

Bipartite functional fractionation within the neural system for social cognition supports the psychological continuity of self versus other

Rocco Chiou ^{1,2,*}, Christopher R. Cox ³, Matthew A. Lambon Ralph ¹

¹MRC Cognition & Brain Science Unit, University of Cambridge, Cambridge, CB2 7EF, United Kingdom.

²Wellcome Centre for Integrative Neuroimaging, University of Oxford, OX3 9DU, Oxford, United Kingdom.

³Department of Psychology, Louisiana State University, LA 70803, Baton Rouge, United States

*Corresponding author: rocco.chiou@ndcn.ox.ac.uk

Abstract

Research of social neuroscience establishes that regions in the brain's default-mode network (DN) and semantic network (SN) are engaged by socio-cognitive tasks. Research of the human connectome shows that DN and SN regions are both situated at the transmodal end of a cortical gradient but differ in their loci along this gradient. Here we integrated these 2 bodies of research, used the psychological continuity of self versus other as a “test-case,” and used functional magnetic resonance imaging to investigate whether these 2 networks would encode social concepts differently. We found a robust dissociation between the DN and SN—while both networks contained sufficient information for decoding broad-stroke distinction of social categories, the DN carried more generalizable information for cross-classifying across social distance and emotive valence than did the SN. We also found that the overarching distinction of self versus other was a principal divider of the representational space while social distance was an auxiliary factor (subdivision, nested within the principal dimension), and this representational landscape was more manifested in the DN than in the SN. Taken together, our findings demonstrate how insights from connectome research can benefit social neuroscience and have implications for clarifying the 2 networks' differential contributions to social cognition.

Key words: default-mode network; cortical organization; semantic network; MVPA; self.

Introduction

Decades of research suggests that humans deduce other individuals' mental states by simulating how oneself would think and feel in similar situations (e.g. [Steinbeis 2016](#)). This ability is rooted in the awareness that self and others are distinct yet relatable social beings. The rudimentary sense of “self vs. other” emerges during infancy, while self-concept (a sophisticated understanding of ourselves in relation to others under social contexts) evolves over lifespan. Neuroimaging evidence (for review, see [Yeshurun et al. 2021](#)) has established that the brain's default-mode network (DN) is robustly engaged by various self-referential processes. This leads to the view that the DN is the key neural substrate that underpins “core self”/“ego” ([Carhart-Harris and Friston 2010](#)). Rather than a homogenous structure, converging results from seed-based connectivity, data-driven parcellation, and task-induced activation have further indicated that the DN comprises functionally distinct and anatomically separable subnetworks (e.g. [Braga and Buckner 2017](#); [Braga et al. 2019](#); [Chiou et al. 2020](#)). The division within the DN into subnetworks begs an important question—whether and how each subnetwork

differentially contributes to the representation of self. In the present study, we investigated this issue by testing how the psychological continuity from “selfness” to “otherness” is encoded in 2 major subnetworks of the DN using multivoxel pattern analysis (MVPA). Our finding unraveled a robust disparity—the neural coding of one subnetwork reliably contained more information about the “self-to-other” continuity than the other subnetwork. Below we first discuss the neural correlates of self-referential processes and the limitation of univariate approach. Next, we discuss evidence for a bipartite fractionation within the DN. We specifically focus on how such a fractionation constrains our scope of scrutiny for neural substrates, and how we can exploit MVPA to investigate the neural coding for “self vs. other” within this bipartite structure.

Self-concept is a collection of beliefs about oneself, embodying the contents of contemplation about “Who am I?” ([Oyserman et al. 2012](#)). It is inextricably intertwined with many critical issues in cognitive and social psychology, including perception of the aesthetics and attractiveness of one's own body (self-image), the capability for reflecting on one's dispositions

Received: May 14, 2021. **Revised:** February 10, 2022. **Accepted:** March 14, 2022

© The Author(s) 2022. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

(self-knowledge), one's recollection of life events (autobiographical memory), etc. Previous research has shown that self-concept depends, in a fundamental way, upon one's capacity to represent self as a psychologically coherent entity persisting through time, whose *past self* is construed as an entity closely related to yet partially separable from the *present self* (Klein et al. 2002). In addition to representing *present self* as a continuation of *past self*, research has shown that self-concept is impacted by the presence of other people in a social environment. Especially, during childhood and adolescence, self-concept is formulated through interaction with significant others, which spawned numerous investigations of familial and peer influences on personality and self-esteem (e.g. Deković and Meeus 1997; Verschueren et al. 2012). The significance of self-related processing is not only seen in the field of social psychology but also in cognitive neuroscience. Researchers have long understood that self-related processing is one of the major neurocognitive dimensions that characterize the functionality of the cortical midline structures, particularly the ventromedial prefrontal cortex (vmPFC; for review, see Lieberman et al. 2019). For example, compared to information regarding another individual, processing information with reference to self amplifies vmPFC activity. This has been robustly observed when one evaluates descriptions with respect to self (Kelley et al. 2002), recalls personal experiences (Spreng and Grady 2010), or even recognizes one's own name/address that is ostensibly unrelated to the task at hand (Moran et al. 2006). In addition, patients with vmPFC lesion fail to show self-referential memory advantage (superior mnemonic performance for information related to self that healthy people reliably show), despite those patients having otherwise intact memory for nonself information (Philippi et al. 2012). Research also indicates that the vmPFC is sensitive to social distance with regard to self, with friends eliciting greater vmPFC activity than strangers, presumably due to one's construal of a "friend" being inextricably linked to oneself (Krienen et al. 2010). In a similar vein, vmPFC activity was found to reflect whether someone's political opinion is in agreement with oneself, with encountering someone holding concordant views eliciting greater vmPFC activity than someone holding opposite views (Mitchell et al. 2006). Together, multiple threads of evidence consistently suggest that the vmPFC is a key neural structure underlying self-related processing and is involved in processing other individuals when they are deemed similar/related to self.

The univariate approach is useful for identifying the loci of maximally responding clusters/regions for a certain cognitive process, giving a bird's eye view about where the "hotspot" is. In the case of self-related processing, the "subtraction" paradigm is often used—the univariate response of each voxel (or averaged response within a region of interest) elicited by self-related processing is compared with the response driven by evaluating, imagining, or reminiscing about someone else. While this approach is useful in identifying

conspicuous targets in the brain (i.e. contiguous voxels that form a cluster surpassing thresholds), it has limited use in unveiling nuances that are jointly encoded by distributed patterns of neural responses across voxels. With MVPA, however, researchers are able to reveal how distributed neural patterns represent the subtle distinction between self- and other-referential processes. For example, the patterns of vmPFC activity have been used to decode if one was entertaining self- or other-referential thoughts (e.g. Chavez et al. 2016; Yankouskaya et al. 2017). In addition to deciphering self versus other using vmPFC patterns, MVPA has been applied to decode various facets of social cognition using neural patterns elsewhere in the brain, such as decoding the identities of fictitious characters (Hassabis et al. 2013), the identities of personally familiar people (Thornton and Mitchell 2017), and whether a behavior shows goodwill or malice (Koster-Hale et al. 2013). Together, these studies demonstrate how MVPA can be used to elucidate our understanding about the neural basis of social cognition, capturing the granularity in collective, distributed codes that often goes unnoticed by the univariate analysis (for review of recent progress, see Wagner et al. 2019).

While the vmPFC is heavily studied, it is not the only region involved in self-related processing. Alongside vmPFC activity, a group of widely distributed brain regions, which collectively form the brain's specialized system for social cognition, have been reliably found to be involved in various self-related/socio-cognitive processes (for review, see Doré et al. 2014). For instance, vmPFC activity usually elevates in tandem with activity of other midline structures—the posterior cingulate cortex (PCC), retrosplenial cortex (RSC), and dorsomedial prefrontal cortex (dmPFC)—when a task requires assessing the personality of people, retrieving autobiographical episodes, or envisaging prospective and counterfactual scenarios that involve human activity. Apart from the brain's medial aspect, self-referential/social-cognitive tasks also recruit a set of regions on the brain's lateral surface (Olson et al. 2013; Doré et al. 2014; Binney et al. 2016; Chiou et al. 2020), including the bilateral anterior temporal lobes (ATL), the left inferior frontal gyrus (IFG), and the inferior parietal lobule (IPL). Resting-state research has shown that all these medial and lateral areas, conjointly as the system for social cognition, tend to stay intrinsically connected during task-free resting moments and overlap substantially with the cortical realm of the DN. Regions implicated in social cognition exhibit a similar functional profile that has been used to characterize the tendencies of DN regions (Raichle et al. 2001; Yeshurun et al. 2021): Akin to the DN, social regions are more active when a situation necessitates integrating internally constructed representations (e.g. memories or schemas) with external signals (e.g. Murphy, Poerio, et al. 2019a). Also resembling the DN, activation of the social system abates when a situation requires externally oriented sensorimotor processes that minimally involve any internal representations

(e.g. Chiou et al. 2020). Importantly, converging evidence from multiple investigations have shown that the extensive DN can be fractionated into (at least) 2 subsystems (Humphreys et al. 2015; Braga and Buckner 2017; Braga et al. 2019; Jackson et al. 2019; Chiou et al. 2020). While the nomenclature used to name the subsystems of DN varies between studies, a bipartite structure has been reliably observed in the profiles of different regions' connectivity alliance and task-driven reaction, partitioning the DN into 2 modules. One subsystem comprises the vmPFC, the PCC/RSC, the posterior part of IPL, and the hippocampal formation. Some nodes of this subsystem are suggested to be the "hub" areas of the DN (Andrews-Hanna et al. 2010), and they are altogether dubbed "Network-A" by more recent studies (e.g. Braga and Buckner 2017). These regions show a strong propensity to be "suppressed" by contexts that require externally focused sensory-motoric processes (and "activated" by introspective processes) and are the quintessential "task-negative" areas by traditional views of the DN (Raichle et al. 2001). The other subsystem consists of the bilateral ATL, the left IFG, and temporoparietal junction. In the field of semantic cognition, these regions are dubbed the "semantic network" (SN; owing to their robust activity during semantic processes), while in the field of connectome research, they are dubbed "Network-B." While many areas of the SN are incorporated under the umbrella of DN regions (Yeo et al. 2011), recent evidence has indicated that the SN is a functionally separable entity from DN. Specifically, DN and SN regions show a similar inclination to deactivate in outwardly oriented sensorimotor tasks (although the extent of "aversion" to sensorimotor processing is significantly more moderate in SN regions, see Chiou et al. 2020). However, the 2 subnetworks diverge on reaction to semantic processing—whereas SN activity intensifies for verbal and nonverbal semantic processing, DN activity attenuates (Humphreys et al. 2015; Chiou et al. 2018; Jackson et al. 2019). The partial dissociation between the 2 subnetworks was found both in resting- and task-state functional magnetic resonance imaging (fMRI) studies. Together, multiple lines of inquiries have consistently demonstrated that the DN and SN are functionally distinct entities (while both are heavily engaged by social cognition). It is noteworthy that there is less consensus on the taxonomy of the dmPFC. The dmPFC shares a similar functional profile to other DN regions during a variety of social-cognitive tasks (Hiser and Koenigs 2018) and is reliably linked with other core nodes of the DN, particularly the vmPFC and PCC/RSC (e.g. Bzdok et al. 2015; Eickhoff et al. 2016; Jackson et al. 2020). However, the dmPFC is also affiliated with the SN—relative to the vmPFC as a seed, the dmPFC has tighter functional coupling with many SN regions in resting state (Bzdok et al. 2013), and is more reactive to tasks that lay emphasis on extracting semantic meaning from lexical stimuli than other DN regions (Chiou et al. 2020).

Pertinent to our main question, while the DN and SN dissociate on their univariate reaction to different contexts, it remains unclear whether the 2 subnetworks carry distinct fine-grained multiple voxel patterns that convey different information about self vs. other.

