# MIT Open Access Articles

## Decoding task and stimulus
## representations in face-responsive cortex

Massachusetts Institute of Technology

DSpace@MIT

# Decoding Task and Stimulus Representations in Face-responsive Cortex:

## Decoding Task and Stimulus from Faces

**Dorit Kliemann**[1,2], **Nir Jacoby**[3], **Stefano Anzellotti**[3], and **Rebecca R. Saxe**[3]

[1]McGovern Institute for Brain Research, Massachusetts Institute of Technology, 43 Vassar Str, Cambridge, MA 02139, USA

[2]Department of Neurology, Massachusetts General Hospital/Harvard Medical School, 149 Thirteenth Street, Charlestown, MA 02129, USA

[3]Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, 43 Vassar Str, Cambridge MA, 02139, USA

## Abstract

Faces provide rich social information about others' stable traits (e.g., age) and fleeting states of mind (e.g., emotional expression). While some of these facial aspects may be processed automatically, observers can also deliberately attend to some features while ignoring others. It remains unclear how internal goals (e.g., task context) influence the representational geometry of variable and stable facial aspects in face-responsive cortex. We investigated neural response patterns related to decoding i) the intention to attend to a facial aspect before its perception, ii) the attended aspect of a face and iii) stimulus properties. We measured neural responses while subjects watched videos of dynamic positive and negative expressions, and judged the age or the expression's valence. Split-half multivoxel pattern analyses (MVPA) showed that (i) the intention to attend to a specific aspect of a face can be decoded from left fronto-lateral, but not face-responsive regions; (ii) during face perception, the attend aspect (age vs emotion) could be robustly decoded from almost all face-responsive regions; and (iii) a stimulus property (valence), was represented in right posterior superior temporal sulcus and medial prefrontal cortices. The effect of deliberately shifting the focus of attention on representations suggest a powerful influence of top-down signals on cortical representation of social information, varying across cortical regions, likely reflecting neural flexibility to optimally integrate internal goals and dynamic perceptual input.

## Keywords

fMRI; Faces; Emotion; MVPA; Social Cognition

---

Corresponding Author: Dorit Kliemann, dorit@mit.edu, 617-324-2891.

## INTRODUCTION

Navigating in our social world requires efficient recognition and processing of information to produce to most optimal behavior (e.g. recognizing a familiar face or identifying an aggressive individual in a crowd). Because our environment is rapidly changing, an immediate integration of multimodal input is necessary. The challenge of effective information processing, however, is not just dependent on the perception of external input. Instead, our internal goals are modulated by context and change our focus of attention and/or the behavioral goal. Effective social functioning thus depends on the flexibility to process perceptual input of the world while optimizing information processing in context.

Among visual social stimuli, faces are one of the most important and salient nonverbal visual sources of information about what others might think and feel. Shortly after birth, humans preferentially orient to face-like stimuli (Goren, Sarty, & Wu, 1975, Johnson, Dziurawiec, Ellis, & Morton, 1991) and over the course of the first few month of life, infants use facial information to interpret external events and to guide their own behavior (Nelson & Dolgin, 1985; Nelson, Morse, & Leavitt, 1979). By adulthood, we are able to extract rich information about another person's stable traits (such as identity, gender, age range) and fleeting states of mind (such as gaze, emotional expression) within 200 milliseconds (Adolphs, 2002). Prominent cognitive models of face perception suggest a division of labor between processing of different stable and variant aspects faces (e.g., facial identity vs emotion recognition) concerted by distinct - but also interconnected - regions within the face network (see, e.g. Bruce & Young et al., 1986). Converging neuroimaging evidence points to the encoding of identity in ventral temporal regions (Nestor et al. 2011, Anzellotti et al. 2013, Anzellotti et al. 2015), and of emotion in ventral temporal, medial prefrontal and posterior lateral temporal regions (Peelen et al 2010, Skerry and Saxe 2014). Some of these aspects of a face may be processed automatically (Critchley et al., 2000), but observers can also deliberately attend to some aspects of facial features, while ignoring others.

Most social neuroscience accounts of face perception focus on investigating bottom-up processing of perceptual information (Riesenhuber & Poggio, 1999, DiCarlo et al., 2012). However, other research points towards feedback loops and recurrent wiring within the face network (e.g., between core and extended face processing regions, Adolphs et al. 2002) and also with other regions of the brain (e.g., Kravitz et al., 2013, Van Essen et al., 1992). The extent to which internal goals (i.e. top down influences) may shape perceptual processing (bottom up) in anticipation or response to facial information is not yet fully understood. Crucial to this investigation is the assumption that objects (and thus faces as well) can be seen as units of attention (O'Craven and Kanwisher, 1999). Of course we would expect different neural representations when attending to different categories of objects (e.g. faces vs houses). However, a task (or a shift in behavioral goals) can also require attending to different aspects of the very same object (i.e., different aspects of a face). Understanding how social information is flexibly represented in neural patterns across the cortex before, during and after stimulus presentation is of high relevance not just for typical but also atypical social-cognitive processing.

Attending to different aspects of a single object can influence the hemodynamic responses to that object in at least two ways over time: by increasing the overall magnitude of response in some cortical regions, and by changing the representation to prioritize relevant dimensions. The effects of attention on the magnitude of response are well-studied. Corbetta and colleagues (1990) showed in one of the earliest PET studies that selective attention influences neural processing of color, velocity and shape of objects in human extrastriate cortex. Since then, a vast body of research including several studies on visual attention (see, e.g., Carrasco, 2011, Kanwisher & Wojciulik 2000, for reviews) gives further insight into how endogenously influenced goals of an agent can modulate neural activity in specific cortical regions. For example, Ganel and colleagues (2005) found that deliberately attending to the emotional expression on a face increased the overall magnitude of response in posterior superior temporal sulcus (pSTS), fusiform face area (FFA) and amygdala, compared to attending to the identity of the same faces (also see, e.g. Fox et al, 2009).

In addition to increases in magnitude, attention can change the representational geometry of neural responses to objects. For example, Harel and colleagues (2014) using multivoxel pattern analyses found that visual processing of objects across the cortex is influenced by behavioral goals. Interestingly, cortical representations of objects were differentially affected by a given task. Ventral-temporal and prefrontal regions showed task-type-dependent representations (physical/conceptual) and the individual tasks used, while neural patterns in early visual cortex showed no significant sensitivity to the type of tasks participants were performing (but to specific physical tasks). Given that the visual input in the different tasks was identical, these results show a striking influence of top-down signals on visual object representation in early stages of object processing. Whether similar effects occur in neural responses to faces is not yet known.

It is also possible that the mere anticipation of processing visual input with differing behavioral goals may already shape neural responses even *before* stimulus onset. Parts of face responsive cortex seem to be affected when imagining faces vs objects without visual input (O'Craven & Kanwisher, 2000). Following this logic, the intention to focus on a certain face aspect may already modulate representational information in neural responses in regions of the face network. Such an effect could be the result of pre-attentive influence of top-down attention regions (see, e.g., Kok et al 2013) on later, more domain specific face-responsive regions.

However, not all aspects of neural responses to a stimulus fluctuate with the context or behavioural goals: some aspects of faces and objects are extracted and recognized automatically. There is some evidence that facial expressions of emotion might be processed automatically: the pSTS, FFA and amygdala all show repetition suppression for repeated emotional expressions – and therefore increased responses when expressions are varied across successive faces – even when participants are attending to facial identity (Ganel et al., 2005). Certain stimulus properties can also facilitate processing of orthogonal facial aspects: is has been repeatedly shown that positive (as compared to negative) facial expression increases magnitude of BOLD response in core face processing regions (potentially via feedback loops, see e.g., Vuillemmier et al, 2003). These results may indicate that some but

not necessarily all of the face network regions extract and represent the emotional expression of a face automatically, even when that feature is not task relevant.

