

# Genetic Algorithms Reveal Identity Independent Representation of Emotional Expressions

Thomas Murray<sup>1</sup>, Nicola Binetti<sup>2</sup>, Christina Carlisi<sup>3</sup>, Vinay Nambodiri<sup>4</sup>, Darren Cosker<sup>4</sup>,  
Essi Viding<sup>3</sup>, and Isabelle Mareschal<sup>5</sup>

<sup>1</sup>Department of Psychology, University of Cambridge

<sup>2</sup>Department of Cognitive Neuroscience, International School for Advanced Studies, Trieste, Italy

<sup>3</sup>Division of Psychology and Language Sciences, University College London

<sup>4</sup>Department of Computer Science, University of Bath

<sup>5</sup>Department of Psychology, Queen Mary University of London

People readily and automatically process facial emotion and identity, and it has been reported that these cues are processed both dependently and independently. However, this question of identity independent encoding of emotions has only been examined using posed, often exaggerated expressions of emotion, that do not account for the substantial individual differences in emotion recognition. In this study, we ask whether people's unique beliefs of how emotions should be reflected in facial expressions depend on the identity of the face. To do this, we employed a genetic algorithm where participants created facial expressions to represent different emotions. Participants generated facial expressions of anger, fear, happiness, and sadness, on two different identities. Facial features were controlled by manipulating a set of weights, allowing us to probe the exact positions of faces in high-dimensional expression space. We found that participants created facial expressions belonging to each identity in a similar space that was unique to the participant, for angry, fearful, and happy expressions, but not sad. However, using a machine learning algorithm that examined the positions of faces in expression space, we also found systematic differences between the two identities' expressions across participants. This suggests that participants' beliefs of how an emotion should be reflected in a facial expression are unique to them and identity independent, although there are also some systematic differences in the facial expressions between two identities that are common across all individuals.

*Keywords:* facial expressions, individual differences, emotion, identity

*Supplemental materials:* <https://doi.org/10.1037/emo0001274.supp>

Faces convey important social information including identity (Bruce & Young, 1998), emotional state (Adolphs, 2002; Michalek et al., 2022), age (Awad et al., 2020; Clifford et al., 2018), and even trustworthiness (Todorov et al., 2008). However, it remains unclear whether these facial cues are processed together as a holistic unit, or separately and recombined at a later stage of processing. For example, traditional models of face perception suggest that identity and expression are processed independently (Bruce & Young, 1986). These models are supported by behavioral evidence that the familiarity of a face does not affect expression matching (Bruce, 1986;

Young et al., 1986), and that expression processing is unaffected in prosopagnosia, a disorder that impairs identity recognition (Duchaine et al., 2003). Using brain imaging, several studies found that a region of the fusiform gyrus (the fusiform face area [FFA]) selectively responded to changes in identity, while a region in the superior temporal sulcus (STS) responded to changes in expression (Andrews & Ewbank, 2004; Eger et al., 2004), consistent with the proposal of a modular system of processing, with identity processed in the FFA, and expression in the STS (Haxby et al., 2000). However, more recent neuroimaging studies using multivariate analysis have challenged the suggestion of

This article was published Online First August 10, 2023.

Thomas Murray  <https://orcid.org/0000-0002-7314-465X>

This study was not preregistered. This study was funded by a Medical Research Council (MRC) Grant, MR/R01177X/1, awarded to Isabelle Mareschal, Essi Viding, and Darren Cosker. The authors have no conflict of interests to declare.

The GA toolkit is publicly available from <https://osf.io/dyfau/>.

The data from this study and code to reproduce analyses are available on the OSF page (<https://osf.io/t7a3w/>).

Thomas Murray served as the lead for formal analysis, visualization, and writing—original draft. Nicola Binetti served in a supporting role for methodology and writing—review and editing. Christina Carlisi served in a supporting role for writing—review and editing. Vinay Nambodiri served in a supporting role

for writing—review and editing. Essi Viding served in a supporting role for writing—review and editing. Isabelle Mareschal served as the lead for conceptualization, funding acquisition, supervision, and writing—review and editing. Vinay Nambodiri and Darren Cosker contributed equally to methodology and software. Darren Cosker and Essi Viding contributed equally to funding acquisition.

Open Access funding provided by University of Cambridge: This work is licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0; <https://creativecommons.org/licenses/by/4.0/>). This license permits copying and redistributing the work in any medium or format, as well as adapting the material for any purpose, even commercially.

Correspondence concerning this article should be addressed to Thomas Murray, Department of Psychology, University of Cambridge, Downing Place, Cambridge CB2 3EB, United Kingdom. Email: [tom29@cam.ac.uk](mailto:tom29@cam.ac.uk)

dissociable processing of identity and expression, with evidence that the FFA also represents expressions (Fox et al., 2009; Kawasaki et al., 2012; Murray et al., 2021; Sormaz et al., 2016; Wegryn et al., 2015; Zhang et al., 2016) and the STS also represents identity (Fox et al., 2009; Tsantani et al., 2019). More recent attempts to reconcile these results have proposed that faces may be represented both along dimensions that code identity and expression independently, as well as along common dimensions that encode both attributes together (Calder & Young, 2005; Rhodes et al., 2015).

Behavioral evidence for an interaction between representations of identity and expression comes from studies using adaptation paradigms. In these paradigms, sensory adaptation to stimuli can bias the perception of new ambiguous test stimuli away from the adapting stimulus. It is suggested that these perceptual changes (aftereffects) reflect changes in the activity of the neural populations that encode the stimuli (e.g., Georgeson, 2004). Such paradigms have been used to probe representations of identity and expression. For example, people perceive a particular identity (e.g., “Jim”) in a test face after they have adapted to faces containing opposite features to the target identities (e.g., “antifaces”; Leopold et al., 2001). Similarly, adaptation to sad faces can result in neutral faces appearing happy (Webster et al., 2004). These results show that adaptation is a powerful tool to examine identity and expression representations.

Rhodes et al. (2015) used adaptation to measure expression and identity aftereffects (i.e., the change in perception of expression and identity following adaptation) and found that the magnitude of the two effects was positively correlated, suggesting that faces may be represented within a common space that encodes both expression and identity. Several studies since have probed this directly by examining whether adaptation to a given expression depends on the adaptor and test being faces of the same identity. The expectation is that if identity and expression are encoded independently of each other, there should still be an aftereffect with different test and adaptor identities. Interestingly, most studies found that expression aftereffects persisted, although the magnitude of the aftereffects was significantly reduced when using different identities between adaptor and test (Campbell & Burke, 2009; Ellamil et al., 2008; Fox & Barton, 2007; Skinner & Benton, 2012; Song et al., 2015). This has been interpreted as evidence for two types of neural populations involved in the coding of facial expressions; one that encodes expressions independently from identity information (i.e., a neural representation of expression that is invariant to identity), and another population that is dependent on identity information (Campbell & Burke, 2009; Fox & Barton, 2007). To account for this, it has been suggested that there may exist an initial stage of visual processing that jointly encodes identity and expression, followed by further stages of (separate) processing of expression and identity (Palermo et al., 2013, 2018; Rhodes et al., 2015). An initial stage of visual processing that encodes multiple facial attributes could also account for the fact that the perception of expression can be affected by other facial attributes, such as age, race, and gender. For example, faces with a darker skin tone, or more masculine features, are more readily perceived as angry than faces with lighter skin or more feminine features (Adams et al., 2015; Becker et al., 2007; Hugenberg & Bodenhausen, 2003). The suggestion of an initial stage of joint processing of facial cues, including identity, is consistent with the seemingly contradicting neuroimaging evidence for both dissociable (Andrews & Ewbank, 2004; Winston et al., 2004) and overlapping (Baseler et al., 2014; Fox et al., 2009; Ganel

et al., 2005; Xu & Biederman, 2010) representations of identity and expression within the face processing network.