In the present study, we investigated the distributed neural coding that underpins the continuity of "self vs. other" representations. Self-concept is known to evolve along one's lifetime. However, although one's beliefs about *present self* might be radically different from his/her opinions about *past self*, the human mind maintains the stability/continuity of self-identity, treating *present self* and *past self* as dissociable yet associated entities that represent the same person (Northoff 2017). Psychological continuity is also applicable to the construal of self versus others, with a close other (or personally familiar person) perceived as more affiliated with (similar to) self than a distant other (personally unrelated person). It remains unknown how such continuity is encoded in the brain. We instigated this issue by gradationally manipulating social distance. We established 4 points of reference: the participant's sense of *present self*, the participant's sense of *past self* 10 years ago, a personally familiar other—the participant's mother, and a personally unfamiliar but well-known other—Queen Elizabeth II. Behavioral rating confirmed that participants rated the 4 persons along a continuum, with *present self* being most distinct from the Queen while *past self* and *mother* situated serially somewhere in between. While undergoing fMRI, participants read descriptions (either positive or negative personality traits) in contexts of each of the 4 individuals and judged whether the description aptly characterized the person. We employed 2 types of MVPA, machine-learning classification (Pereira et al. 2009) and representational similarity analysis (Nili et al. 2014), to investigate how the 2 subnetworks encoded the psychological continuity of self versus other and the possible divergence between DN and SN.

Replicated across multiple analyses, we found (i) the broad-stroke division of self versus other was the principal axis of representational space while social distance was an auxiliary axis, nested within the principal dimension; (ii) the DN dissociated from the SN, with the former carrying more information about personal identities. In the section of Discussion, we expound on how recent progress in understanding the topography of cortical mantle (particularly studies that demonstrated that the DN is the apex of a cortical hierarchy; Margulies et al. 2016) could be leveraged to explain the dyadic split between the DN and SN and their differential contributions to cognition.

Materials and methods

Participants

Twenty-four volunteers gave informed consent before the fMRI experiment. The sample consisted of a 10/14

female-to-male ratio, with average age = 33 years and standard deviation (SD) = 11. All volunteers are right-handed and speak English as their mother tongue. All of them completed the MRI safety screening questionnaire before the experiment and reported not having any neurological or psychiatric condition. This study was reviewed and approved by the local research ethics committee.

Experimental design

Participants completed 2 experiments while undergoing fMRI in a single session. In the main experiment, participants read short phrases describing various traits of personality or temperament, either positive or negative, and made a binary button response to answer whether they reckoned the phrase rightly described the characteristics of the specific person that was under consideration. In the localizer experiment, participants read short narratives describing either events that involved human activity or events that involved alterations of the physical environment and made a binary button response to answer a comprehension question related to the narrative they just read. A session began with the acquisition of each participant's anatomical scan, followed by the main experiment that consisted of 8 functional runs of scanning, and ended with the localizer experiment that contained only a single run.

The main experiment had a 4×2 factorial design (4 individuals: *Present Self*, *Past Self*, *Mother*, and *Queen Elizabeth II*; 2 sides of emotive valence: positive vs. negative traits). We adopted and modified a well-established fMRI paradigm that has been widely used to assess the neural substrates of self- and other-referential processing (e.g. Kelley et al. 2002; Heatherton et al. 2006; Meyer and Lieberman 2018). The task was to read a short description in each trial and to evaluate whether or not it appropriately depicted a particular individual. When performing the task, participants were presented with a fixation dot (0.5 s) in each trial, followed by words (3.3 s). They were required to make a response within the 3.3-s time limit. The target of assessment (*Present Self*, *Past Self*, *Mother*, and *Queen*) was shown above the fixation dot, and a short phrase describing a certain personality or temperament trait was shown below (e.g. "Sincere to friends" or "Anxious about uncertainty"). When the target was *Past Self*, participants reflected on themselves specifically 10 years ago with respect to the phrase. Stimuli were presented using a block design, controlled with E-Prime (Psychology Software Tools). Each run consisted of 16 blocks of trials, with each of the 8 conditions having 2 blocks. Across the runs and participants, the order in which the 8 task conditions were presented was counterbalanced so that each task condition was equally likely to appear in each of the 128 possible slots of the 8 runs (i.e. each condition was equally probable to precede or succeed any other condition), with stimuli randomly drawn from a designated stimuli set for a given run and shuffled across blocks. The stimuli sets were

also counterbalanced across participants; thus, each set was equally likely to be presented in each run. This fully counterbalanced design is vital for the subsequent leave-one-run-out decoding analysis, ensuring that every task condition and stimuli set appeared in each fold of the cross-validation. Each block was 19-s long, containing 5 trials and no inter-trial interval. Each run of scanning was 380-s long, containing 16 task blocks, 15 inter-block intervals (blank screen, 5-s each), and a 1-s blank at the end. All text stimuli were white in color, Arial typeface, 28-point in font size displayed on a black background. Participants reacted to the questions by pressing 1 of the 2 designated buttons on a MR-compatible response pad with their right index or middle finger. All visual stimuli were displayed using high-resolution LCD goggles (NordicNeuroLab) mounted on top of the head coil.

A total of 160 short phrases were used for evaluating personality traits. Each phrase contained 3 or 4 words. A half of the phrases were designed to convey positive meaning, whereas the remaining half conveyed negative meaning (the complete set of stimuli are reported in [Supplementary Material](#)). To ascertain that the phrases express the emotive meaning as intended, we asked 6 volunteers (none participated in the later fMRI study) to rate the emotive valence of the phrases using a 5-point scale (1 being most negative, 5 being most positive). Results of rating support the adequacy of our stimuli: By-subject analysis showed that phrases designed to convey negative meaning (average \pm SEM: 1.7 ± 0.1) were rated significantly lower than those designed to convey positive meaning (average \pm SEM: 4.4 ± 0.1 ; $t_{(5)} = 14.5$, $P < 0.0001$). This difference was also seen in by-item analysis for every volunteer (all P s $< 10^{-10}$). The length of text stimuli was also equated: No difference was found between the letter counts between negative (average \pm SD: 22 ± 3.4) and positive phrases (average \pm SD: 21 ± 3.8 , $P > 0.21$, n.s.). Crucially, each of the 160 phrases was equiprobable to be assessed with reference to any of the 4 target individuals. Namely, every phrase appeared 4 times in the experiment but referred to a different individual each time it was presented. This circumvented the potential biasing effect of specific stimuli by equating their frequency in every condition and ensured that none of the stimuli was repeated—each combination of description and person was encountered only once during the experiment and appeared "novel" from a participant's perspective. This set-up also ensured that identical stimuli (descriptions of personality) were used in each condition of the 4 target individuals, with the only difference being the cue word that prompted the current target.

The localizer experiment was based on a well-established paradigm that has been repeatedly used to assess the brain regions associated with the processes of mentalizing/theory of mind (e.g. Saxe and Kanwisher 2003; Saxe and Wexler 2005; Dodell-Feder et al. 2011). There were 2 conditions in this localizer run: In the Social condition, participants read a narrative describing

human interactions; afterwards, they were presented with a statement about the beliefs or feelings that a person in the narrative might have, and were asked to verify whether the statement is true or false. Answering such questions entailed changing perspectives and making inferences about someone's mental states. In the Nonsocial condition, participants read a narrative describing the physical state of the world (e.g. "a large oak tree stood in front of the city hall from the time the building was built...") and answered a comprehension question related to the nonsocial narrative. Participants had 10 s to read the narrative and were probed with an ensuing question that was shown for 4 s. The localizer run was 520-s long, containing 20 task blocks (10 Social blocks, 10 Nonsocial blocks) and 20 blank intervals (each 12-s long) that followed each task block. The 2 conditions were presented in an alternating order and counterbalanced across participants (i.e. a half of them started with the Social condition; the other half started with the Nonsocial). We used the same story materials previously used by [Dodell-Feder et al. \(2011\)](#). We adopted this localizer task due to the fact that, in the original study, this paradigm had proved effective in reliably activating a set of widely distributed brain regions known to be sensitive to mentalization in specific (and social cognition in general), which incorporates all of our regions of interest (ROIs) in the DN, as well as in the SN (see the details of how we defined the voxels in the section of Regions of Interest).

MRI acquisition

Some of the regions of our primary interest are situated in the rostroventral aspects of the brain (e.g. the ATL and the vmPFC), which are known to be particularly susceptible to signal-dropout issues ([Visser et al. 2010](#)). To combat signal-dropout in these areas, we adopted a dual-echo EPI sequence, which has been demonstrated to effectively improve signal-to-noise ratio in dropout-prone regions, compared to other conventional imaging protocols (for precedents using this dual-echo acquisition protocol, see [Halai et al. 2014](#); [Jackson et al. 2015](#); [Chiou and Lambon Ralph 2019](#)). Scans were acquired using a 3T Phillips Achieva scanner equipped with a 32-channel coil and a SENSE factor of 2.5. Using this protocol, each scan consisted of 2 images acquired simultaneously with 2 echo times: a short echo optimized to obtain maximum signal from the ventral parts and a long echo optimized for whole-brain coverage. The sequence included 31 slices covering the whole brain with repetition time (TR)=2.8 s, short/long echo times (TE)=12/35 ms, flip angle=85°, field of view (FOV)=240 × 240 mm, resolution matrix=80 × 80, slice thickness=4 mm, and voxel dimension=3 × 3 mm on the x-axis and y-axis. To reduce ghosting artifacts in the temporal lobes, all functional scans were acquired using a tilted angle, upward 45° off the AC-PC line. For the main experiment, the EPIs were collected over 8 runs; each run was 380-s long during which 136 dynamic volumes

were acquired (alongside 2 dummy scans, discarded). For the localizer experiment, the EPIs were collected from a single run (520 s); 186 dynamic volumes (and 2 dummies) were acquired. To tackle field inhomogeneity, a B_0 field map was acquired using identical parameters to the EPIs except for the following: TR=599 ms, short/long TEs = 5.19/6.65 ms. Total B_0 scan time was 1.6 min. A high-resolution T1-weighted structural scan was acquired for spatial normalization (260 slices covering the whole brain with TR=8.4 ms, TE=3.9 ms, flip angle=8°, FOV=240 × 191 mm, resolution matrix=256 × 163, and voxel size=0.9 × 1.7 × 0.9 mm).

Preprocessing and GLM

An established procedure was used to combine the 2 volumes from the dual-echo dataset. Using SPM8 (Wellcome Department of Imaging Neuroscience), we integrated the standard preprocessing procedure (realignment, slice-time correction, co-registration, and the linear integration of the long- and short-echo images) with B_0 field-map correction to prevent distortion due to inhomogeneity. The linear averaging approach has been well established in previous studies (e.g. [Poser et al. 2006](#); [Halai et al. 2014](#); [Chiou et al. 2018](#)). The combined images were realigned using rigid body transformation (correction for motion-induced artifacts) and unwrapped using B_0 field map (correction for field inhomogeneity). The integrated EPIs were then co-registered with each participant's T_1 anatomical image. For the first-level individual analysis, the β -weight of each experimental regressor was estimated by convolving each task block with a canonical hemodynamic response function. Six motion parameters were added into the model as nuisance covariates in the general linear model (GLM). Behavioral reaction times were also modeled as parametric modulators to account for the influence of fluctuating reaction times within a condition. For the main experiment, each of the 8 experimental conditions was modeled explicitly as a separate regressor, while resting baseline was modeled implicitly. For the localizer experiment, the entire 14-s duration of "story" and "question" intervals was convolved with a canonical hemodynamic response function, as per previous studies using the same localizer paradigm (e.g. [Dodell-Feder et al. 2011](#); [Skerry and Saxe 2014, 2015](#)). Low-frequency drifts were removed using a high-pass filter of 128 s. Normalized beta-estimates associated with each voxel and each regressor were subsequently submitted to MVPA.

Decoding analysis was performed on each participant's brain in native space, without spatial normalization and smoothing. Normalization [Montreal Neurological Institute (MNI)] and smoothing [full-width at half-maximum (FWHM)=8 mm] were only done on the individual whole-brain searchlight outcomes (accuracy maps), prior to the second-level random-effect analysis. The decoding accuracy maps of whole-brain searchlight analysis were normalized into the MNI standard space using the DARTEL Toolbox of SPM ([Ashburner 2007](#)),

which has been shown to produce highly accurate inter-subject alignment (Klein et al. 2009). Specifically, the T_1 -weighted image of each subject was partitioned into gray matter, white matter, and CSF tissues using SPM8's "Segmentation" function; afterwards, the DARTEL toolbox was used to create an average template combining all participants of the group. The gray matter component of this template was registered into the SPM's gray matter probability map (in MNI) using affine transformation. In the process of creating the group's template using individual T_1 , for each individual DARTEL estimated "flow fields" that contained the parameters for contorting native T_1 -weighted images to the group template. SPM8 deformation utility was then applied to combine group-to-MNI affine parameters with each participant's "flow fields" to enable tailored warping into the standard MNI space. At the end of the procedure, voxel-size of the whole-brain map was resampled to $3 \times 3 \times 3$ mm. Smoothing on the normalized accuracy maps was then applied using an 8-mm Gaussian FWHM kernel, consistent with prior studies (e.g. Halai et al. 2014; Jackson et al. 2015).