We sought to test the sensitivity of neural patterns to internal goals when processing faces. The main focus of the current study is to investigate how shifting attention between two aspects of facial information modulates the neural representation of faces in independently localized regions of face-responsive cortex (see methods section for further details on a priori selection of regions and supplementary material for further regions). Specifically, we asked whether we can decode i) the intention to attend to a specific facial aspect before its actual perception, ii) the attended aspect of a face, independent of the stimulus and iii) stimulus properties, independent of the attended aspect. The stimulus property we targeted was the emotional valence of a dynamic facial expression. In prior studies, emotional valence of a facial expression could be decoded from regions in posterior STS and MPFC (Said et al 2010, Peelen et al 2010, Skerry and Saxe 2014). In addition, these responses seem to be fairly abstract: facial emotional valence could be decoded using a model trained on neural responses to positive versus negative emotion in voices, body movements (in superior temporal gyrus, Peelen et al 2010) or animated cartoons (in MPFC, Skerry and Saxe 2014). However, in all of these prior studies, participants were instructed to attend to the character's emotion. Prior evidence, using only univariate analyses, provides hints both that attention affects processing of emotional expression in these regions, and that the valence of the face may be represented automatically. Therefore we tested whether neural patterns would be robust to changing behavioral goals, or whether the robust and abstract response to emotional valence observed in prior studies depends on participants deliberately attending to emotions.

We manipulated the observer's internal goals by instructing participants to discriminate either the target aspect of the face (emotion: positive versus negative) or an orthogonal distractor aspect (age: over versus under 40 years old), in a dynamic naturalistic movie clip (Skerry and Saxe 2014). In order to identify the intention to attend to one of these aspects, we separated the instructions from the stimulus by a long and jittered delay, and used two physically dissimilar cues to instruct each task. During the dynamic movie clips, information about both invariant and changeable aspects of the faces was presented simultaneously, and relied on the same facial features. For example, both age and emotional valence are conveyed disproportionately by the eye and mouth regions (Gamer et al., 2010, Kwart et al., 2012). Nevertheless, attending to the person's age versus emotion could lead to a change in the representation of the face, that would be reflected in different patterns of response across cortex.

## METHODS

### Participants

Twenty-eight right-handed adult participants (11 female, aged 21–33 years (mean (SD) = 26.6 (4.2)) with no history of neurological or psychiatric disorders and normal or corrected-to-normal vision participated in the study. We excluded three participants' data (1 female) from further analyses due to poor task performance (see results section for details on exclusion criteria) resulting in a final dataset of 25 participants. Participants were paid for

participation and gave written informed consent prior to participating, in accordance with the Committee on the Use of Human Experimental Subjects at the Massachusetts Institute of Technology (MIT).

## Procedure

Participants completed two fMRI tasks in the scanner (a task to individually localize brain regions involved in emotional face processing and an emotion/age attribution task).

## fMRI tasks

**Localizer Task—**To identify a broad spectrum of brain regions involved in processing faces or emotion, we presented 45 unique triplets of emotional faces versus colored shapes in a block-design (Hariri et al., 2000). Shapes consisted of colored geometrical shapes (e.g., cylinders, triangles, rectangles). Triplets of faces were happy and angry emotional expressions taken from a standardized database (Tottenham et al., 2002). In each trial, participants were asked to indicate via button press which face (or shape) from a pair at the bottom of the screen most closely resembles the target face (or shape) at the top of the screen according to emotional expression (for faces) or geometrical shape characteristic (for shapes). The task consisted of one run, with 6 blocks (3 blocks per condition, no inter-block or inter-trial delays), starting with the presentation of a blank screen (8s) before the first block. Each block consisted of 15 trials (2s each) resulting in a total experiment time of 184s. Participants responded via button press (left versus right button) during each trial. All participants completed a standardized instruction prior to scanning.

**Emotion/Age Attribution Task—**Participants watched short movie clips of dynamic positive and negative facial expressions (for further details on stimuli and emotional valence validation, see Skerry & Saxe, 2014). Faces were close-ups of different individuals, taken from TV-shows and movies, thus representing relatively uncontrolled but naturalistic visual stimuli (compared to highly controlled but less naturalistic stimuli, such as face morphs). We chose to use more naturalistic stimuli to elicit neural representations with high ecological validity (see, e.g. Zaki & Ochsner 2009).

The complete stimulus set comprised 192 unique stimuli (96 positive, 96 negative, within each valence 48 males and 48 females) presented over 8 runs. The experiment followed a jittered event-related 2×2 design of two task ("age" versus "emotion" tasks) and two stimulus conditions (positive versus negative emotion). Participants judged the valence of the emotional expression (emotion task: positive versus negative) or, to direct attention away from emotions, judged the individual's age (age task: over versus under 40 years old). Information relevant for both tasks was available immediately after movie onset. Each trial started with a prompt screen (1s) indicating the task for this trial. After a jittered delay (4 – 12s, mean = 8s) a movie was presented for 4s, followed by a short 250ms delay and the response screen (1.75s). The next trial started immediately after the response screen. The last trial in each run ended with the presentation of a blank screen for 12s, resulting in a total run time of 372s. Prompts were presented in two formats: three letters (emotion task: "EMO"; age task: "AGE") or iconic symbols (emotion task: smiling and sad emoticon; age task: small and bigger neutral emoticon; see Fig. 1). Response screens were identical for both task

conditions, consisting of a plus and a minus symbol (emotion task: plus = positive, minus = negative; age task: plus = 'over 40', minus = 'under 40'), and their position was randomized across trials. Participants responded by pressing the left or right button.

To optimize the presentation order of the four main conditions within each and over all runs, we created 8 schedules using Optsep2 (http://surfer.nrm.mgh.harvard.edu/optseq) with a first-order counterbalancing constraint. The order of items within a scheduled condition was then pseudo-randomized across runs, with the constraint that each movie clip was presented once in each task condition over runs. The orders of response option arrangement, gender of the face, and task prompt format were balanced within runs (i.e. each run had the same number of females, symbol prompts, etc.). Participants were trained on the tasks and completed one practice run before the scan, with different movie clips.

## fMRI Acquisition

Data were acquired on a 3-Tesla Tim Trio scanner (Siemens; Erlangen, Germany) at the Athinoula A. Martinos Imaging Center at the McGovern Institute for Brain Research at MIT, using a Siemens 32-channel phased-array head coil. To register functional data to individual and standard anatomy we first collected high-resolution T1-weighted anatomical images (MPRAGE, voxel size = 1×1×1 mm, TR = 2530 ms, slices = 176, FoV = 256 mm) with whole brain coverage. We then collected functional images acquired with a gradient-echo EPI sequence sensitive to Blood Oxygen Dependent (BOLD) contrast (voxel size = 3×3×3 mm, TR = 2000 ms, TE = 30 ms, flip angle = 90 degrees, FoV = 192 mm). Slices were aligned with the anterior/posterior commissure and provided near whole-brain coverage (excluding the cerebellum).

## fMRI Data Analyses

**Preprocessing**—MRI data were analyzed using SPM8 (http://www.fil.ion.ucl.ac.uk/spm/) and custom software written in Matlab (www.mathworks.com; Nattick, MA, USA). Each participant's data were registered to the first image of the first run. All functional runs were co-registered with the participant's anatomical scan and all images (functional and anatomical) were normalized to a common (Montreal Neurological Institute, EPI template) brain space. Functional images were smoothed using a Gaussian kernel filter (5mm FWHM (full-width-half-maximum)). Note that smoothing does not substantially affect decoding performance (Zhang et al., 2010; de Beek 2010). Data were high-pass filtered (cut-off 128s) to remove low-frequency noise. Functional data were further corrected for motion artifacts, defined as timepoints during which motion exceeded 2mm in any direction relative to the previous timepoint or a change in global signal exceeded a threshold of three standard deviations form the mean global signal. Time points with motion artifacts were removed during modeling with artifact timepoint regressors. For all analyses – except the bold pattern analyses in detail described below –, we performed whole-brain first level analyses on each participant's functional data by applying a general linear model (GLM) with SPM modeled as a boxcar function using a standard hemodynamic response function (HRF) matching the onset and duration of experiment specific regressors: For the *localizer task*, data were modeled with the two condition regressors (faces/shapes). For the *age/emotion attribution* task we modeled four prompt types (word/symbol × emotion/age task), four stimulus

conditions (age/emotion task × positive/negative stimulus) and the response. Nuisance covariates were added to the model i) for timepoints with head motion artifacts, ii) to correct for run effects, and (iii) reaction time, using a parametric regressor for each trial, with an amplitude on each trial corresponding to the mean-centered reaction time. To further investigate the influence of reaction time effects on neural patterns, we conducted two control analyses (see supplementary material). In short, the control analyses replicated all of the key effects reported in the main analyses.

