

FFA and OFA Encode Distinct Types of Face Identity Information

Maria Tsantani,¹ Nikolaus Kriegeskorte,² Katherine Storrs,³ Adrian Lloyd Williams,¹ Carolyn McGettigan,⁴ and Lúcia Garrido¹

¹Division of Psychology, Department of Life Sciences, Brunel University London, Uxbridge, UB8 3PH, United Kingdom, ²Zuckerman Mind Brain Behavior Institute, Columbia University, New York, New York 10027, ³Department of Experimental Psychology, Justus Liebig University, Giessen, 35390, Germany, and ⁴Speech Hearing and Phonetic Sciences, University College London, London WC1N 1PF, United Kingdom

Faces of different people elicit distinct fMRI patterns in several face-selective regions of the human brain. Here we used representational similarity analysis to investigate what type of identity-distinguishing information is encoded in three face-selective regions: fusiform face area (FFA), occipital face area (OFA), and posterior superior temporal sulcus (pSTS). In a sample of 30 human participants (22 females, 8 males), we used fMRI to measure brain activity patterns elicited by naturalistic videos of famous face identities, and compared their representational distances in each region with models of the differences between identities. We built diverse candidate models, ranging from low-level image-computable properties (pixel-wise, GIST, and Gabor-Jet dissimilarities), through higher-level image-computable descriptions (OpenFace deep neural network, trained to cluster faces by identity), to complex human-rated properties (perceived similarity, social traits, and gender). We found marked differences in the information represented by the FFA and OFA. Dissimilarities between face identities in FFA were accounted for by differences in perceived similarity, Social Traits, Gender, and by the OpenFace network. In contrast, representational distances in OFA were mainly driven by differences in low-level image-based properties (pixel-wise and Gabor-Jet dissimilarities). Our results suggest that, although FFA and OFA can both discriminate between identities, the FFA representation is further removed from the image, encoding higher-level perceptual and social face information.

Key words: face identity; face processing; FFA; OFA; representational similarity analysis

Significance Statement

Recent studies using fMRI have shown that several face-responsive brain regions can distinguish between different face identities. It is however unclear whether these different face-responsive regions distinguish between identities in similar or different ways. We used representational similarity analysis to investigate the computations within three brain regions in response to naturalistically varying videos of face identities. Our results revealed that two regions, the fusiform face area and the occipital face area, encode distinct identity information about faces. Although identity can be decoded from both regions, identity representations in fusiform face area primarily contained information about social traits, gender, and high-level visual features, whereas occipital face area primarily represented lower-level image features.

Received June 8, 2020; revised Dec. 18, 2020; accepted Dec. 22, 2020.

Author contributions: M.T., N.K., C.M., and L.G. designed research; M.T., N.K., C.M., and L.G. performed research; M.T., N.K., K.S., A.L.W., and L.G. contributed unpublished reagents/analytic tools; M.T. and L.G. analyzed data; M.T. and L.G. wrote the first draft of the paper; M.T., N.K., K.S., A.L.W., C.M., and L.G. edited the paper; M.T., N.K., K.S., and L.G. wrote the paper.

This work was supported by Leverhulme Trust Research Grant RPG-2014-392 to L.G., N.K., and C.M. We thank Tiana Rakotonombana, Roxanne Zamyadi, Rasanat Nawaz, and Natasha Baxter for help with stimulus preparation and with testing.

M. Tsantani's present address: Department of Psychological Sciences, Birkbeck, University of London, London WC1E 7HX, United Kingdom.

L. Garrido's present address: Department of Psychology, City, University of London, London EC1V 0HB, United Kingdom.

The authors declare no competing financial interests.

Correspondence should be addressed to Maria Tsantani at maria.tsantani@gmail.com or Lúcia Garrido at lucia.garrido@city.ac.uk.

<https://doi.org/10.1523/JNEUROSCI.1449-20.2020>

Copyright © 2021 the authors

Introduction

The human brain contains several face-selective regions that consistently respond more to faces than other visual stimuli (Kanwisher et al., 1997; Pitcher et al., 2011; Rossion et al., 2012; Khuvis et al., 2021; Axelrod et al., 2019). fMRI has revealed that some of these regions represent different face identities with distinct brain patterns. Specifically, studies using fMRI multivariate pattern analysis have shown that face identities can be distinguished based on their elicited response patterns in the fusiform face area (FFA), occipital face area (OFA), posterior superior temporal sulcus (pSTS), and anterior inferior temporal lobe (Nestor et al., 2011; Goesart and Op de Beeck, 2013; Verosky et al., 2013; Anzellotti et al., 2014; Axelrod and Yovel, 2015; Zhang et al., 2016; Anzellotti and Caramazza, 2017; Guntupalli et al.,

2017; di Oleggio Castello et al., 2017; Tsantani et al., 2019; for results using intracranial EEG [iEEG], see also Davidesco et al., 2014; Ghuman et al., 2014; Khuvis et al., 2021). But do these regions represent the same information and, if not, what information is explicitly encoded in each of these face-selective regions?

Behaviorally, we distinguish between different faces using the surface appearance of the face, the shape of face features, and their spacing or configuration (e.g., Rhodes, 1988; Calder et al., 2001; Yovel and Duchaine, 2006; Russell and Sinha, 2007; Russell et al., 2007; Tardif et al., 2019). In particular, Abudarham and Yovel (2016) recently showed that features, such as lip thickness, hair color, eye color, eye shape, and eyebrow thickness, were crucial in distinguishing between individuals (see also Abudarham et al., 2019). Additionally, we perceive a vast amount of socially relevant information from faces that can be used to distinguish between different individuals, such as gender, age, ethnicity, and social traits (Oosterhof and Todorov, 2008; Sutherland et al., 2013), and even relationships and social network position (Parkinson et al., 2014, 2017). Therefore, if the response patterns in a certain brain region distinguish between two individuals, that region could be representing any one, or a combination of, these dimensions.

Like several other studies (see above), Goesaert and Op de Beeck (2013) demonstrated that the FFA, OFA, and a face-selective region in the anterior inferior temporal lobe could all decode between different face identities based on fMRI response patterns. Importantly, the authors further tested what type of face information was encoded in these different regions. The authors found that all three regions could distinguish between faces using both configural and featural face information; therefore, all regions seemed to represent similar information. Goesaert and Op de Beeck (2013) also showed that representational distances between different faces in face-selective regions did not correlate with low-level pixel-based information. This study, however, used one single image for each person's face, making it difficult to disentangle whether representations in a certain brain region are related to identity *per se* or related to the specific images used.

To determine whether brain response patterns represent face identity *per se*, it is necessary to show that patterns generalize across different images of the same person's face, in addition to distinguishing that person's face from the faces of other people. Anzellotti et al. (2014) showed that classifiers trained to decode face identities in the FFA, OFA, anterior temporal lobe, and pSTS (later analyzed in Anzellotti and Caramazza, 2017) could also decode the same faces from novel viewpoints. Guntupalli et al. (2017) additionally showed a hierarchical organization of the functions of face-selective regions, with the OFA decoding viewpoint of face independently of the face identity, the anterior inferior temporal lobe (and a region in the inferior frontal cortex) decoding face identity independently of the viewpoint, and the FFA decoding both viewpoint and identity information (see also Dubois et al., 2015). Extending these findings and using iEEG in epilepsy patients, Ghuman et al. (2014) showed invariant decoding in the FFA across different facial expressions. In contrast, Grossman et al. (2019) have recently shown that representational distances between different face identities (computed from brain response patterns recorded from implanted electrodes) were very similar across the OFA and the FFA (in the left hemisphere). Crucially, the representational geometries in both regions were associated with differences in image-level descriptions computed from a deep neural network (VGG-Face), which were not

generalizable across different viewpoints of the same person's face. These results thus suggest that the OFA and FFA both represent complex configurations of image-based information and not face identity *per se*.

Also using iEEG, Davidesco et al. (2014) further showed that representational distances between face images in the FFA (and to a lesser extent in the OFA) were associated with perceived similarity and characteristics of facial features (e.g., face area and mouth width), but not with low-level features related to pixel-based information (see also Ghuman et al., 2014). Some fMRI studies have shown that even lower-level stimulus-based properties of face images, such as those computed by Gabor filters, explain significant variance in the representational geometries in the FFA (Carlin and Kriegeskorte, 2017) as well as OFA and pSTS (Weibert et al., 2018). On the other hand, other studies have shown that more high-level information, such as biographical information and social context, affects the similarity of response patterns to different faces in the FFA (Verosky et al., 2013; Collins et al., 2016).

There is thus mixed evidence regarding whether different face-selective regions rely on similar or distinct information to distinguish between face identities, and what type of information may be encoded in different regions. In the present study, we used representational similarity analysis (RSA) (Kriegeskorte et al., 2008a,b) to investigate what type of identity-distinguishing information is encoded in different face-selective regions. In our previous work (Tsantani et al., 2019), we showed that famous face identities could be distinguished in the right FFA, OFA, and pSTS based on their elicited fMRI response patterns. Here, for the same set of famous identities and using the same data as in Tsantani et al. (2019), we compared the representational distances between identity-elicited fMRI patterns in these regions with diverse candidate models of face properties that could potentially be used to distinguish between identities.

Importantly, we used multiple naturalistically varying videos for each identity that varied freely in terms of viewpoint, lighting, head motion, and general appearance. In addition, our representational distances were cross-validated across different videos, to deconfound identity from incidental image properties. By using a large, diverse set of candidate models, based on image properties of the stimuli (Image-computable models) and on human-rated properties (Perceived-property models), we were able to determine what types of identity-distinguishing information are encoded in different face-selective regions.

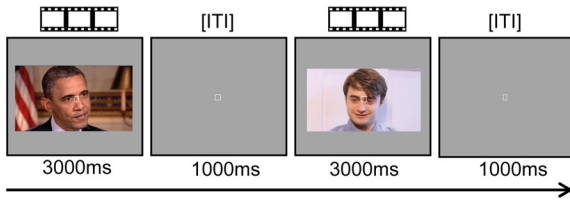
Materials and Methods

This study involved an fMRI component, in which we measured brain representations of faces and voices, and a behavioral component, in which we collected ratings of the same faces and voices on social traits and perceived similarity. The fMRI part corresponds to the same experiment and data described in Tsantani et al. (2019), and the behavioral part is reported here for the first time. In the present study, we analyzed the data related to faces only.

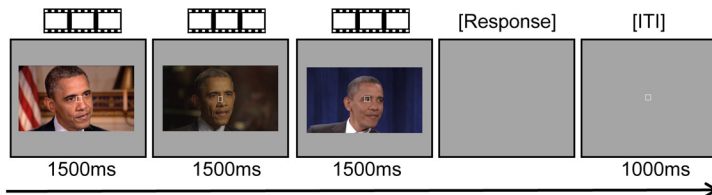
Participants

We recruited 31 healthy right-handed adult participants to take part in two fMRI sessions and a behavioral session (all on separate days, resulting in at least 6 h of testing per participant). We did not conduct a formal power analysis as there were no previous studies at the time of the study design that had investigated the main effect described by Tsantani et al. (2019). Our sample size was determined based on similar fMRI studies within the field and on available funding. To ensure adequate exposure to our stimulus set of famous people, participants were required to be native English speakers between 18 and 30 years of age, and to have

fMRI



Social Trait Judgements



Perceived Similarity



Figure 1. Examples of face trials in the fMRI and behavioral experiments. All experiments presented the same videos of moving, nonspeaking, faces of 12 famous people. For each famous person, we presented six naturalistically varying videos of their face. In an event-related fMRI task, each trial presented a single face video. This task also contained trials of the same length featuring voice clips (excluded from the present analysis), stimuli relating to the anomaly detection task, and fixation (null events). In each trial of the Social Trait Judgments Tasks (separate tasks for Trustworthiness, Dominance, Attractiveness, and Valence), participants viewed three videos of the face of the same identity and judged the intensity of the target trait (on a scale from 1 to 7). In each trial of the Perceived Similarity Task, participants viewed three videos of one identity followed by three videos of a different identity and rated their visual similarity (from 1 to 7). Face videos were presented for their full duration of 3000 ms in the fMRI experiment, whereas only the first 1500 ms were presented in the behavioral experiments.

been resident in the United Kingdom for at least 10 years. We also independently verified that all participants knew the famous people used in the experiment (see Tsantani et al., 2019). No inclusion or exclusion criteria were applied based on race or ethnicity, and we did not formally record this information. It has been shown that the other-race effect does not apply to familiar faces (McKone et al., 2007; Zhou and Mondloch, 2016). Participants were recruited at Royal Holloway, University of London, and Brunel University London. One participant was excluded because of excessive head movement in the scanner. The final sample consisted of 30 participants (22 females, 8 males) with a mean age of 21.2 years ($SD = 2.37$ years, range = 19–27 years). Participants reported normal or corrected-to-normal vision and normal hearing, provided written informed consent, and were reimbursed for their participation. The study was approved by the Ethics Committee of Brunel University London.

Stimuli

The same stimuli were used in the fMRI and behavioral testing, and consisted of videos of the faces and sound recordings of 12 famous individuals, including actors, comedians, TV personalities, pop stars, and politicians: Alan Carr, Daniel Radcliffe, Emma Watson, Arnold Schwarzenegger, Sharon Osbourne, Graham Norton, Beyonce Knowles, Barbara Windsor, Kylie Minogue, Barack Obama, Jonathan Ross, and Cheryl Cole. These individuals were selected based on pilot studies that showed that participants (18–30 years of age and living in the United Kingdom) could recognize them easily from their faces and voices.

