It has previously been suggested that faster stimulus presentation leads to
538 more transient representations (Mohsenzadeh, Qin, Cichy, & Pantazis, 2018); however, the
539 backward masking procedure used here disrupts the formation of a stable representation
540 by entering the visual stream, and it is unclear whether different methods of preventing
541 awareness would lead to the same results.

542 Information supporting face decoding outside of subjective awareness was localized
543 to occipitotemporal cortex in our searchlight source-space decoding analysis (Figure 4C).
544 Given the disruption of recurrent processing in backward masking (Lamme, Zipser, &
545 Spekreijse, 2002; Boehler, Schoenfeld, Heinze, & Hopf, 2008), the early stages of this
546 response can be attributed to either purely feedforward activity, or to feedback connections
547 targeting V1 at early processing stages (Wyatte, Jilk, & O'Reilly, 2014; Mohsenzadeh et
548 al., 2018). Furthermore, the fact that we detect a response to faces, and not to expression,
549 suggests that the different tasks of identification and categorization are supported by
550 qualitatively different mechanisms. However, the spatial resolution of MEG, together
551 with recent observations of information spreading in searchlight source-space MVPA (Sato,
552 Yamashita, Sato, & Miyawaki, 2018), prevent us from drawing strong conclusions about

553 the origin of this response to faces. To minimize such concerns, we restricted our source-
554 space decoding analysis to localizing effects identified at the sensor level, and we applied
555 randomization testing with an omnibus threshold in order to avoid spurious effects.

556 All expressions presented for at least 30 ms were decodable from MEG data starting
557 at ~100 ms. Since all analyses were performed across facial identity and stimuli were
558 matched for low-level properties, this suggests that expression categorization begins at
559 the early stages of visual perception (Aguado et al., 2012; Dima, Perry, Messaritaki,
560 Zhang, & Singh, 2018), in line with behavioural goals. However, in terms of non-conscious
561 expression processing, the results are mixed. Despite the absence of a subliminal expression
562 effect in MEG responses, behavioural data suggest that expression (specifically, a model
563 differentiating between emotional and neutral stimuli) explains approximately one quarter
564 of the variance in behavioural responses given to faces presented for 10 ms. This effect is
565 not revealed by the more traditional accuracy-based behavioural analysis, suggesting that
566 model-based approaches to the analysis of behavioural responses can provide additional
567 information. With the caveat that low numbers of trials were included in this analysis,
568 the fact that cross-subject patterns of response reflected shared variance between the
569 models based on expression, facial features and facial configuration points to a certain
570 degree of expression processing taking place outside of subjective awareness. The absence
571 of a subliminal expression effect in the neural data may be explained by several factors,
572 including the limited ROI used in RSA, the study design minimizing residual awareness,
573 and challenges in the detection of a potential subcortical response.

574 Representational similarity analysis results linked stages in time-resolved decoding to
575 stages in feature extraction and to behavioural responses. Ventral stream areas encoded
576 sustained and extensive behavioural representations as early as 120 ms after stimulus onset
577 (Figure 7), suggesting that the extraction of features essential in behavioural decision-
578 making is a rapid process accomplished in face-responsive cortex. This is in line with
579 evidence found in higher-level object and scene perception (Walther et al., 2009; Bankson,
580 Hebart, Groen, & Baker, 2018; Groen et al., 2018) and with previous studies showing that
581 the perceptual similarity of faces is represented in neural patterns (Said, Moore, Engell,
582 & Haxby, 2018; Furl et al., 2017).

583 Furthermore, ventral stream areas encoded facial features prior to facial configuration
584 when faces were presented for 150 ms. This adds to evidence suggesting that emotional
585 face perception is supported by the processing of diagnostic features, such as the eyes
586 and mouth (Wegrzyn, Vogt, Kireclioglu, Schneider, & Kissler, 2017). What is more,

587 configural representations explain behaviour and overlap with behavioural representations,
588 suggesting that it is face configuration that drives expression-selective responses in ventral
589 stream areas and guides behaviour.

590 Previous studies have shown differential modulation of ERP components by first-order
591 and second-order face configuration. Some studies have shown early components (P1,
592 N170) to encode the former only (Mercure, Dick, & Johnson, 2008; Zion-Golumbic &
593 Bentin, 2007), while others have also shown effects of second-order configuration at N170
594 latencies (Eimer, Gosling, Nicholas, & Kiss, 2011). Furthermore, fMRI studies have re-
595 ported a division of labour in the face-selective network, with the FFA thought to play
596 a special role in representing both types of configural information (Golarai, Ghahremani,
597 Eberhardt, & Gabrieli, 2015). Recently, it has been suggested that featural and configural
598 processing of even non-face objects elicit face-like responses in the OFA and FFA (Zachar-
599 iou, Safiullah, & Ungerleider, 2018). Here, we combined the strengths of source-localized
600 MEG data and the RSA framework to tease apart the two models using a single stimulus
601 set. The searchlight RSA analysis revealed that the two models overlap spatially in a right
602 ventral stream area corresponding to the FFA, but are dissociated temporally: for 150 ms
603 faces, representations switch from first-order to second-order at ~300 ms after stimulus
604 onset, bringing together previous fMRI and electrophysiological findings.