We defined 8 face-responsive regions of interest (ROIs) based on prior findings about face-selective regions (e.g. Haxby et al., 2000) and regions selective for stimulus representations independently (Skerry & Saxe, 2014; Peelen et al. 2010) with the localizer task's data: bilateral anterior and posterior STS (aSTS, pSTS), right fusiform face area (rFFA), right occipital face area (rOFA) as well as dorsal and ventral medial prefrontal cortex (d/vMPFC) (see supplementary section for additional regions). We first applied a univariate whole-brain random effects analysis submitting all individuals' contrast images derived from the first-level analysis to a second-level analysis for the contrast faces > objects. Second, we defined a hypothesis space based on the peak voxel of the resulting group effect's t-map for each ROI (k > 9, p < .001, see table 1 for details on each ROI) for the face-responsive regions. Third, individual ROIs were defined per participant based on first-level t-maps as the top 80 most activated voxels within the hypothesis space. Participants showing a smaller number of activated voxels (t > 0) within each ROI were excluded for the respective ROI MVPA analyses (3 participants in the dMPFC-ROI, see table 1).

As a control region within early visual cortex (EVC) - not specifically related to emotional face processing - we first created an anatomical ROI along the calcarine sulcus (with the WFU_pickatlas, Maldijan et al., 2003, based on the Talairach Labels (Lancaster et al., 1997)) that served as hypothesis space to then select the 80 most active voxels for the contrast objects > rest in the face localizer tasks' t-maps for each participant.

We selected a fixed number of voxels to minimize differences in the number of voxels across regions and participants. We could not reliably identify a sufficiently sized left FFA group ROI (total number of voxels < 80) and hence did not include this region in our analyses. Information about additional regions (bilateral amygdala, IFG, early visual cortex) are reported in the supplementary material.

**Multivariate Pattern Analyses—**We asked three main questions: First, can we decode the face aspect that participants are preparing to attend, at the time of the task prompt? Second, can we decode the face aspect that participants are actually attending, at the time of stimulus? And third, can we decode the stimulus (i.e. emotional valence) independent of the face aspect participants were attending? To address these questions, we used split-half Multi-Voxel Pattern Analyses (MVPA; Haxby et al., 2001).

**Pattern analyses: beta values—**Each participant's data were binned into odd (1,3,5,7) versus even (2,4,6,8) runs and the mean response (beta value) for every voxel in a defined region. For each participant, we computed the correlation of beta values across voxels and compared the averaged correlation within versus across condition comparisons. Correlations

were Fisher Z transformed to allow statistical comparisons with parametrical tests. Voxel-wise within and across condition data was not normalized prior to comparing correlations (see Garrido et al., 2013 for a general discussion on the topic of normalizing; see supplementary material for main results when normalizing). If the *within* correlation is significantly greater than the *across* correlation, as assessed by a student's T complementary cumulative distribution function, the neural pattern for these two conditions, and therefore some aspect of the way these conditions are represented neurally, is distinct.

**Pattern analyses: BOLD values—**In addition to calculating beta values in response to the main regressors, we also assessed difference in split-half correlations between conditions per timepoint over the course of trials. This analysis allowed us to ask what is the earliest time the neural patterns represent information about the task in a given trial and in particular, are neural responses task-dependent in face responsive ROIs *before* stimulus onset? First, for every voxel in the respective regions we extracted the full timecourse, then applied temporal filtering (with the same filter that was applied to the model estimating beta values, i.e. 128s) and normalized the timecourse per voxel (subtracting the mean and dividing by the standard deviation). Next, all relevant timepoints' z-scored values were averaged over 2 TRs (i.e. 4 sec) in relation to the onset of an event (i.e. the i) prompt or the ii) stimulus (video)). For the *prompt-locked analysis* we calculated difference scores (within minus between correlations as described above) from the time of the prompt up to three timepoints after the prompt. For the *video-locked analyses*, we calculated difference scores for three timepoints before and four during/after the video. We only analyzed trials in which the delay between the prompt and the movie was at least 6s long, i.e. excluding 4 trials per run that had a 4s delay. The focus of this analysis was to identify the earliest time in a trial, at which task information is decoded (see Figure 2 and 3). Note that this is a completely different analysis approach than the beta value based analyses described above. Instead of modeling betas per condition regressors, we extracted bold per time point.

**Whole-brain searchlight pattern analyses—**We conducted whole-brain searchlight analyses to ask whether there are regions in the brain (other than the predefined face-sensitive ROIs) that contain a distinct neural pattern for i) the task subjects are intending to perform on a given trial (prompt content) and ii) the valence of the emotional expression in the movie. The spatial correlations across and within conditions were computed in voxels selected by a Gaussian searchlight sphere (9mm FWHM) moving iteratively across the brain. By using a Gaussian kernel, the influence of voxels at increasing distances from the reference voxel is de-emphasized (Dehaene & Cohen 2007). Resulting whole-brain maps (Fisher Z transformed to allow statistical comparisons with parametrical tests) for each participant were next submitted to second-level analyses using one-sampled t-tests (corrected for multiple comparisons at $p < .05$ using Monte Carlo permutation tests to establish empirical null distributions for the peak T and cluster size with $\theta = 0.5$ (Statistical non-Parametric Mapping (SnPM), www2.Warwick.ac.uk/snpm; Nichols and Holmes, 2002; Hayasaka and Nichols, 2004), if not specified otherwise).

When the searchlight identified regions that were not among our initial regions of interest, we used an iterative leave-one-participant-out procedure to define independent ROIs for

further analysis. We ran a whole brain searchlight analysis leaving out data for one participant at a time and defined n-1 ROIs at p<.001, k > 10. We repeated this process for each participant, and extracted BOLD responses and/or beta values from the resulting ROIs defined from the analysis that left out each participant's data.

## RESULTS

### Behavioral Results

**Emotion/Age Attribution Task—**To ensure participants were effectively attending to the cues facial aspect, runs were excluded if less than 83% of trials in the emotion task condition (i.e., less than 20 of 24 trials) were answered correctly. Only the emotion task was used to exclude runs, because the correct answer was unambiguous (see Mechanical Turk ratings in Skerry & Saxe, 2014). Participants were excluded if more than one run was excluded. These a priori exclusion criteria led to exclusion of three participants and one run of a fourth participant, leading to a final sample of n=25.

In the remaining data, participants showed overall high accuracy rates when judging emotional valence (mean (SD) = 97.9 (2.1)), and were equally accurate at detecting positive and negative emotions (positive: 97.55, negative: 98.14, t(24): −0.74, p = .47). A 2×2 repeated measurement analysis of variance (ANOVA) on RT with the factors *task* (emotion × age) and *stimulus* (positive × negative) yielded main effects of task, F(1, 24) = 23.49, p < .001, $\eta_p^2$ = .49, and stimulus, F(1, 24) = 7.27, p = .013, $\eta_p^2$ = .23. Participants were significantly faster in responding to emotion as compared to age trials (emotion: 593ms, age: 637ms), and when responding to positive as compared to negative faces (positive: 608ms, negative: 628ms). There was no significant interaction of the two factors.

### fMRI Results

In this section, we outline the results with regards to the three main questions of this study: What are the neural responses for i) the intention to attend to a face aspect, ii) attending to a face aspect and iii) distinguishing the features of the stimuli themselves?