For each identity, six silent, nonspeaking video clips of their moving face were obtained from videos on www.YouTube.com (Fig. 1). The six clips were obtained from different original videos. In total, we obtained 72 face stimuli. Face videos were selected so that the background did not

provide any cues to the identity of the person. The face videos were primarily front-facing and did not feature any speech but were otherwise unconstrained in terms of facial motion. Head movements included nodding, smiling, and rotating the head. Videos were edited so that they were 3 s long, 640×360 pixels, and centered on the bridge of the nose, using Final Cut Pro X (Apple).

For purposes not related to this study, we also presented 72 voice stimuli, which consisted of recordings of the voices of the same 12 famous individuals (6 clips per identity) obtained from videos on www.YouTube.com. Speech clips were selected so that the speech content, which was different for every recording, did not reveal the identity of the speaker. Recordings were edited so that they contained 3 s of speech after removing long periods of silence using Audacity 2.0.5 recording and editing software (RRID:SCR_007198). The recordings were converted to mono with a sampling rate of 44,100, low-pass filtered at 10 kHz, and root mean square normalized using Praat (version 5.3.80; www.praat.org) (Boersma and Weenink, 2014).

Participants were familiarized with all stimuli via one exposure to each clip immediately before the first scanning session.

MRI data acquisition and preprocessing

Participants completed two MRI sessions: in each session, participants completed a structural scan, three runs of the main experiment, and functional localizer scans (for face and voice areas, but below we only describe the localizer of face-selective regions). Participants were scanned using a 3.0 Tesla Tim Trio MRI scanner (Siemens) with a 32 channel head coil. Scanning took place at the Combined Universities Brain Imaging Center at Royal Holloway, University of London. We acquired whole-brain T1-weighted anatomic scans using MPRAGE (1.0×1.0 in-plane resolution; slice thickness, 1.0 mm; 176 axial

interleaved slices; PAT, factor 2; PAT mode, GeneRalized Autocalibrating Partially Parallel Acquisitions; TR, 1900 ms; TE, 3.03 ms; flip angle, 11°; matrix, 256 × 256; FOV, 256 mm).

For the functional runs, we acquired T2*-weighted functional scans using EPI [3.0 × 3.0 in-plane resolution; slice thickness, 3.0 mm; PAT, factor 2; PAT mode, GeneRalized Autocalibrating Partially Parallel Acquisitions; 34 sequential (descending) slices; TR, 2000 ms; TE, 30 ms; flip angle, 78°; matrix, 64 × 64; FOV, 192 mm]. Slices were positioned at an oblique angle to cover the entire brain, except for the most dorsal part of the parietal cortex. Each run of the main experiment comprised 293 brain volumes, and each run of the face localizer had 227 brain volumes.

Functional images were preprocessed using Statistical Parametric Mapping (SPM12; Wellcome Department of Imaging Science, London; RRID:SCR_007037; <http://www.fil.ion.ucl.ac.uk/spm>) operating in MATLAB (version R2013b, The MathWorks; RRID:SCR_001622). The first three EPI images in each run served as dummy scans to allow for T1-equilibration effects and were discarded before preprocessing. Data from each of the two scanning sessions, which took place on different days, were first preprocessed independently with the following steps for each session. Images within each brain volume were slice-time corrected using the middle slice as a reference, and were then realigned to correct for head movements using the first image as a reference. The participants' structural image in native space was coregistered to the realigned mean functional image, and was segmented into gray matter, white matter, and cerebrospinal fluid. Functional images from the main experimental runs were not smoothed, whereas images from the localizer runs were smoothed with a 4 mm Gaussian kernel (FWHM). To align the functional images from the two scanning sessions, the structural image from the first session was used as a template, and the structural image from the second session was coregistered to this template; we then applied the resulting transformation to all the functional images from the second session.

Functional localizers and definition of ROIs

Face-selective regions were defined using a dynamic face localizer that presented famous and nonfamous faces, along with a control condition consisting of objects and scenes. The stimuli were silent, nonspeaking videos of moving faces, and silent videos of objects and scenes, presented in an event-related design. Participants completed between one and two runs of the localizer across the two scanning sessions. The localizer presented different stimuli in each of two runs. For full details of the localizer, see Tsantani et al. (2019).

Functional ROIs were defined using the Group-Constrained Subject-Specific method (Fedorenko et al., 2010; Julian et al., 2012), which has the advantage of being reproducible and reducing experimenter bias by providing an objective means of defining ROI boundaries. Briefly, subject-specific ROIs were defined by intersecting subject-specific localizer contrast images with group-level masks for each ROI obtained from an independent dataset. In this study, we obtained group masks of face-selective regions (right FFA [rFFA], right OFA [rOFA], and right pSTS [rpSTS]) from a separate group of participants who completed the same localizer (for details, see Tsantani et al., 2019). We focused on face-selective regions from the right hemisphere because they have been shown to be more consistent and larger compared with the left hemisphere (e.g., Rossion et al., 2012). Our masks are publicly available at <https://doi.org/10.17633/rd.brunel.6429200.v1>.

Contrast images were defined for each individual participant. Face selectivity was defined by contrasting activation to faces versus nonface stimuli using *t* tests. We then intersected these subject-specific contrasts with the group masks, and extracted all significantly activated voxels at $p < 0.001$ (uncorrected) that fell within the boundaries of each mask. In cases where the resulting ROI included fewer than 30 voxels, the threshold was lowered to $p < 0.01$ or $p < 0.05$. ROIs that included fewer than 30 voxels at the lowest threshold were not included, and this occurred for the rFFA in 2 participants and for the rOFA in 1 participant. For full details of size and location of all ROIs, see Tsantani et al. (2019).

Experimental design and statistical analysis

Main experimental fMRI runs

In the main experimental runs, face stimuli were presented intermixed with voice stimuli within each run in an event-related design. The

experiment was programmed using the Psychophysics Toolbox (version 3; RRID:SCR_002881) (Brainard, 1997; Pelli, 1997) in MATLAB and was displayed through a computer interface inside the scanner. Participants were instructed to fixate on a small square shape that was constantly present in the center of the screen. From a distance of 85 cm, visual stimuli subtended 20.83 × 12.27 degrees of visual angle on the 1024 × 768 pixel screen.

The experiment was presented in two scanning sessions, with three runs in each session. Each run featured two unique videos of the face of each of the 12 identities, presented twice. Each run therefore contained 48 face trials (12 identities × 2 videos × 2 presentations), intermixed with 48 voice trials (96 experimental trials in total). In other words, across all three runs within a session, each of the 12 face identities appeared in 12 trials, featuring six unique videos of their face. Stimuli were presented in a pseudorandom order that prohibited the succeeding repetition of the same stimulus and ensured that each identity could not be preceded or succeeded by another identity more than once within the same modality. Each trial presented a stimulus for 3000 ms and was followed by a 1000 ms intertrial interval (Fig. 1).

To maintain attention to stimulus identity in the scanner, participants performed an anomaly detection task in which they indicated via button press when they were presented with a famous face or voice that did not belong to one of the 12 famous individuals that they had been familiarized with before the experiment. Therefore, each run also included 12 randomly presented task trials (six faces and six voices). Finally, each run contained 36 randomly interspersed null fixation trials, resulting in a total of 144 trials in each run lasting ~10 min.

The three experimental runs that were completed in the first scanning session were repeated in the second session with the same stimuli, but in a new pseudorandom order. The task stimuli, however, were always novel for each run. The three runs, which had different face videos, were presented in counterbalanced order across participants in both sessions.

Behavioral session

All participants completed a behavioral session in a laboratory, which took place on a separate day and always after the fMRI sessions had been completed. In this session, participants rated the same faces with which they had been presented in the scanner on perceived social traits and on perceived pairwise visual similarity. Participants also rated voices (the order of tasks was counterbalanced across modality), but these results are not presented here. All tasks and stimuli were presented using the Psychophysics Toolbox and MATLAB.

Social Trait Judgment Tasks. In the Social Trait Judgment Tasks, participants were asked to make judgments about the perceived Trustworthiness, Dominance, Attractiveness, and positive-negative Valence of the face identities. There were four blocks, one for each judgment, and their order was counterbalanced across participants. Face stimuli were presented in the center of the screen. In contrast to the fMRI runs, in which stimuli were presented for the full 3 s of their duration, here all stimuli were only presented for the first 1500 ms of their duration, to reduce testing time.

All blocks followed the same trial structure (Fig. 1). In each trial, a face identity was presented with three videos: these were presented successively with no gap in between them (total of 4500 ms). Participants were then asked to rate how trustworthy/dominant/attractive/negative-positive the face was, and they were asked to base their judgment on all three videos of the face. The rating scale ranged from 1 (very untrustworthy/nondominant/unattractive/negative) to 7 (very trustworthy/dominant/attractive/positive), and participants responded using the corresponding keys on the keyboard. There was a 1000 ms intertrial interval following the response.

Each identity was presented in two trials: one trial presented three face videos randomly selected from the six available, and the other trial presented the remaining three videos. This resulted in 24 trials in each block (12 identities × 2 presentations). The videos within each trial were presented in a random order, and the trial order was also randomized. Trustworthiness was defined as “able to be relied on as honest and truthful.” Dominance was defined as “having power and influence over other people.” No definition was deemed necessary for valence or

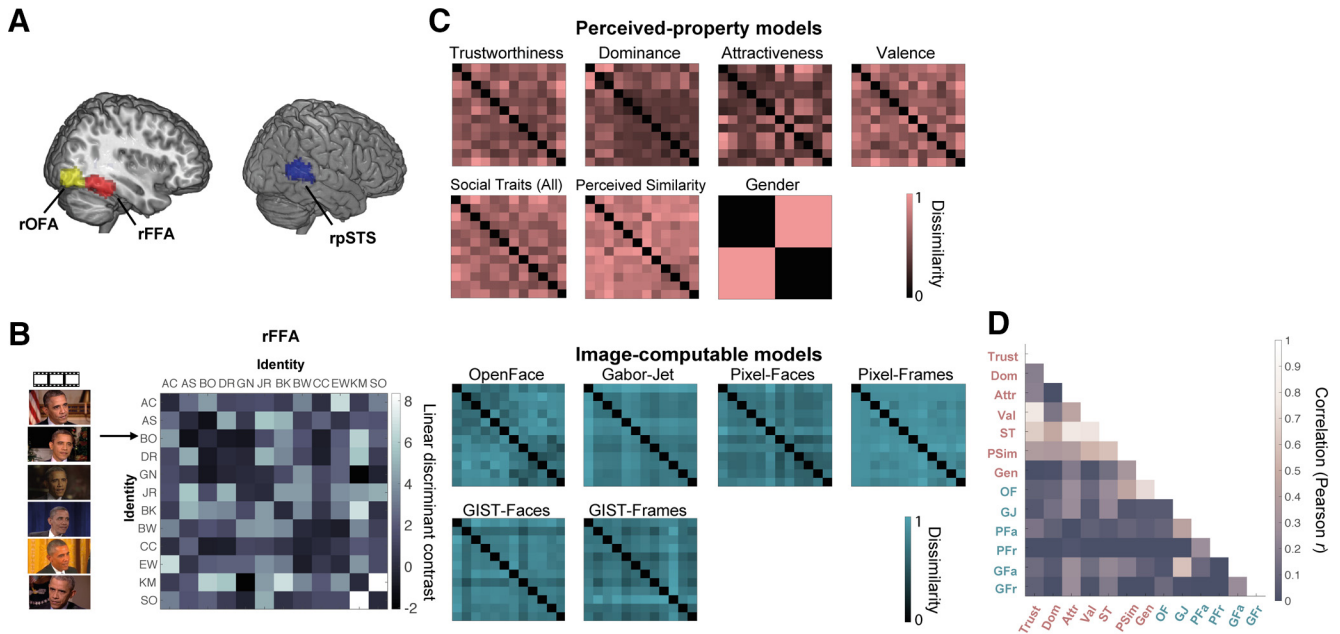


Figure 2. Brain and model representational dissimilarity matrices (RDMs). **A**, Location in MNI space of the three face-selective regions localized in our participants: rOFA, rFFA, and rpSTS (all regions in the right hemisphere). These probabilistic maps were created for illustration purposes (in our analyses, we only used subject-specific ROIs) and show all voxels that were present in at least 20% of participants. **B**, Example brain RDM for the rFFA. For each ROI and each participant, we computed RDMs showing the dissimilarity of the brain response patterns between all pairs of identities. Each row and column represent one identity, and response patterns are based on all six presented videos of that identity. Each cell represents the LDC distance between the response patterns of two identities (higher values indicate higher dissimilarity), cross-validated across runs presenting different videos of the face of each identity. The matrix is symmetric around a diagonal of zeros. **C**, Model RDMs for Image-computable properties (blue) and Perceived properties (pink). These models are in the same format as the brain RDMs and show the dissimilarity between two identities on each property (see Materials and Methods). Image-computable models include a neural network trained to distinguish between face identities (OpenFace), a Gabor-Jet model, Pixel Dissimilarity (both for faces only, Pixel-Faces; and the whole frames, Pixel-Frames), and a GIST Descriptor model (both for faces, GIST-Faces; and the whole frames, GIST-Frames). The RDMs computed per image (before averaging across identity) are shown in Extended Data Figure 2-1, although those 72×72 RDMs were not used in any analysis. Perceived-property models include perceived social traits [Trustworthiness, Dominance, Attractiveness, Valence, Social Traits (All)], Perceived Similarity, and Gender. Models based on participant ratings were averaged across participants. All models were built based on multiple images (Image-computable models) or videos (Perceived-property models) of the face of each identity. For visualization purposes, all model RDMs were scaled to a range between 0 (no dissimilarity) and 1 (maximum dissimilarity). **D**, Correlations (Pearson) between the different model RDMs. The different candidate models were compared with each other using Pearson correlation. Extended Data Figure 2-2 shows this same matrix with added correlation values.

attractiveness. Participants were advised that there was no time limit to their responses and that they should follow their first judgment. The duration of each block was ~3 min.