Regions of interest

As discussed earlier, we adopted a localizer paradigm that had previously proved effective in detecting various target regions in the DN and SN. In the original study that reported this fMRI paradigm, Dodell-Feder et al. (2011) tested a sizable sample of 62 participants and contrasted the Social condition against the Nonsocial one. Using this contrast, they detected robust activation in 5 areas of the DN—the dmPFC, vmPFC, PCC, and left/right IPL, as well as in 3 areas of the SN—the left/right ATL and IFG. Because the Dodell-Feder et al. results provided a useful exemplar regarding the loci of neural activity, their group-level peak coordinates (in the MNI space) could be used to guide the localization of ROIs in our participants' native space. For each participant's localizer data, we applied a t-test, voxel-wise thresholded at $P < 0.001$, to generate a whole-brain map of t-values to identify voxels that responded more intensely to the Social than Nonsocial condition. Using each participant's reversed-normalization parameters computed by SPM, we first identified 8 "landmark" points (5 DN regions, 3 SN regions) in each person's native brain that corresponded to the coordinates from the Dodell-Feder et al. study. Next, we identified the local maxima nearest to each of the 8 "landmark" points and created a spherical ROI (radius = 10 mm) centered at the local maximal coordinate. If no activity was detected at $P < 0.001$, we repeated the procedure using $P < 0.005$ and $P < 0.01$. This procedure constrained the localization of ROIs using the group-level results from Dodell-Feder et al. (2011) while allowing subject-specific variation in functional activation. For every participant, we were able to identify 8 ROIs—5 regions of the DN (the dmPFC, vmPFC, PCC, and left/right IPL) and 3 regions of the SN (the left/right ATL, and IFG). Meticulous care was taken to ensure that all

ROIs were spatially mutually exclusive—namely, there was no overlap in the voxels contained in each functionally defined ROI (particularly for the dmPFC and vmPFC). It is worth emphasizing that, with an independent localizer task to select the ROIs, our subsequent MVPA analyses were devoid of the statistically "double-dipping" issue (Button 2019).

Multivoxel pattern analysis

Prior to multivariate decoding, we fitted the main experiment's data of each participant using a standard GLM, implemented in SPM8, to compute the β -estimates of voxel-wise activation elicited by each experimental regressor. Beta-estimates were computed for each of the 8 experimental conditions (4 individuals: *Present Self*, *Past Self*, *Mother*, and *Queen* by 2 valence qualities: Positive vs. Negative) and for the 8 runs of scanning, yielding 64 β -weights. It is important to note that, for machine learning classification, there is a trade-off between having many noisy examples (e.g. one β -value per trial) and having fewer but cleaner examples (e.g. one β for each task condition per run—obtained by averaging examples of the same class to subdue their variability; for discussion on this issue, see Pereira et al. 2009; Haynes 2015). Because our investigation concerned the very subtle representational differences between high-order concepts (e.g. one's sense of self at present vs. in the past), the β -estimate of a single trial (or even the β -value computed for a single block) might introduce stimuli-specific noise and fail to reflect the most critical facet. Therefore, in order to denoise, we prioritized having somewhat fewer but robust examples over many but flimsy examples by estimating a β -value per regressor per run. However, it is important to emphasize that, due to the factorial design that we employed, we still obtained sufficient examples for the training and testing—for instance, in the 4-way cross-classification wherein we decoded the 4 persons, trained using positive trials, and tested using negative ones (and vice versa), 2 separate sets of 32 examples were used for training and testing. Our dataset is sufficiently large for training and validation by the conventional practice of decoding research (Pereira et al. 2009). We performed MVPA with the Decoding Toolbox (TDT; Hebart et al. 2015), which employed a linear support vector machine (SVM) for classification with a "C = 1" cost parameter (Chang and Lin 2011). Using standard approaches of cross-validation and cross-classification, we carried out 7 analyses to decode different aspects of "selfness vs. otherness." For each classification, a leave-one-run-out 8-fold splitter was used whereby the algorithm learnt to distinguish between relevant categories using the data from 7 of 8 runs; its ability to correctly predict was tested using the unseen data from the remaining "held-out" run. This procedure was iterated over all possible combinations of runs used for training and testing. By partitioning the datasets based on different runs of scanning, we ensured that there was no contamination

of information leaking from the training sets to testing sets. The accuracy scores were then averaged across folds to produce a mean accuracy score for further statistical analysis. This was done separately for each participant, each ROI, and each of 7 decoding analyses. The analyses included: (1) an aggregate sense of “self” (incorporating *Present Self* and *Past Self*) versus an aggregate sense of “others” (incorporating *Mother* and *Queen*); (2) finer-grained differentiation within the self-concepts—*Present Self* versus *Past Self*; (3) finer-grained differentiation within the concepts about others—*Mother* versus *Queen*; (4) 4-way differentiation among the 4 individuals—*Present Self* versus *Past Self* versus *Mother* versus *Queen*; (5) cross-classification of self versus other across near and far social distances—the classifier was trained to tell apart self versus other based on samples with closer social distance (i.e. *Present Self* vs. *Mother*) and tested using samples with farther distance (i.e. *Past Self* vs. *Queen*); this was repeated with the reverse mapping (i.e. using the “distant” pair for the algorithm to learn and the “close” pair to test its generalizability); (6) cross-classification of the abstract sense of social distance across the self and other domain—the classifier was trained to distinguish *Present Self* from *Past Self* and tested using *Mother* versus *Queen* (and vice versa). Successful cross-classification in this case indicated the acquisition of neural patterns that encoded the abstract information about near versus far social relationship, applicable both to the self and other domain; (7) cross-classification of 4 individuals across the 2 sides of emotive valence; the classifier was trained to perform 4-way classification using the dataset wherein personality judgments were based on positive traits, and then was tested whether it could generalize to the unseen data based on negative traits. This cross-classifying was repeated with the inverse mapping (trained on the negative trials, tested on the positive ones). Bonferroni correction was applied to adjust multiple comparisons based on the number of ROIs in each analysis (α -level: $0.05/8 = 0.006$).

As an exploratory analysis, we used the roaming whole-brain “searchlight” method to test whether brain regions outside our selected ROIs also carried task-relevant information that allowed successful multivoxel pattern decoding (Kriegeskorte et al. 2006). Each “searchlight” was composed of a multivoxel pattern in a local neighborhood (a sphere of 10 mm radius) that surrounded each voxel of the brain. For each sphere, a linear SVM classifier was trained to decode relevant information and tested using the same leave-one-run-out procedure as described above; the accuracy score was assigned to the centroid voxel. This process was repeated each time using a different voxel as the center as the searchlight roamed across the brain. Prior to group-level analysis on the decoding result, each participant’s classification accuracy map was normalized into the MNI space using the DARTEL toolbox (see the Preprocessing section) and smoothed using a Gaussian kernel of 8-mm FWHM. Voxel-wise decoding accuracy was tested against the chance level (50% for binary classifications; 25% for

4-way classifications) using a one-sample t-test. Multiple comparisons were constrained using SPM’s standard procedure (based on the random-field theory), keeping familywise error below $P < 0.05$ for each voxel.

As a complementary approach to machine learning decoding, representational similarity analysis (RSA) was employed to investigate whether there was any systematic structure that underlay the neural representations as revealed by the pair-wise resemblance between different task conditions (e.g. whether the multivoxel patterns of *Present Self* are more similar to *Past Self* than to *Queen*). Based on the established methods (Kriegeskorte et al. 2008; Nili et al. 2014), we calculated the neural similarity between each pair of experimental conditions as the Pearson correlation of their vectorized patterns of voxel-wise activity. Note that we used similarity (Pearson’s r between patterns) as the metric, rather than their distance/dissimilarity ($1 - \text{Pearson’s } r$), to characterize the resemblance between contexts, for the sake of more intuitive and straightforward interpretations (the 2 approaches generated exactly the same conclusion as by mathematical definition they were 2 sides of the same coin; for discussion, see Dimsdale-Zucker and Ranganath 2018). Hierarchical clustering analysis was performed on the distance measures ($1 - r$ or $1 - \tau_a / \rho$) to visualize the categorical grouping of neural representations, embodied in the structure of a dendrogram tree. As per the standard approach of RSA research (e.g. Nili et al. 2014), we used the rank-based correlation indices—Kendall’s Tau (τ_a) and Spearman’s Rho (ρ)—to assess the second-order relationship between representational similarity matrices, and considered only the off-diagonal lower-triangular elements of each matrix to prevent inflation of correlation size. Given the nonparametric nature of τ_a and ρ , signed-rank test was used to assess whether the correlation between 2 representational similarity matrices was significantly greater than chance (one-tailed), as well as whether the correlation for a certain pair of matrices significantly differed from those of other matrices (2-tailed). Multiple comparisons were adjusted based on the number of ROIs in each analysis (Bonferroni correction). Finally, to evaluate the magnitude of correlation with reference to a hypothetically “true” model’s optimal performance (given the amount of noise in the present data), the noise-ceiling was derived by (i) calculating the averaged correlation between the group’s matrix and every individual matrix (the upper bound) and (ii) performing an iterative leave-one-subject-out procedure that correlated each individual’s matrix with all other participants’ averaged matrix (the lower-bound; cf. Nili et al. 2014).

Results

Machine learning classification analysis

We sought to clarify (i) whether there is a spectrum-like change in neural coding in relation to social distance and (ii) whether the 2 subnetworks were equally capable of representing this “psychological continuity.”

To answer the 2 questions, we deciphered the multivoxel patterns of DN and SN to know how the continuity and boundary of “self vs. other” representations were encoded. In [Supplemental Results 1](#), we report behavioral data and a confirmatory analysis of the motor cortex that refutes performance-related explanations. Using the subject-specific peak coordinates derived from the localizer data, for each individual, we defined 8 spherical ROIs within the DN and SN—5 ROIs are typically affiliated with the DN (the dmPFC, the vmPFC, the PCC, the left IPL, the right IPL), while 3 are associated with the SN (the left ATL, the right ATL, the left IFG). Here we used a well-established localizer paradigm ([Dodell-Feder et al. 2011](#)) to identify ROIs, pinpointing the peaks in the DN and SN that had greater activation for the Social than Nonsocial condition. To ascertain the robustness of our findings, we conducted 3 additional analyses^{1,2,3}. SVMs were trained on the neural pattern of each ROI. Using a supervised-learning cross-validation procedure, the algorithms were tested using “quarantined” data to evaluate whether it was able to predict the mental content and whether the chance of successful prediction varied systematically with network membership.

We began by verifying whether the regional multivoxel pattern of each ROI allowed deciphering the broad-stroke information about the aggregate sense of selfness (*Present Self* and *Past Self*) versus the aggregate sense of otherness (*Mother* and *Queen*). Significantly above-chance decoding was achieved in every ROI ([Fig. 1A](#)), suggesting that this coarse-grained, binary distinction between self and other is discernible using the patterns of local DN/SN activity. Motivated by this result, we tested whether it would be possible to decode nuance within the orbit of self-referential ideation (*Present Self* vs. *Past Self*) and other-referential ideation (*Mother* vs. *Queen*). As [Fig. 1A](#) shows, statistically reliable decoding was achieved in nearly all ROIs (all, except for the left IFG) for the differentiation between *Present Self* and *Past Self* and in every ROI for the differentiation between *Mother* and *Queen*. This indicates that the regional pattern of DN and SN activity enabled accurate classification not only “between” the domains of self versus other but also the finer-grained subgroups “within” the realm of self and other. Successful decoding of a socially proximal

(*Present Self*) versus a distant concept (*Past Self*) within the “self” domain suggests that information about “social distance” intersected with “self vs. other.” To untangle this intricacy, we performed 4-way classification among the 4 individuals to understand how decoding fared when the algorithm had to consider “self vs. other” and “social distance” simultaneously. As shown in [Fig. 1A](#), while the chance-level dropped from 50% (binary) to 25% (4-way), reliable above-chance decoding was achieved in every ROI (also see [Supplemental Results 6](#) for the confusion matrices), indicating that the neural pattern that encoded each individual’s specific identity can be reliably differentiated from all other identities. Critically, a systematic trend was replicated across all of the 4 analyses: Although significantly above-chance decoding was obtained in every ROI, classification accuracy was reliably higher in the 5 ROIs that belong to the DN (with the PCC reliably reaching the highest accuracy), compared to the 3 ROIs that belong to the SN. This “imbalance” between the 2 subnetworks implies that the DN, as a whole, carried greater amount of information about personal identities that could be extracted by the classifier to enable correct predictions. We also examined whether the univariate amplitude of each ROI could be used to differentiate individuals. As shown in [Fig. 1B](#), results revealed that univariate contrasts did not reliably differ between conditions, suggesting that encoding personal identity relied on multivariate patterns rather than differential amplitudes at the univariate-level (also see [Supplemental Results 5](#) for the data of whole-brain interrogation for univariate effects).