**The intention to attend to a face aspect—**At the start of each trial, participants saw a prompt (either letters or symbols) indicating which face aspect (age or emotion) they would attend in the upcoming movie. We asked whether there is information about the intention to attend a specific facial aspect in the pattern of neural response (i.e. before onset of the actual face video). We found no significant decoding of information about the task in any face-responsive region in response to the prompt (*beta pattern analyses*, for details on statistics see supplementary material). In addition, the earliest time (*bold pattern analyses*, see Figure 3) at which the classic face-responsive regions showed above chance decoding of the task was much later, following onset of the face stimulus (see next section). The control EVC region also showed no effect of task in response to the prompt.

By contrast, a whole brain searchlight (SnPM corrected, p < .05) identified two regions containing information about the task at the time of the prompt (*beta pattern analyses*): left precentral gyrus (lPCG) and left inferior frontal gyrus (lIFG). We created regions of interest in these two regions, using a leave-one-subject-out iterated analysis, so that in each fold, the

extracted responses were independent of the data used to select the ROIs. In both of these regions, we found that the intended task could be decoded from 4 to 8 seconds after the prompt was presented, i.e. in response to the prompt (see Fig. 2, p < .05, blue background). Furthermore, these regions both represented the intended task, and not the visual image of the prompt: we could decode participants' intended task, even when requiring generalization across the two prompt formats (e.g. correlating within tasks and testing between letters and symbols, *beta pattern analyses*).

In addition, a whole brain searchlight (SnPM corrected, p < .05, *beta pattern analyses*) revealed information about the prompt format (letters vs numbers) independent of task, in both bilateral occipital regions (bilateral middle occipital gyrus) and in distinct frontal regions (bilateral inferior frontal gyrus, anterior cingulate gyrus, left precentral gyrus).

**Attending to a face aspect—**After a jittered delay, participants saw a naturalistic dynamic video of a single person, whose emotional expression was positively or negatively valenced. We asked whether we could decode the facial aspect participants were attending in the video, i.e. age or emotion, generalizing across distinct stimuli. All of the a priori face-sensitive ROIs –except the rOFA - contained information about the attended aspect at the time of the video (*beta pattern analyses*, see Table 2, Supplementary Figure S2). There were also some regional differences in strength of the effect, with raSTS and rFFA showing the weakest representation of task in its neural pattern. Furthermore, the pattern of BOLD response in the face-responsive ROIs was similar when participants attended the same aspect of the face, even when the stimuli (and therefore participants' responses) were different. The task effect was robustly represented in neural patterns even when generalizing across valence, i.e., when correlating within tasks but testing across positive and negative expressions.

The visual control region (EVC) showed no decoding of the task participants were performing. To statistically test for regional differences we first averaged the decoding accuracies for ventral and medial MPFC), as well bilateral posterior and anterior STS (pSTS, aSTS), respectively. We then conducted a repeated measures ANOVA with the within-subjects factors *ROI* [MPFC, aSTS, pSTS, rFFA, rOFA, EVC] and *condition-comparisons* [averaged versus generalized across the stimulus valence condition]. This analyses (see Figure 4) showed significant main effects of ROI ($F_{(1,21)} = 3.3$, p = .00001, ŋ = .213) and condition-comparison ($F_{(1,21)} = 15.1$, p = .001, ŋ = .4).

In all of these face-responsive ROIs that show a task effect, the time-locked analyses (*bold pattern analyses*) indicated that information about the attended face aspect emerged around 4 seconds after stimulus presentation (see, Fig. 3, red background).

**The stimulus property: valence of the facial expression—**The facial expressions in the videos were all unambiguously perceived as positively or negatively valenced; and in prior research, we found that the valence of a facial emotional expression could be decoded from the pattern of BOLD response in pSTS and MPFC. We therefore sought to (a) replicate this result, when emotional expression was attended, and (b) ask how the representation of emotional valence was affected by attending to a different face aspect, age.

When participants were attending to emotion, a whole brain searchlight revealed decoding of emotional valence in both pSTS and MPFC, consistent with prior finding, though only when using a relatively lenient statistical threshold ($p < .001$, voxelwise, $k > 30$, uncorrected, see Fig. 5). When participants were attending to age, the searchlight revealed no regions with significant classification of emotional valence in the stimulus (except early visual cortex, see Table 3); however this difference between the two tasks was not significant (i.e. in whole brain analyses, we did not observe a task × decoding interaction). None of the a priori face-sensitive ROIs (or the control EVC region) showed significant classification of the stimulus valence, during either task; although our a priori ROIs included regions in both pSTS and MPFC, the whole brain analyses indicated that the a priori ROIs did not overlap spatially with the regions that showed successful stimulus classification in the whole brain searchlight (see Peelen et al., 2010).

## DISCUSSION

In our everyday life, we are presented with stable and changing aspects of objects and other social agents in our environment. Ideally, internal goals may shape perceptual processing towards optimized representation of information (e.g. O'Craven & Kanwisher, 1999). In the case of face processing, it remains unclear when, where and how representations of faces are affected by changing the internal goals of an agent, or whether they remain independent. To this end, we asked in the current study how a modulation of internal goals affects patterns in brain activity representing information about the task subjects were performing and about the stimulus itself before and during stimulus perception. The design of our task allowed us to identify the effects, on face representations, of: the task prompt (the initial intention to attend to a face aspect), the prompt format (e.g. letters versus numbers), the attended aspect of the face, and the features of the face stimulus.

### Representation of the intention to attend to an aspect of a face

The earliest time at which the independently localized face-responsive ROIs could decode the task was in response to the stimulus, *not* in response to the prompt, before the face was presented. Thus, preparing to attend to a facial aspect appears not to elicit a task-specific pattern of response in face-responsive regions cortex. In other words, the presentation of the specific task (or goal) did not set the relevant face-responsive brain regions into a 'process-ready' state before stimulus onset. The shift in representational geometry in face-responsive cortex only occurred while attending to the face itself (see Figs. 3 and 4). Note, however, that null results in MVPA must be interpreted with caution and it is possible that the information was present but not decodable.

In contrast, two fronto-lateral regions successfully decoded whether participants intended to attend to a specific facial aspect before the presentation of the face, most likely reflecting domain-general processes for task preparation. The neural patterns successfully decoded a participant's internal goal independent of the respective prompt type (letters versus symbols). These regions are broadly consistent with brain regions previously implicated in task preparation (see, e.g., Ruge, Jamadar, Zimmermann & Karayanidis 2013). Task preparation is a complex process involving several steps including (but not limited to)

encoding of task cue, retrieval of relevant task set rules, inhibition of concurring/previous task rules, intention to allocate attention to specific features/aspects of the stimulus and preparation of behavioral response to actual stimulus (see, e.g. Roger & Monsell 1995). Our experimental design allowed us to distinguish the representation of the task (what to attend and what to ignore) from perception of the cue (using two different prompt formats), from preparation of the response (because the response time and mapping was unknown), and from deploying attention to the stimulus. In particular, at the time of the prompt the onset of the stimulus was unpredictable, the value of the stimulus on the attended dimension was unpredictable, and the mapping of the task to the response buttons was unpredictable. Thus, the pattern of activity in fronto-lateral regions likely reflects a representation of the task rules themselves.

Because the fronto-lateral regions we found in the whole brain searchlight analyses contained task information already at the time of the cue, and face-responsive regions only contained task information at the time of the video (and not before stimulus onset), we hypothesize that the fronto-lateral regions support shaping the change in stimulus processing implemented in the face network. This is in line with suggestions from the cognitive control literature indicating that selective representations in the fronto-lateral attention network may guide subsequent brain regions and networks that process stimulus related information (Desimone & Duncan 1995, Miller and Cohen 2001, Kanwisher & Wojciulik 2000). Additionally this interpretation is consistent with the more general idea that recurrent connections can contribute to optimize perceptual processing for context specific goals.