605 Furthermore, this two-stage process appears to depend on the amount of information
606 available to the visual system. For 150 ms faces, local and global configuration models
607 make unique, temporally distinct contributions to explaining the data, as shown in the
608 partial correlation analysis. For 30 ms faces, no unique variance is explained by the
609 two models; furthermore, representations are temporally overlapping in the correlation
610 analysis and occur after 300 ms (Figure 8). This complements our sensor-level temporal
611 generalization findings: 30 ms faces are processed through a series of transient coding steps
612 at early stages and a stable representation is formed after 300 ms, when both first-order
613 and second-order features are represented. On the other hand, for 150 ms faces, a two-stage
614 process takes place, with an initial stable representation emerging at M170 latencies and
615 supported mainly by first-order features, and a later representation after 300 ms encoding
616 second-order configuration. Feature representations thus appear to be linked to the late
617 emergence of stable representations, thought to be reflective of recurrent processing and
618 categorization. Importantly, this idea is supported by spatially and temporally overlapping
619 behavioural representations in ventral stream areas.

620 The findings we present here constitute a stepping stone towards a better understand-

621 ing of high-level representations in face perception. While binary categorical models can
622 estimate high-level representations and task-related processing, the code supporting visual
623 perception is likely to be better understood in terms of behavioural goals and the visual
624 features supporting them. We show that face-responsive cortex dynamically encodes fa-
625 cial configuration starting with first-order features, and that this supports behavioural
626 representations when participants are performing an expression discrimination task. Fur-
627 thermore, we show that the cascade of processing stages changes with stimulus duration,
628 pointing to the adaptability of the face processing system in achieving goals with lim-
629 ited visual input. This highlights the importance of investigating neural computations
630 in a spatiotemporally resolved fashion; furthermore, when employing rapid presentation
631 paradigms, it is important to consider the changes in neural dynamics and stimulus repre-
632 sentations induced by relatively small changes in stimulus duration. Together, our results
633 bridge findings from previous fMRI and electrophysiological research, revealing the spa-
634 tiotemporal structure of face representations in human occipitotemporal cortex.

635 Acknowledgements

636 The authors would like to thank Lorenzo Magazzini and Gavin Perry for advice on the
637 study design, and Dimitrios Pantazis for helpful comments on the manuscript. The study
638 was supported by the UK MEG Partnership Grant (MRC/EP SRC, MR/K005464/1),
639 CUBRIC and the School of Psychology at Cardiff University.

640 Conflict of interest

641 The authors declare no competing interests.

642 References

- 643 Aguado, L., Valdés-Conroy, B., Rodríguez, S., Román, F. J., Diéguez-Risco, T., & Fernández-
644 Cahill, M. (2012). Modulation of early perceptual processing by emotional expression
645 and acquired valence of faces: An ERP study. *Journal of Psychophysiology*, *26*(1),
646 29–41. doi:10.1027/0269-8803/a000065
- 647 Axelrod, V. (2010). The Fusiform Face Area: In Quest of Holistic Face Processing. *Journal*
648 *of Neuroscience*, *30*(26), 8699–8701. doi:10.1523/JNEUROSCI.1921-10.2010
- 649 Axelrod, V., Bar, M., & Rees, G. (2015). Exploring the unconscious using faces. *Trends*
650 *in Cognitive Sciences*, *19*(1), 35–45. doi:10.1016/j.tics.2014.11.003
- 651 Baltrusaitis, T., Robinson, P., & Morency, L.-P. (2016). OpenFace: an open source fa-
652 cial behaviour analysis toolkit. *2016 IEEE Winter Conference on Applications of*
653 *Computer Vision (WACV)*, 1–10.