Although identical stimuli were used in each condition (i.e. the same personality descriptions were assessed with reference with different individuals), participants saw different cue words that reminded them of the current target person. Thus, an alternative explanation is that our decoding was driven by word length—e.g. 1 word (*Mother*) versus 2 words (*Present Self*). This alternative is unlikely given the fact that robust decoding was still achieved when the lengths of cue words were matched between conditions—see [Fig. 2A](#) for the whole-brain searchlight interrogation of *Present Self* versus *Past Self* and [Fig. 2B](#) for *Mother* versus *Queen*. To further rule out the potential effects of low-level visual features (i.e. vertical/horizontal lines and curves that constituted the cue words), we trained the classifier to tell apart *Present Self* versus *Past Self* and tested whether it could distinguish *Mother* versus *Queen* (and vice versa). This cross-classification was a stringent test to assay whether the algorithm truly deciphered high-level information (in this case, social distance that was generalizable across self and other), independent of sensory factors (there was minimal visual resemblance between the pair of “*Present Self* vs. *Past Self*” and “*Mother* vs. *Queen*”). Results showed that cross-classification was detected in the default network—models trained to discriminate *Mother* from *Queen* could also reliably discern *Present Self* from *Past Self* (and vice versa). This suggests that the

¹ First, we used the coordinates from activation-likelihood estimation (ALE) meta-analyses on the literature of semantic processing and social cognition to relocalize ROIs and performed multivoxel decoding at these literature-defined locations ([Supplemental Results 2](#)). The decoding results based on ALE are highly consistent with those based on the Dodell-Feder localizer, demonstrating the robustness of our findings that they are not reliant on a particular localizer paradigm.

² Second, we used NeuroSynth to identify the voxels robustly engaged by semantic tasks (based on 40,030 activations from 1,031 fMRI studies; [Supplemental Results 3](#)). We found that all of the 3 semantic ROIs defined by our localizer contrast are situated within the semantic clusters of NeuroSynth, suggesting concordance between our semantic ROIs and the sites found in the fMRI literature of semantics.

³ Third, we used NeuroSynth to identify the voxels robustly engaged by self-related processing (based on 4,728 activations from 166 neuroimaging studies; [Supplemental Results 4](#)). We found that all of the 5 default-mode ROIs defined by our localizer contrast are situated within the self-related clusters of NeuroSynth, suggesting agreement between our default-mode ROIs and the sites found in the fMRI literature of self-concept.

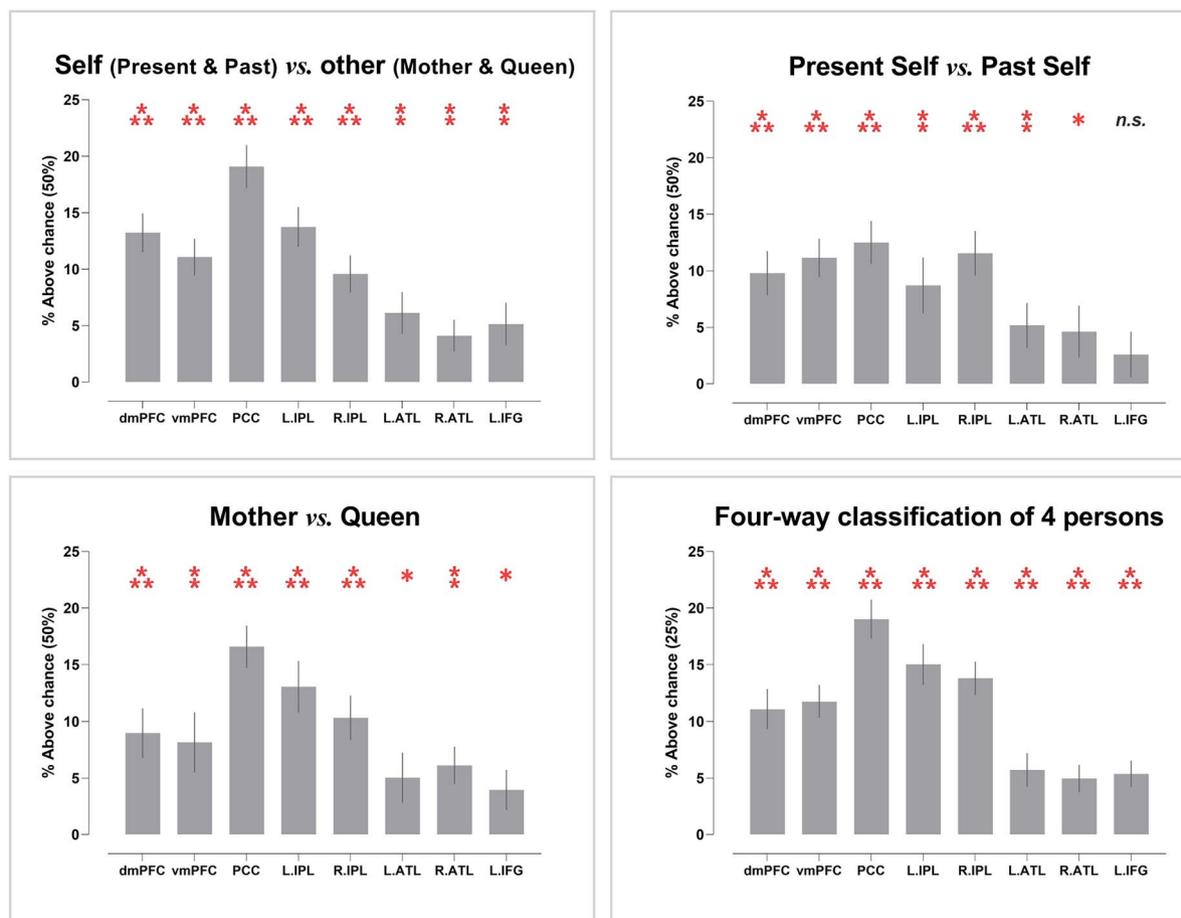
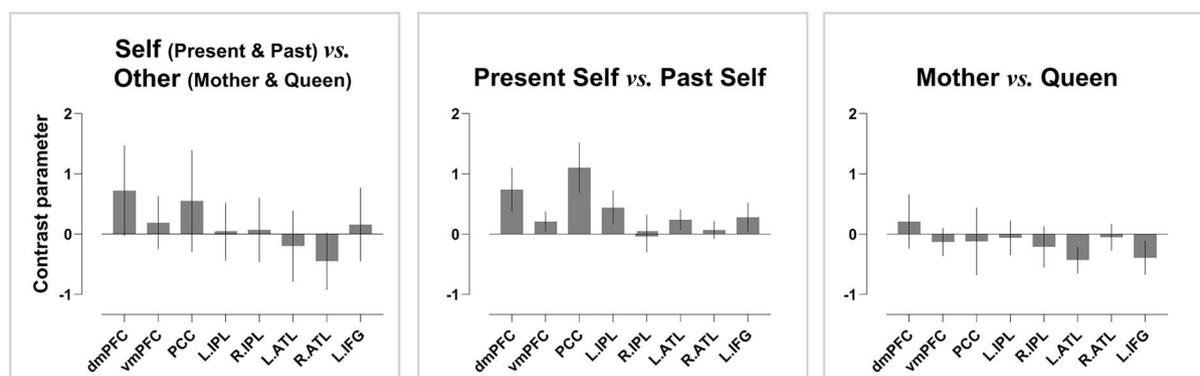
(A) Multivariate decoding at DN & SN regions**(B) Univariate contrasts at DN & SN regions**

Fig. 1. A) Average accuracy of 4 classification analyses in the 8 ROIs; asterisk indicate significantly above-chance decoding (Bonferroni-corrected for multiple comparisons). B) Univariate contrast parameters in the same 8 ROIs. * $P < 0.01$ (marginal); ** $P < 0.005$; *** $P < 0.001$.

multivariate classifiers were, at least in part, overlooking sensory differences of letters and discovering neural patterns that enciphered social distance. As illustrated in Fig. 2C, successful cross-classification was achieved using the patterns of 3 DN regions in the midline structures (with the vmPFC being the most informative area), while decoding was at chance in 2 control regions (the primary visual and motor cortices). This

rigorous verification suggests that DN regions carried information about personal identity and social proximity and abstracted such information away from the sheer appearance of visual stimuli.

While the standard cross-validation approach unraveled whether the information of self versus other existed in the multivoxel pattern of an ROI (e.g. whether the classifier succeeded in predicting if one was reflecting

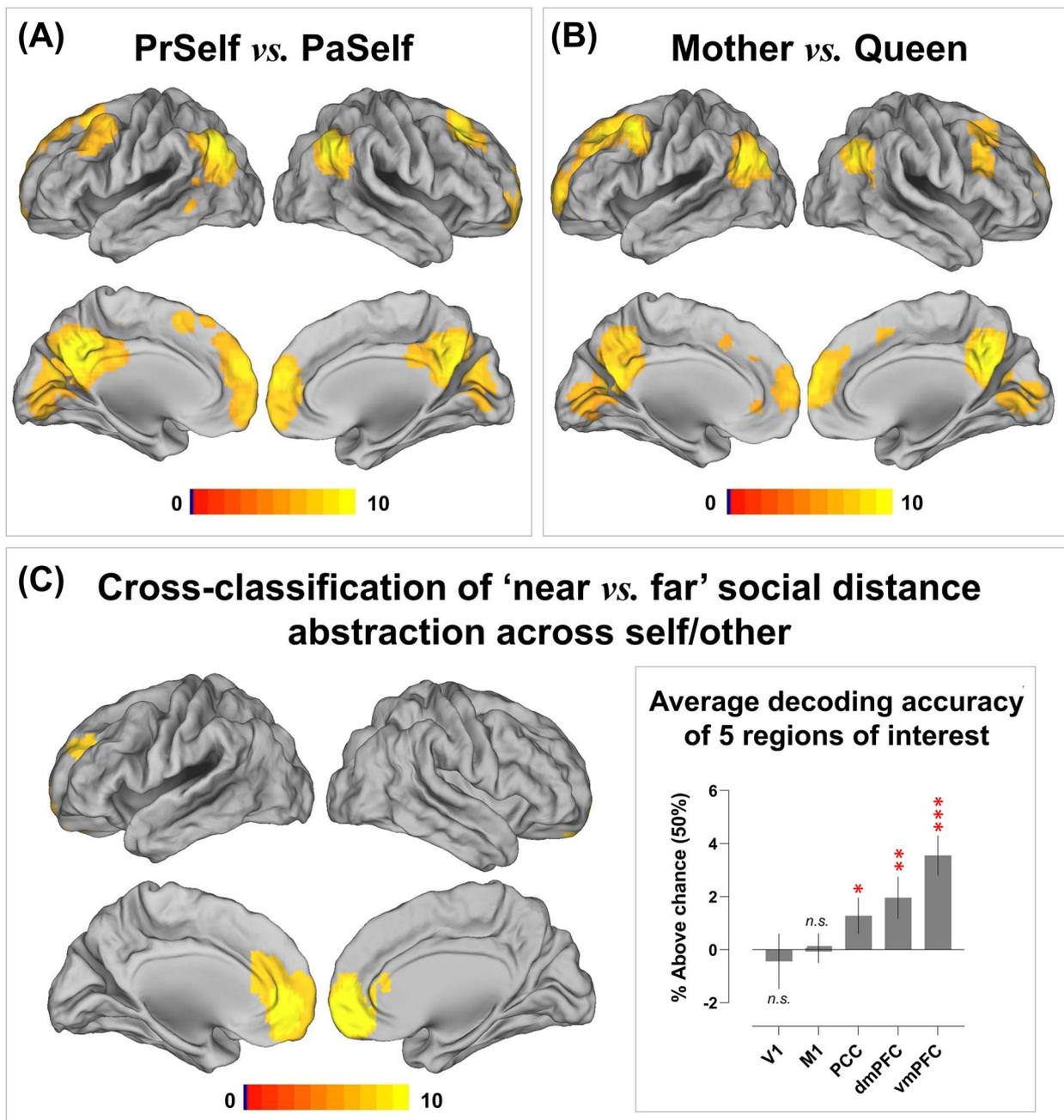


Fig. 2. Whole-brain searchlight interrogation results for (A) Present Self versus Past Self, (B) Mother versus Queen, and (C) Cross-classification of “near vs. far” social distance across the self and other domain. Shown in the inset box is the extracted decoding accuracy from 5 target anatomically defined ROIs (2 control areas: V1/the primary visual cortex, M1/the primary motor cortex; 3 midline structures of the default-mode system: the PCC, dmPFC, vmPFC; these ROIs were selected *a priori* and defined using the anatomical masks of the Wake Forest PickAtlas Toolbox). Correction for multiple comparisons was via constraining voxel-wise FWE under $P < 0.05$.