One difficulty in reconciling some of the above studies is that most use stimuli of posed, stereotypical expressions that do not reflect the types of facial expressions we naturally encounter. We recently found considerable individual differences in people’s representations of emotions which reflect how they think an emotion should be reflected in a facial expression. Importantly, these individual differences in emotion representations accounted for differences in their performance on subsequent emotion recognition tasks (Binetti et al., 2022). It is possible therefore that the results of the previous behavioral and neuroimaging experiments may have probed common representations of *posed, stereotypical* expressions which may not often match the unique representations held by the perceiver. Therefore, it is important to account for individual differences in expression recognition when examining whether identity and expression are processed independently.

We recently developed and validated, a genetic algorithm (GA) toolkit that allows users to create photorealistic facial expressions that they believe correspond to a given target emotion (Binetti et al., 2022; Carlisi et al., 2021; Roubtsova et al., 2021). The GA controls changes in facial stimuli by manipulating a set of 149 weights on facial features of a computer avatar, which are evolved over several iterations to converge on the set of weights that represents an expression tailored to the individual. These expressions, therefore, exist as points within an expression space (where each dimension corresponds to the weight given to facial features), and, like previous research using an expression space (e.g., Jack et al., 2012; Skerry & Saxe, 2015), can serve as proxies for how the participant internally represents expressions of emotion. In the present study, participants evolved facial expressions rendered onto two different identity avatars, to determine whether their representations of emotions differed when generated on the two different identity faces. Since the GA toolkit results in expressions that are defined by 149 weights, we can quantify the extent of independence by directly comparing the weights between the two identity faces created for the same emotion.

## Method

### Participants

Eighty-five participants took part in this study, either online or in-person ( $N$  male = 34,  $N$  female = 50;  $N$  online = 54,  $N$  in-person = 31;  $M_{\text{age}} = 26.20$  [ $SD = 9.60$ ]; mean years of education = 14.77 [ $SD = 3.19$ ]). Based on simulations in Binetti et al. (2022), we established that a sample of 40 participants was sufficient to replicate (with respect to the emotional representations) the distribution of 350 participants. That is, 40 participants capture the variability within the population for emotion representations. However, as we were unsure of the variability there might be within a participant for different faces, we aimed to test 100. No exclusions were applied. We have previously demonstrated that the data acquired using the GA tool do not differ if collected online or in-person (Binetti et al., 2022; also in the [online supplemental materials](#)). Demographic data is missing for one (online) participant. Online participants were recruited via Prolific and completed the GA task through Google’s remote access client while liaising with the experimenter through the Prolific messaging system, and were paid £7.44 p/h. In-person participants were recruited from the QMUL participant pool (via the SONA system);

[www.sona-systems.com](http://www.sona-systems.com)) and took part in exchange for course credits. This study was approved by the QMUL ethics board (QMERC2019/81), and participants provided written consent. We report how we determined our sample size, all data exclusions (if any), all manipulations, and all measures in the study.

## GA Task

Participants evolved facial expressions representing the emotions of anger, fear, happiness, and sadness with our GA toolkit, using two different identity avatars (one male and one female) in a counter-balanced order. Details of the task procedure can be found in Binetti et al. (2022), and technical details for the underlying algorithm can be found in Roubtsova et al. (2021). Briefly, participants were given a target emotion, and used the toolkit to evolve expressions to resemble what they think the expression for the target emotion should look like (their “emotion representation”). Participants evolved the faces using a GA procedure over eight iterations (i.e., with seven evolutions in the algorithm) per emotion and identity, to converge on a final expression that the participant believes best displays the target emotion.

On the first iteration, participants viewed 10 random expressions (displayed by the same identity) on the screen, and were instructed to select the faces that they believed displayed the expressions resembling the target emotion, with no constraints on the number of selections (other than they had to choose at least one face). We note that the final expression that a participant converges on does not depend on the expressions shown in the first iteration (Binetti et al., 2022). After their selection, the participant indicated which of the selected faces best matched the target emotion. After each iteration, the selected expressions were evolved by the GA via processes of selection, crossbreeding, mutation, and replacement (Roubtsova et al., 2021), so that on the next iteration participants viewed new expressions that resulted from evolution of the selected expressions on the previous iterations.

The expressions displayed on the faces were controlled by manipulating a set of 149 “blendshape” weights. Manipulation of the blendshapes controls the position of the vertices of the three-dimensional mesh on which the texture of the avatar is displayed, thereby manipulating the position and shape of the facial features in the displayed expression. Importantly, while there were differences in the morphology of the two identities, the manipulation of the blendshape weights had a consistent effect on the expressions displayed by the two identities (Roubtsova et al., 2021). The GA takes the vectors of blendshape weights for the selected expressions on a given iteration and “breeds” them to generate the vectors of blendshapes that characterize the expressions displayed in subsequent iterations.

Forty-one of the blendshapes approximately correspond to facial movements defined by the Facial Action Coding system (FACs), and these 41 were the blendshapes used within the subsequent analyses. The remaining blendshapes controlled the symmetrical movement of the core 41 blendshapes, acted as “corrective” blendshapes controlling combinations of movements to ensure realistic facial actions, and were blendshapes that controlled head movement and gaze direction. For the purposes of the current study, the task enforced symmetry between the blendshapes on the left and right sides of the face, and the head and gaze position were fixed. On the final iteration, participants indicated which face was the best match to the target

emotion they were creating, and we called this their “preferred expression.” As such, we acquired a point estimate for each of the four expressions displayed by the two identities within the latent expression-space, for each participant. This space encompasses a range of plausible expressions, where each of the 41 dimensions corresponds to the weight of a particular facial action (as controlled by a single blendshape)—a point at the origin within this space represents a neutral face. Importantly, this space is not constrained to prototypical expressions, allowing participants to generate realistically plausible expressions that may not match the posed expressions commonly found within stimulus sets.

After completing the task for a given identity and emotion category, participants were asked to rate on a scale of 1–7 how satisfied they were that their final preferred expression matched the target emotion they were trying to create. For participants that completed the task in-person, screen size was 310 × 174 mm and faces were displayed in a window of size 165 × 123 mm, where each face was 25 × 52 mm (subtending 3.97° at a viewing distance of 75 cm). An example trial is provided in Figure S1 in the online supplemental materials.

## Analysis

We analyzed the data at four different levels to examine how emotion categories are represented on the different identity faces. Firstly, we used machine learning to see if a classifier could determine the identity of the face, within each of the four different emotion categories, based on the (final) preferred expressions created by participants. This analysis was conducted across participants, and so examined the presence of any systematic differences in the representation of the two identities. Secondly, we examined within participants, whether their preferred expressions were independent of identity, using representational similarity analysis (RSA) to measure whether the (within-participant) representations of faces in expression space were clustered according to identity or emotion. Thirdly, we examined whether the two identities for each emotion category were represented in approximately the same parts of expression space for each participant, by comparing the within- and between-participant distances between the two identities. Finally, we examined systematic differences in the weights given to individual blendshapes for the two identities using paired sample comparisons. We used cosine distance of the weights of the 41 core blendshapes to measure the pairwise similarity of faces as it provides a reliable metric for high-dimensional sparse vectors such as the blendshape representation (Roubtsova et al., 2021).