- 654 Bankson, B. B., Hebart, M. N., Groen, I. I., & Baker, C. I. (2018). The temporal evo-
655 lution of conceptual object representations revealed through models of behavior,
656 semantics and deep neural networks. *NeuroImage*, *178*, 172–182. doi:10.1016/J.
657 NEUROIMAGE.2018.05.037
- 658 Boehler, C. N., Schoenfeld, M. A., Heinze, H.-J., & Hopf, J.-M. (2008). Rapid recurrent
659 processing gates awareness in primary visual cortex. *Proceedings of the National
660 Academy of Sciences*, *105*(25), 8742–8747.
- 661 Calder, A. J., Young, A. W., Keane, J., & Dean, M. (2000). Configural information in facial
662 expression perception. *Journal of Experimental Psychology: Human Perception and
663 Performance*, *26*(2), 527–551. doi:10.1037/0096-1523.26.2.527
- 664 Chang, L. & Tsao, D. Y. (2017). The Code for Facial Identity in the Primate Brain. *Cell*,
665 *169*(6), 1013–1028. doi:10.1016/j.cell.2017.05.011
- 666 Dehaene, S. (2016). Decoding the Dynamics of Conscious Perception : The Temporal
667 Generalization Method. In B. G & C. Y (Eds.), *Micro-, meso- and macro-dynamics
668 of the brain* (pp. 85–97). New York: Springer. doi:10.1007/978-3-319-28802-4
- 669 Diamond, R. & Carey, S. (1986). Why faces are and are not special: an effect of expertise.
670 *Journal of experimental psychology*, *115*(2), 107–117.
- 671 Dima, D. C., Perry, G., Messaritaki, E., Zhang, J., & Singh, K. D. (2018). Spatiotemporal
672 dynamics in human visual cortex rapidly encode the emotional content of faces.
673 *Human Brain Mapping*, *39*(10), 3993–4006. doi:10.1002/hbm.24226
- 674 Dima, D. C., Perry, G., & Singh, K. D. (2018). Spatial frequency supports the emer-
675 gence of categorical representations in visual cortex during natural scene perception.
676 *NeuroImage*, *179*, 102–116. doi:10.1016/j.neuroimage.2018.06.033
- 677 Eimer, M., Gosling, A., Nicholas, S., & Kiss, M. (2011). The N170 component and its
678 links to configural face processing: A rapid neural adaptation study. *Brain Research*,
679 *1376*, 76–87. doi:10.1016/J.BRAINRES.2010.12.046
- 680 Ekman, P. & Friesen, W. (1977). *Facial action coding system: a technique for the mea-
681 surement of facial movement*. Palo Alto, CA: Consulting Psychologists Press.
- 682 Fan, R.-E., Chang, K.-W., Hsieh, C.-J., Wang, X.-R., & Lin, C.-J. (2008). LIBLINEAR:
683 A Library for Large Linear Classification. *Journal of Machine Learning Research*,
684 *9*(2008), 1871–1874. doi:10.1038/oby.2011.351
- 685 Farah, M. J., Wilson, K. D., & Tanaka, J. N. (1998). What Is “ Special ” About Face
686 Perception ? *Psychological review*, *105*(3), 482–498. doi:10.1037//0033-295X.105.3.
687 482
- 688 Freiwald, W., Duchaine, B., & Yovel, G. (2016). Face Processing Systems: From Neurons
689 to Real-World Social Perception. *Annual Review of Neuroscience*, *39*(1), 325–346.
690 doi:10.1146/annurev-neuro-070815-013934
- 691 Furl, N., Lohse, M., & Pizzorni-Ferrarese, F. (2017). Low-frequency oscillations employ a
692 general coding of the spatio-temporal similarity of dynamic faces. *NeuroImage*, *157*,
693 486–499. doi:10.1016/j.neuroimage.2017.06.023
- 694 Gohel, B., Lim, S., Kim, M.-Y., Kwon, H., & Kim, K. (2018). Dynamic pattern decoding of
695 source-reconstructed MEG or EEG data: Perspective of multivariate pattern analysis
696 and signal leakage. *Computers in Biology and Medicine*, *93*, 106–116. doi:10.1016/j.
697 complbiomed.2017.12.020
- 698 Golarai, G., Ghahremani, D. G., Eberhardt, J. L., & Gabrieli, J. D. E. (2015). Distinct rep-
699 resentations of configural and part information across multiple face-selective regions
700 of the human brain. *Frontiers in Psychology*, *6*, 1710. doi:10.3389/fpsyg.2015.01710
- 701 Greene, M. R., Baldassano, C., Esteva, A., Beck, D. M., & Fei-fei, L. (2016). Visual Scenes
702 are Categorized by Function. *Journal of Experimental Psychology: General*, *145*(1),
703 82–94. doi:10.1037/xge0000129.Visual