One open question from our study concerns the representation of the task during the delay between the prompt and the stimulus movie. A priori, one might expect a task representation to remain activated in fronto-lateral regions until the change in face processing could be implemented; that is, there should be a temporal overlap between the representation of the task in fronto-lateral regions, and the effects of the task on face processing regions. By contrast, we found that in the two fronto-lateral regions, the task-specific pattern of activity decayed over time, and was not detectable in the seconds just before the movie onset, leaving a large gap before the task effect in face-responsive cortex emerged. One interesting conjecture is that during the delay, task information can be maintained in a sparse or weak format that is not detectable using MVPA.

### Representation of the attended aspect of a face

Once the face movie was presented, patterns in almost all of the tested face-responsive regions robustly discriminated the attended aspect. Anterior and posterior STS, FFA, as well as ventral and dorsal MPFC decoded which task participants were performing when watching the videos averaging across stimulus aspect (positive vs negative emotion). For most of these regions, these patterns were robust enough to show strong within task correlations even when generalizing across the face's valence (see Fig. 4). The current results thus suggest that the representation of a face, in most regions of face-responsive cortex, is sensitive to the observer's internal goals. In other words, the task influences the representation of the faces themselves in these regions (rather then representing the task per se). Our results are consistent with many prior demonstrations that the *magnitude* of

response in these regions is affected by attention: for example, the magnitude of response is higher in FFA and pSTS when attending to a face's emotion (Vuillemier et al., 2001, Ganel et al, 2005). The current results show that in addition to changing which regions are more recruited on average, attention to a specific aspect of a face can shift the pattern of internal representation of that face in much of face-responsive cortex.

Although the task participants were performing had effects on multiple (and widespread) regions of the face network, we also found differences across regions. We found no effect in regions involved in early aspects of face processing (rOFA and EVC, see supplementary), smaller effects in rFFA (not generalizing across stimulus properties) and most reliable effects in higher-order face processing regions (pSTS, MPFC). These findings are consistent with prominent face processing models (e.g., Bruce & Young, 1986, Haxby et al., 2000) that suggest at least partially distinct cognitive and underlying neural mechanisms for processing different facial aspects within the face network. For instance, insights from congenital prosopagnosia strongly suggest separate mechanisms for emotion vs identity recognition from faces (e.g., Bate et al., 2009, Duchaine et al., 2003). Additionally, research in typical development suggests one-directional (asymmetric) routes of influence from early (e.g. EVC, OFA, FFA) to later processing stages (e.g. STS, MPFC) of processing different facial aspects (e.g. identity, emotion, ethnicity) tested with behavioral (e.g., Atkinson et al., 2005, Karnadewi & Lipp et al., 2011) or neural measures (see, e.g. Alonso-Petro et al., 2015). Our results are consistent with the idea that representations at the later stages contain more relevant information for recognizing emotional expressions, but suggest that these representations are also more flexible, potentially carrying relevant information for multiple different deliberate tasks.

### Representation of a specific stimulus property

To measure the representation of the stimulus itself, independent of the task, we focused on the valence of the emotional expression. We chose this aspect of faces because multiple prior studies suggest that emotional valence of faces is represented in (and can be decoded from) pSTS and MPFC (Said et al 2010, Peelen et al 2010, Skerry & Saxe 2014). Replicating these prior studies, emotional valence could be decoded from regions of pSTS and MPFC in our study. However, this response was weaker and less robust than the task effect, and could only be observed when participants were instructed to attend to the emotional expression (though we did not find a significant task by stimulus interaction). In contrast to the hypothesis of automatically computed stimulus property representations, these results hint that representation of emotional valence of faces, in pSTS and MPFC, is context dependent: there was stronger evidence for representation of valence when participants attended to emotion, and overall the evidence for valence information was weaker in the current study (in which participants switched between tasks) than in prior studies (when participants attended only to emotion).

Distinct patterns for positive and negative valence were found near, but not in, face-responsive regions of pSTS and MPFC, again replicating prior findings (e.g. Peelen et al 2010, Skerry & Saxe, 2015). Emotional valence appears to be represented within distinct functional subregions within pSTS and MPFC from overall responses to faces (see, e.g.

Deen et al., 2015). As a result emotional valence was not represented in any of our a priori ROIs. To directly test whether these representations are task specific, future studies will therefore need to use a different strategy to identify regions of interest (i.e. not a face-localizer).

**Implications for the Metaphor of Representational Geometry**—By looking at spatial patterns of response within face-responsive regions, rather than the average magnitude of response, we hoped to make inferences about how attention influences the representation of faces. A currently popular metaphor in cognitive science considers representations to be points in a multi-dimensional space defined by the population code of activity within a region. We anticipated that attention could expand representational space along the attended dimension, thus decreasing neural similarity along the attended dimension, while potentially contracting representation space (increasing neural similarity) along unattended dimensions (Çukur et al., 2013, Kriegeskorte eta l., 2008, Reddy et al. 2009). The empirical signature of this mechanism would be that regions containing information about facial emotion (i.e. regions in which emotion-relevant features are salient dimensions of the representational space) would show better decoding of valence when attending to emotion. While we found evidence consistent with this prediction in the whole brain analysis, the region of interest analyses revealed an unexpected pattern: large shifts in the pattern of response, even in regions that did not represent emotional valence. That is, attention seemed to change the representation of faces not only by specifically enhancing the representation of the attended dimension, but also by affecting representations along many other (as yet unknown) dimensions.

Because the effects of task were so pervasive, especially in higher-level face-processing regions, it is plausible that many other experiments using multi-voxel pattern analyses to resolve the structure of cognitive representations are also revealing the neural similarity structure *within a specific task context*. A similar insight in cognitive science is that participants' explicit similarity judgments characterize the similarity structure of a conceptual domain, only *with respect* to some task or context (Goldstone 1994). In other words, the similarity of two concepts (zebra, horse; zebra, newspaper) depends on the relevance of different attributes (animacy, color). As another example, "to say that surgeons are like butchers means something different than to say butchers are like surgeons" (Medin et al 1993). In future research, it will be critical to combine cognitive and neural approaches to characterize how attention changes the similarity structure of concepts, especially beyond just highlighting or expanding the focal dimension.

**Conclusions and Future work**—In sum, our results suggest that participants' deliberate focus of attention dramatically shapes the information represented about faces in face-responsive cortical regions: information about the intended task is 1) endogenously represented prior to stimulus onset in fronto-lateral regions and 2) at the time of the stimulus in face-responsive regions, while 3) we only found weak stimulus representation in previously reported regions. These results illustrate the powerful influence of top-down signals on cortical representations of faces at the time of stimulus processing. Future designs could for instance be adapted to include i) eye movement measurements to account for