### *Between-Participant Analysis: Support Vector Machine (SVM) Classification*

We used a SVM to perform binary classification of the identity of the face from the vectors of core blendshape weights, separately for each emotion, using five-fold cross-validation to assess classification accuracy. In five-fold cross-validation, data is randomly split into five equal subsets, and classification accuracy for each (and every) subset is tested after training the model using the four remaining subsets. Overall classification accuracy of the model was determined by calculating the mean accuracy for classifying the identity of the faces across the five subsets of data. As SVMs are not scale invariant, the training and testing data were independently normalized so that the vector of weights for each blendshape had a mean of 0 and *SD* of 1.

Hyperparameters (the kernel type [“radial basis function,” “polynomial,” or “sigmoid”], the associated gamma, and regularization parameter;  $C$ ) for the final model were determined using five-fold cross-validation to assess the classification accuracy for every combination of parameters. Candidate values for gamma and  $C$  were 30 log-linear values ranging from  $10^{-6}$  to  $10^6$ . Permutation testing was used to assess whether the classification accuracy was above chance by randomly shuffling the class (i.e., identity) labels 1,000 times and repeating the cross-validation and classification procedure on each permutation to simulate a distribution of classification accuracies under the null hypothesis (i.e., that there is no dependency between the data and the labels). The proportion of the permutations with classification accuracies higher than the accuracy obtained by the original data was taken as the  $p$ -value, and we accepted an  $\alpha$  level of .05. This analysis was conducted using the “SVC,” “GridSearchCV,” and “permutation\_test\_score” functions in the scikit-learn library for python ([scikit-learn.org](http://scikit-learn.org)). The full procedure was conducted separately for each emotion.

As SVMs treat each data point as an independent observation, the classification accuracy (and associated  $p$ -value) reflects the ability of a model to detect any systematic differences between the vectors of blendshape weights for the two identities across participants, while ignoring any within-participant idiosyncrasies in emotion representation. If identity is classified above chance, it would suggest that participants create faces that are systematically different depending on the identity used.

### Within-Participant Analysis: RSA

We examined whether the within-participant representations of the different faces were clustered according to identity or emotion, using RSA. RSA is a method for examining how the structure of representations within one multidimensional space relates to the structure of representations in another space, and can be used to compare representational geometries with conceptual models (Kriegeskorte et al., 2008). In RSA, matrices are constructed where the value in each cell describes the similarity of two representations, either as observed in the data or as suggested by a conceptual model. These matrices are then statistically compared to assess how well the pairwise similarities as suggested by conceptual models can explain the similarities observed in the data. In our case, we examined the extent to which the representational structure of the different face categories can be explained by a model where faces belonging to the same identity are clustered together (identity-similarity), or a model where faces belonging to the same emotion are clustered together (emotion-similarity). To do this, for each participant, we computed the cosine distance between the core blendshapes for every pair of faces and created a representational dissimilarity matrix (RDM). Next, we constructed the two conceptual models (identity-similarity and emotion-similarity) where pairs of faces with the same identity or emotion are represented in the same space (cosine distance = 0) and pairs of faces with a different identity or emotion are represented in a different space (cosine distance = 1). Then, for each participant we performed multiple linear regression using the two conceptual models as predictors of the pairwise cosine distances, using the unique lower-triangle cells in each matrix to avoid repeated pairwise cosine distances. We used multiple linear regression to assess the extent to which each model explains the variance in the cosine distances after controlling for any variance shared with the other model. This allowed us to examine

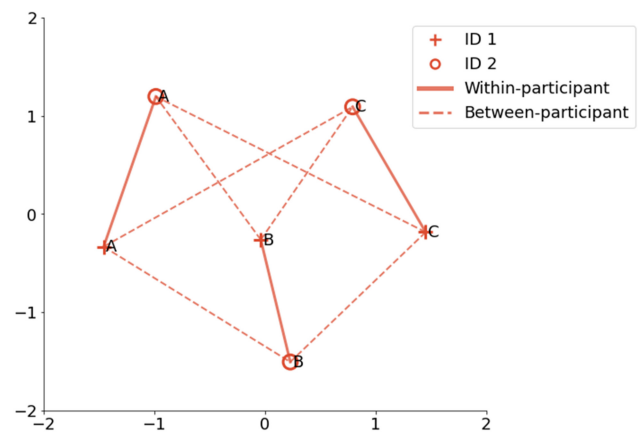
whether any variance in the cosine distances is associated with the similarity structure proposed by the identity model after controlling for variance associated with the emotion model. Finally, one-sample  $t$  tests were conducted using the regression coefficients across participants, against the test value of 0 (i.e., the coefficient for a model that accounts for no variance).

### Within-Participant Analysis: Between-Identity Representational Distances

To examine whether the representation of specific emotions depends on face identity, we examined the difference between the within- and between-participant distances between the two identities. If emotions are represented independently of identity, we expect that the cosine distance between one person’s preferred expressions using the two identities (for the same emotion) should be smaller than the cosine distance between any two people’s preferred expressions using the two identities (for the same emotion). That is, we expect that a participant’s two faces for a given emotion are closer to each other in expression space, than they are to another participant’s faces for the same emotion. We calculated every between-identity cosine distance of the core blendshapes, both within- and between-participants ( $N_{\text{within}} = 85$ ;  $N_{\text{between}} = 7,140$ ), separately for each emotion. Figure 1 shows a two-dimensional (2D) visualization of the representations of faces for three example participants, with the within- and between-participant distances between the representations.

We calculated the true difference between the means for these two sets of distances (i.e., how much smaller the mean of the within-participant distances is than the mean of the between-participant distances), then performed a bootstrapping procedure to model the null hypothesis that there is no difference between the within- and between-participant distances. We did this by resampling with replacement two new distributions of equivalent sizes ( $N_{\text{within}} = 85$ ;

**Figure 1**  
Representations of Faces From Three Example Participants (A, B, and C), Visualized in Two-Dimensional Space Using Multidimensional Scaling



*Note.* Circles and + signs represent each identity. Solid lines represent within-participant (between-identity) distances, and dashed lines represent between-participant (between-identity) distances. For visualization purposes, the axes represent Euclidean distance. See the online article for the color version of the figure.



$N_{\text{between}} = 7,140$ ) 100,000 times from the concatenated array of within and between-participant distances. On each iteration, the mean difference between the resampled distributions was recorded. To assess whether the true mean difference was larger than would be expected by chance, we found the proportion of resampled mean differences that were larger than the true difference, accepting an  $\alpha$  level of .05.

### **Within-Participant Analysis: Identity Blendshape Weight Differences**

To examine any differences between the identities in the specific blendshape weights in the preferred expressions, we used paired-samples Wilcoxon signed-rank tests to compare the vectors of weights between the two identities for each of the core blendshapes, across the four emotion categories. False discovery rate (FDR) was controlled for with the [Benjamini and Hochberg \(1995\)](#) procedure.

### **Data Availability and Preregistration Statement**

The GA toolkit is publicly available (<https://osf.io/dyfau/>), and data from this study are available on the Open Science Framework (OSF, <https://osf.io/t7a3w/>). This study was not preregistered.

## **Results**

[Figure 2](#) shows examples of five participants' preferred expressions for the emotion "angry," displayed on the two identities. This visualization suggests clear individual differences in preferred expressions between participants that appear consistent within the two different identities.

Participants reported high satisfaction with the expressions they created (median rating = 6/7 for all emotions and identities; see [Table S1 in the online supplemental materials](#)). These ratings (after eight iterations of the GA) are comparable to the satisfaction ratings provided by participants after 10 iterations ([Binetti et al., 2022](#)), suggesting that the GA toolkit converged onto expressions that participants believed represented the target emotion by the eighth iteration.