on *Mother* or *Queen* based on distributed activities), it lacked the ability to test whether the brain reinstates a reproducible and generalizable neural code across varying situations. A generalizable code indicates *abstraction* despite contextual variation. A hypothetical example that epitomizes this generalizability and invariance would be that the brain summons reproducible neural patterns to represent one’s identity no matter whether that person is seen, heard, or recalled. To overcome the limitation of standard cross-validation, we conducted cross-classification—the classifier was trained on data

from one cognitive context and tested on another; with this procedure, we tested whether there was any commonality in neural codes invariant to contextual changes. Two analyses were performed: First, we tested whether there is generalizable neural coding of self versus other, irrespective of social distance. This was achieved by training the algorithm to classify *Present Self* from *Mother* (closer distance) and testing whether the codes were transferrable to distinguish *Past Self* from *Queen* (farther distance), and vice versa. Second, we tested whether generalizable patterns of neural

activities were used to represent the 4 individuals, irrespective of the emotive valence of descriptions that were used to probe person-related concepts. By virtue of our factorial design that crossed 4 identities with valence, this cross-classification was achieved by training the algorithm to classify the 4 persons using data from the positive-valence context and testing the generalizability using the negative-valence data (and vice versa). Cross-classification provides a rigorous test to assay whether invariant neural codes were reliably summoned to represent the mental concepts about particular individuals, unaltered by social distance or emotive valence.

As illustrated in Fig. 3A, the analysis showed that, based on the neural patterns of 4 regions of the DN (the dmPFC, vmPFC, PCC, and left IPL), the classifier was able to extrapolate information (e.g. neural codes elicited by the pair of near distance—*Present Self* vs. *Mother*) from one context and successfully applied it to cross-classify in another context with a different degree of social distance (e.g. far: *Past Self* vs. *Queen*). This suggests context-invariant patterns that encode the essence of “self vs. other,” generalizable across close and distant social relationships. The representational content of these regions is instantiated in the 3 inset boxes of Fig. 2A. Here we saw a consistent pattern across the 3 ROIs with highest cross-classification accuracies—when the ground truth was “self,” the classifier was more inclined to predict “self” than “other,” and vice versa when the ground truth was “other.” A coherent but more striking pattern was found in the 4-way cross-classification which generalized across positive versus negative emotive valence. As illustrated in Fig. 3B, based on the regional activity of an ROI, the algorithm was capable of extrapolating personal identities from one context and leveraging the codes to make predictions in another context of different emotive valence. Significantly above-chance cross-classification was achieved in every ROI of the DN and SN. The representational contents of such cross-valence neural coding were exemplified by the pattern of PCC activity: As illustrated in the inset box of Fig. 3B, the percentage of the classifier’s prediction clearly followed a sequence-like pattern (particularly conspicuous in the polar “extreme” cases of *Present Self* and *Queen*). For instance, while the classifier correctly predicted *Present Self* most frequently when the truth was indeed *Present Self*, erroneous response was affected by social distance in a sequential way—*Present Self* was most confused with *Past Self*, followed by *Mother*, and least confused with *Queen*. A similar pattern (with the opposite order) was observed for *Queen*. Together, these indicated that the underlying neural representations that were employed to represent the distinction among personal identities were reproducible across contexts, invariant to changes of social distance and emotive valence, and were structured in a gradational manner that followed interpersonal distance.

By visual inspection on cross-validation and cross-classification analyses, we noticed a disparity that decoding accuracy was obviously better in DN regions than in SN regions. This implies that DN regions might contain more person-related information than SN regions. To further investigate this gap in predictive power between the 2 subnetworks, we employed the combinatorial-ROI decoding analysis (for precedents, see Clithero et al. 2009 ; Smith et al. 2013 ; Wang et al. 2017). This analysis entailed a stepwise procedure wherein decoding was conducted on a “combined” ROI—each combined ROI incorporated the pattern of an original ROI “plus” another ROI, followed by decoding analysis on the combined pattern; this was iterated for every pairwise combination of the 8 ROIs. The outcome of joint-ROI decoding was subsequently compared with the “baseline” where the decoding was conducted using the pattern from the original ROI alone. This procedure revealed whether adding a given ROI improved or impaired the classifier’s performance through serially coupling this ROI with all other ROIs and assessing how the joint-decoding altered relative to baseline. As has been demonstrated by previous research (Clithero et al. 2009; Smith et al. 2013; Wang et al. 2017), this method bypassed the difficulty of directly comparing between ROIs on the amount of information they carried (i.e. in the serial joint-decoding, the number of voxels between combined patterns were exactly matched; this way, we were able to quantitatively assess if an ROI was a reliable “contributor” or “beneficiary” when it was coupled with another ROI). We found that (i) adding any of the 5 DN regions to the combinatorial decoding robustly boosted accuracy whereas adding any of the 3 SN regions had little impact on decoding; (ii) the 3 SN regions reliably benefited more from the addition of another ROI whereas the DN regions benefited less. Shown in Fig. 4 are the example results of joint decoding for “Self vs. Other” (Fig. 4A), “Present Self vs. Past Self” (Fig. 4B), and “Mother vs. Queen” (Fig. 4C). A reliable pattern was seen in these results (as well as in other joint decoding results). Dovetailing our earlier data, the joint decoding revealed a robust “imbalance” between the amount of information carried by the 2 subnetworks. Whereas DN regions carried more information about personal identities (making them “givers” that reinforced classification accuracy), SN regions carried less information (making them reliable “takers” in the combinatorial-ROI decoding).

Representational similarity analysis

Results of the classification analysis gave testable hypotheses about the configuration of neural coding: The brain might encode personal identities as continuous progression along a single-dimension spectrum of social distance, from close to distant, without any categorical cut-off. Alternatively, the neural codes could be configured with a boundary that categorically

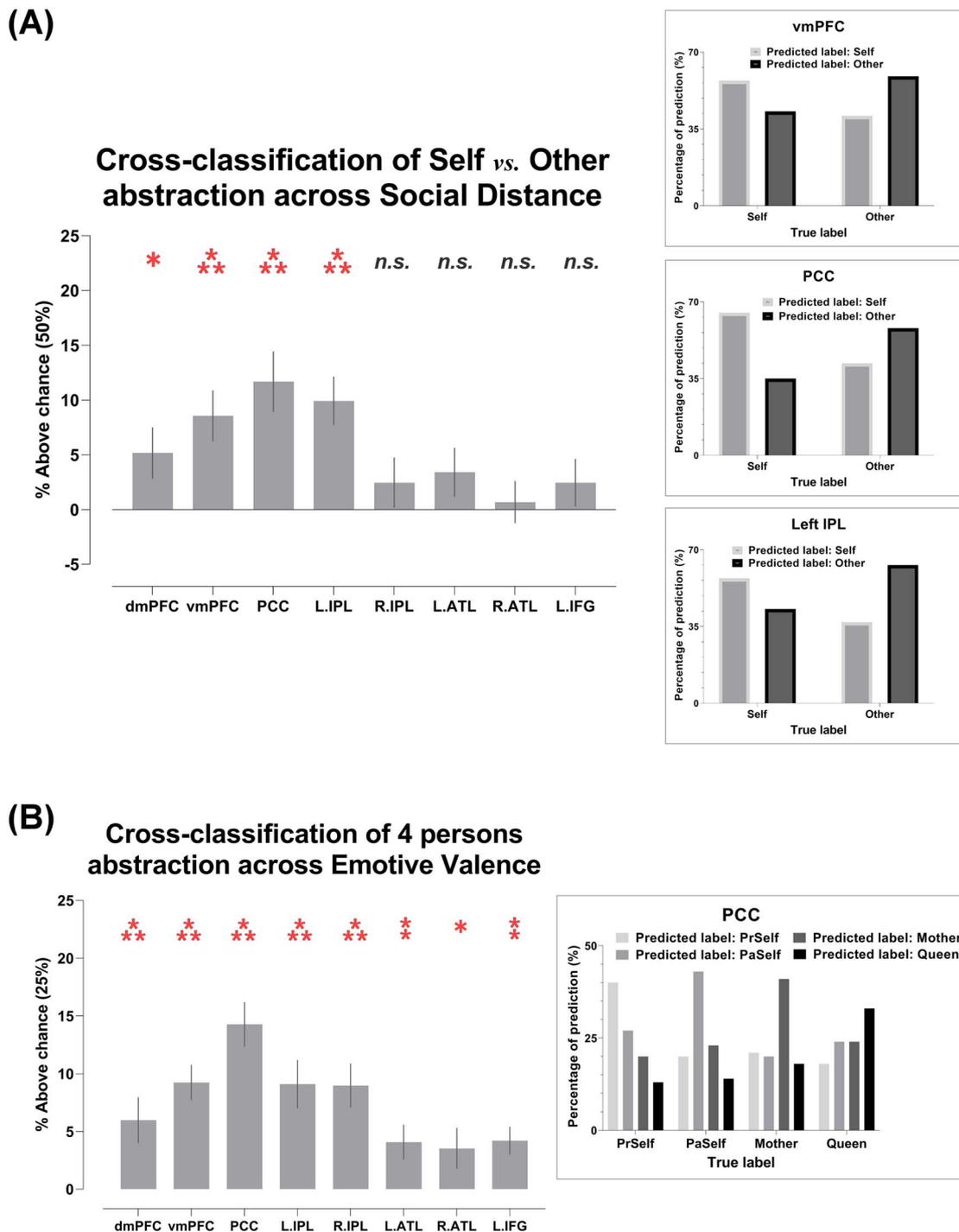


Fig. 3. A) Cross-classification of self versus other between near and far interpersonal distances. B) Cross-classification of 4 individual identities between positive- and negative-valence stimuli. Inset boxes show the confusion pattern of ROIs with highest cross-classification accuracy. Bonferroni-correction was applied to constrain multiple comparisons. * $P < 0.01$ (marginal); ** $P < 0.005$; *** $P < 0.001$.

separates selfness from otherness, with finer-grained differentiation nested within the “self” or “other” domain. The confusion matrices of classification analysis (see [Supplemental Results 6](#)) only provided an equivocal answer—although the patterns of confusion showed that within-domain confusion was more frequent than between-domain (e.g. *Queen* was more confused with *Mother* than with the 2 variants of self, which

implied a bipartite structure), such analyses did not directly quantify the extent of the similarity between 2 representations. Unlike classification-style analysis that by nature discretizes different categories via imposing a decision boundary, RSA quantifies the extent of similarity using continuous measures (e.g. correlation coefficient). Therefore, we exploited RSA to further investigate the representational “geometry” that