- 704 Grill-Spector, K., Weiner, K. S., Gomez, J., Stigliani, A., & Natu, V. S. (2018). The func-
705 tional neuroanatomy of face perception: From brain measurements to deep neural
706 networks. *Interface Focus*, 8. doi:10.1098/rsfs.2018.0013
- 707 Groen, I. I., Greene, M. R., Baldassano, C., Fei-Fei, L., Beck, D. M., & Baker, C. I.
708 (2018). Distinct contributions of functional and deep neural network features to
709 representational similarity of scenes in human brain and behavior. *eLife*, 7, e32962.
710 doi:10.7554/eLife.32962
- 711 Grootswagers, T., Wardle, S. G., & Carlson, T. A. (2017). Decoding Dynamic Brain Pat-
712 terns from Evoked Responses: A Tutorial on Multivariate Pattern Analysis Applied
713 to Time Series Neuroimaging Data. *Journal of Cognitive Neuroscience*, 29(4), 677–
714 697. doi:10.1162/jocn.1001068
- 715 Guggenmos, M., Sterzer, P., & Cichy, R. M. (2018). Multivariate pattern analysis for
716 MEG: A comparison of dissimilarity measures. *NeuroImage*, 173, 434–447. doi:10.
717 1016/J.NEUROIMAGE.2018.02.044
- 718 Harris, A. M. & Aguirre, G. K. (2008). The effects of parts, wholes, and familiarity on
719 face-selective responses in MEG. *Journal of Vision*, 8(10), 4–4. doi:10.1167/8.10.4
- 720 Hedger, N., Gray, K. L. H., Garner, M., & Adams, W. J. (2016). Are visual threats priori-
721 tised without awareness? A critical review and meta analysis involving 3 behavioural
722 paradigms and 2696 observers. *Psychological Bulletin*, 142(9), 934–968.
- 723 Henriksson, L., Mur, M., & Kriegeskorte, N. (2015). Faciotopy-A face-feature map with
724 face-like topology in the human occipital face area. *Cortex*, 72, 156–167. doi:10.1016/
725 j.cortex.2015.06.030
- 726 Hillebrand, A., Barnes, G. R., Bosboom, J. L., Berendse, H. W., & Stam, C. J. (2012).
727 Frequency-dependent functional connectivity within resting-state networks : An atlas-
728 based MEG beamformer solution. *NeuroImage*, 59(4), 3909–3921. doi:10.1016/j.
729 neuroimage.2011.11.005
- 730 Hillebrand, A., Singh, K. D., Holliday, I. E., Furlong, P. L., & Barnes, G. R. (2005). A new
731 approach to neuroimaging with magnetoencephalography. *Human Brain Mapping*,
732 25(2), 199–211. doi:10.1002/hbm.20102
- 733 King, J.-R. & Dehaene, S. (2014). Characterizing the dynamics of mental representations
734 : the temporal generalization method. *Trends in Cognitive Sciences*, 18(4), 203–210.
735 doi:10.1016/j.tics.2014.01.002
- 736 Koster, E. H. W., Verschuere, B., Burssens, B., Custers, R., & Crombez, G. (2007). At-
737 tention for Emotional Faces Under Restricted Awareness Revisited : Do Emotional
738 Faces Automatically Attract Attention ? *Emotion*, 7(2), 285–295. doi:10.1037/1528-
739 3542.7.2.285
- 740 Lamme, V. A. F., Zipser, K., & Spekreijse, H. (2002). Masking Interrupts Figure-Ground
741 Signals in V1. *Journal of Cognitive Neuroscience*, 14(7), 1044–1053. doi:10.1162/
742 089892902320474490
- 743 Leopold, D. A., O’Toole, A. J., Vetter, T., & Blanz, V. (2001). Prototype-referenced
744 shape encoding revealed by high-level aftereffects. *Nature Neuroscience*, 4(1), 89–94.
745 doi:10.1038/82947
- 746 Longmore, C. A., Liu, C. H., & Young, A. W. (2015). The importance of internal facial
747 features in learning new faces. *Quarterly Journal of Experimental Psychology*, 68(2),
748 249–260. doi:10.1080/17470218.2014.939666
- 749 Maurer, D., Grand, R. L., & Mondloch, C. J. (2002). The many faces of configural process-
750 ing. *Trends in Cognitive Sciences*, 6(6), 255–260. doi:10.1016/S1364-6613(02)01903-
751 4
- 752 Mercure, E., Dick, F., & Johnson, M. H. (2008). Featural and configural face processing
753 differentially modulate ERP components. *Brain Research*, 1239, 162–170. doi:10.
754 1016/J.BRAINRES.2008.07.098