[Figure S2 in the online supplemental materials](#) shows the averages of the preferred expressions, created by averaging the weights of 149 blendshapes and rendering them onto the avatars with which they were created. [Figure S3 in the online supplemental materials](#) shows the preferred expressions of an example participant, with each rendered onto both avatars to visualize how a given set of blendshape weights appears on each identity.

To visualize any clustering of the faces in expression space, [Figure 3](#) plots the similarity structure of all faces using multidimensional scaling (MDS), where the distances in 2D space approximate the distances between the points in high-dimensional space. For the visualization, we used Euclidean distance between vectors of core blendshape weights instead of cosine distance, as it allows for the distance of each face from a vector of zeros (i.e., a neutral face) to be calculated, in addition to the distances between face pairs. The resulting 2D  $X$ - $Y$  coordinates of the neutral face were subtracted from the  $X$ - $Y$  coordinates of the other faces, to center the origin of the figure on the neutral face (where all blendshapes are set to 0). While there are individual differences in the preferred expressions, there is also clear clustering of the different emotion categories, with happy faces particularly distinct from the negative emotion clusters, consistent with [Binetti et al. \(2022\)](#).

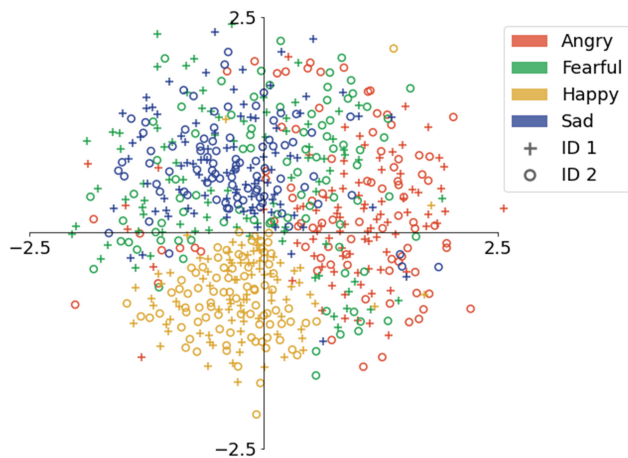
**Figure 2**

*Examples From Five Participants of Their Preferred Expressions for Anger Displayed by the Two Identities*



*Note.* See the online article for the color version of the figure.

**Figure 3**  
*Visualization of Preferred Expressions in Two-Dimensional Space, Using Multidimensional Scaling*



*Note.*  $N = 85$  participants. The origin of the figure represents the position of a neutral face relative to the positions of emotional faces. Scales on  $X$  and  $Y$  axes represent Euclidean distance. See the online article for the color version of the figure.

### Between-Participant Analysis: SVM Classification

We used SVMs to classify the identity of the faces from the vectors of blendshape weights across participants, separately for each emotion. Classification accuracy was calculated, and permutation testing was used to assess whether this accuracy was higher than could be expected by chance (accuracy of 0.25). Results showed that SVMs classified the two identities above chance for all emotion categories (angry: mean accuracy = 0.588,  $p = .021$ ; fear: mean accuracy = 0.624,  $p = .006$ ; happy: mean accuracy = 0.606,  $p = .015$ ; sad: mean accuracy = 0.600,  $p = .004$ ), suggesting that there is some consistent between-participant signal in the blendshape weights that distinguishes the two identities.

### Within-Participant Analysis: RSA

We used RSA to assess whether the within-participant positions of the different face categories in expression space were clustered according to identity or emotion. Two conceptual models were constructed: in the first, faces belonging to the same identity are represented in the same space, independent of emotion (identity model), and in the second, faces belonging to the same emotion are represented in the same space, independent of identity (emotion model). These models were used as predictors of the pairwise cosine distances for each participant, and a one-sample  $t$  test (one-tailed) was used to examine whether the mean of the distributions of regression coefficients for each of the models was greater than zero. Figure 4 shows the mean RDM across participants, and the two conceptual models. The coefficients for the emotion model were greater than zero, mean  $\beta = 0.152$ ,  $t(84) = 16.60$ ,  $p < .001$ , whereas they were not greater than zero for the identity model, mean  $\beta < 0.001$ ,  $t(84) = 0.05$ ,  $p = .519$ . This suggests that the similarity of pairs of faces as defined by the emotion category predicts the representational similarity (after controlling for the similarity as defined by the identity), and that the position of faces in representational space is determined more by emotion than identity.

### Within-Participant Analysis: Between-Identity Distances in Expression Space

We examined whether participants represent the two identity faces in a similar expression space, relative to other participants, by assessing whether the within-participant distances were smaller than the between-participant distances. We used bootstrapping to assess whether the difference between the means of these distances was larger than could be expected by chance. We found that the true mean difference was larger than could be expected by chance (i.e., that the within-participant distances were smaller than the between-participant distances) for angry faces (within = 0.623, between = 0.669,  $p = .003$ ), fearful faces (within = 0.622, between = 0.661,  $p = .011$ ), happy faces (within = 0.575, between = 0.612,  $p = .017$ ), but not sad faces (within = 0.616, between = 0.638,  $p = .101$ ). These results suggest that, for angry, fearful, and happy faces, the two faces created by a given participant are closer to each other than they are to the faces created by other participants. The average within- and between-participant distances, and the distributions of mean differences between the resampled arrays of within- and between-participant distances are presented in Figure 5 for each emotion, with the orange line representing the true within/between mean difference.

### Within-Participant Analysis: Identity Blendshape Weight Differences

Figure 6 plots the weights of the 41 core blendshapes used to create the preferred expressions for each identity and emotion category, shown for individual participants (top row) and averaged across participants for each emotion category (bottom row). Initial visual inspection supports the idea that participants used similar combinations of blendshapes on the two identities for happy and sad expressions, with more between-identity differences in the weights for angry and fearful faces.

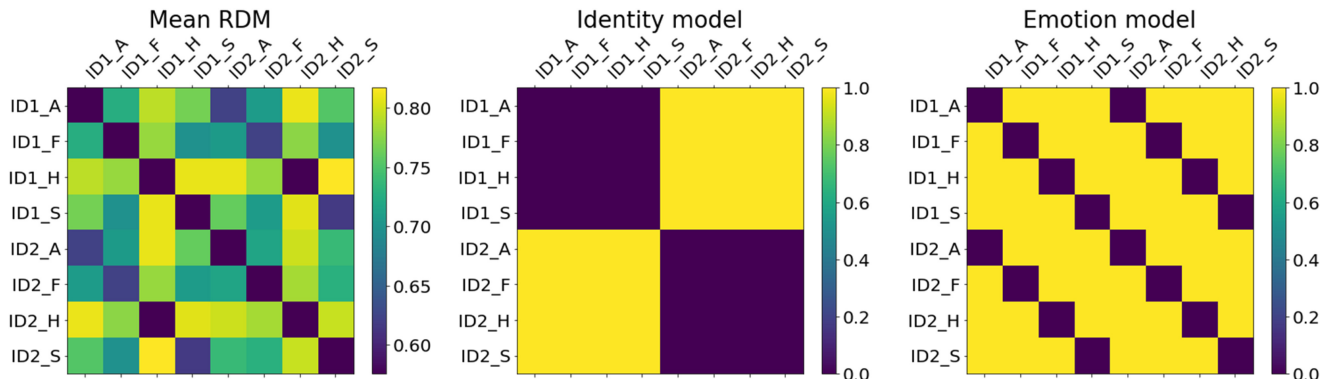
We then examined whether the weight assigned to each blendshape differed between the two identities. After correcting for multiple comparisons with FDR correction separately for each emotion, only one blendshape showed a significant weight difference between the two identities, for happy faces only. The “lower funneler” blendshape (approximately corresponding to the “lip funneler” action unit in the FACS) was larger for the first identity (male) faces (0.077) than the second identity (female) faces (0.012);  $W(84) = 245$ ;  $p(\text{uncorrected}) < .001$ ;  $p(\text{FDR}) = .006$ . All significant comparisons (before correction) are reported in Table S2 in the online supplemental materials.