potential different scan path on faces for task or stimulus variations or ii) different task modulations with regards to variant and invariant facial aspects (e.g. sex, trustworthiness, ethnicity (see, e.g. Karnadewi & Lipp, 2011)). We believe that the current study and the specific task design provide promising opportunities to identify group differences in stimulus and task representations between typical and atypical social cognition. Differences in of top-down and bottom-up effects on neural response patterns in face processing regions between groups might lead to important insights into their specific alterations. For instance, one possible explanation for the striking absence of clear group differences in recent studies in social cognition in Autism (e.g. Dufour et al, 2013) may be that experimental paradigms rarely capture the rapidly changing internal and external factors that are a prerequisite for effective social functioning. Investigating neural representations of flexibility in social information processing could provide a new fruitful approach to study subtle differences in social cognition in the laboratory to identify quantifiable biomarkers of atypical social information processing.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.
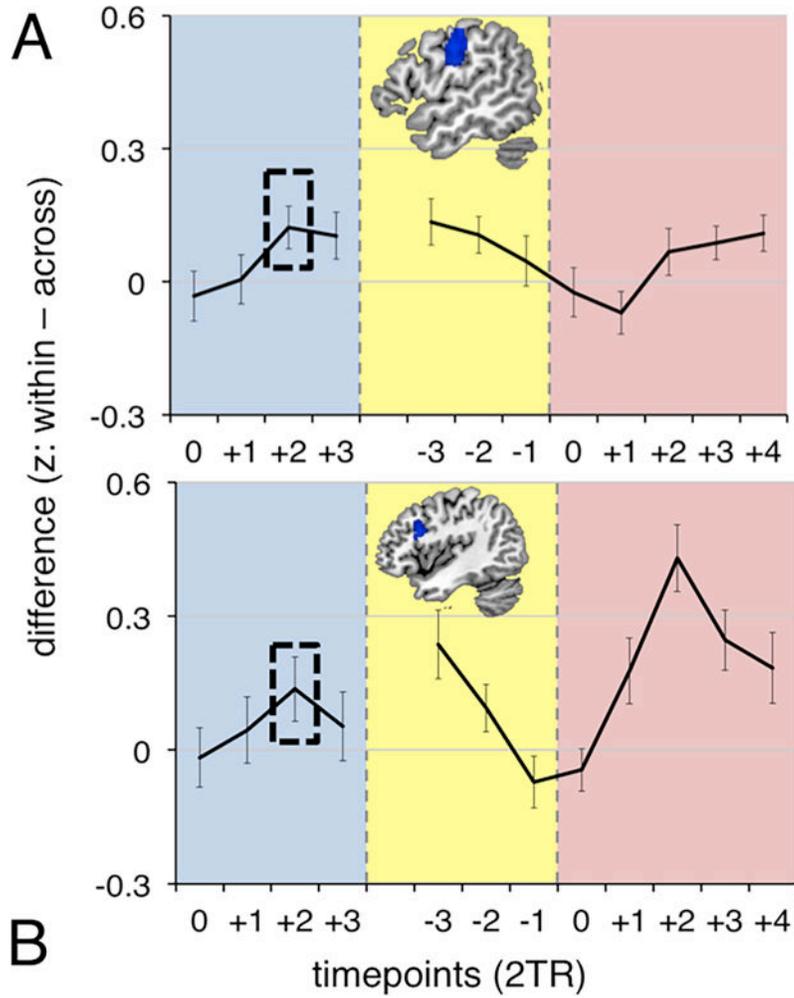
## Acknowledgments

## References

Adolphs R. Recognizing emotion from facial expressions: Psychological and neurological mechanisms. Behavioral and cognitive neuroscience reviews. 2002; 1(1):21–62. [PubMed: 17715585]

Alonso-Prieto E, Oruç I, Rubino C, Zhu M, Handy T, Barton JJ. Interactions between the perception of age and ethnicity in faces: an event-related potential study. Cognitive neuropsychology. 2015; 32(6): 368–384. [PubMed: 26226051]

Alink A, Schwiedrzik CM, Kohler A, Singer W, Muckli L. Stimulus predictability reduces responses in primary visual cortex. The Journal of Neuroscience. 2010; 30(8):2960–2966. [PubMed: 20181593]

Atkinson AP, Tipples J, Burt DM, Young AW. Asymmetric interference between sex and emotion in face perception. Perception & Psychophysics. 2005; 67(7):1199–1213. [PubMed: 16502842]

Bruce V, Young A. Understanding face recognition. British journal of psychology. 1986; 77(3):305–327. [PubMed: 3756376]

Corbetta M, Miezin FM, Dobmeyer S, Shulman GL, Petersen SE. Attentional modulation of neural processing of shape, color, and velocity in humans. Science. 1990; 248(4962):1556. [PubMed: 2360050]

Carrasco M. Visual attention: The past 25 years. Vision research. 2011; 51(13):1484–1525. [PubMed: 21549742]

Critchley HD, Daly EM, Bullmore ET, Williams SC, Van Amelsvoort T, Robertson DM, Rowe A, Phillips M, McAlonan G, Howlin P, Murphy DG. The functional neuroanatomy of social behaviour. Brain. 2000; 123(11):2203–2212. [PubMed: 11050021]

Çukur T, Nishimoto S, Huth AG, Gallant JL. Attention during natural vision warps semantic representation across the human brain. Nature neuroscience. 2013; 16(6):763–770. [PubMed: 23603707]

de Beeck HPO. Against hyperacuity in brain reading: spatial smoothing does not hurt multivariate fMRI analyses? Neuroimage. 2010; 49(3):1943–1948. [PubMed: 19285144]

Çukur T, Nishimoto S, Huth AG, Gallant JL. Attention during natural vision warps semantic representation across the human brain. Nature neuroscience. 2013; 16(6):763–770. [PubMed: 23603707]

Deen B, Koldewyn K, Kanwisher N, Saxe R. Functional organization of social perception and cognition in the superior temporal sulcus. Cerebral Cortex. 2015:bhv111.

Dehaene S, Cohen L. Cultural recycling of cortical maps. Neuron. 2007; 56(2):384–398. [PubMed: 17964253]

Delorme A, Thorpe SJ. Face identification using one spike per neuron: resistance to image degradations. Neural Networks. 2001; 14(6):795–803. [PubMed: 11665771]

Den Ouden HE, Friston KJ, Daw ND, McIntosh AR, Stephan KE. A dual role for prediction error in associative learning. Cerebral Cortex. 2009; 19(5):1175–1185. [PubMed: 18820290]

Desimone R, Duncan J. Neural mechanisms of selective visual attention. Annual review of neuroscience. 1995; 18(1):193–222.

DiCarlo JJ, Zoccolan D, Rust NC. How does the brain solve visual object recognition? Neuron. 2012; 73(3):415–434. [PubMed: 22325196]

Dubois J, de Berker AO, Tsao DY. Single-unit recordings in the macaque face patch system reveal limitations of fMRI MVPA. The Journal of Neuroscience. 2015; 35(6):2791–2802. [PubMed: 25673866]

Dufour N, Redcay E, Young L, Mavros PL, Moran JM, Triantafyllou C, Gabrieli JD, Saxe R. Similar brain activation during false belief tasks in a large sample of adults with and without autism. PloS one. 2013; 8(9):e75468. [PubMed: 24073267]

Fox CJ, Moon SY, Iaria G, Barton JJ. The correlates of subjective perception of identity and expression in the face network: an fMRI adaptation study. Neuroimage. 2009; 44(2):569–580. [PubMed: 18852053]

Ganel T, Valyear KF, Goshen-Gottstein Y, Goodale MA. The involvement of the "fusiform face area" in processing facial expression. Neuropsychologia. 2005; 43(11):1645–1654. [PubMed: 16009246]

Goldstone RL. The role of similarity in categorization: Providing a groundwork. Cognition. 1994; 52(2):125–157. [PubMed: 7924201]

Goren CC, Sarty M, Wu PY. Visual following and pattern discrimination of face-like stimuli by newborn infants. Pediatrics. 1975; 56(4):544–549. [PubMed: 1165958]

Harel A, Kravitz DJ, Baker CI. Task context impacts visual object processing differentially across the cortex. Proceedings of the National Academy of Sciences. 2014; 111(10):E962–E971.

Hariri AR, Bookheimer SY, Mazziotta JC. Modulating emotional responses: effects of a neocortical network on the limbic system. Neuroreport. 2000; 11(1):43–48. [PubMed: 10683827]

Hayasaka S, Nichols TE. Combining voxel intensity and cluster extent with permutation test framework. Neuroimage. 2004; 23(1):54–63. [PubMed: 15325352]

Haxby JV, Hoffman EA, Gobbini MI. The distributed human neural system for face perception. Trends in cognitive sciences. 2000; 4(6):223–233. [PubMed: 10827445]

Haxby JV, Gobbini MI, Furey ML, Ishai A, Schouten JL, Pietrini P. Distributed and overlapping representations of faces and objects in ventral temporal cortex. Science. 2001; 293(5539):2425–2430. [PubMed: 11577229]

Johnson MH, Dziurawiec S, Ellis H, Morton J. Newborns' preferential tracking of face-like stimuli and its subsequent decline. Cognition. 1991; 40(1):1–19. [PubMed: 1786670]

Kanwisher N, Wojciulik E. Visual attention: insights from brain imaging. Nature Reviews Neuroscience. 2000; 1(2):91–100. [PubMed: 11252779]

Karnadewi F, Lipp OV. The processing of invariant and variant face cues in the Garner Paradigm. Emotion. 2011; 11(3):563. [PubMed: 21668107]

Kok P, Jehee JF, de Lange FP. Less is more: expectation sharpens representations in the primary visual cortex. Neuron. 2012; 75(2):265–270. [PubMed: 22841311]

Kok P, Brouwer GJ, van Gerven MA, de Lange FP. Prior expectations bias sensory representations in visual cortex. The Journal of Neuroscience. 2013; 33(41):16275–16284. [PubMed: 24107959]

Kriegeskorte N, Mur M, Ruff DA, Kiani R, Bodurka J, Esteky H, Bandettini PA. Matching categorical object representations in inferior temporal cortex of man and monkey. Neuron. 2008; 60(6):1126–1141. [PubMed: 19109916]

Kravitz DJ, Saleem KS, Baker CI, Ungerleider LG, Mishkin M. The ventral visual pathway: an expanded neural framework for the processing of object quality. Trends in cognitive sciences. 2013; 17(1):26–49. [PubMed: 23265839]

Medin DL, Goldstone RL, Gentner D. Respects for similarity. Psychological review. 1993; 100(2):254.