- 755 Micallef, L. & Rodgers, P. (2014). eulerAPE: Drawing Area-Proportional 3-Venn Diagrams
756 Using Ellipses. *PLoS ONE*, *9*(7), e101717. doi:10.1371/journal.pone.0101717
- 757 Mohsenzadeh, Y., Qin, S., Cichy, R. M., & Pantazis, D. (2018). Ultra-Rapid serial visual
758 presentation reveals dynamics of feedforward and feedback processes in the ventral
759 visual pathway. *eLife*, *7*(e36329). doi:10.7554/eLife.36329.001
- 760 Nichols, T. E. & Holmes, A. P. (2001). Nonparametric Permutation Tests For Functional
761 Neuroimaging : A Primer with Examples. *Human Brain Mapping*, *25*(15), 1–25.
762 doi:10.1002/hbm.1058
- 763 Nili, H., Wingfield, C., Walther, A., Su, L., Marslen-Wilson, W., & Kriegeskorte, N. (2014).
764 A Toolbox for Representational Similarity Analysis. *PLoS Computational Biology*,
765 *10*(4), e1003553. doi:10.1371/journal.pcbi.1003553
- 766 Oliva, A. & Torralba, A. (2001). Modeling the Shape of the Scene : A Holistic Repre-
767 sentation of the Spatial Envelope. *International Journal of Computer Vision*, *42*(3),
768 145–175. doi:10.1023/A:1011139631724
- 769 Oostenveld, R., Fries, P., Maris, E., & Schoffelen, J. M. (2011). FieldTrip: Open source
770 software for advanced analysis of MEG, EEG, and invasive electrophysiological data.
771 *Computational Intelligence and Neuroscience*, *2011*, 156869. doi:10.1155/2011/
772 156869
- 773 Perry, G. & Singh, K. D. (2014). Localizing evoked and induced responses to faces us-
774 ing magnetoencephalography. *European Journal of Neuroscience*, *39*(9), 1517–1527.
775 doi:10.1111/ejn.12520
- 776 Pessoa, L., Japee, S., & Sturman, D. (2006). Target Visibility and Visual Awareness Mod-
777 ulate Amygdala Responses to Fearful Faces. *Cerebral Cortex*, *16*(March), 366–375.
778 doi:10.1093/cercor/bhi115
- 779 Piepers, D. W. & Robbins, R. A. (2012). A review and clarification of the terms "holistic,"
780 "configural," and "relational" in the face perception literature. *Frontiers in Psychol-
781 ogy*, *3*, 1–11. doi:10.3389/fpsyg.2012.00559
- 782 Said, C. P., Moore, C. D., Engell, A. D., & Haxby, J. V. (2018). Distributed representations
783 of dynamic facial expressions in the superior temporal sulcus. *Journal of Vision*,
784 *10*(5), 1–12. doi:10.1167/10.5.11.Introduction
- 785 Sato, M., Yamashita, O., Sato, M.-a., & Miyawaki, Y. (2018). Information spreading by
786 a combination of MEG source estimation and multivariate pattern classification.
787 *PLOS ONE*, *13*(6), e0198806. doi:10.1371/journal.pone.0198806
- 788 Schlossmacher, I., Junghöfer, M., Straube, T., & Bruchmann, M. (2017). No differential
789 effects to facial expressions under continuous flash suppression: An event-related
790 potentials study. *NeuroImage*, *163*, 276–285. doi:10.1016/j.neuroimage.2017.09.034
- 791 Singh, K. D., Barnes, G. R., & Hillebrand, A. (2003). Group imaging of task-related
792 changes in cortical synchronisation using nonparametric permutation testing. *Neu-
793 roImage*, *19*, 1589–1601. doi:10.1016/S1053-8119(03)00249-0
- 794 Stelzer, J., Chen, Y., & Turner, R. (2013). Statistical inference and multiple testing correc-
795 tion in classification-based multi-voxel pattern analysis (MVPA): Random permu-
796 tations and cluster size control. *NeuroImage*, *65*, 69–82. doi:10.1016/j.neuroimage.
797 2012.09.063
- 798 Studebaker, G. (1985). A "rationalized" arcsine transform. *Journal of speech and hearing
799 research*, *28*, 455–462. doi:10.1044/jshr.2803.455
- 800 Su, L., Fonteneau, E., Marslen-wilson, W., & Kriegeskorte, N. (2012). Spatiotemporal
801 Searchlight Representational Similarity Analysis in EMEG Source Space. In *Sec-
802 ond international workshop on pattern recognition in neuroimaging spatiotemporal*.
803 doi:10.1109/PRNI.2012.26
- 804 Tottenham, N., Tanaka, J. W., Leon, A. C., McCarry, T., Nurse, M., Hare, T. a., ...
805 Nelson, C. (2009). The NimStim set of facial expressions: judgments from untrained

- 806 research participants. *Psychiatry research*, *168*(3), 242–9. doi:10.1016/j.psychres.
807 2008.05.006
- 808 Tzourio-Mazoyer, N., Landeau, B., Papathanassiou, D., Crivello, F., Etard, O., Delcroix,
809 N., . . . Joliot, M. (2002). Automated Anatomical Labeling of Activations in SPM
810 Using a Macroscopic Anatomical Parcellation of the MNI MRI Single-Subject Brain.
811 *NeuroImage*, *15*, 273–289. doi:10.1006/nimg.2001.0978
- 812 Van Veen, B., van Drongelen, W., Yuchtman, M., & Suzuki, A. (1997). Localization of brain
813 electrical activity via linearly constrained minimum variance spatial filtering. *IEEE*
814 *Transactions on Biomedical engineering*, *44*(9), 867–880. doi:10.1109/10.623056
- 815 Vidaurre, D., Myers, N., Stokes, M., Nobre, A. C., & Woolrich, M. W. (2018). Tempo-
816 rally unconstrained decoding reveals consistent but time-varying stages of stimulus
817 processing. *bioRxiv*, 260943. doi:10.1101/260943
- 818 Visconti Di Oleggio Castello, M., Wheeler, K. G., Cipolli, C., & Gobbini, I. (2017). Famil-
819 iarity facilitates feature-based face processing. *PLoS One*, *12*(6), e0178895. doi:10.
820 1371/journal.pone.0178895
- 821 Vrba, J. & Robinson, S. E. (2001). Signal processing in magnetoencephalography. *Methods*,
822 *25*(2), 249–271. doi:10.1006/meth.2001.1238
- 823 Walther, A., Nili, H., Ejaz, N., Alink, A., Kriegeskorte, N., & Diedrichsen, J. (2016).
824 Reliability of dissimilarity measures for multi-voxel pattern analysis. *NeuroImage*,
825 *137*, 188–200. doi:10.1016/j.neuroimage.2015.12.012
- 826 Walther, D. B., Caddigan, E., Fei-Fei, L., & Beck, D. M. (2009). Natural Scene Cate-
827 gories Revealed in Distributed Patterns of Activity in the Human Brain. *Journal of*
828 *Neuroscience*, *29*(34), 10573–10581. doi:10.1523/JNEUROSCI.0559-09.2009
- 829 Wegrzyn, M., Vogt, M., Kireclioglu, B., Schneider, J., & Kissler, J. (2017). Mapping the
830 emotional face . How individual face parts contribute to successful emotion recogni-
831 tion. *PLoS ONE*, *12*(5), 1–15.
- 832 Willenbockel, V., Sadr, J., Fiset, D., Horne, G. O., Gosselin, F., & Tanaka, J. W. (2010).
833 Controlling low-level image properties: The SHINE toolbox. *Behavior Research Meth-*
834 *ods*, *42*(3), 671–684. doi:10.3758/BRM.42.3.671
- 835 Wyatte, D., Jilk, D. J., & O’Reilly, R. C. (2014). Early recurrent feedback facilitates
836 visual object recognition under challenging conditions. *Frontiers in Psychology*, *5*,
837 674. doi:10.3389/fpsyg.2014.00674
- 838 Zachariou, V., Safiullah, Z. N., & Ungerleider, L. G. (2018). The Fusiform and Occipi-
839 tal Face Areas Can Process a Nonface Category Equivalently to Faces. *Journal of*
840 *Cognitive Neuroscience*, *30*(10), 1499–1516. doi:10.1162/jocn.1162/jocn.1162.01288
- 841 Zion-Golumbic, E. & Bentin, S. (2007). Dissociated Neural Mechanisms for Face Detection
842 and Configural Encoding: Evidence from N170 and Induced Gamma-Band Oscilla-
843 tion Effects. *Cerebral Cortex*, *17*(8), 1741–1749. doi:10.1093/cercor/bhl100