### Discussion

Our aim was to examine whether individuals’ representations of emotions (as reflected by their preferred expressions) were independent of facial identity. To do this, we used our GA toolkit which allows precise quantification of facial expressions to analyze people’s preferred expressions in three different ways. We found that while a classifier was able to broadly distinguish the two identities between participants, individuals’ representations of emotions did not depend on identity. Firstly, we used RSA to examine how the within-subject structure of preferred expressions in expression-space relates to the structure of representations proposed by two conceptual models. This analysis showed that representations of faces belonging

**Figure 4**

Mean RDM Across Participants (Left) and Two Conceptual RDMs Modeling the Similarities From Identity (Middle) and Emotion (Right)



Note. A = angry; F = fearful; H = happy; S = sad; RDM = representational dissimilarity matrix. See the online article for the color version of the figure.

to the same emotion category are clustered together, but faces belonging to the same identity are not. Additionally, the dissimilarity matrix showed that, of all pairwise comparisons, the closest distances were the four within-emotion between-identity distances. Secondly, comparing the difference between the within-participant and between-participant distances showed that expressions of anger, fear, and happiness displayed by the different identities were represented in a space that is unique to the individual, although expressions of sadness were not. As identity was classified above chance for sad faces, it may be that representations of the two identities displaying sadness lie in two separate spaces that are common across participants. Finally, we examined whether there were differences in the weights of specific blendshapes between the two identities. We found only one comparison across the four emotion categories with significant between-identity differences, where the weight for the lower funneler (approximately corresponding to the Lip funneler action unit) was larger for happy expressions displayed by the male face than the female face. No other comparisons showed significant between-identity differences, suggesting that the weights of the specific blendshapes used to generate the expressions did not differ between the two identities. The significant difference in the weight of the lower funneler for happy faces suggests a consistent difference in the representations of happy expressions displayed by the two identities, but is not inconsistent with the result that participants represent happy expressions displayed by both identities in a space that is unique to them.

One point to consider is that our results might not reflect identity-independent representations of emotion, but instead sex-independent representations, since our two identity faces were of different sexes. Similar to the research surrounding the representation of emotion and identity, some studies using adaptation paradigms have reported representations of sex that are both dependent and independent from emotion (Bestelmeyer et al., 2010; Jaquet & Rhodes, 2008; Little et al., 2005). Considering the lack of a significant difference between the within- and between-subject distances for sadness, it might be the case that the manipulation of sex has a larger effect on people's unique representations of sadness than the other emotions, and that solely manipulating identity (e.g., faces of the same sex), might produce a similar pattern of results to the other emotions. Without testing multiple identities of the

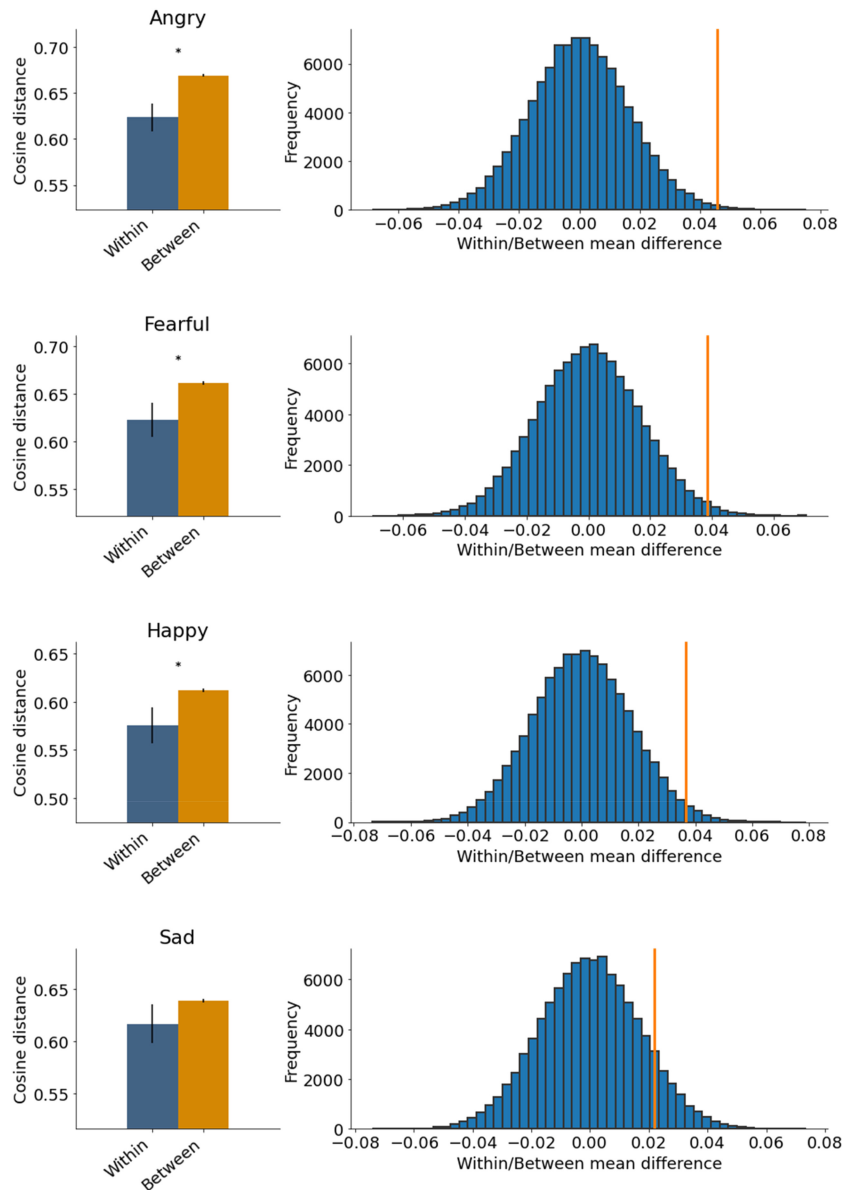
same sex, we cannot disentangle the role of sex and identity. One could replicate the study using avatars with multiple male and female identities; if the classification of identity is above chance within the same sex, it might suggest that the broad clustering of identity found in this study was due to the manipulation of morphological features rather than any involvement of higher-level social category representations. Alternatively, one could manipulate other social categories, such as race and age, to investigate how the representation of these categories might affect perceptual representations of expressions.

Another issue to consider is that these results might have occurred due to the nature of the task itself. The GA task necessarily requires participants to attend to the emotion rather than the identity. Given the evidence that neural representations of identity and expression within the face-processing network may depend on the task (Cohen Kadosh et al., 2010), it may be possible that representations of identity and expression were not probed to the same extent within the GA task. If the task was to manipulate identities rather than expression, we may have found expression-independent representation of identity.