Miller EK, Cohen JD. An integrative theory of prefrontal cortex function. Annual review of neuroscience. 2001; 24(1):167–202.

Nelson CA, Dolgin KG. The generalized discrimination of facial expressions by seven-month-old infants. Child development. 1985:58–61. [PubMed: 3987408]

Nelson CA, Morse PA, Leavitt LA. Recognition of facial expressions by seven-month-old infants. Child development. 1979:1239–1242. [PubMed: 535438]

O'Craven KM, Downing PE, Kanwisher N. fMRI evidence for objects as the units of attentional selection. Nature. 1999; 401(6753):584–587. [PubMed: 10524624]

O'Craven KM, Kanwisher N. Mental imagery of faces and places activates corresponding stimulus-specific brain regions. Journal of cognitive neuroscience. 2000; 12(6):1013–1023. [PubMed: 11177421]

Peelen MV, Atkinson AP, Vuilleumier P. Supramodal representations of perceived emotions in the human brain. The Journal of neuroscience. 2010; 30(30):10127–10134. [PubMed: 20668196]

Reddy L, Kanwisher NG, VanRullen R. Attention and biased competition in multi-voxel object representations. Proceedings of the National Academy of Sciences. 2009; 106(50):21447–21452.

Riesenhuber M, Poggio T. Hierarchical models of object recognition in cortex. Nature neuroscience. 1999; 2(11):1019–1025. [PubMed: 10526343]

Rogers RD, Monsell S. Costs of a predictable switch between simple cognitive tasks. Journal of experimental psychology: General. 1995; 124(2):207.

Ruge H, Jamadar S, Zimmermann U, Karayanidis F. The many faces of preparatory control in task switching: reviewing a decade of fMRI research. Human brain mapping. 2013; 34(1):12–35. [PubMed: 21998090]

Said CP, Moore CD, Norman KA, Haxby JV, Todorov A. Graded representations of emotional expressions in the left superior temporal sulcus. Frontiers in systems neuroscience. 2010; 4:6. [PubMed: 20305753]

Skerry AE, Saxe R. A common neural code for perceived and inferred emotion. The Journal of Neuroscience. 2014; 34(48):15997–16008. [PubMed: 25429141]

Summerfield C, Koechlin E. A neural representation of prior information during perceptual inference. Neuron. 2008; 59(2):336–347. [PubMed: 18667160]

Tottenham, N., Borscheid, A., Ellertsen, K., Marcus, D., Nelson, CA. The NimStim face set. 2002. Retrieved from http://www.macbrain.org/faces/index.htm

Vuilleumier P, Armony JL, Driver J, Dolan RJ. Effects of attention and emotion on face processing in the human brain: an event-related fMRI study. Neuron. 2001; 30(3):829–841. [PubMed: 11430815]

Van Essen DC, Anderson CH, Felleman DJ. Information processing in the primate visual system: an integrated systems perspective. Science. 1992; 255(5043):419–423. [PubMed: 1734518]

Zaki J, Ochsner K. The need for a cognitive neuroscience of naturalistic social cognition. Annals of the New York Academy of Sciences. 2009; 1167(1):16–30. [PubMed: 19580548]

Zhang J, Meeson A, Welchman AE, Kourtzi Z. Learning alters the tuning of functional magnetic resonance imaging patterns for visual forms. The Journal of Neuroscience. 2010; 30(42):14127–14133. [PubMed: 20962233]
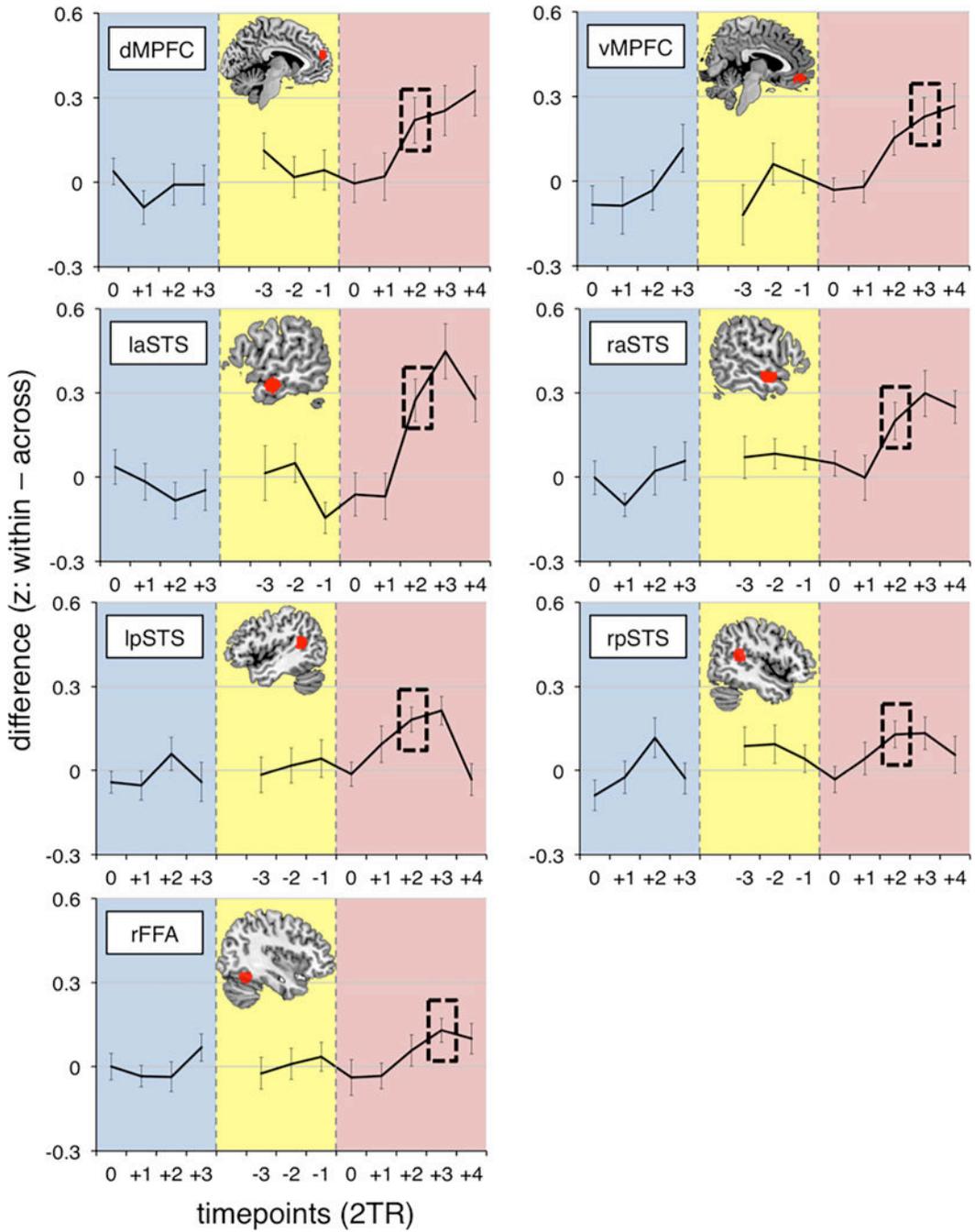
**Figure 1.**
Emotion/age attribution task. Trials start with a prompt indicating whether participants have to judge emotion or age (blue). Two prompt formats were used: face symbols or letters. Then there was a jittered delay, with a blank screen (yellow). Next participants saw a naturalistic video, of a dynamic facial emotional expression (red). Then participants saw the response mapping, and made their response (grey). To account for reaction time (RT) each trial is modeled as a nuisance regressor per participant.