844

Appendix

	Sensor-space			Source-space
	150 ms	30 ms	10 ms	10 ms
Max % accuracy	82.3	76.8	56.8	59.62
SD (%)	13.6	14.18	9.3	8.35
Decoding onset (ms)	100	100	147	120-150

Supplementary Table 1: Face decoding results.

	Stimulus duration											
	150 ms			30 ms			10 ms			30 + 150 ms		
	A-N	H-N	A-H	A-N	H-N	A-H	A-N	H-N	A-H	A-N	H-N	A-H
Max % accuracy	61.9	63.1	60.76	57.79	58.49	58.12	56.62	55.86	55.87	60.48	60.21	59.74
SD (%)	8.57	6.78	9.34	10.91	9.92	10.38	10.88	9.11	13.66	9.04	10.52	13.41
Decoding onset (ms)	180	113	220	437	120	633	N/A	N/A	N/A	107	113	117
	Perceptual rating											
	2			1			0			2 + 1		
	A-N	H-N	A-H	A-N	H-N	A-H	A-N	H-N	A-H	A-N	H-N	A-H
Max % accuracy	59.55	62.54	64.03	56.56	56.88	56.63	57.64	55.32	56.01	60.43	62.25	60.24
SD (%)	12.24	11.6	10.82	12.1	13.63	13.21	14.46	10.24	12.47	11.95	12.07	12.25
Decoding onset (ms)	230	113	523	307	120	130	N/A	N/A	N/A	220	113	127

Supplementary Table 2: Expression decoding results.

<i>Model</i>	Behavioural	Expression		Spatial Envelope		Global Config		Local Config	
150 ms	ρ	ρ	ρ_{part}	ρ	ρ_{part}	ρ	ρ_{part}	ρ	ρ_{part}
Max rho	0.23	0.14	0.14	0.17	0.17	0.18	0.17	0.18	0.16
SD	0.12	0.03	0.03	0.06	0.06	0.09	0.08	0.07	0.06
Onset (ms)	120	N/A	N/A	420	390	300	330	120	120
Offset (ms)	390	N/A	N/A	540	540	450	390	300	180
30 ms	ρ	ρ	ρ_{part}	ρ	ρ_{part}	ρ	ρ_{part}	ρ	ρ_{part}
Max rho	0.17	0.14	0.14	0.13	0.12	0.15	0.13	0.14	0.12
SD	0.07	0.03	0.03	0.05	0.04	0.07	0.03	0.05	0.04
Onset (ms)	210	300	300	N/A	N/A	390	N/A	360	N/A
Offset (ms)	540	360	360	N/A	N/A	450	N/A	510	N/A
10 ms	ρ	ρ	ρ_{part}	ρ	ρ_{part}	ρ	ρ_{part}	ρ	ρ_{part}
Max rho	0.18	0.09	0.1	0.13	0.14	0.17	0.16	0.11	0.12
SD	0.04	0.03	0.03	0.04	0.05	0.04	0.04	0.03	0.04
Onset (ms)	120	N/A	N/A	N/A	N/A	N/A	N/A	N/A	180
Offset (ms)	660	N/A	N/A	N/A	N/A	N/A	N/A	N/A	240

Supplementary Table 3: RSA results for the 5 models achieving significant correlations.