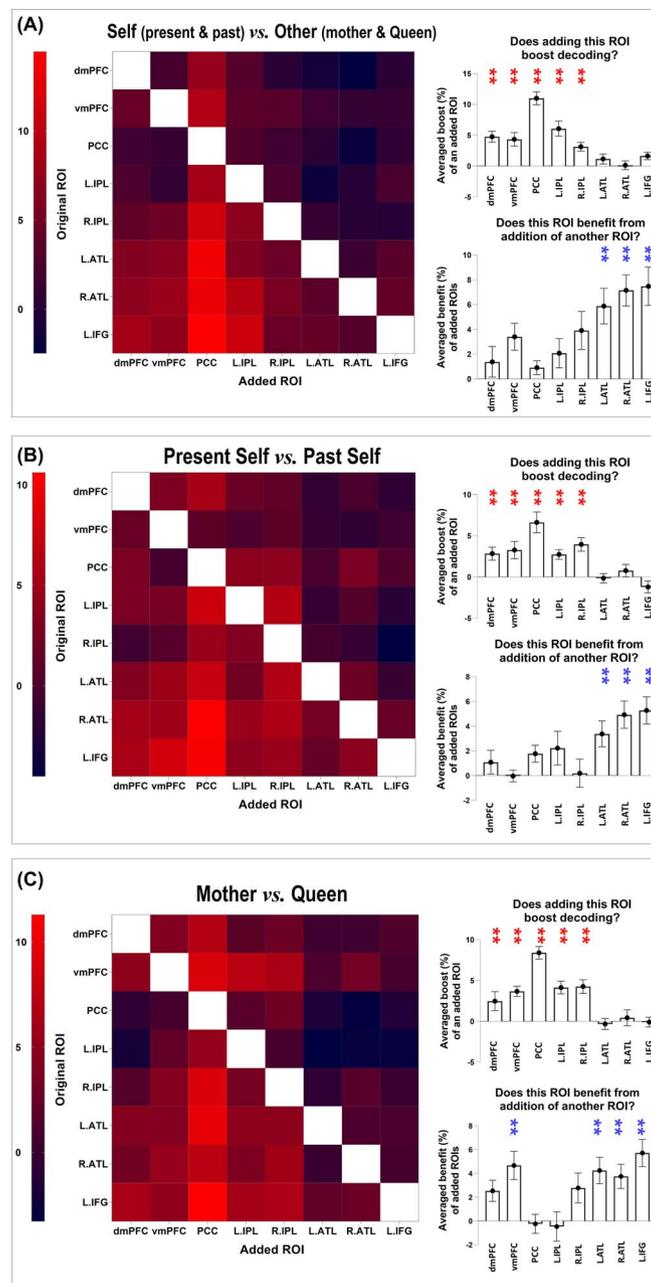


Fig. 4. Combinatorial decoding was based on an iteratively procedure that examined how decoding performance varied when the neural pattern of an original ROI was joined by the pattern of another ROI. The values on the color scale indicates changes in decoding accuracy when an ROI was included, compared to the decoding result based on the original ROI alone, from dark purple (below zero, indicating “decrement”) to bright red (indicating most “improvement”). A reliable pattern was found across all analyses—as shown by the 3 examples here: A) Present Self + Past Self versus Mother + Queen; B) Present Self versus Past Self; C) Mother versus Queen. The bar graphs indicate that, across analyses and across original ROIs, the addition of an ROI that belongs to the default network led to significant improvement of decoding results, whereas ROIs of the SN reliably benefited most from the addition of another ROI.

underlies the similarity between different categories and different brain regions.

We first verified whether there was good concordance between the neural representations of positive- and negative-valence conditions, given the fact that the SVM algorithm successfully cross-classified individual identities across valence. As illustrated in Fig. 5, in every ROI, a statistically robust correlation was found between the representational matrices of positive- and negative-valence contexts, suggesting a coherent motif that the classifier could extrapolate from one situation

and apply to another. Moreover, the outcomes of RSA and pairwise classification concurred with one another. An example result from one ROI is illustrated in Fig. 6: When the representational similarity was high between 2 conditions (which indicated a higher degree of overlap in neural coding), there was a corresponding decline in the accuracy scores of classification (the classifier was more prone to confusion). This resulted in negative correlations, reliably seen in both positive- and negative-valence contexts. Next, hierarchical clustering was used to visualize the representational distance between

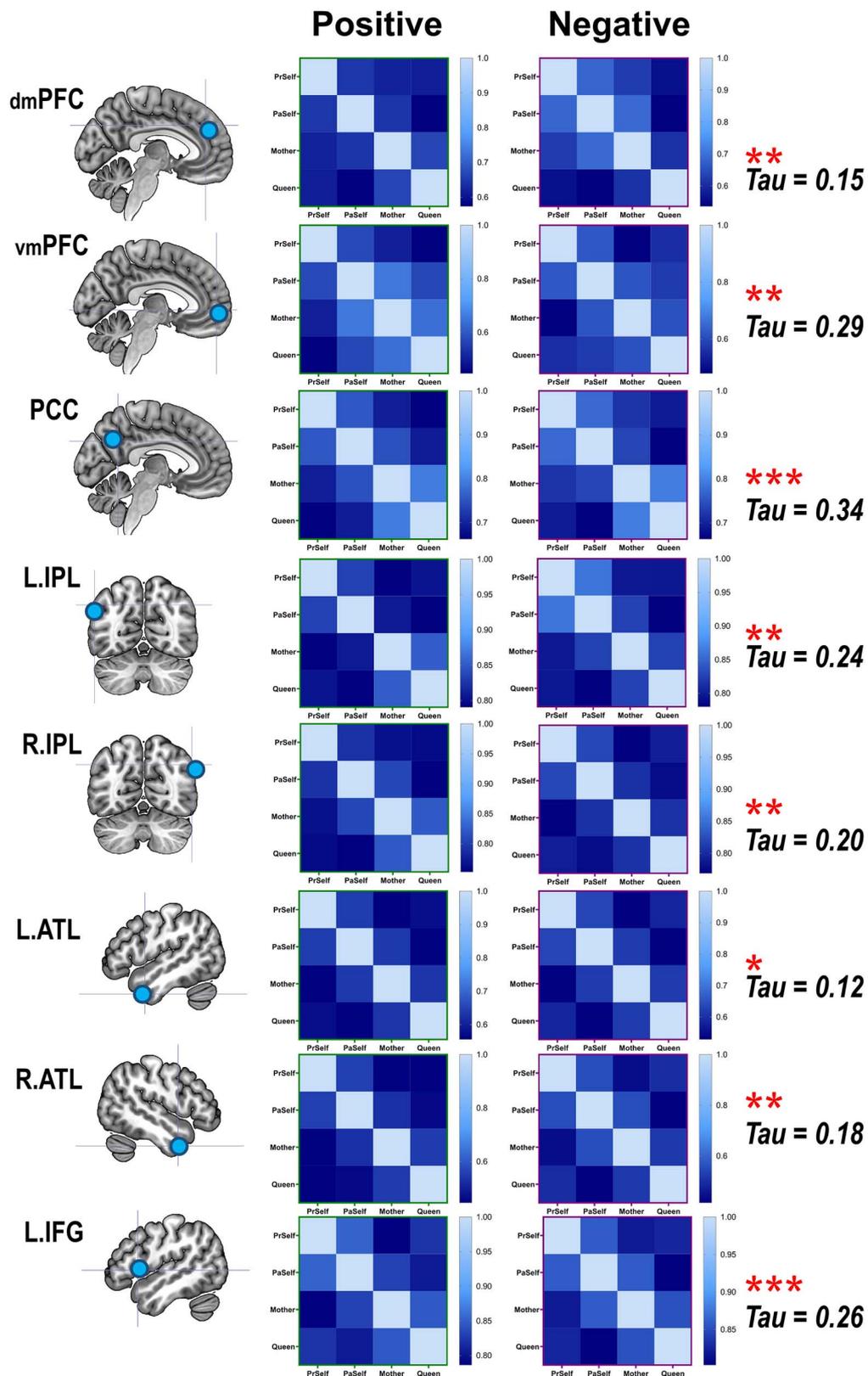


Fig. 5. Corroborating the results of cross-classification between positive- and negative-valence stimuli, the results of representational similarity analysis showed that the neural patterns elicited by the 4 personal identities are significantly correlated between the contexts of positive and negative valence. The correlation was found in every ROI of the DN and SN.

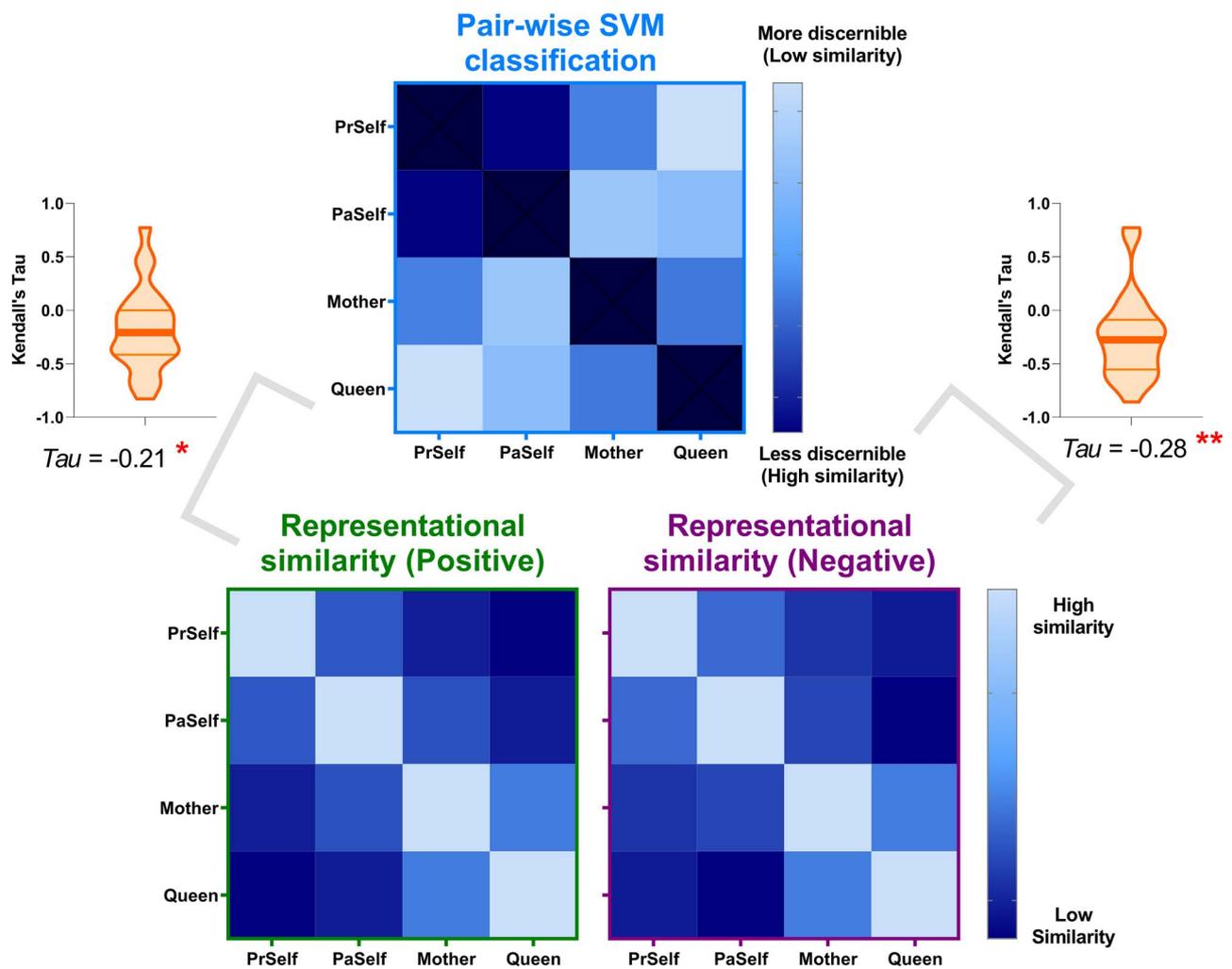


Fig. 6. The results of representational similarity and machine learning classification analyses dovetailed with each other. Illustrated here are example results based on the neural patterns of the PCC (coherent results were observed for all other regions)—when 2 experimental conditions were representationally more similar to each other (hence more confusable), the classifier was less able to predict the correct category label, resulting in significant negative correlations found in both positive- and negative-valence conditions.

conditions, quantified by the branching of a dendrogram tree. As Fig. 7 illustrates, in these 4 “core” regions of the DN (in which we saw generally higher classification accuracy), there is a clearly bipartite structure in which neural representations were stratified into 2 branches by “selfness” versus “otherness.” One branch comprised *Present Self* and *Past Self*, while the other branch comprised *Mother* and *Queen*. Interestingly, in the vmPFC, the representations of *Present Self* (encompassing both positive and negative) formed a distinct cluster that was separable from the remainder (even detached from *Past Self*). This is consistent with the literature that the vmPFC has a unique role in representing the “essence” of self-concept. While the clustering was dominated by the difference of self versus other, emotive valence had minimal impact on how the neural codes of different persons were arranged. Moreover, whereas a clear bipartite split between self and other was found in the representations of DN regions, the demarcating dimension was not as clear-cut in the SN

regions, which offers explanation as to why classification accuracy was generally lower in the SN regions. Taken together, these results emphasize the convergence of evidence that we obtained from the RSA/correlation and SVM/classification approaches.