Pairwise Visual Similarity Task. In the Pairwise Visual Similarity Task, participants rated the perceived visual similarity of pairs of face identities. Each of the 12 identities was paired with the other 11 identities creating 66 identity pairs. Each identity was presented by three videos, randomly selected from the six available videos. Each identity pair was presented in two trials, counterbalancing the presentation order of each identity in the pair. There were therefore 132 trials in each task (66 identity pairs \times 2 presentations). The presentation order of the Pairwise Similarity Task in relation to the Social Trait Judgment Tasks was also counterbalanced across participants.

Participants were instructed to rate the similarity between the visual appearance of the two face identities in each pair, focusing on the facial features. Participants were asked to rate how similar the two faces looked on a scale from 1 (very dissimilar) to 7 (very similar). Participants were advised that there was no time limit to their responses and that they should follow their first instinct. Participants were told to ignore similarities between people that were related to biographical or semantic information (e.g., if both identities were actors). Furthermore, to encourage participants to base their judgments on perceptual information, participants were advised to consider to what extent two identities could potentially be related to each other (i.e., be part of the same family) based on how they looked.

In each trial, participants were first presented with the three videos of the face of one identity (Fig. 1). Following a 500 ms fixation screen, they were presented with the three videos of the face of the second identity. Videos for each identity were presented successively with no gap in between. Each video was presented for 1500 ms, and there was a 1000 ms

intertrial interval following the response. The presentation order of the trials was randomized. The duration of each task was ~30 min.

Brain representational dissimilarity matrices (RDMs)

RDMs showing the discriminability of the brain response patterns elicited by the 12 face identities (during the fMRI experimental runs) were created for each individual participant and for each ROI.

First, to obtain brain responses at each voxel for each of the 12 face identities, mass univariate time-series models were computed for each participant using a high-pass filter cutoff of 128 s and autoregressive AR(1) modeling to account for serial correlation. Regressors modeled the BOLD response at stimulus onset and were convolved with a canonical HRF. We defined a model for each run separately, and for every possible pair of runs within a scanning session (by concatenating the two runs), to create data partitions for cross-validation (described below). Each model contained a regressor for the face of each of the 12 identities, which incorporated the different videos of their face (two per run) and the repetitions of those videos. The model also included regressors for each of the 12 voice identities, task trials, and the six motion parameters obtained during the image realignment preprocessing stage (included as regressors of no interest).

Second, within each ROI, we extracted the β estimates at each voxel for each of the 12 face identities. This resulted in 12 vectors of β values per ROI that described the response patterns (across voxels) elicited by the 12 face identities.

Third, these vectors of β estimates were used to compute 12×12 Face RDMs in face-selective ROIs, in which each cell showed the distance between the response patterns of two identities (Fig. 2B). RDMs were computed using the linear discriminant contrast (LDC), a cross-validated distance measure (Nili et al., 2014; Walther et al., 2016), which

we implemented using in-house MATLAB code and the RSA toolbox (Nili et al., 2014). Two RDMs were created for each ROI, one for each scanning session. Each RDM was computed using leave-one-run-out cross-validation across the three runs, which presented different stimuli for each identity. Therefore, RDMs showed the dissimilarities between face identities, rather than specific face videos. In each cross-validation fold, concatenated data from two runs formed Partition A, and data from the left-out run formed Partition B. For each pair of identities (e.g., ID1 and ID2), Partition A was used to obtain a linear discriminant, which was then applied to Partition B to test the degree to which ID1 and ID2 could be discriminated. Under the null hypothesis, LDC values are distributed ~ 0 when two patterns cannot be discriminated. Values > 0 indicate higher discriminability of the two response patterns (Walther et al., 2016).

The discriminability of face identities in each ROI was computed by calculating the mean LDC across all cells of each participant's RDM, and comparing the mean LDC distances against 0 (Tsantani et al., 2019).

Full details of this analysis are presented in Tsantani et al. (2019), and the data to compute brain RDMs are available at <https://doi.org/10.17633/rd.brunel.6429200.v1>. Here, we used the RDMs for three face-selective regions (rFFA, rOFA, and rpSTS). All three of these regions showed significant discriminability of face identities.

RDMs based on Image-computable properties

We computed dissimilarities between the 12 face identities based on visual descriptions of their faces obtained using the models described below. We did not use the full videos as input to these models, but instead extracted one still frame from each face video used in the experiment (typically the first frame in which the full face was visible and the image was not blurred). Thus, we obtained six different images of the face of each identity, taken from the six different videos in which the identity was presented, resulting in 72 images in total.

OpenFace model. The OpenFace model RDM was computed from low-dimensional face representations obtained from OpenFace (Amos et al., 2016) (<http://cmusatyalab.github.io/openface/>). Briefly, OpenFace uses a deep neural network that has been pretrained (using 500,000 faces) to learn the best features or measurements that can group two pictures of the same identity together and distinguish them from a picture of a different identity. We used this pretrained neural network to generate measurements for each of our face pictures and to compare these measurements between each pair of pictures. OpenFace first performs face detection, identifies prespecified landmarks, and does an affine transformation so that the eyes, nose, and mouth appear in approximately the same location. The faces are then passed on to the pretrained neural network to generate 128 descriptor measurements for each face. To create an RDM, we used the program's calculated distances between the measurements for each pair of faces images. A value of 0 indicates that two images are identical, and values between 0 and 1 suggest that two different images likely show the same person's face. Values > 1 indicate that the two images show the faces of two different people. We found that OpenFace performed well at grouping different images of the same person's face compared with images of different people's faces in our image set (Extended Data Fig. 2-1 includes full 72×72 matrices showing distances between all images, but these full matrices were not used in any analysis). To obtain a 12×12 RDM for the 12 identities, which would be comparable to the brain RDMs, we computed the mean of all cells that showed images of the same identity pair (Fig. 2C). The 12×12 RDMs were used in all analyses.

Gabor-Jet model. The Gabor-Jet model RDM was computed from visual descriptors of face images obtained using the Gabor-Jet model (Biederman and Kalocsai, 1997; Yue et al., 2012; Margalit et al., 2016). This model was designed to simulate response properties of cells in area V1, and has been found to correlate with psychophysical measures of facial similarity (Yue et al., 2012). In addition, Carlin and Kriegeskorte (2017) showed that the dissimilarity of response patterns to different faces in the FFA was predicted by image properties based on Gabor filters. First, we used OpenFace 2.0 (Baltrusaitis et al., 2018) to automatically detect the faces in each image, and the pictures were grayscaled. The MATLAB script provided in http://www.geon.usc.edu/GWTgrid_

simple.m was then used to create a 100×40 Gabor descriptor for each face. After transforming these matrices into vectors, we computed the Euclidean distance between the vectors from each pair of faces (Extended Data Fig. 2-1), and then averaged the distances across all pairs of stimuli that showed the same two identities, resulting in a 12×12 RDM (Fig. 2C).

GIST model (faces only and whole frames). The GIST model RDMs were computed from visual descriptors of pictures obtained using the GIST model (Oliva and Torralba, 2001). The GIST model estimates information about the spatial envelope of scenes, and it is related to perceived dimensions of naturalness, openness, roughness, expansion, and ruggedness. Weibert et al. (2018) showed that the similarity between the representations of different faces in the FFA, OFA, and posterior STS was predicted by the similarity of the different pictures computed using the GIST descriptor model. We extracted GIST descriptors both from the full picture (whole Frames) and just from the face (Faces only; we used the same stimuli as in the Gabor-Jet model). We then used the MATLAB script provided in <http://people.csail.mit.edu/torralba/code/spatialenvelope> to compute GIST descriptors for each picture, and computed Euclidean distances between each pair of pictures (Extended Data Fig. 2-1). We finally averaged the distances across all pairs of stimuli that showed the same two identities, resulting in 12×12 RDMs (Fig. 2C).

Pixel model (faces only and whole frames). Finally, we computed model RDMs based on pixel dissimilarity between each pair of pictures. As for the GIST model, we computed this model both for the full picture (whole Frames) and just for the face (Faces only). We extracted pixel grayscale values for each image, computed Pearson correlations between the vectors of each pair of images, and used correlation distance as the output measure ($1 - r$) (Extended Data Fig. 2-1). We finally averaged the distances across all pairs of stimuli that showed the same two identities, resulting in 12×12 RDMs (Fig. 2C).

RDMs based on Perceived properties

Social trait models: Trustworthiness, Dominance, Attractiveness, Valence, Social Traits (All). RDMs for ratings of the 12 face identities on Trustworthiness, Dominance, Attractiveness, and positive-negative Valence were computed using Euclidean distances. For each participant and each social trait, the Euclidean distance between the ratings of each pair of identities was calculated (ratings were averaged across the two trials in which the same identity was presented), resulting in a 12×12 RDM per trait. We then averaged the matrices for the same trait across participants (Fig. 2C).

We also created Social Traits (All) RDMs combining all four Social Traits, by calculating the Euclidean distance between all trait ratings for each pair of identities, resulting in a 12×12 trait RDM per participant. We then computed the mean matrix for all Social Traits across participants (Fig. 2C).

To get estimates of the intersubject reliability of these models, we computed the correlations between each participant's RDM and the average RDMs across all participants (i.e., the RDMs that we used as models), and then averaged the correlations across participants. The reliabilities were $r = 0.34$ for Trustworthiness, $r = 0.48$ for Dominance, $r = 0.67$ for Attractiveness, $r = 0.31$ for Valence, and $r = 0.48$ for Social Traits (All). We also computed the average correlations between each participant's RDM and the average RDM of all remaining participants. These reliabilities were $r = 0.24$ for Trustworthiness, $r = 0.42$ for Dominance, $r = 0.63$ for Attractiveness, $r = 0.20$ for Valence, and $r = 0.42$ for Social Traits (All).

Perceived Similarity model. The judgments in the Pairwise Visual Similarity Task indicated the degree of visual similarity between all possible pairs of identities. These ratings were averaged across the two trials in which each identity-pair was presented, and were reverse-coded to match the LDC and Euclidean distance measures, where a higher value indicates higher dissimilarity. The resulting values were arranged into a 12×12 face RDM for each participant and were then averaged across participants (Fig. 2C).

Intersubject reliability, estimated by computing the average correlation between each participant's RDM and the average RDM across all participants, was $r = 0.65$. Reliability computed as the average correlation

between each participant's RDM and the average RDM of all remaining participants was $r = 0.61$.

Gender model. Finally, a 12×12 RDM for Gender was constructed by assigning a value of 0 to same Gender identity pairs, and a value of 1 to different-Gender identity pairs (Fig. 2C).

Correlations between all 13 models are presented in Figure 2D and Extended Data Figure 2-2.

Individual model analysis: RSA comparing brain RDMs to candidate model RDMs using correlation

For each individual participant and each ROI, we compared the brain RDM for faces with each of the candidate model RDMs defined above using Pearson correlation (Fig. 3A). We then tested whether the correlations across participants for each ROI were significantly >0 , using two-sided one-sample Wilcoxon signed-rank tests (Nili et al., 2014). p values were corrected for multiple comparisons using FDR correction ($q = 0.05$) across all 13 comparisons for each ROI. We also compared the correlations across all pairs of models within each ROI, to test which model was the best predictor of the variance in brain RDMs in each ROI. For these pairwise comparisons, we used two-sided Wilcoxon signed-rank tests and only significant FDR-corrected values (for 78 comparisons) are reported.

An estimate of the noise ceiling was calculated for each ROI, to estimate the maximum correlation that any model could have with the brain RDMs in each ROI given the existing noise in the data. We estimated the noise ceiling using the procedures described by Nili et al. (2014). The lower bound of the noise ceiling was estimated by calculating the Pearson correlation of the brain RDM for each participant with the average brain RDM across all other participants (after z scoring the brain RDM for each participant). The upper bound of the noise ceiling was estimated by computing the Pearson correlation of the brain RDM for each participant with the average brain RDM across all participants (after z scoring the brain RDM for each participant).

Weighted model-combination analysis: weighted representational modeling

We also used weighted representational modeling (Khaligh-Razavi and Kriegeskorte, 2014; Jozwik et al., 2016, 2017) to combine individual models via reweighting and thus investigate whether combinations of different model RDMs could explain more variance in representational geometries than any single model. For each combined model, we used linear non-negative least squares regression (lsqnonneg algorithm in MATLAB) to estimate a weight for each component of the combined model. We fitted the weights and tested the performance of the reweighted (combined) model on nonoverlapping groups of both participants and stimulus conditions within a cross-validation procedure, and used bootstrapping to estimate the distribution of the combined model's performance (Storrs et al., 2020).

We used six different combinations of component models: *Image-computable* properties (OpenFace, GIST, GaborJet, and Pixel), *Social Traits* (comprising a weighted combination of the Trustworthiness, Dominance, Attractiveness, and Valence properties), *Perceived properties* (Trustworthiness, Dominance, Attractiveness, Valence, Perceived Similarity, and Gender), *Low-Level* properties (GIST, GaborJet, and Pixel), *High-Level* properties (Trustworthiness, Dominance, Attractiveness, Valence, Perceived Similarity, Gender, and OpenFace), and *All properties*.