AU Code	Facial Action Coding System Name
AU01	Inner brow raiser
AU02	Outer brow raiser
AU04	Brow lowerer
AU06	Cheek raiser
AU09	Nose wrinkler
AU10	Upper lip raiser
AU12	Lip corner puller
AU14	Dimpler
AU15	Lip corner depressor
AU17	Chin raiser
AU20	Lip stretcher
AU25	Lips part

Supplementary Table 4: Action Units (AU) used to create the Action Unit model RDM.

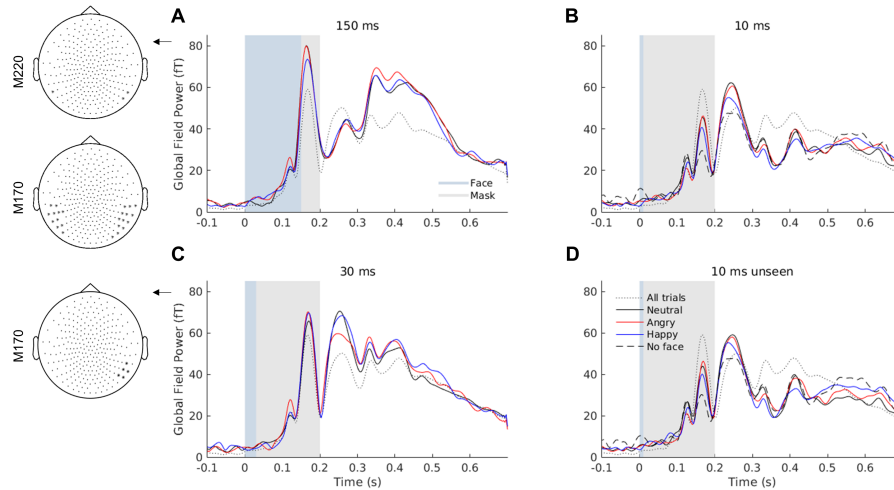
845 **Supplementary Analysis 1: Event-related field (ERF) analysis**

846 We assessed the presence of difference between conditions in event-related fields (ERF).
847 For the purposes of this analysis, MEG data were bandpass-filtered between 0.1 and 30 Hz
848 and axial gradiometer event-related fields were averaged across subjects to calculate the
849 global field power across all trials and conditions. This allowed us to determine three time
850 windows of interest for evoked response component analysis: 63-137 ms (M100), 137-203
851 ms (M170), and 203 – 306 ms (M220).

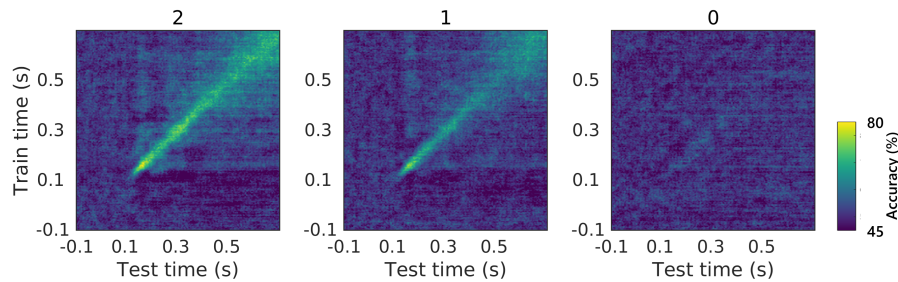
852 Next, we averaged evoked response fields for each condition and subject within the
853 three time windows. We tested for differences between responses to faces and scrambled
854 stimuli, and between responses to different emotional expressions, using paired t-tests and
855 repeated-measures ANOVAs respectively at each sensor and time window. Significant
856 sensors were determined using randomization testing (5000 iterations) and corrected for
857 multiple comparisons using the maximal statistic distribution ($\alpha = 0.001$).

858 We assessed the presence of a response to faces by contrasting neutral faces with
859 scrambled stimuli at each stimulus duration. For 150 ms faces, we found significant
860 differences at M170 latencies and M220 latencies ($P < 0.0007, t(24) > 6.07$), but no
861 significant effects at M100 latencies surviving our alpha of 0.001 (only one occipital sen-
862 sor showed a non-significant effect with $P = 0.0059, t(24) = 4.89$). A significant, but
863 smaller, cluster of right temporal sensors was also found for 30 ms faces at M170 la-
864 tencies ($P < 0.0004, t(24) > 5.99$). No conclusive effects were found when contrasting
865 faces presented for 10 ms with their scrambled counterparts, regardless of whether trials
866 where a face was perceived were excluded or not ($P > 0.015, t(24) < 4.66$ across compar-
867 isons), and no effect of emotional expression was found at any of the stimulus durations