Despite strong evidence of identity-independent representations within individuals, there were some systematic differences in the representations of emotions displayed by male and female faces across participants, as revealed by the SVM analysis and direct comparison of the blendshape weights. One possibility for these systematic differences could be the role of feedback processes from higher-level conceptual representations. Rhodes et al. (2015) highlight that it may be difficult to disentangle whether the representations of emotions that have been probed in adaptation paradigms and imaging studies reflect purely perceptual processes, or whether they are shaped by postperceptual feedback from higher-level emotion processing regions. Recent models of social perception put emphasis on feedback processing from higher-level regions of the brain that encode learned cultural knowledge of stereotype associations during the processing of faces (Freeman & Johnson, 2016). Additionally, some have reported a perceptual similarity of faces belonging to different sex and emotion categories (e.g., the perceived similarity of angry faces and male faces), and the similarity of the neural representations of those faces in the fusiform gyrus is explained by the subjectively rated conceptual association of those categories (Barnett et al., 2021; Stolier & Freeman, 2016). Other research showed that

**Figure 5**

*Histograms Showing the Mean Within- and Between-Participant Cosine Distances, With Error Bars Representing Standard Error (Left), and Distributions of Mean Differences From the Bootstrapping Analysis, for Each Emotion (Right)*



*Note.* Vertical reference lines within each distribution indicate the true within/between difference. See the online article for the color version of the figure.

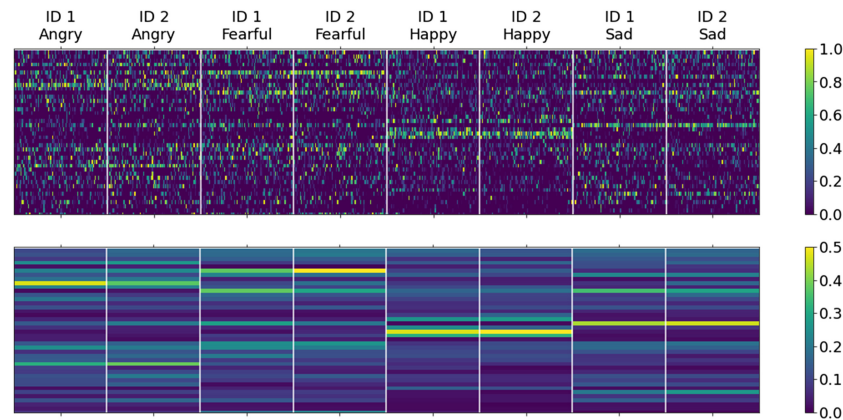
people more readily perceive anger displayed by male faces and happiness displayed by female faces (Becker, 2017; Becker et al., 2007; Hess et al., 2004), suggesting a possible interaction between gender stereotypes and emotion perception, that is measurable across participants. It is possible, therefore, that the above-chance classification of our two identity avatars across participants could reflect the involvement of higher-level representations of learned stereotypes surrounding sex and emotion on the unique representations of emotions displayed by male and female faces. An interesting question for future research, therefore, could be to examine whether there is any

relationship between the structure of an individual's representations of emotional male and female faces, and the strength of the conceptual associations surrounding sex and emotion held by the individual. Should such an association exist, this would demonstrate the role of higher-level social category representations on unique representations of emotional faces held by individuals.

Alternatively, the systematic across-participant differences in the preferred expressions in male and female faces might reflect some compensation for sex differences in morphological features of the face. One explanation for the facilitated perception of anger in



**Figure 6**  
*Visualizations of the Weights of the 41 Core Blendshapes in the Preferred Expressions Across the Different Identity and Emotion Categories*



*Note.* The top row displays the weights as 85 columns (one per participant) for each identity and emotion category. The bottom row displays weights averaged across participants. Each row is one of the 41 core blendshapes. See the online article for the color version of the figure.

male faces and happiness in female faces is that there are morphological cues shared by the sex and emotion categories that could affect facial displays of emotion (Adams et al., 2015; Becker, 2017). For example, both anger and masculinity are associated with greater “angularity” of the face, and the manipulation of this cue can facilitate the perception of both attributes (Becker et al., 2007). According to this account, the systematic across-participant differences in representations might simply reflect a consistent adjustment of facial features in a given sex to compensate for morphological cues that are not shared with the target emotion (e.g., adjusting features of female faces to increase “angularity” for angry expressions). It is worth noting that these two accounts are not mutually exclusive, and given the proposed interaction between feedforward and feedback processes in person perception (Freeman & Johnson, 2016), the consistent across-participant differences in emotion representations likely reflect both the matching of morphological features and the involvement of higher-level social category representations.

To account for conflicting results, it has been suggested that there is an initial stage of visual processing where identity and expression are jointly processed, followed by independent systems in which each attribute is encoded separately (Palermo et al., 2013, 2018; Rhodes et al., 2015). Given that we report systematic across-participant differences in the representation of emotions in male and female faces, this could suggest a common identity-dependent processing system, followed by identity-independent representations of emotion within a system dedicated to processing expression alone are unique to the individual. Future research could examine to what extent these shared joint representations are shaped by learned higher-level social category representations or by morphological cues.

There is growing evidence that conceptual knowledge influences the perception of facial expressions (Brooks & Freeman, 2018; Brooks et al., 2019; Murray et al., 2021). Our task involves participants associating a label (e.g., “fear”) with their conceptual knowledge of the emotion, to guide their creation of facial expressions. It is worth considering that these results might therefore reflect individual differences in the conceptual representations of different emotions.

The emotion category labels used in this experiment cover a wide range of emotional experiences; for example, “anger” encompasses emotions ranging from annoyance to rage, and it is likely that individuals have unique interpretations of “anger” as an emotion. One study found that people reliably report emotional experiences that are captured by 27 distinct categories (e.g., “amusement,” “nostalgia,” and “awe”) when viewing emotionally evocative videos, which are more subtle than the basic emotions used throughout emotion research (Cowen & Keltner, 2017). It is possible, therefore, that the proximity of the male and female preferred expressions belonging to each individual might not reflect an identity-independent representation of emotion, but rather representations of faces within a more subtle emotion space. Future research could examine representations of expressions relating to more fine-grained emotional experiences.

Relatedly, it is also worth noting that our design limits the data to a single-point estimate per identity and emotion category. Constructivist theories of emotion propose that a single emotion can be associated with a range of expressions (Barrett et al., 2019), and people judge a range of expressions as belonging to a given emotion category (Kohler et al., 2004). It is worth considering that it may be more appropriate to model unique representations of emotion as distributions within an expression-space to encompass the range of expressions that an individual associates with a given emotion, and future research could examine whether these distributions are identity-independent.

We conclude that people’s unique representations of emotion are likely independent of the identity of the face, although there may be some systematic effects of identity on (population level) emotion representations across participants. Further research is needed to examine whether this clustering is the result of between-identity morphological differences or due to any involvement of higher-level conceptual representations of social categories.