Within each cross-validation fold, data from 8 participants for four stimulus identity conditions were assigned to serve as test data, and the remainder were used to fit the weights for each component of each of the six combined models. Because the cross-validation was performed within a participant-resampling bootstrap procedure, the number of participant data RDMs present in each cross-validation fold was sometimes smaller than eight (when a participant was not present in the bootstrap) or larger than eight (when a participant was sampled multiple times in the bootstrap). All data from the same participant were always assigned only to either the training or test split. A reweighting target RDM was constructed by averaging the training-split participants' RDMs for training-split stimulus conditions, and weights were fitted to the components of each combined model to best predict this target RDM. The six

resulting combined models, as well as the 13 individual models, were then correlated separately with each of the brain RDMs from test participants for test conditions, using Pearson correlation. The noise ceiling was also computed within every cross-validation fold using the same procedure as for the main analysis. In other words, we correlated (Pearson correlation) each test participant's RDM with the average of all other test RDMs, excluding their own (for the lower bound of the noise ceiling) and with the average of all test participants' RDMs including their own (for the upper bound of the noise ceiling). This procedure was repeated for 30 participant cross-validation folds within 30 stimulus-condition cross-validation folds to provide a stabilized estimate of the noise ceiling and the performance of each model (Storrs, et al., 2020).

The cross-validation procedure was repeated for 1000 bootstrap resamplings of participants for each face-selective ROI. From the resulting bootstrap distribution, we computed the mean estimate of the lower bound of the noise ceiling, as well as the mean of each model's correlation with human data for both individual models and combined models (Fig. 3B). Correlations between model and brain RDMs were considered significantly >0 if the 95% CI of the bootstrap distribution did not include 0. Bonferroni correction was applied to correct for multiple comparisons. Finally, we compared each pair of models by testing whether the distributions of the differences between each pair of models contained 0. We only report pairwise differences that were significant after Bonferroni correction. Code for this analysis was adapted from the following: https://github.com/tinyrobots/reweighted_model_comparison.

Data and code accessibility

Data and code for main analysis are available as follows: <https://doi.org/10.25383/city.11890509.v1>.

Results

We tested 30 participants in an fMRI experiment, in which they were presented with faces of 12 famous people (same fMRI data as in Tsantani et al., 2019), and in a separate behavioral experiment, in which participants rated the faces of the same people on perceived similarity and social traits (Fig. 1). We then computed RDMs showing the representational distances between the brain response patterns elicited by the face identities in the face-selective rFFA, rOFA, and rpSTS. The distance measure that we used to compute the RDMs was the LDC, which is a cross-validated estimate of the Mahalanobis distance (Walther et al., 2016). The mean LDC across each RDM showed that response patterns to different face identities were discriminable in all three regions (Tsantani et al., 2019). To investigate the informational content of brain representations of the face identities in each face-selective region, we used RSA (Kriegeskorte et al., 2008a,b) to compare the brain RDMs with a diverse set of candidate model RDMs (Fig. 2). We used candidate models based on the physical properties of the stimuli (Image-computable models), including low-level stimulus properties (based on Pixel-wise, GIST [Oliva and Torralba, 2001], and Gabor-jet [Biederman and Kalocsa, 1997] dissimilarities) and higher-level image-computable descriptions obtained from a deep neural network trained to cluster faces according to identity (OpenFace; see Materials and Methods) (Amos et al., 2016). Additionally, we used candidate models based on perceived higher-level properties (perceived-property models), including Gender and participants' ratings of the face identities on Perceived Similarity and social traits (Trustworthiness, Dominance, Attractiveness, Valence, and Social Traits (All), which corresponds to all traits combined) in a behavioral experiment.

Individual model analysis

In our main analysis, we computed Pearson's correlations between RDMs in the rFFA, rOFA, and rpSTS, and each

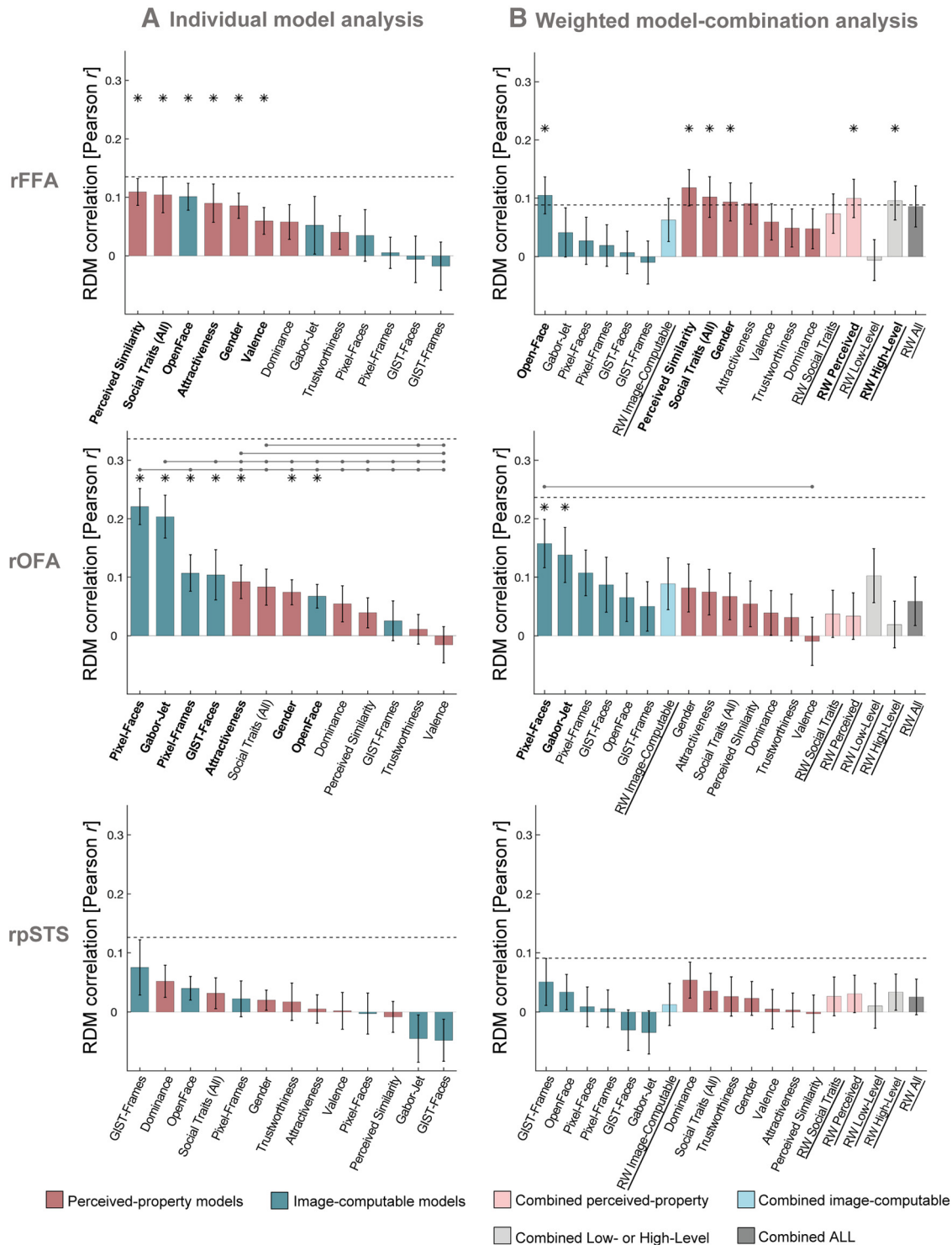


Figure 3. rFFA and rOFA show distinct representational profiles of face identity information. **A**, Similarity (Pearson correlations) between brain RDMs (in rFFA, rOFA, and rpSTS) and each of the individual candidate models. Bars represent mean correlations across participants. Error bars indicate SE. Blue represents correlations with Image-computable models. Pink represents Perceived-property models. Horizontal dashed lines indicate the lower bound of the noise ceiling. An asterisk above a bar and the name of the model in bold indicate that correlations with that model were significantly >0 . Correlations with individual models are sorted from highest to lowest. Horizontal lines above bars indicate significant differences between the correlations of the first marked column with the subsequent marked columns (FDR-corrected for multiple comparisons). Full results are given in Table 1, and single-subject data are shown in Figure 4. **B**, Similarity (Pearson correlations) between brain RDMs (in rFFA, rOFA, and rpSTS) and each of the candidate models in the weighted representational modeling analysis. Bars represent mean correlations. Error bars indicate SE across 1000 bootstrap samples. Horizontal dashed lines indicate the lower bound of the noise ceiling, averaged across bootstrap samples. An asterisk above a bar and the name of the model in bold indicate that correlations with that model were significantly >0 . Correlations with individual models are blocked by type of model (Image-computable models followed by Perceived-property models) and sorted from highest to lowest. RW refers to the combined and reweighted models. Light blue represents models that combine Image-computable properties. Light pink represents models that combine Perceived properties. Gray represents models that combine both types of properties. None of the combined models outperformed individual models. Full results are reported in Table 2. The results of both analyses show that in the rFFA, the models that explained most of the variance are related to high-level properties, such as perceived properties of the stimuli and the Image-computable OpenFace model of face recognition. In contrast, brain RDMs in rOFA correlated mainly with low-level Image-computable properties, such as pixel dissimilarity and the Gabor-Jet model. No significant correlations were found in rpSTS.

Table 1. Results of individual model analysis^a

| | Pearson correlation between RDMS | | | | Noise ceiling (lower bound, upper bound) |
|----------------------|----------------------------------|-------|----------|---------------------------------|--|
| | Mean <i>r</i> | SE | <i>Z</i> | <i>p</i> < 0.05 (FDR-corrected) | |
| rFFA | | | | | 0.135, 0.262 |
| Perceived Similarity | 0.109 | 0.023 | 3.689 | Yes | |
| Social Traits (All) | 0.104 | 0.031 | 2.710 | Yes | |
| OpenFace | 0.101 | 0.023 | 3.461 | Yes | |
| Attractiveness | 0.090 | 0.033 | 2.687 | Yes | |
| Gender | 0.086 | 0.021 | 3.302 | Yes | |
| Valence | 0.060 | 0.023 | 2.391 | Yes | |
| Dominance | 0.058 | 0.030 | 1.640 | No | |
| Gabor-Jet | 0.052 | 0.049 | 0.956 | No | |
| Trustworthiness | 0.040 | 0.029 | 1.594 | No | |
| Pixel-Faces | 0.035 | 0.044 | 0.865 | No | |
| Pixel-Frames | 0.005 | 0.027 | 0.159 | No | |
| GIST-Faces | −0.006 | 0.040 | 0.114 | No | |
| Pixel-Frames | −0.018 | 0.041 | −0.478 | No | |
| rOFA | | | | | 0.337, 0.408 |
| Pixel-Faces | 0.221 | 0.031 | 4.357 | Yes | |
| Gabor-Jet | 0.204 | 0.037 | 3.968 | Yes | |
| Pixel-Frames | 0.107 | 0.031 | 3.016 | Yes | |
| GIST-Faces | 0.104 | 0.043 | 2.216 | Yes | |
| Attractiveness | 0.092 | 0.029 | 2.843 | Yes | |
| Social Traits (All) | 0.083 | 0.031 | 1.979 | No | |
| Gender | 0.074 | 0.021 | 2.757 | Yes | |
| OpenFace | 0.067 | 0.020 | 2.952 | Yes | |
| Dominance | 0.055 | 0.031 | 1.546 | No | |
| Perceived Similarity | 0.039 | 0.026 | 1.416 | No | |
| GIST-Frames | 0.025 | 0.034 | 0.746 | No | |
| Trustworthiness | 0.011 | 0.025 | 0.400 | No | |
| Valence | −0.016 | 0.031 | −0.573 | No | |
| rpSTS | | | | | 0.126, 0.252 |
| GIST-Frames | 0.075 | 0.047 | 1.800 | No | |
| Dominance | 0.052 | 0.027 | 1.800 | No | |
| OpenFace | 0.040 | 0.020 | 2.129 | No | |
| Social Traits (All) | 0.032 | 0.026 | 1.018 | No | |
| Pixel-Frames | 0.022 | 0.030 | 0.956 | No | |
| Gender | 0.020 | 0.017 | 0.956 | No | |
| Trustworthiness | 0.017 | 0.032 | 0.524 | No | |
| Attractiveness | 0.005 | 0.024 | 0.134 | No | |
| Valence | 0.002 | 0.031 | 0.051 | No | |
| Pixel-Faces | −0.003 | 0.035 | −0.113 | No | |
| Perceived Similarity | −0.008 | 0.026 | −0.072 | No | |
| Gabor-Jet | −0.045 | 0.040 | −1.100 | No | |
| GIST-Faces | −0.048 | 0.036 | −1.368 | No | |

^aValues correspond to the results presented in Figure 3A. For each ROI, we show the mean correlations between brain RDMS with each model, SE, *Z* statistics from two-sided one-sample Wilcoxon signed-rank tests, and whether correlations were significantly >0. We also show the estimated lower and upper bounds of the noise ceiling for each ROI. Models are ordered by effect size.

candidate model RDM. Correlations were computed for each individual participant, and then correlations across participants for each model were compared against 0 using two-sided one-sample Wilcoxon signed-rank tests. For each ROI and each model that showed significant correlations with participants' brain RDMS, we report below the mean correlation across participants, and the *Z* statistic and *p* value obtained from the signed-rank test, corrected for multiple comparisons using FDR correction. Full results are presented in Figure 3A and Table 1, and individual-subject correlations are presented in Figure 4. We also compared the correlations across all pairs of models using two-sided Wilcoxon signed-rank tests.