Replicated across multiple classification analyses, decoding accuracy was found to be robustly higher in DN regions compared to SN regions. Furthermore, combinatorial-ROI decoding also showed that the addition of DN regions reliably boosted accuracy whereas adding SN regions caused no improvement, suggesting an asymmetry of information quantity between the networks. To investigate this issue further, we used RSA to examine whether this asymmetry translates into representational distance between the 2 networks. Hierarchical clustering showed that the principal factor that characterized the heterogeneity between neural representations was the separation between DN and SN. As the configuration of the dendrogram illustrates (Fig. 8A), the initial bifurcation between regions was

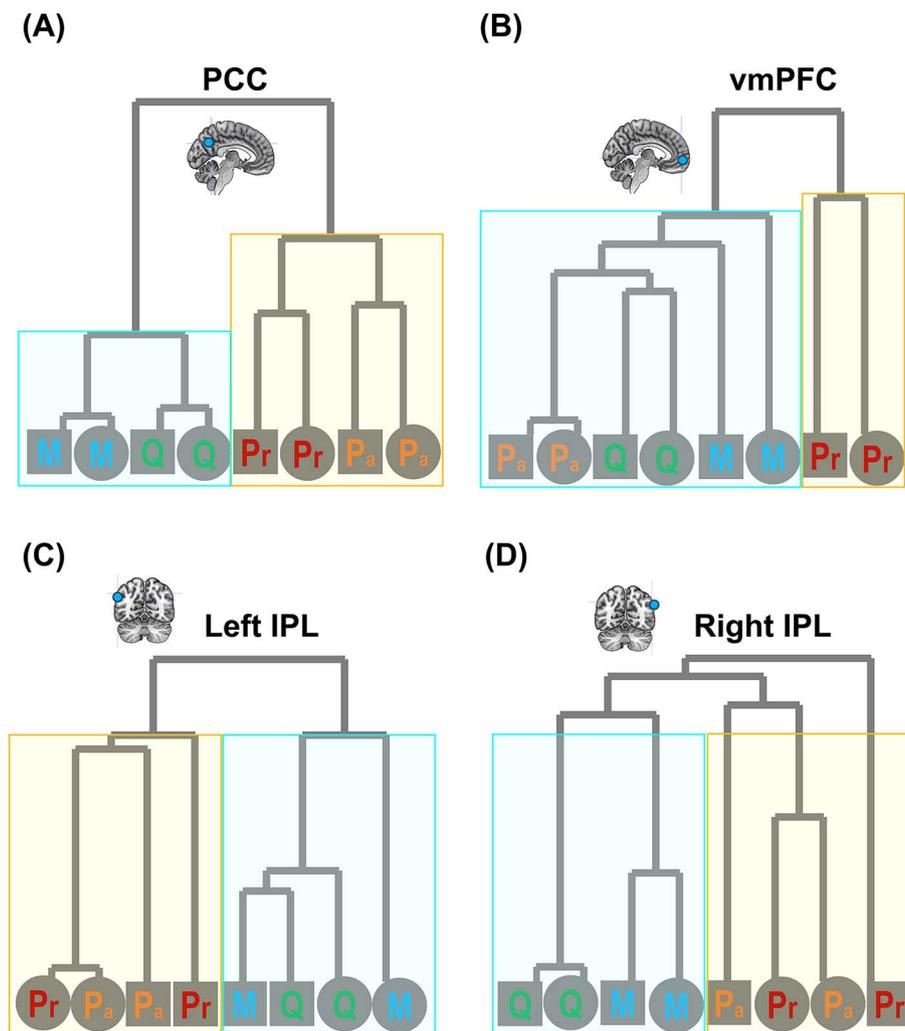


Fig. 7. The configurations of each dendrogram tree here show a clear separation between an aggregate class of *selfness* (*Present Self* and *Past Self*; color-coded using red and orange, respectively) and aggregate sense of otherness (*Mother* and *Queen*; coded using blue and green, respectively). Corroborating the cross-valence classification, the structure of representational dissimilarity ($1-r$) showed that emotive valence had little impact on clustering arrangement (which indicates the extent of neural similarity). Social distance (coded using different colors) indeed impacted on the clustering, but it was couched within the “self/other” separation. Emotive valence is coded by geometric shapes: circle—positive, square—negative. This reliable broad-brush segregation between self and other is found in core DN regions of (A) PCC, (C) left IPL, and (D) right IPL. The vmPFC—in Panel (B)—is somewhat different that it treats *Present Self* as a unique category different from all other categories.

driven by the network membership an area belongs to. Within the cluster of DN, there were 2 subclusters—the 2 medial prefrontal regions (the dmPFC and vmPFC) formed a subgroup separable from the posterior regions (medial-parietal: the PCC; lateral-parietal: the left/right IPL). Within the cluster of SN, the clustering was consistent with contemporary theories and findings of semantic cognition (e.g. Lambon Ralph et al. 2017)—the areas representing semantic meaning *per se* (the left/right ATL) formed a subcluster separable from the area that controls the retrieval of semantic meaning (the IFG). This structure is also evident in the similarity matrices whereby we color-coded the strength of correlation/similarity (Fig. 8B): Close inspection of the layout of correlation matrix revealed that the 3 SN regions congregated to form a cluster that was separable from all of the remaining DN regions; within the 5 DN

regions, the dmPFC and vmPFC formed a subcluster while other posterior regions formed another subcluster. These 2 visualization methods provide complementary yet consistent insight into the separation of the 2 networks. Finally, we investigated whether the size of correlation between regions was modulated by whether it was correlating 2 areas within the same network or between the 2 networks, and whether positive-/negative-valence contexts made a difference. Results showed a significant effect of network membership (Fig. 8C)—correlation size was significantly greater within the same network than between networks (Kendall’s τ_{-a} : $F_{(1,23)} = 6.36$, $P = 0.01$, $\eta_p^2 = 0.22$), indicating greater representational similarity between areas of the same clan. By contrast, valence had no effect (indicating a coherent representational structure across valence), nor did the valence \times network interaction (both $P_s > 0.55$). To

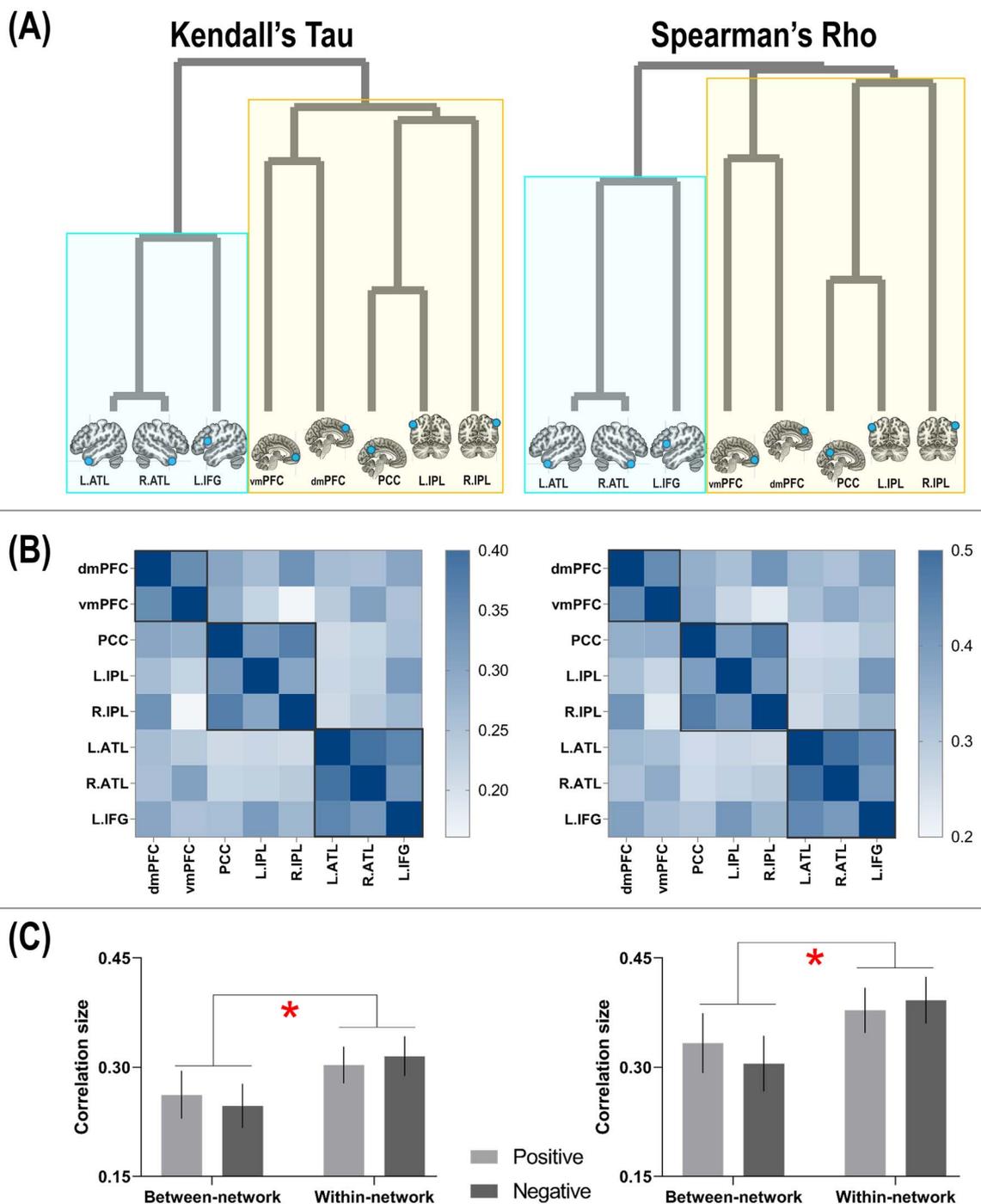


Fig. 8. A) The configurations of dendrogram show that representational distance ($1 - \tau_a$ or ρ) dissociated between the SN and DN; this was seen both using Kendall's Tau-a and Spearman's Rho. B) Here representational similarity is delineated using correlation matrix. A consistent pattern with the dendrogram structure is highlighted using black boxes. C) Within-network representational similarity is significantly higher than between-network similarity, reliably found in both the contexts of positive and negative valence.

ascertain robustness, we repeated all these analyses using a different type of rank-correlation measure (Spearman's ρ) and obtained entirely consistent results (see Fig. 8). Together, these results complement our observations of the classification-based analysis and further highlight the subtlety of subdivision between networks (and within a network).