## References

- Adams, R. B., Hess, U., & Kleck, R. E. (2015). The intersection of gender-related facial appearance and facial displays of emotion. *Emotion Review*, 7(1), 5–13. <https://doi.org/10.1177/1754073914544407>

- Adolphs, R. (2002). Recognizing emotion from facial expressions: Psychological and neurological mechanisms. *Behavioral and Cognitive Neuroscience Reviews*, 1(1), 21–62. <https://doi.org/10.1177/1534582302001001003>
- Andrews, T. J., & Ewbank, M. P. (2004). Distinct representations for facial identity and changeable aspects of faces in the human temporal lobe. *NeuroImage*, 23(3), 905–913. <https://doi.org/10.1016/j.neuroimage.2004.07.060>
- Awad, D., Clifford, C. W. G., White, D., & Mareschal, I. (2020). Asymmetric contextual effects in age perception. *Royal Society Open Science*, 7(12), Article 200936. <https://doi.org/10.1098/rsos.200936>
- Barnett, B. O., Brooks, J. A., & Freeman, J. B. (2021). Stereotypes bias face perception via orbitofrontal-fusiform cortical interaction. *Social Cognitive and Affective Neuroscience*, 16(3), 302–314. <https://doi.org/10.1093/scan/nsaa165>
- Barrett, L. F., Adolphs, R., Marsella, S., Martinez, A. M., & Pollak, S. D. (2019). Emotional expressions reconsidered: Challenges to inferring emotion from human facial movements. *Psychological Science in the Public Interest*, 20(1), 1–68. <https://doi.org/10.1177/1529100619832930>
- Baseler, H. A., Harris, R. J., Young, A. W., & Andrews, T. J. (2014). Neural responses to expression and gaze in the posterior superior temporal sulcus interact with facial identity. *Cerebral Cortex*, 24(3), 737–744. <https://doi.org/10.1093/cercor/bhs360>
- Becker, D. V. (2017). Facial gender interferes with decisions about facial expressions of anger and happiness. *Journal of Experimental Psychology: General*, 146(4), 457–463. <https://doi.org/10.1037/xge0000279>
- Becker, D. V., Kenrick, D. T., Neuberg, S. L., Blackwell, K. C., & Smith, D. M. (2007). The confounded nature of angry men and happy women. *Journal of Personality and Social Psychology*, 92(2), 179–190. <https://doi.org/10.1037/0022-3514.92.2.179>
- Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(1), 289–300. <https://doi.org/10.1111/j.2517-6161.1995.tb02031.x>
- Bestelmeyer, P. E. G., Jones, B. C., DeBruine, L. M., Little, A. C., & Welling, L. L. M. (2010). Face aftereffects suggest interdependent processing of expression and sex and of expression and race. *Visual Cognition*, 18(2), 255–274. <https://doi.org/10.1080/13506280802708024>
- Binetti, N., Roubtsova, N., Carlisi, C., Cosker, D., Viding, E., & Mareschal, I. (2022). Genetic algorithms reveal profound individual differences in emotion recognition. *Proceedings of the National Academy of Sciences*, 119(45), Article e2201380119. <https://doi.org/10.1073/pnas.2201380119>
- Brooks, J. A., Chikazoe, J., Sadato, N., & Freeman, J. B. (2019). The neural representation of facial-emotion categories reveals conceptual structure. *Proceedings of the National Academy of Sciences*, 116(32), 15861–15870. <https://doi.org/10.1073/pnas.1816408116>
- Brooks, J. A., & Freeman, J. B. (2018). Conceptual knowledge predicts the representational structure of facial emotion perception. *Nature Human Behaviour*, 2(8), 581–591. <https://doi.org/10.1038/s41562-018-0376-6>
- Bruce, V. (1986). Influences of familiarity on the processing of faces. *Perception*, 15(4), 387–397. <https://doi.org/10.1068/p150387>
- Bruce, V., & Young, A. (1986). Understanding face recognition. *British Journal of Psychology*, 77(3), 305–327. <https://doi.org/10.1111/j.2044-8295.1986.tb02199.x>
- Bruce, V., & Young, A. (1998). *In the eye of the beholder: The science of face perception*. Oxford University Press.
- Calder, A. J., & Young, A. W. (2005). Understanding the recognition of facial identity and facial expression. *Nature Reviews Neuroscience*, 6(8), 641–651. <https://doi.org/10.1038/nrn1724>
- Campbell, J., & Burke, D. (2009). Evidence that identity-dependent and identity-independent neural populations are recruited in the perception of five basic emotional facial expressions. *Vision Research*, 49(12), 1532–1540. <https://doi.org/10.1016/j.visres.2009.03.009>
- Carlisi, C. O., Reed, K., Helmink, F. G. L., Lachlan, R., Cosker, D. P., Viding, E., & Mareschal, I. (2021). Using genetic algorithms to uncover individual differences in how humans represent facial emotion. *Royal Society Open Science*, 8(10), Article 202251. <https://doi.org/10.1098/rsos.202251>
- Clifford, C. W. G., Watson, T. L., & White, D. (2018). Two sources of bias explain errors in facial age estimation. *Royal Society Open Science*, 5(10), Article 180841. <https://doi.org/10.1098/rsos.180841>
- Cohen Kadosh, K., Henson, R. N. A., Cohen Kadosh, R., Johnson, M. H., & Dick, F. (2010). Task-dependent activation of face-sensitive cortex: An fMRI adaptation study. *Journal of Cognitive Neuroscience*, 22(5), 903–917. <https://doi.org/10.1162/jocn.2009.21224>
- Cowen, A. S., & Keltner, D. (2017). Self-report captures 27 distinct categories of emotion bridged by continuous gradients. *Proceedings of the National Academy of Sciences*, 114(38), E7900–E7909. <https://doi.org/10.1073/pnas.1702247114>
- Duchaine, B. C., Parker, H., & Nakayama, K. (2003). Normal recognition of emotion in a prosopagnosic. *Perception*, 32(7), 827–838. <https://doi.org/10.1068/p5067>
- Eger, E., Schyns, P. G., & Kleinschmidt, A. (2004). Scale invariant adaptation in fusiform face-responsive regions. *NeuroImage*, 22(1), 232–242. <https://doi.org/10.1016/j.neuroimage.2003.12.028>
- Ellamil, M., Susskind, J. M., & Anderson, A. K. (2008). Examinations of identity invariance in facial expression adaptation. *Cognitive, Affective, & Behavioral Neuroscience*, 8(3), 273–281. <https://doi.org/10.3758/CABN.8.3.273>
- Fox, C. J., & Barton, J. J. S. (2007). What is adapted in face adaptation? The neural representations of expression in the human visual system. *Brain Research*, 1127, 80–89. <https://doi.org/10.1016/j.brainres.2006.09.104>
- Fox, C. J., Moon, S., Iaria, G., & Barton, J. (2009). The correlates of subjective perception of identity and expression in the face network: An fMRI adaptation study. *NeuroImage*, 44(2), 569–580. <https://doi.org/10.1016/j.neuroimage.2008.09.011>
- Freeman, J. B., & Johnson, K. L. (2016). More than meets the eye: Split-second social perception. *Trends in Cognitive Sciences*, 20(5), 362–374. <https://doi.org/10.1016/j.tics.2016.03.003>
- Ganel, T., Valyear, K. F., Goshen-Gottstein, Y., & Goodale, M. A. (2005). The involvement of the “fusiform face area” in processing facial expression. *Neuropsychologia*, 43(11), 1645–1654. <https://doi.org/10.1016/j.neuropsychologia.2005.01.012>
- Georgeson, M. (2004). Visual aftereffects: Cortical neurons change their tune. *Current Biology*, 14(18), R751–R753. <https://doi.org/10.1016/j.cub.2004.09.011>
- Haxby, J. V., Hoffman, E. A., & Gobbini, M. I. (2000). The distributed human neural system for face perception. *Trends in Cognitive Sciences*, 4(6), 223–233. [https://doi.org/10.1016/S1364-6613\(00\)01482-0](https://doi.org/10.1016/S1364-6613(00)01482-0)
- Hess, U., Adams, R. B., & Kleck, R. E. (2004). Facial appearance, gender, and emotion expression. *Emotion*, 4(4), 378–388. <https://doi.org/10.1037/1528-3542.4.4.378>
- Hugenberg, K., & Bodenhausen, G. V. (2003). Facing prejudice: Implicit prejudice and the perception of facial threat. *Psychological Science*, 14(6), 640–643. [https://doi.org/10.1046/j.0956-7976.2003.psci\\_1478.x](https://doi.org/10.1046/j.0956-7976.2003.psci_1478.x)
- Jack, R. E., Garrod, O. G. B., Yu, H., Caldara, R., & Schyns, P. G. (2012). Facial expressions of emotion are not culturally universal. *Proceedings of the National Academy of Sciences*, 109(19), 7241–7244. <https://doi.org/10.1073/pnas.1200155109>
- Jaquet, E., & Rhodes, G. (2008). Face aftereffects indicate dissociable, but not distinct, coding of male and female faces. *Journal of Experimental Psychology: Human Perception and Performance*, 34(1), 101–112. <https://doi.org/10.1037/0096-1523.34.1.101>
- Kawasaki, H., Tsuchiya, N., Kovach, C. K., Nourski, K. V., Oya, H., Howard, M. A., & Adolphs, R. (2012). Processing of facial emotion in the human fusiform gyrus. *Journal of Cognitive Neuroscience*, 24(6), 1358–1370. [https://doi.org/10.1162/jocn\\_a\\_00175](https://doi.org/10.1162/jocn_a_00175)