Brain RDMS in the rFFA had the highest mean correlation with the Perceived Similarity model (mean *r* = 0.11, *Z* = 3.69,

p = 0.0002), followed by perceived Social Traits (All) (mean *r* = 0.10, *Z* = 2.71, *p* = 0.0067), the Image-computable neural network OpenFace (mean *r* = 0.10, *Z* = 3.46, *p* = 0.0005), perceived Attractiveness (mean *r* = 0.09, *Z* = 2.69, *p* = 0.0072), Gender (mean *r* = 0.09, *Z* = 3.30, *p* = 0.0010), and Valence (mean *r* = 0.06, *Z* = 2.39, *p* = 0.0168) (Fig. 3A). We estimated the lower bound of the noise ceiling as the mean correlation between each participant's rFFA RDM and the average of all other participants' rFFA RDMS (Nili et al., 2014). This estimates the non-noise variance in the data, and is not overfit to the present data. None of the mean correlations reached the lower bound of the noise ceiling for the rFFA (*r* = 0.14); this suggests that there could be models outside those tested here that would better explain the representational distances in rFFA. Pairwise comparisons showed no significant differences between the correlations of any pairs of models (all *p* > 0.0041; no significant results after FDR correction).

In contrast with the rFFA, the brain RDMS in the rOFA had the highest mean correlations with low-level Image-computable models. The highest mean correlation was observed with the Pixel-Faces model (mean *r* = 0.22, *Z* = 4.36, *p* < 0.0001) (Fig. 3A), followed by the Gabor-Jet (mean *r* = 0.20, *Z* = 3.97, *p* < 0.0001), Pixel-Frames (mean *r* = 0.11, *Z* = 3.02, *p* = 0.0026), GIST-Faces (mean *r* = 0.10, *Z* = 2.22, *p* = 0.0267), perceived Attractiveness (mean *r* = 0.09, *Z* = 2.84, *p* = 0.0045), Gender (mean *r* = 0.07, *Z* = 2.76, *p* = 0.0058), and the OpenFace model (mean *r* = 0.07, *Z* = 2.95, *p* = 0.0032). None of the mean correlations reached the lower bound of the noise ceiling (*r* = 0.34). Pairwise comparisons between model correlations revealed that the Pixel-Faces model had significantly higher correlations with the rOFA RDMS than all other models (all *p* < 0.0058, FDR-corrected), except for the Gabor-Jet model and the GIST-Faces model. The Gabor-Jet model also had significantly higher correlations with the brain RDMS in rOFA than all other models (all *p* < 0.0058, FDR-corrected), except the Pixel-Faces and Pixel-Frames models. Perceived Attractiveness had significantly higher correlations with the rOFA RDMS than perceived Valence (*p* = 0.0051), and Social traits (All) was significantly higher than Trustworthiness and Valence (both *p* < 0.0018).

Finally, we investigated which model best explained the variance in representational distances in the rpSTS. We found no significant correlations between any of the candidate models and the brain RDMS in this region (all *p* > 0.0333; no significant results after FDR correction) (Fig. 3A). None of the models reached the lower bound of the noise ceiling (*r* = 0.13), and there were no significant differences between models (all *p* > 0.0140; no significant results after FDR correction).

These results show a clear distinction between the types of models that were associated with the representational geometries of face identities in the rFFA and rOFA. Representational distances of face identities in the rFFA were most associated with high-level Perceived Similarity, Gender, and Social Traits, as well as a high-level model of Image-computable properties (OpenFace), whereas representations in rOFA were most associated with low-level Image-computable properties. To test this directly, we compared the correlation profiles between the two regions. We first averaged all correlations per participant (after Fisher's transformation) for the same type of model (all Perceived-property models and all Image-computable models) for each ROI (rFFA and rOFA). In the rFFA, the mean correlation with Perceived-property models was 0.08 (SD = 0.095) and 0.03 (SD = 0.109) with Image-

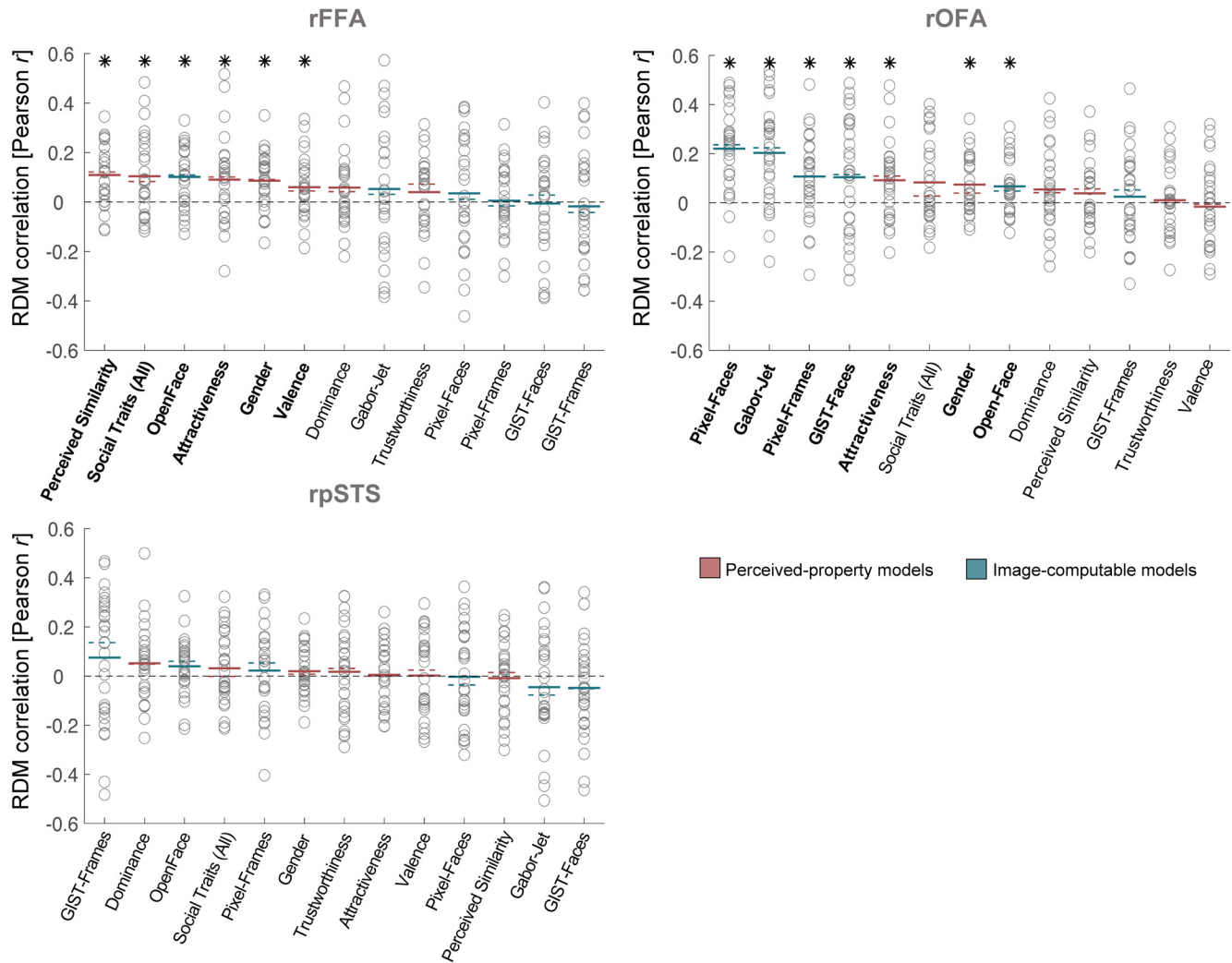


Figure 4. Similarity between brain RDMs (in rFFA, rOFA, and rpSTS) and each of the candidate models, showing individual participant data. This figure shows the same data as in Figure 3A, but with added individual data. Circles represent correlations for individual participants. Colored lines indicate mean (full lines) and median (dotted lines) correlations across participants. Pink represents correlations with models based on Perceived-property models. Blue represents correlations with Image-computable models. Horizontal black dashed lines indicate the 0 correlation point. An asterisk above a bar and the name of the model in bold indicate correlations that were significantly >0 . Correlations with individual models are sorted from highest to lowest based on the mean correlation across participants to match the format of Figure 3A.

computable models. In the rOFA, the mean correlation with Perceived-property models was 0.05 (SD = 0.108) and 0.13 (SD = 0.102) with Image-computable models. We then conducted a 2×2 repeated-measures ANOVA with ROI and type of model as variables. There was no main effect of ROI ($F_{(1,27)} = 3.37$, $p = 0.0773$) or type of model ($F_{(1,27)} = 0.36$, $p = 0.5519$), but there was a significant interaction between the two variables ($F_{(1,27)} = 23.75$, $p < 0.0001$). Pairwise comparisons (using two-sided Wilcoxon signed-rank tests) showed that, in the rFFA, the correlations with Perceived-property models were significantly higher than correlations with Image-computable models ($Z = 2.25$, $p = 0.0242$), whereas in the rOFA, correlations with Perceived-property models were significantly lower than correlations with Image-computable models ($Z = -3.17$, $p = 0.0015$). We also divided the models into low-level properties (GIST, Gabor-Jet, and Pixel) and high-level properties (Trustworthiness, Dominance, Attractiveness, Valence, Perceived Similarity, Gender, and OpenFace), and computed means per participant and per ROI for each of these types of models. In the rFFA, there was a mean correlation of 0.08 (SD = 0.090) with high-level properties, and of 0.02 (SD = 0.157) with low-level properties. In the rOFA,

there was a mean correlation of 0.05 (SD = 0.102) with high-level properties, and of 0.16 (SD = 0.141) with low-level properties. A 2×2 repeated-measures ANOVA showed a significant effect of ROI ($F_{(1,27)} = 5.44$, $p = 0.0274$), no significant effect of model ($F_{(1,27)} = 0.43$, $p = 0.5201$), and a significant interaction between the two variables ($F_{(1,27)} = 21.64$, $p < 0.0001$). Pairwise comparisons showed that in the rFFA, the correlations with high-level models were significantly higher than correlations with low-level models ($Z = 2.21$, $p = 0.0272$), whereas in the rOFA, correlations with high-level models were significantly lower than correlations with low-level models ($Z = -3.25$, $p = 0.0011$). These results demonstrate the clear distinct patterns of correlations for the rFFA and rOFA.

Our Image-computable models used a single image from each video clip. We recomputed all models using 72 images per clip, and averaged the features across all images of the same clip. We then computed distances between video clips in the same manner as before, and averaged distances for each pair of identities, resulting in 12×12 RDMs for each model. The results were very similar when using 72 images per clip compared with one image per clip (Fig. 5A). We additionally showed that we