**Figure 2.**
Difference in average within versus across task (emotion vs age) correlation over time in two latero-frontal regions in left medial/inferior frontal gyrus (A, upper) and precentral gyrus (B, lower) over the course of a trial. The onset on the left (blue background) on the x-axis reflects the timepoints in response to the prompt. The onset on the right (red background) reflects the timepoints in response to the actual video, whereas the middle area reflects the timepoints before the onset of the video (yellow background). The box (dashed line) indicates the earliest significant decoding (p < .05). Error bars reflect standard error of the mean, across subjects.

**Figure 3.**
Difference in average within versus across task (emotion vs age) correlation over time in 7 independently identified face-responsive regions over the course of a trial. The onset on the left (blue background) on the x-axis reflects the timepoints in response to the prompt. The onset on the right (red background) reflects the timepoints in response to the actual video, whereas the middle area reflects the timepoints before the onset of the video (yellow background). The box (dashed line) represents the earliest significant decoding (corrected for number of face-regions tested.05/8, p < .00625). The rOFA ROI here is not shown

because of the absence of a task effect (see Table 2). Error bars reflect standard error of the mean, across subjects.

**Figure 4.**
Differences in within vs across across task correlations (z-scored) per ROI at the time of the prompt (blue background) for *averaged* (training and testing on both positive and negative emotions) and *generalized* (training on videos depicting one emotion (positive vs negative valence) stimulus condition comparisons. Decoding accuracies differed between regions (main effect of region) and were stronger when averaging as compared to generalizing (main effect of condition-comparisons). Error bars reflect standard error of the mean, across subjects.

**Figure 5.**
Whole-brain searchlight for negative vs positive emotions averaged over the age and emotion task (upper) and only in the emotion task (lower).

**Table 1**

**Number of voxels in hypothesis space, peak voxel coordinates and number of subjects passing the 80 voxels criterion per ROI**

Abbreviations: rpSTS, right posterior superior temporal sulcus; lpSTS, left posterior superior temporal sulcus; raSTS, right anterior superior temporal sulcus; laSTS, left anterior superior temporal sulcus; rFFA, right fusiform face area; dMPFC, dorsal medial prefrontal cortex; vMPFC, ventral medial prefrontal cortex; n, number of.

|  | n(voxels) | peak voxel | | | n(subjects) |
|---|---|---|---|---|---|
|  |  | x | y | z |  |
| **rpSTS** | 365 | 48 | −42 | 16 | 25 |
| **lpSTS** | 340 | −46 | −48 | 16 | 25 |
| **raSTS** | 268 | 58 | −4 | 268 | 25 |
| **laSTS** | 269 | −62 | −12 | −18 | 25 |
| **rFFA** | 199 | 40 | −56 | −18 | 25 |
| **rOFA** | 304 | 36 | −82 | −18 | 25 |
| **dMPFC** | 169 | −8 | 54 | 20 | 22 |
| **vMPFC** | 272 | 4 | 48 | −14 | 25 |

**Table 2**

**Differences in within vs across condition correlations comparing task type (emotion vs age) per ROI based on beta values averaged and generalized across stimulus aspect (emotional expression). Bold values indicate significance after correcting for multiple comparison correction**

Abbreviations: rpSTS, right posterior superior temporal sulcus; lpSTS, left posterior superior temporal sulcus; raSTS, right anterior superior temporal sulcus; laSTS, left anterior superior temporal sulcus; rFFA, right fusiform face area; dMPFC, dorsal medial prefrontal cortex; vMPFC, ventral medial prefrontal cortex.

| Region | | Within | Across | Significance |
|---|---|---|---|---|
| rpSTS | *averaged* | 2.1(.09) | 1.9(.09) | t(24) = 2.8, **p = .0053** |
| | *generalized* | 1.9(.09) | 1.8(.09) | t(24) = 2.9, **p = .0041** |
| raSTS | *averaged* | 1.6(.09) | 1.4(.07) | t(24) = 2.7, p = .007 |
| | *generalized* | 1.4(.08) | 1.3(.07) | t(24) = 2.5, p = .0088 |
| lpSTS | *averaged* | 1.9(.09) | 1.7(.09) | t(24) = 4.9, **p = 2.6e–05** |
| | *generalized* | 1.7(.09) | 1.6(.09) | t(24) = 4.8, **p = 3.5e–05** |
| laSTS | *averaged* | 1.2(.09) | .99(.09) | t(24) = 4.2, **p = .00015** |
| | *generalized* | 1(.08) | .88(.08) | t(24) = 3.2, **p = 0.0021** |
| dMPFC | *averaged* | 1.4(.09) | 1.2(.09) | t(21) = 4.8, **p = 4.9e–05** |
| | *generalized* | 1.2(.09) | 1(.09) | t(21) = 3.6, **p = .00076** |
| vMPFC | *averaged* | 1.4(.10) | 1.3(.09) | t(24) = 3.2, **p = .002** |
| | *generalized* | 1.2(.09) | 1.1(.07) | t(24) = 3.4, **p = .0011** |
| rFFA | *averaged* | 2.1(.10) | 2(.09) | t(24) = 3.1, **p = .0023** |
| | *generalized* | 1.9(.10) | 1.9(.09) | t(24) = .96, p = .17 |
| rOFA | *averaged* | 2.2(.11) | 2.1(.12) | t(24) = .97, p = .17 |
| | *generalized* | 2.0(.12) | 2.0(.13) | t(24) = .05, p = .48 |
| EVC | *averaged* | 2.3(.08) | 2.4(.08) | t(24) = −.49, p = .69 |
| | *generalized* | 2.2(.08) | 2.2(.08) | t(24) = −.3, p = .62 |

**Table 3**

**Results of the whole-brain searchlight analysis (uncorrected p < .001, k > 30) on the influence of stimulus aspect showing brain regions, cluster extent, local peaks in MNI space, peak (pseudo) t-value**

Abbreviations: rpSTS, right posterior superior temporal sulcus; lpSTS, left posterior superior temporal sulcus; raSTS, right anterior superior temporal sulcus; laSTS, left anterior superior temporal sulcus; rFFA, right fusiform face area; dMPFC, dorsal medial prefrontal cortex; vMPFC, ventral medial prefrontal cortex; n, number of.

| cluster | region | n voxel | x | y | z | peak |
|---|---|---|---|---|---|---|
| | *negative vs positive (averaged over tasks)* | | | | | |
| 1 | *right STS* | 118 | 60 | −48 | 6 | 3.87 |
| 2 | right middle occipital gyrus | 101 | 28 | −90 | 0 | 3.82 |
| 3 | right lateral inferior fusiform gyrus | 37 | −30 | −24 | −32 | 3.49 |
| | *emotion task: negative vs positive* | | | | | |
| 1 | right middle occipital gyrus | 329 | 22 | −92 | 2 | 4.1 |
| 2 | *right STS* | 44 | 46 | −46 | 4 | 3.82 |
| 3 | MPFC | 56 | 16 | 58 | 2 | 3.81 |
| 4 | right superior occipital gyrus | 31 | 42 | −84 | 26 | 3.67 |
| | *age task: negative vs positive* | | | | | |
| 1 | right middle occipital gyrus | 53 | −26 | −90 | 0 | 3.95 |
| 2 | left middle occipital gyrus | 51 | 28 | −90 | −2 | 3.58 |