- Kohler, C. G., Turner, T., Stolar, N. M., Bilker, W. B., Brensing, C. M., Gur, R. E., & Gur, R. C. (2004). Differences in facial expressions of four universal emotions. *Psychiatry Research*, *128*(3), 235–244. <https://doi.org/10.1016/j.psychres.2004.07.003>
- Kriegeskorte, N., Mur, M., & Bandettini, P. A. (2008). Representational similarity analysis—Connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, *2*(4), 1–28. <https://doi.org/10.3389/neuro.06.004.2008>
- Leopold, D. A., O'Toole, A. J., Vetter, T., & Blanz, V. (2001). Prototype-referenced shape encoding revealed by high-level aftereffects. *Nature Neuroscience*, *4*(1), 89–94. <https://doi.org/10.1038/82947>
- Little, A. C., DeBruine, L. M., & Jones, B. C. (2005). Sex-contingent face after-effects suggest distinct neural populations code male and female faces. *Proceedings of the Royal Society B: Biological Sciences*, *272*(1578), 2283–2287. <https://doi.org/10.1098/rspb.2005.3220>
- Michalek, J., Lisi, M., Binetti, N., Ozkaya, S., Hadfield, K., Dajani, R., & Mareschal, I. (2022). War-related trauma linked to increased sustained attention to threat in children. *Child Development*, *93*(4), 900–909. <https://doi.org/10.1111/cdev.13739>
- Murray, T., O'Brien, J., Sagiv, N., & Garrido, L. (2021). The role of stimulus-based cues and conceptual information in processing facial expressions of emotion. *Cortex*, *144*, 109–132. <https://doi.org/10.1016/j.cortex.2021.08.007>
- Palermo, R., Jeffery, L., Lewandowsky, J., Fiorentini, C., Irons, J. L., Dawel, A., Burton, N., McKone, E., & Rhodes, G. (2018). Adaptive face coding contributes to individual differences in facial expression recognition independently of affective factors. *Journal of Experimental Psychology: Human Perception and Performance*, *44*(4), 503–517. <https://doi.org/10.1037/xhp0000463>
- Palermo, R., O'Connor, K. B., Davis, J. M., Irons, J., & McKone, E. (2013). New tests to measure individual differences in matching and labelling facial expressions of emotion, and their association with ability to recognise vocal emotions and facial identity. *PLoS ONE*, *8*(6), Article e68126. <https://doi.org/10.1371/journal.pone.0068126>
- Rhodes, G., Pond, S., Burton, N., Kloth, N., Jeffery, L., Bell, J., Ewing, L., Calder, A. J., & Palermo, R. (2015). How distinct is the coding of face identity and expression? Evidence for some common dimensions in face space. *Cognition*, *142*, 123–137. <https://doi.org/10.1016/j.cognition.2015.05.012>
- Roubtsova, N., Parsons, M., Binetti, N., Mareschal, I., Viding, E., & Cosker, D. (2021). *EmoGen: Quantifiable emotion generation and analysis for experimental psychology* (pp. 1–21). Retrieved from <http://arxiv.org/abs/2107.00480>
- Skerry, A. E., & Saxe, R. (2015). Neural representations of emotion are organized around abstract event features. *Current Biology*, *25*(15), 1945–1954. <https://doi.org/10.1016/j.cub.2015.06.009>
- Skinner, A. L., & Benton, C. P. (2012). The expressions of strangers: Our identity-independent representation of facial expression. *Journal of Vision*, *12*(2), Article 12. <https://doi.org/10.1167/12.2.12>
- Song, M., Shinomori, K., Qian, Q., Yin, J., & Zeng, W. (2015). The change of expression configuration affects identity-dependent expression aftereffect but not identity-independent expression aftereffect. *Frontiers in Psychology*, *6*, 1–12. <https://doi.org/10.3389/fpsyg.2015.01937>
- Sormaz, M., Watson, D. M., Smith, W. A. P., Young, A. W., & Andrews, T. J. (2016). Modelling the perceptual similarity of facial expressions from image statistics and neural responses. *NeuroImage*, *129*, 64–71. <https://doi.org/10.1016/j.neuroimage.2016.01.041>
- Stolier, R. M., & Freeman, J. B. (2016). Neural pattern similarity reveals the inherent intersection of social categories. *Nature Neuroscience*, *19*(6), 795–797. <https://doi.org/10.1038/nn.4296>
- Todorov, A., Baron, S. G., & Oosterhof, N. N. (2008). Evaluating face trustworthiness: A model based approach. *Social Cognitive and Affective Neuroscience*, *3*(2), 119–127. <https://doi.org/10.1093/scan/nsn009>
- Tsantani, M., Kriegeskorte, N., McGettigan, C., & Garrido, L. (2019). Faces and voices in the brain: A modality-general person-identity representation in superior temporal sulcus. *NeuroImage*, *201*, Article 116004. <https://doi.org/10.1016/j.neuroimage.2019.07.017>
- Webster, M. A., Kaping, D., Mizokami, Y., & Duhamel, P. (2004). Adaptation to natural facial categories. *Nature*, *428*(6982), 557–561. <https://doi.org/10.1038/nature02420>
- Wegrzyn, M., Riehle, M., Labudda, K., Woermann, F., Baumgartner, F., Pollmann, S., Bien, C. G., & Kissler, J. (2015). Investigating the brain basis of facial expression perception using multi-voxel pattern analysis. *Cortex*, *69*, 131–140. <https://doi.org/10.1016/j.cortex.2015.05.003>
- Winston, J. S., Henson, R. N. A., Fine-Goulden, M. R., & Dolan, R. J. (2004). fMRI-adaptation reveals dissociable neural representations of identity and expression in face perception. *Journal of Neurophysiology*, *92*(3), 1830–1839. <https://doi.org/10.1152/jn.00155.2004>
- Xu, X., & Biederman, I. (2010). Loci of the release from fMRI adaptation for changes in facial expression, identity, and viewpoint. *Journal of Vision*, *10*(14), Article 36. <https://doi.org/10.1167/10.14.36>
- Young, A. W., McWeeny, K. H., Hay, D. C., & Ellis, A. W. (1986). Matching familiar and unfamiliar faces on identity and expression. *Psychological Research*, *48*(2), 63–68. <https://doi.org/10.1007/BF00309318>
- Zhang, H., Japee, S., Nolan, R., Chu, C., Liu, N., & Ungerleider, L. G. (2016). Face-selective regions differ in their ability to classify facial expressions. *NeuroImage*, *130*, 77–90. <https://doi.org/10.1016/j.neuroimage.2016.01.045>

Received November 17, 2022

Revision received May 23, 2023

Accepted June 6, 2023 ■