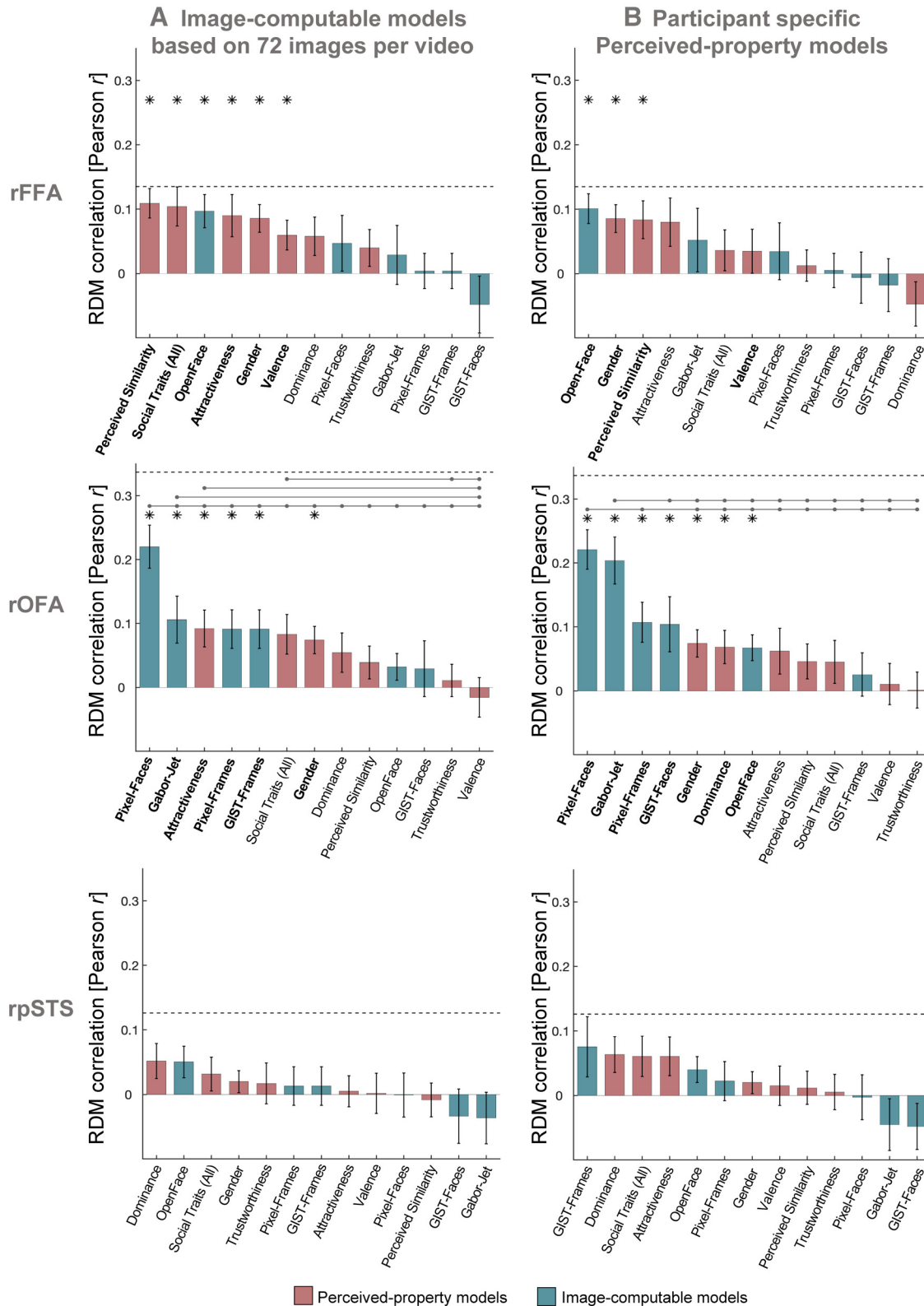


Figure 5. Control analyses with modified model RDMs. **A**, Similarity between brain RDMs (in rFFA, rOFA, and rpSTS) and each of the candidate models, using Image-computable models derived from 72 images per video. Our main analysis in Figure 3A used a single image per video to compute Image-computable models. Here, we repeated all analyses of Image-computable models using 72 frames for each video. We extracted 72 image frames for each video, and applied each model to each image. For each model, after extracting the features of each image of each video, we averaged the values for all images belonging to the same video. We then computed distances between videos in the same manner as before, and averaged distances for each pair of identities. We note that these results were very similar to the ones using just with one image per video, but some correlations were lower. **B**, Similarity between brain RDMs (in rFFA, rOFA, and rpSTS) and each of the individual candidate models, using behavioral models based on individual participant ratings. The analysis was the same as in Figure 3A; but instead of using average behavioral RDMs, each participant’s brain RDM was correlated to their own behavioral RDMs for Perceived Similarity, Trustworthiness, Dominance, Attractiveness, Valence, and Social Traits (All). The pattern of results looked very similar to the ones in Figure 3A, but correlations with Perceived-property models were overall lower when using each participant’s own model RDMs. Bars show mean correlations across participants and error bars show standard error. Horizontal dashed lines show the lower bound of the noise ceiling. An asterisk above a bar and the name of the model in bold indicate that correlations with that model were significantly higher than zero. Horizontal lines above bars show significant differences between the correlations of the first marked column with the subsequent marked columns (FDR corrected for multiple comparisons).

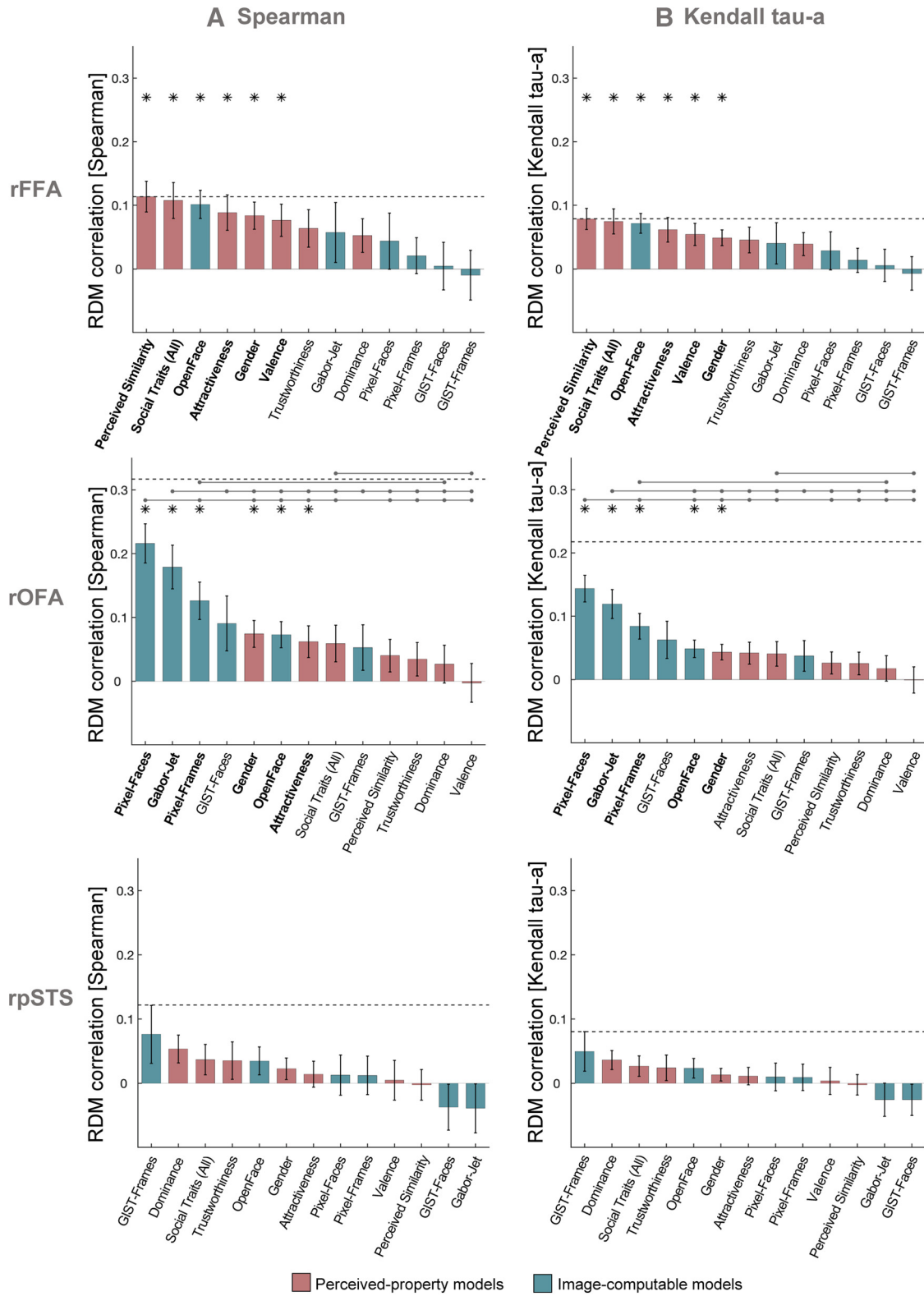


Figure 6. Control analyses using other similarity measures between RDMs. Similarity between brain RDMs (in rFFA, rOFA, and rpSTS) and each of the candidate models using Spearman correlation (A) and Kendall tau-a (B). These analyses were identical to the analysis using Pearson correlations (Fig. 3A), with the exception that noise ceiling was computed after rank-transforming the RDMs (Nili et al., 2014). The pattern of results was similar across all three correlation measures. Bars show mean correlations across participants and error bars show standard error. Horizontal dashed lines show the lower bound of the noise ceiling. An asterisk above a bar and the name of the model in bold indicate that correlations with that model were significantly higher than zero. Horizontal lines above bars show significant differences between the correlations of the first marked column with the subsequent marked columns (FDR corrected for multiple comparisons).

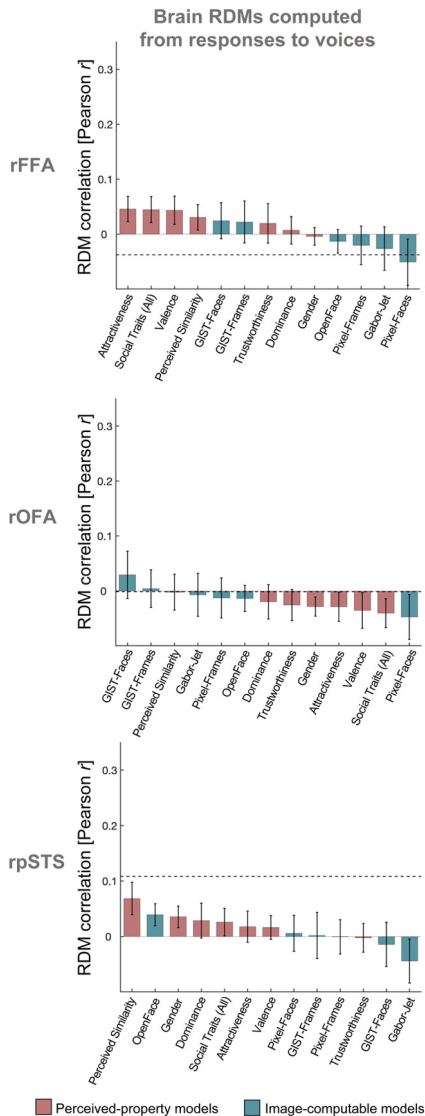


Figure 7. Control analysis with modified brain RDMs. Similarity between brain RDMs for voices (in rFFA, rOFA, and rpSTS) and each of the candidate models for faces. We computed RDMs from response patterns to voices in the rFFA, rOFA, and rpSTS, and compared them with our model RDMs for faces (same models as in Fig. 2). The voice stimuli belonged to the same 12 identities as the face stimuli and were presented interspersed among the face videos in the same runs (see Materials and Methods). RDMs for voice identities were computed using the same procedure as for face identities (see Materials and Methods) and were compared with model RDMs for faces using Pearson correlation. Bars show mean correlations across participants and error bars show standard error. Horizontal dashed lines show the lower bound of the noise ceiling. Correlations with individual models are sorted from highest to lowest. None of the correlations was significantly >0 after correction for multiple comparisons. Pairwise comparisons showed no significant differences between the correlations of any pairs of models.

obtained similar results to those in Figure 3A when using other similarity measures between RDMs (Spearman correlation, Kendall tau-a), demonstrating that these results are not dependent on using Pearson correlation (Fig. 6). Finally, we conducted an additional control analysis using brain RDMs in the same ROIs but built from response patterns to voices of the same individuals, instead of brain responses to faces. There were no significant correlations between any of the model RDMs for faces and brain RDMs for voices after correcting for multiple comparisons in the rFFA (all $p > 0.040$), rOFA (all $p > 0.103$), or rpSTS (all $p > 0.063$) (Fig. 7). Pairwise comparisons showed no significant

differences between the correlations of any pairs of models (all $p > 0.034$). The estimated lower bounds of noise ceilings for the voices brain RDMs were very low for rFFA ($r = -0.038$) and rOFA ($r = -0.001$), and higher for rpSTS ($r = 0.108$). This control analysis demonstrates that the above results for rFFA and rOFA are specific to visual stimuli (faces). To conclude, we find that the structure of the model correlations is reliable and is systematically different between the rFFA and rOFA.

Weighted model-combination analysis

Although our models accounted for a large portion of the explainable variance (based on the noise ceiling) in brain representations in the rFFA and rOFA, none of the mean correlations reached the lower bound of the noise ceiling. It could be that each individual model captured only a portion of the information represented in each brain region, in which case we may be able to fully explain the brain representations by combining multiple models. We thus used weighted representational modeling (Khaligh-Razavi and Kriegeskorte, 2014; Jozwik et al., 2016, 2017) to combine sets of models into weighted combinations via cross-validated fitting on the human data, and to investigate whether these combined models resulted in better predictions of the brain dissimilarities in each brain region (see Materials and Methods). We considered six different combined models: *Image-computable* properties (OpenFace, GIST, GaborJet, and Pixel), *Social Traits* (comprising a weighted combination of the Trustworthiness, Dominance, Attractiveness, and Valence properties), *Perceived* properties (Trustworthiness, Dominance, Attractiveness, Valence, Perceived Similarity, and Gender), *Low-Level* properties (GIST, GaborJet, and Pixel), *High-Level* properties (Trustworthiness, Dominance, Attractiveness, Valence, Perceived Similarity, Gender, and OpenFace), and *All properties*.

We used linear non-negative least squares regression to estimate a weight for each component of each combined model. We fitted the weights and tested the performance of the reweighted (combined) model on nonoverlapping groups of both participants and stimulus conditions within a cross-validation procedure, and used bootstrapping to estimate the distribution of the combined model's performance (Storrs et al., 2020). Figure 3B shows the results of this analysis. p values were corrected for multiple comparisons using Bonferroni correction. For the rFFA, the combined models for Perceived properties and High-Level properties had the highest mean correlations with the brain RDMs, and the individual-subject correlations were significantly >0 . For the rOFA, the combined model of all Low-Level properties and that of all Image-computable properties had the highest mean correlations with the brain RDMs, although the individual-subject correlations were not significantly >0 after correcting for multiple comparisons. Importantly, however, none of the combined models performed better than the best of the individual models (see full results in Table 2). Instead, the models with best performance in the previous (main) analysis also showed the highest correlations in this analysis. These results suggest that the models that best explained representational distances in each face-selective region share overlapping variance, given that combining them did not improve model performance. Last, replicating the findings of the previous analysis using more stringent statistical methods (cross-validation across stimuli and participants) provides further evidence of a reliable pattern of model correlations in rFFA and rOFA that reveals a distinction between the type of information encoded in these two regions.