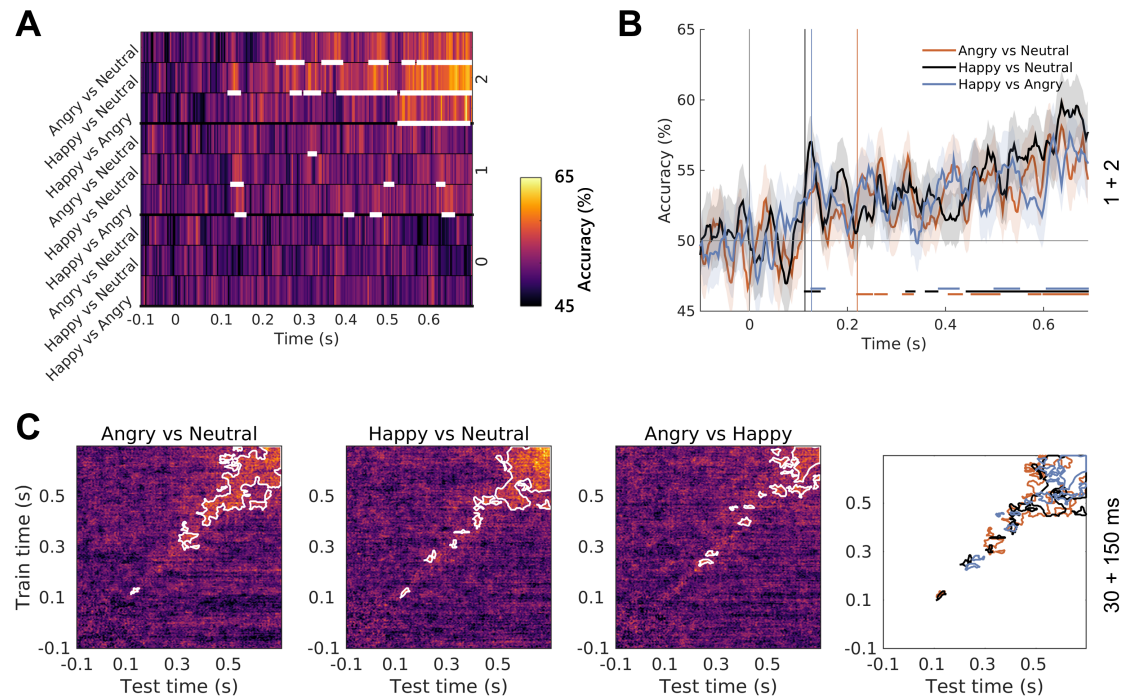
868 ($P > 0.06$, $F(2, 48) < 8.59$). Several factors could explain the absence of emotional expres-
869 sion effects in our ERF data: (1) stimuli were highly controlled for low-level properties,
870 minimizing visually-driven differences in early time windows; (2) our time windows of in-
871 terest did not include late stages dominated by task-related processing of expression; (3) we
872 performed a whole-brain analysis with a conservative correction for multiple comparisons.



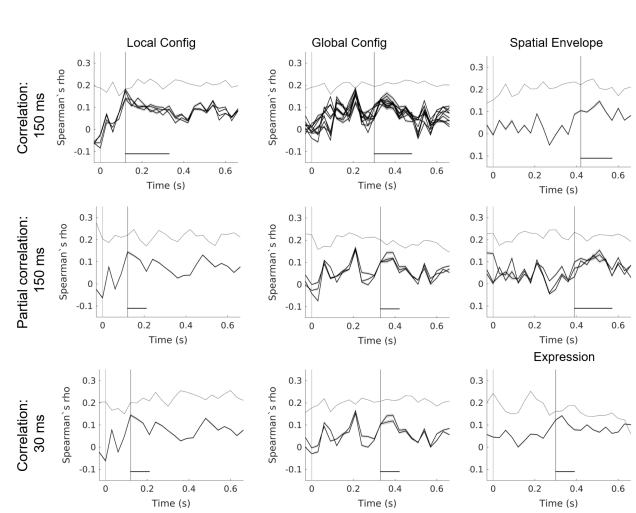
Supplementary Figure 1: ERF analysis results. **A-D**. Global field power averaged across participants and trials for each stimulus duration condition. Note decreasing M170 amplitudes with stimulus duration. **Left**. Significant sensors in the face vs scrambled contrast at M170 (137-203 ms) and M220 (203-306 ms) latencies ($P < 0.001$ corrected).



Supplementary Figure 2: Face vs scrambled temporal generalization decoding for each perceptual rating category. The same progression from stable to transient representations is observed as when datasets were split according to stimulus duration.



Supplementary Figure 3: Expression decoding. **A**. Time-resolved decoding accuracy for each pair of expressions and perceptual rating, with above-chance time-windows highlighted in white ($P < 0.05$ corrected). **B**. Accuracy time-courses obtained using pooled datasets (awareness ratings of 1 + 2). **C**. Temporal generalization accuracy and significant clusters (white contours; $P < 0.05$, corrected) for the three decoding problems using the pooled datasets (duration of 30 + 150 ms). The last panel shows significant temporal generalization clusters for all three decoding problems. Angry vs neutral decoding leads to earlier stable representations.



Supplementary Figure 4: Correlation time-courses obtained in the RSA analysis. All significant searchlights are plotted separately against a noise ceiling averaged across significant searchlights.