Table 2. Results of weighted representational modeling analysis^a

| | | Pearson correlation between RDMS | | | Noise ceiling | |
|-------|----------------------|----------------------------------|-------|--|----------------------------|--|
| | | Mean <i>r</i> | SE | <i>p</i> < 0.05 (Bonferroni-corrected) | (lower bound, upper bound) | <i>p</i> < 0.05 (Bonferroni-corrected) |
| rFFA | | | | | 0.089, 0.286 | |
| | OpenFace | 0.105 | 0.032 | Yes | | No |
| | Gabor-Jet | 0.041 | 0.042 | No | | No |
| | Pixel-Faces | 0.027 | 0.040 | No | | No |
| | Pixel-Frames | 0.019 | 0.036 | No | | No |
| | GIST-Faces | 0.007 | 0.037 | No | | No |
| | GIST-Frames | −0.010 | 0.037 | No | | No |
| | RW Image-computable | 0.063 | 0.037 | No | | No |
| | Perceived Similarity | 0.118 | 0.031 | Yes | | No |
| | Social Traits (All) | 0.102 | 0.035 | Yes | | No |
| | Gender | 0.094 | 0.033 | Yes | | No |
| | Attractiveness | 0.091 | 0.035 | No | | No |
| | Valence | 0.059 | 0.031 | No | | No |
| | Trustworthiness | 0.049 | 0.033 | No | | No |
| | Dominance | 0.048 | 0.034 | No | | No |
| | RW Social Traits | 0.074 | 0.034 | No | | No |
| | RW Perceived | 0.100 | 0.033 | Yes | | No |
| | RW Low-Level | −0.006 | 0.035 | No | | No |
| | RW High-Level | 0.096 | 0.033 | Yes | | No |
| | RW ALL | 0.086 | 0.035 | No | | No |
| rOFA | | | | | 0.237, 0.372 | |
| | Pixel-Faces | 0.158 | 0.041 | Yes | | No |
| | Gabor-Jet | 0.138 | 0.047 | Yes | | No |
| | Pixel-Frames | 0.108 | 0.039 | No | | Yes |
| | GIST-Faces | 0.087 | 0.047 | No | | No |
| | OpenFace | 0.066 | 0.041 | No | | Yes |
| | GIST-Frames | 0.050 | 0.042 | No | | Yes |
| | RW Image-computable | 0.089 | 0.044 | No | | No |
| | Gender | 0.082 | 0.041 | No | | No |
| | Attractiveness | 0.075 | 0.039 | No | | Yes |
| | Social Traits (All) | 0.067 | 0.040 | No | | Yes |
| | Perceived Similarity | 0.055 | 0.039 | No | | Yes |
| | Dominance | 0.039 | 0.038 | No | | Yes |
| | Trustworthiness | 0.031 | 0.040 | No | | Yes |
| | Valence | −0.010 | 0.041 | No | | Yes |
| | RW Social Traits | 0.037 | 0.040 | No | | Yes |
| | RW Perceived | 0.033 | 0.040 | No | | Yes |
| | RW Low-Level | 0.103 | 0.046 | No | | No |
| | RW High-Level | 0.019 | 0.040 | No | | Yes |
| | RW ALL | 0.059 | 0.041 | No | | Yes |
| rpSTS | | | | | 0.091, 0.277 | |
| | GIST-Frames | 0.051 | 0.040 | No | | No |
| | OpenFace | 0.034 | 0.030 | No | | No |
| | Pixel-Faces | 0.009 | 0.034 | No | | No |
| | Pixel-Frames | 0.006 | 0.032 | No | | No |
| | GIST-Faces | −0.031 | 0.034 | No | | No |
| | Gabor-Jet | −0.038 | 0.037 | No | | No |
| | RW Image-computable | 0.013 | 0.036 | No | | No |
| | Dominance | 0.054 | 0.030 | No | | No |
| | Social Traits (All) | 0.035 | 0.030 | No | | No |
| | Trustworthiness | 0.026 | 0.033 | No | | No |
| | Gender | 0.023 | 0.029 | No | | No |
| | Valence | 0.005 | 0.033 | No | | No |
| | Attractiveness | 0.003 | 0.029 | No | | No |
| | Perceived Similarity | −0.003 | 0.032 | No | | No |
| | RW Social Traits | 0.026 | 0.033 | No | | No |
| | RW Perceived | 0.031 | 0.032 | No | | No |
| | RW Low-Level | 0.010 | 0.038 | No | | No |
| | RW High-Level | 0.033 | 0.031 | No | | No |
| | RW ALL | 0.025 | 0.030 | No | | No |

^a Values correspond to the results presented in Figure 3B. Within each ROI, we show the mean correlations between brain RDMS with each model (individual models and combined models), and whether correlations were significantly >0. We also show the estimated lower and upper bounds of the noise ceiling for each ROI, and whether correlations were significantly below the noise ceiling. Models are ordered by effect size and grouped first by Image-computable models, then Perceived-property models, and then models that combined both types of properties. RW refers to combined and reweighted models.

Individual differences and idiosyncratic representations

It is possible that there were substantial individual differences in face identity representations that limit the magnitude of the correlations between brain and model RDMs in our analyses. Brain and behavioral representations of face identities could be idiosyncratic and thus characteristic of each individual. We considered below three ways in which we could test this hypothesis.

First, we considered whether there were substantial individual differences in brain RDMs. To estimate the lower-bound of the noise ceiling, we had computed intersubject reliabilities of brain RDMs. If, however, there were substantial individual differences in the brain RDMs, we would expect that representational distances in each of the face-selective ROIs could be highly reliable within each participant but not across participants. We thus computed intrasubject reliabilities of brain RDMs by correlating the brain RDMs calculated independently from two separate testing sessions for each participant, and then averaging the correlations across participants. We note that, in all other analyses in the present manuscript, the brain RDMs for each participant corresponded to the average of these two sessions. For all three face-selective ROIs, we observed intrasubject reliabilities (rFFA: $r = 0.063$; rOFA: $r = 0.079$; rpSTS: $r = 0.094$) that were on average lower than the intersubject reliabilities (rFFA: $r = 0.135$; rOFA: $r = 0.337$; rpSTS: $r = 0.126$; see Table 1), suggesting that in fact, in this case, the brain RDMs were not more reliable within each individual. It is important to note, however, that there was much less data to compute intrasubject reliabilities than intersubject reliabilities.

Second, idiosyncratic brain representations could also result in higher correlations between each participant's brain RDM and behavioral RDMs based on their own ratings, compared with the average behavioral RDMs that we used in the main analyses. We thus repeated the main analysis using each individual's own RDMs for the rating-based Perceived-property models, namely, Perceived Similarity, Trustworthiness, Dominance, Attractiveness, Valence, and Social Traits (All). The results, however, did not reveal higher correlations when using these participant-specific behavioral models (Fig. 5B). In contrast, correlations with the participants' individual behavioral models were slightly lower than when using average behavioral models.

A third possibility is that idiosyncratic representational geometries could result in the variance of each participant's brain RDMs being best explained by a uniquely weighted combination of candidate models (even if no set of weightings would perform well for all participants). However, we did not have sufficient data per participant to test this possibility here.

Discussion

We aimed to investigate what information is explicitly encoded in the face-selective rFFA, rOFA, and rpSTS. We extracted fMRI patterns elicited by famous face identities in these regions, and computed face identity RDMs, which showed that face identities could be distinguished based on their elicited response patterns in all three regions. Using RSA, we compared the brain RDMs for the rFFA, rOFA, and rpSTS with multiple model RDMs ranging from low-level image-computable properties (pixel-wise, GIST, and Gabor-Jet dissimilarities), through higher-level image-computable descriptions (OpenFace deep neural network, trained to cluster faces by identity), to complex human-rated face properties (perceived visual similarity, social traits, and gender). We found that the rFFA and rOFA encode face identities in

a different manner, suggesting distinct representations in these two regions. The representational geometries of face identities in the rFFA were most associated with high-level properties, such as perceived visual similarity, social traits, gender, and high-level image features extracted with a deep neural network (OpenFace) (Amos et al., 2016). In contrast, the representational geometries of faces in the rOFA were most associated with low-level image-based properties, such as pixel similarity and features extracted with Gabor filters that simulate functioning of early visual cortex. While previous studies had shown that low-level properties of images extracted with Gabor filters were associated with representational distances of faces in rFFA (Carlin and Kriegeskorte, 2017; Weibert et al., 2018), our results suggest that representations in rFFA use more complex combinations of stimulus-based features and relate to higher-level perceived and social properties (see also Davidesco et al., 2014). These results inform existing neurocognitive models of face processing (Haxby et al., 2000; Duchaine and Yovel, 2015) by shedding light on the much-debated computations of face-responsive regions, and providing new evidence to support a hierarchical organization of these regions from the processing of low-level image-computable properties in the rOFA to higher-level visual features and social information in the rFFA.

Our initial prediction was that, by combining and reweighting different candidate models, we would be better able to explain the brain RDMs. However, we did not find evidence for this in any of our face-selective ROIs. These results suggest that, when more than one model was significantly correlated with the brain RDMs for a certain brain region, they tended to explain overlapping variance in the brain RDMs. For example, while Perceived Similarity and OpenFace both explained the representational geometries in rFFA, their combination did not explain more variance than each model individually. However, our pattern of results suggests a clear distinction between the types of models that are associated with representations in the rFFA and rOFA, with higher-level properties explaining more variance in the rFFA and lower-level image-based properties explaining more variance in the rOFA.

One crucial aspect of our study is that we used naturalistically varying video stimuli and multiple depictions for each identity. Brain RDMs were built by cross-validating the response patterns across runs featuring different videos of the face of each identity, and behavioral models were based on averages of ratings of multiple videos for each identity. Image-based models were built by calculating dissimilarities between image frames taken from multiple videos of the face of each identity, and then computing the mean dissimilarity across different image pairs featuring the same identity pair. Behavioral studies have demonstrated that participants make more mistakes in “telling together” (i.e., grouping multiple images of the same identity, which is different process from “telling apart,” or distinguishing, between different identities) different photographs of the same person when those photographs were taken with different cameras, on different days, or with different lighting conditions, compared with when photographs were taken on the same day and with the same camera (Bruce et al., 1999; Jenkins et al., 2011). Most previous fMRI studies, however, used very visually similar images, or even just a single image, for each identity, making it difficult to determine whether a brain region represents different face images or different face identities. Here, by having multiple videos for each person, we can be more confident that we are capturing representations of specific identities rather than specific stimuli.

Related to the previous point, Abudarham and Yovel (2016) have recently shown that humans are more sensitive in perceiving changes in some face features (e.g., lip thickness, hair, eye color, eye shape, and eyebrow thickness) compared with others (e.g., mouth size, eye distance, face proportion, skin color). Changes in the former type of features (also known as critical features) are perceived as changes in identity and those features tend to be invariant for different images of the same identity. Interestingly, Abudarham et al. (2019) showed that the OpenFace algorithm that we used in the present study also seemed to be capturing those same critical features. Given our results in rFFA, it would be interesting to see whether representations in this region can also distinguish between the processing of the critical and noncritical face features as described by Abudarham et al. (2016, 2019).

Grossman et al. (2019) have also recently shown that representations in the FFA relate to image-computable descriptors from a deep neural network. There are two main differences, however, between our results and those of Grossman et al. (2019). First, Grossman et al. (2019) found similar representational geometries across all face-selective ventral temporal cortex, and no differentiation between OFA and FFA. One possible reason for this difference is that the authors were only able to define OFA and FFA in the left hemisphere, whereas our face-selective regions were defined in the right hemisphere. Face-selective regions are more consistent and larger in the right hemisphere (e.g., Rossion et al., 2012). A second main difference between our results and those of Grossman et al. (2019) is that the deep neural network that we used here showed high generalization across different images of the same person. OpenFace (Amos et al., 2016) was trained specifically to group together images of the same person and distinguish images of different people, and it performed very well in doing this in our set of stimuli (see Extended Data Fig. 2-1), where it showed high generalization across very variable pictures of the same person. This was not the case with the VGG-Face network used by Grossman et al. (2019). Future studies should focus on describing and comparing the image-level descriptions of different types of neural networks.

Previous studies have demonstrated that face-selective regions are sensitive to the viewpoint from which faces are presented (Grill-Spector et al., 1999; Axelrod and Yovel, 2012; Kietzmann et al., 2012; Ramírez et al., 2014; Dubois et al., 2015; Guntupalli et al., 2017). However, there is also evidence that the FFA, OFA, anterior temporal lobe, and pSTS represent face identity across different viewpoints (Anzellotti et al., 2014; Anzellotti and Caramazza, 2017; Guntupalli et al., 2017). In our video stimuli, the faces were mostly front-facing, but were free to vary in terms of changes in viewpoint (e.g., turning the head to the side during the video). Given that our patterns for each identity were estimated across multiple different videos of their face, it is unlikely that viewpoint alone could explain the differences between identities. Therefore, our results suggest that the FFA and OFA encode information that relate to face identity, beyond viewpoint.

We note that the lower bounds on the noise ceiling in our analyses were consistently quite low, especially for rFFA and rpSTS. However, these values are similar to the lower bounds of the noise ceiling in other studies using RSA (e.g., Jozwik et al., 2016; Carlin and Kriegeskorte, 2017; Thornton and Mitchell, 2017, 2018). We considered whether the low correlations could reflect substantial individual differences in face identity brain representations, but our results did not support this possibility. The low noise ceilings in our study likely reflect the fact that the

differences between brain-activity patterns associated with faces of different people are small compared with the differences between patterns associated with different visual categories (e.g., faces and places). Moreover, we used identity-based, rather than image-based, patterns (by cross-validating across runs presenting different videos for each identity), and this is likely to have introduced additional variability to the pattern estimates. It is also possible that we needed more data per participant, and future studies should consider ways to increase the amount of explainable variance. A related issue is that the Perceived-property models had intersubject reliabilities that varied between 0.2 and 0.6; thus, correlations between these models and brain RDMs would be affected by these low reliabilities.

None of the models that we considered here explained the representational geometry of responses in the face-selective rpSTS. It is likely that the pSTS, as defined in the present study, contains overlapping and interspersed groups of voxels that respond to faces only, voices only, or both faces and voices (Beauchamp et al., 2004) that make the overlapping representational geometry difficult to explain. On the other hand, it is possible that the pSTS represents information about people that we did not consider here, such as idiosyncratic facial movements (Yovel and O'Toole, 2016), emotional and mental states (Thornton et al., 2019), biographical knowledge (Verosky et al., 2013; Collins et al., 2016; Thornton et al., 2019), social distance or network position (Parkinson et al., 2014, 2017), or type of social interactions (Walbrin and Koldewyn, 2019). Future studies may need to explore an even richer set of social, perceptual, and stimulus-based models to better characterize responses in the pSTS (and investigate representations beyond face-selective regions).

A limitation of our study was the lack of diversity of our face identities in terms of race and ethnicity (10 identities were White and 2 were Black), which limits the generalizability of our results to faces of different ethnicities. It was essential to our study that our set of celebrities were highly familiar to our sample of young British participants, and they were chosen based on their recognizability (of both faces and voices) (see Tsantani et al., 2019). Future work will need to incorporate more diversity in the face stimuli. This is also crucial when considering the Image-computable models. In particular, OpenFace has been developed, trained, and evaluated on databases that contain large proportions of White faces compared with other ethnicities. Future work using larger samples of identities should evaluate the biases caused by these procedures, and develop models trained on more representative and diverse databases.

In conclusion, our study highlights the importance of using multiple and diverse representational models to characterize how face identities are represented in different face-selective regions. Although similar levels of identity decodability were observed in both OFA and FFA (Tsantani et al., 2019), the information explicitly encoded in these two regions is indeed distinct, suggesting that the two regions serve quite different computational roles. Future work attempting to define the computations of cortical regions that appear to serve the same function (e.g., discriminating between identities) would benefit from comparing representations in those regions with multiple and diverse candidate models to reveal the type of information that is encoded.

References

Abudarham N, Yovel G (2016) Reverse engineering the face space: discovering the critical features for face identification. *J Vis* 16:40.

- Abudarham N, Shkiller L, Yovel G (2019) Critical features for face recognition. *Cognition* 182:73–83.
- Amos B, Ludwiczuk B, Satyanarayanan M (2016) Openface: a general-purpose face recognition library with mobile applications. *CMU School Comput Sci* 6:2.
- Anzellotti S, Caramazza A (2017) Multimodal representations of person identity individuated with fMRI. *Cortex* 89:85–97.
- Anzellotti S, Fairhall SL, Caramazza A (2014) Decoding representations of face identity that are tolerant to rotation. *Cereb Cortex* 24:1988–1995.
- Axelrod V, Yovel G (2012) Hierarchical processing of face viewpoint in human visual cortex. *J Neurosci* 32:2442–2452.
- Axelrod V, Yovel G (2015) Successful decoding of famous faces in the fusiform face area. *PLoS One* 10:e0117126.
- Axelrod V, Rozier C, Malkinson TS, Lehongre K, Adam C, Lambrecq V, Navarro V, Naccache L (2019) Face-selective neurons in the vicinity of the human fusiform face area. *Neurology* 92:197–198.
- Baltrusaitis T, Zadeh A, Lim YC, Morency LP (2018). Openface 2.0: facial behavior analysis toolkit. In: 2018 13th IEEE International Conference on Automatic Face and Gesture Recognition, pp 59–66.
- Beauchamp MS, Argall BD, Bodurka J, Duyn JH, Martin A (2004) Unraveling multisensory integration: patchy organization within human STS multisensory cortex. *Nat Neurosci* 7:1190–1192.
- Biederman I, Kalocsai P (1997) Neurocomputational bases of object and face recognition. *Philos Trans R Soc Lond B Biol Sci* 352:1203–1219.
- Boersma P, Weenink D (2014) Praat: Doing Phonetics by Computer [Computer Software]. Version 5.3.8.0.
- Brainard DH (1997) The psychophysics toolbox. *Spat vis* 10:433–436.
- Bruce V, Henderson Z, Greenwood K, Hancock PJ, Burton AM, Miller P (1999) Verification of face identities from images captured on video. *J Expl Psychol Appl* 5:339–360.
- Calder AJ, Burton AM, Miller P, Young AW, Akamatsu S (2001) A principal component analysis of facial expressions. *Vision Res* 41:1179–1208.
- Carlin JD, Kriegeskorte N (2017) Adjudicating between face-coding models with individual-face fMRI responses. *PLoS Comput Biol* 13:e1005604.
- Collins JA, Koski JE, Olson IR (2016) More than meets the eye: the merging of perceptual and conceptual knowledge in the anterior temporal face area. *Front Hum Neurosci* 10:189.
- Davidesco I, Zion-Golumbic E, Bickel S, Harel M, Groppe DM, Keller CJ, Schevon CA, McKhann GM, Goodman RR, Goelman G, Schroeder CE, Mehta AD, Malach R (2014) Exemplar selectivity reflects perceptual similarities in the human fusiform cortex. *Cereb Cortex* 24:1879–1893.
- di Oleggio Castello MV, Halchenko YO, Guntupalli JS, Gors JD, Gobbini MI (2017) The neural representation of personally familiar and unfamiliar faces in the distributed system for face perception. *Sci Rep* 7:1–14.
- Dubois J, de Berker AO, Tsao DY (2015) Single-unit recordings in the macaque face patch system reveal limitations of fMRI MVPA. *J Neurosci* 35:2791–2802.
- Duchaine B, Yovel G (2015) A revised neural framework for face processing. *Annu Rev Vis Sci* 1:393–416.
- Fedorenko E, Hsieh PJ, Nieto-Castañón A, Whitfield-Gabrieli S, Kanwisher N (2010) New method for fMRI investigations of language: defining ROIs functionally in individual subjects. *J Neurophysiol* 104:1177–1194.
- Ghuman AS, Brunet NM, Li Y, Konecky RO, Pyles JA, Walls SA, Destefino V, Wang W, Richardson RM (2014) Dynamic encoding of face information in the human fusiform gyrus. *Nat Commun* 5:1–10.
- Goesaert E, Op de Beeck HP (2013) Representations of facial identity information in the ventral visual stream investigated with multivoxel pattern analyses. *J Neurosci* 33:8549–8558.
- Grill-Spector K, Kushnir T, Edelman S, Avidan G, Itzhak Y, Malach R (1999) Differential processing of objects under various viewing conditions in the human lateral occipital complex. *Neuron* 24:187–203.
- Grossman S, Gaziv G, Yeagle EM, Harel M, Mégevand P, Groppe DM, Khuvis S, Herrero JL, Irani M, Mehta AD, Malach R (2019) Convergent evolution of face spaces across human face-selective neuronal groups and deep convolutional networks. *Nat Commun* 10:1–13.
- Guntupalli JS, Wheeler KG, Gobbini MI (2017) Disentangling the representation of identity from head view along the human face processing pathway. *Cereb Cortex* 27:46–53.
- Haxby JV, Hoffman EA, Gobbini MI (2000) The distributed human neural system for face perception. *Trends Cogn Sci* 4:223–233.
- Jenkins R, White D, Van Montfort X, Burton AM (2011) Variability in photos of the same face. *Cognition* 121:313–323.
- Julian JB, Fedorenko E, Webster J, Kanwisher N (2012) An algorithmic method for functionally defining regions of interest in the ventral visual pathway. *Neuroimage* 60:2357–2364.
- Jozwik KM, Kriegeskorte N, Mur M (2016) Visual features as stepping stones toward semantics: explaining object similarity in IT and perception with non-negative least squares. *Neuropsychologia* 83:201–226.
- Jozwik KM, Kriegeskorte N, Storrs KR, Mur M (2017) Deep convolutional neural networks outperform feature-based but not categorical models in explaining object similarity judgments. *Front Psychol* 8:1726.
- Kanwisher N, McDermott J, Chun MM (1997) The fusiform face area: a module in human extrastriate cortex specialized for face perception. *J Neurosci* 17:4302–4311.
- Khaligh-Razavi SM, Kriegeskorte N (2014) Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS Comput Biol* 10:e1003915.
- Khuvis S, Yeagle EM, Norman Y, Grossman S, Malach R, Mehta AD (2021) Face-selective units in human ventral temporal cortex reactivate during free recall. *J Neurosci*. Advance online publication. Retrieved Jan 11, 2021. doi: 10.1523/JNEUROSCI.2918-19.2020.
- Kietzmann TC, Swisher JD, König P, Tong F (2012) Prevalence of selectivity for mirror-symmetric views of faces in the ventral and dorsal visual pathways. *J Neurosci* 32:11763–11772.
- Kriegeskorte N, Mur M, Ruff DA, Kiani R, Bodurka J, Esteky H, Tanaka K, Bandettini PA (2008a) Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron* 60:1126–1141.
- Kriegeskorte N, Mur M, Bandettini P (2008b) Representational similarity analysis: connecting the branches of systems neuroscience. *Front Syst Neurosci* 2:4–28.
- Margalit E, Biederman I, Herald SB, Yue X, von der Malsburg C (2016) An applet for the Gabor similarity scaling of the differences between complex stimuli. *Atten Percept Psychophys* 78:2298–2306.
- McKone E, Brewer JL, MacPherson S, Rhodes G, Hayward WG (2007) Familiar other-race faces show normal holistic processing and are robust to perceptual stress. *Perception* 36:224–248.
- Nestor A, Plaut DC, Behrmann M (2011) Unraveling the distributed neural code of facial identity through spatiotemporal pattern analysis. *Proc Natl Acad Sci USA* 108:9998–10003.
- Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N (2014) A toolbox for representational similarity analysis. *PLoS Comput Biol* 10:e1003553.
- Oliva A, Torralba A (2001) Modeling the shape of the scene: a holistic representation of the spatial envelope. *Int J Comput Vis* 42:145–175.
- Oosterhof NN, Todorov A (2008) The functional basis of face evaluation. *Proc Natl Acad Sci USA* 105:11087–11092.
- Parkinson C, Liu S, Wheatley T (2014) A common cortical metric for spatial, temporal, and social distance. *J Neurosci* 34:1979–1987.
- Parkinson C, Kleinbaum AM, Wheatley T (2017) Spontaneous neural encoding of social network position. *Nat Hum Behav* 1:1–7.
- Pelli DG (1997) The VideoToolbox software for visual psychophysics: transforming numbers into movies. *Spat vis* 10:437–442.
- Pitcher D, Dilks DD, Saxe RR, Triantafyllou C, Kanwisher N (2011) Differential selectivity for dynamic versus static information in face-selective cortical regions. *Neuroimage* 56:2356–2363.
- Ramírez FM, Cichy RM, Allefeld C, Haynes JD (2014) The neural code for face orientation in the human fusiform face area. *J Neurosci* 34:12155–12167.
- Rhodes G (1988) Looking at faces: first-order and second-order features as determinants of facial appearance. *Perception* 17:43–63.
- Rossion B, Hanseeuw B, Dricot L (2012) Defining face perception areas in the human brain: a large-scale factorial fMRI face localizer analysis. *Brain Cogn* 79:138–157.

- Russell R, Biederman I, Naderhouser M, Sinha P (2007) The utility of surface reflectance for the recognition of upright and inverted faces. *Vision Res* 47:157–165.
- Russell R, Sinha P (2007) Real-world face recognition: the importance of surface reflectance properties. *Perception* 36:1368–1374.
- Storrs K, Khaligh-Razavi S, Kriegeskorte N (2020) Noise ceiling on the cross-validated performance of reweighted models of representational dissimilarity: addendum to Khaligh-Razavi and Kriegeskorte (2014). *bioRxiv* 003046. doi: 10.1101/2020.03.23.003046.
- Sutherland CA, Oldmeadow JA, Santos IM, Towler J, Burt DM, Young AW (2013) Social inferences from faces: ambient images generate a three-dimensional model. *Cognition* 127:105–118.
- Tardif J, Morin Duchesne X, Cohan S, Royer J, Blais C, Fiset D, Duchaine B, Gosselin F (2019) Use of face information varies systematically from developmental prosopagnosics to super-recognizers. *Psychol Sci* 30:300–308.
- Thornton MA, Mitchell JP (2017) Consistent neural activity patterns represent personally familiar people. *J Cogn Neurosci* 29:1583–1594.
- Thornton MA, Mitchell JP (2018) Theories of person perception predict patterns of neural activity during mentalizing. *Cereb Cortex* 28:3505–3520.
- Thornton M, Weaverdyck M, Tamir D (2019) The brain represents people as the mental states they habitually experience. *Nat Commun* 10:2291.
- 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* 201:116004.
- Verosky SC, Todorov A, Turk-Browne NB (2013) Representations of individuals in ventral temporal cortex defined by faces and biographies. *Neuropsychologia* 51:2100–2108.
- Walbrin J, Koldewyn K (2019) Dyadic interaction processing in the posterior temporal cortex. *Neuroimage* 198:296–302.
- Walther A, Nili H, Ejaz N, Alink A, Kriegeskorte N, Diedrichsen J (2016) Reliability of dissimilarity measures for multi-voxel pattern analysis. *Neuroimage* 137:188–200.
- Weibert K, Flack TR, Young AW, Andrews TJ (2018) Patterns of neural response in face regions are predicted by low-level image properties. *Cortex* 103:199–210.
- Yovel G, Duchaine B (2006) Specialized face perception mechanisms extract both part and spacing information: evidence from developmental prosopagnosia. *J Cogn Neurosci* 18:580–593.
- Yovel G, O'Toole AJ (2016) Recognizing people in motion. *Trends Cogn Sci* 20:383–395.
- Yue X, Biederman I, Mangini MC, Malsburg C, von der Amir O (2012) Predicting the psychophysical similarity of faces and non-face complex shapes by image-based measures. *Vision Res* 55:41–46.
- Zhang H, Japee S, Nolan R, Chu C, Liu N, Ungerleider L (2016) Face-selective regions differ in their ability to classify facial expressions. *Neuroimage* 130:77–90.
- Zhou X, Mondloch CJ (2016) Recognizing “Bella Swan” and “Hermione Granger”: no own-race advantage in recognizing photos of famous faces. *Perception* 45:1426–1429.