Finally, we examined 3 theoretical models, testing their explanatory power to account for the neural data (see Fig. 9). The first model was derived from each participant's subjective rating of the pairwise similarity between the personality traits of 2 individuals using a continuous scale (most dissimilar: 0—most similar: 100). The second model was a hypothetical model

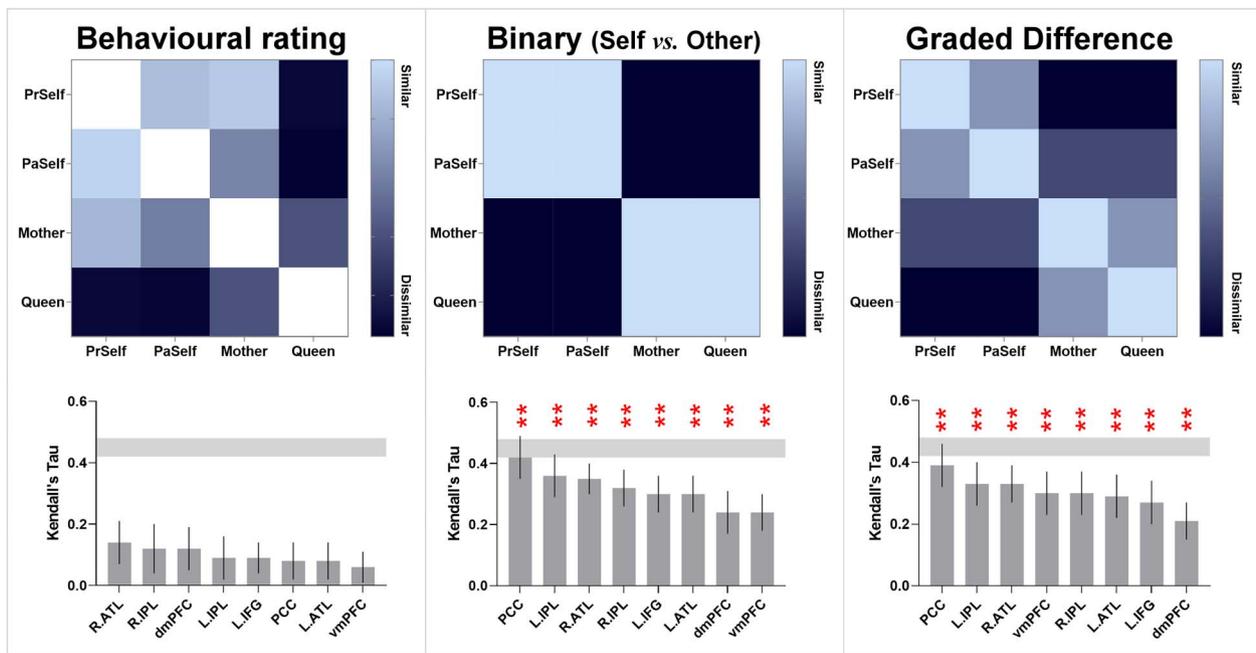


Fig. 9. The hypothetical models that assumed a binary difference between self and other (middle) and gradational differences of the 4 individuals (right) are significantly correlated with the neural pattern of every region. $**P < 0.005$. Multiple comparisons are controlled by Bonferroni correction. The noise ceiling is averaged across all regions. By contrast with the other 2 models, the outcomes of behavioral rating (left) were not correlated with the neural pattern of any region.

based on a binary distinction between self and other—*Present Self* and *Past Self* were collapsed under the “self” umbrella while *Mother* and *Queen* were under the “other” umbrella. The third model was a hypothetical model based on gradual changes, with each of the 4 individuals assumed to be equidistant from one another on the spectrum. We correlated the neural similarity matrix of each brain region with these theoretical models and statistically examined its reliability. We found that the neural pattern of every brain region significantly correlated with the binary model and the graded model (all $P_s < 0.005$). As illustrated in Fig 9, for both of the binary and graded models, the neural pattern of the PCC was most correlated with the theoretical models and approached the lower bound of noise ceiling, which indicates these models’ near-optimal capability to explain neural data. The correlation sizes with the theoretical model did not reliably differ between regions, nor did the comparison between the binary model and graded model (all $P_s > 0.05$ by Bonferroni correction). By contrast, the behavioral rating matrix was not reliably correlated with any region. Critically, both the binary and graded models outstripped the model of behavioral rating in terms of their explanatory power on the neural data (binary vs. behavioral: $P = 0.0008$; graded vs. behavioral: $P = 0.001$). While subjective rating best characterized how a participant perceived and compared the resemblance of 2 characters, such high-dimensional, multi-faceted understanding (quantified as their rating) was not fully reflected in the local neural patterns of DN and SN regions. Instead, while the binary

model was less elaborated (or more impoverished in the dimensions to characterize the difference of personality) than behavioral rating, it was better able to capture the “representational landscape” of neural data. This concurs with earlier results that the broad-brush separation of “self vs. other” explained most variance. By contrast, while socially “close vs. distant” was decodable from the cross-classification between self and other (implying that it could be an orthogonal dimension to self/other), it was a less dominant factor in shaping the landscape of neural representations.

Discussion

In the present study, we used a series of multivoxel decoding to unravel the neural representations of self-versus other-referential concepts, a pivotal psychological construct that permeates all aspects of social life and has been argued as a core function of a specialized system for social cognition. Across multiple analyses, we found a robust dyadic fractionation between the 2 subnetworks within the social system—multivoxel patterns reliably dissociated between regions that are canonically affiliated with the DN and those with the SN, evident both in the outcomes of supervised classification and representational similarity. Moreover, the representational geometry of neural responses to self- and other-referential thoughts mirrored the psychological continuum of interpersonal distance, from an integral sense of *Present Self* to a distant other (*Queen*), with the broad-stroke segregation of “selfness” from “otherness” being the chief factor that

divided the representational space, and social distance⁴ being an auxiliary factor that subdivided concepts within the self/other domain. Below we discuss the implications of the present results.

The fusion and fission of DN and SN

In the literature of social neuroscience, various regions of the DN and SN have been demonstrated to represent diverse mental states (e.g. Tamir et al. 2016) and personality traits (e.g. Thornton and Mitchell 2018), which leads to an agglomeration of the 2 networks as the social system. In 2 separate bodies of literature, however, neurolinguistics and human-connectome research have demonstrated that the DN and SN are 2 distinct entities that unite and disjoin in various resting and task states, despite them both favoring social tasks. Data-driven parcellation methods have demonstrated that DN and SN regions form a coherent network at rest (e.g. Yeo et al. 2011). Seed-based connectivity methods have also shown that key nodes of the SN—the ATL and IFG—are connected with DN nodes at rest (Jackson et al. 2016; Humphreys and Lambon Ralph 2017). Their association during resting state has also been found during task states: DN and SN activities are both enhanced by mnemonic retrieval (e.g. autobiographical memory) and inhibited by externally driven processes (sensory input or motoric output; Chiou et al. 2020). However, the 2 networks have also been found to reliably dissociate: Chiou et al. (2020) demonstrated that, while both networks showed heightened activation for socio-cognitive tasks, DN regions preferred tasks emphasizing the retrieval of episodic details whereas SN regions preferred tasks emphasizing the semantic interpretation of perceptual input. Jackson et al. (2019) and Humphreys et al. (2015) both found that typical semantic tasks drove a dissociation—they activated SN regions but suppressed DN regions. Given the mixed picture, the findings reported in the present study have important implications for clarifying the relationship between SN and DN—we demonstrated that both networks contained information for decoding social categories; however, compared to the SN, the DN possessed more abstract/generalizable information that allowed cross-classifying across distance and valence and exhibited sharper representational partitions that individualized each social category. The disparity between the 2 networks has implications for constraining the interpretation on the activity of SN regions during various socio-cognitive tasks (e.g. Olson et al. 2013; Binney et al. 2016; Wang et al. 2017)—for instance, if some SN regions prefer social knowledge to other semantic contents, there might be more exchange of information between these SN regions and DN regions that could be unraveled using

connectivity analysis, compared to SN regions that prefer nonsocial contents. More broadly, human-connectome research has shown that both the DN and SN are situated at the high-order, transmodal end of macroscale cortical hierarchy (Margulies et al. 2016). In this regard, our finding may pave the way for future research to clarify how their positions in this macroscale architecture affect the information they carry. We elaborate on this in the following section.

The bipartite split in multivoxel patterns between networks

Decades of research have accumulated a wealth of data about how inter-connected brain regions form large-scale networks, as well as how network architecture can be mapped onto cognitive functions (for review, see Uddin et al. 2019). Of particular interest is the functionality of the expansive DN. This widely dispersed group of brain areas was originally identified based on their “deactivation” (relative to task state) during active engagement in a task state and heightened activation during wakeful resting periods (Raichle et al. 2001). Owing to its initial “task-negative” definition, for a long time this system was assumed to play little role in goal-oriented behavior. Later research has identified its contribution in a panoply of goal-directed cognitive tasks that depend on internally constructed representations (e.g. memory, schema, etc.), such as social cognition (e.g. Skerry and Saxe 2015), retrieval of episodic/autobiographical memory (e.g. Spreng and Grady 2010), skilful application of schema to solve a task (Vatansver et al. 2017), and visual memory (as compared to visual perception; see Murphy, Wang, et al. 2019b). Despite mounting evidence, the “misnomer” of default-mode system lingers (which is rooted in its original task-negative definition as the brain’s metabolic default). Recent research has uncovered a bipartite structure within the default system (e.g. Braga and Buckner 2017; Braga et al. 2019; Chiou et al. 2020; DiNicola et al. 2020; also see Andrews-Hanna et al. 2010 for a tripartite fractionation of the DN). The bipartite structure has been discussed under the framework of a macroscale “gradient” spanning the entirety of cerebrum (Margulies et al. 2016; Huntenburg et al. 2017). This gradient reflects a recursive process of information convergence, occurring in multiple places of the brain, from sensory-motoric representations that encode the “here-and-now” of perceptible entities to multisensory representations that are stored in the DN/SN and encode memories and concepts. Core regions of the DN (e.g. the PCC and vmPFC) sit atop this gradient, while core regions of the SN (e.g. the ATL and IFG) are suspected to be situated at tiers just below the apex positions (i.e. the DN’s cores). The functional dissociation between DN and SN might result from their differential positions on the cortical gradient. Our results lend some support to such interpretations—owing to its position at the gradient apex of multisensory convergence, the DN may contain

⁴ It is noteworthy that temporal distance, as an auxiliary factor, was one of a multitude of factors that could determine the partition of social categories. For instance, in the differentiation of Present Self versus Past Self, temporal distance was a factor closely entangled with social distance and could be an alternative/additional reason that configured neural representations.

more information about person identities, which are abstract concepts distilled from multiple modalities and invariant to varying contexts. Thus, DN regions could successfully cross-classify across valence and social distance, unaffected by stimuli appearance, whereas SN regions exhibited much less capability to cross-classify (implying that perceptual changes might have detrimental impact on the SN's ability to generalize across situations). These findings are compatible with the interpretation that the DN contains more abstract information for social cognition, despite the fact that DN and SN (univariate) activities both intensified in response to the demand of social tasks.

Significantly above-chance decoding was found in both networks for coarser- and finer-grained classifications of self versus other, indicating that both of them carried sufficient information that allowed the algorithm to delimit a margin to distinguish between 2/4 classes of data points. Among all ROIs, the PCC contained most information about personal identities. This is observed across multiple analyses—decoding accuracy was highest based on the patterns of PCC; when the PCC was added to the combinatorial decoding, it led to greatest boost to the outcome of decoding; the neural pattern of PCC was most correlated with the psychological continuity of self versus other, leading to correlations that approached the noise ceiling (indicating near optimum). These data are consistent with the view that the PCC serves as one of the cores of the DN (Andrews-Hanna et al. 2010) such that its representational pattern “echoes” the activities elsewhere in the brain through long-range connections (Leech et al. 2012) and its dysfunction causes memory-related ailments (Leech and Sharp 2014). Moreover, it is noteworthy that the vmPFC represents *Present Self* as a distinct entity from all other categories, as evident in how the dendrogram initially bifurcates in Fig. 7B. This result is consistent with the specialized role of the vmPFC in representing the essence of self-concept, adding to a large body of evidence (for review, Wagner et al. 2019). Recently, the “self-in-context” model about vmPFC function has been proposed (Koban et al. 2021)—according to this view, the vmPFC represents “self” in a compressed low-dimensional space that captures relevant features of a context (e.g. human interaction, social norm) to construct “self.”

It is important to note that any individual brain region may participate in multiple brain networks and subserve multiple cognitive functions even though a region has a primary network affiliation or has certain types of functions that it is frequently associated with (for discussion, see Pessoa 2014). The “network membership” of a brain region can be affected by various factors, such as its position in the network (i.e. “connector” regions between 2 networks tend to have more fluid network affiliation, relative to regions within the “heartland” of a module; see Spreng et al. 2010; Spreng et al. 2013), the cognitive state one is under (i.e. contexts can drive a

region to fluidly couple with different networks), as well as the methods used to probe the relationship between regions (for discussion, see Petersen and Sporns 2015). These factors may particularly affect several “connector” regions that bridge between 2 networks, making them exhibit characteristics of both systems (such as the dmPFC that bridges core DN nodes with SN nodes; see Spreng et al. 2013). We speculate that this may offer explanations regarding the fluid functional profile of dmPFC: Previously, using univariate analysis, we found that the dmPFC responded preferentially to semantic tasks over episodic tasks, which makes it more akin to the behavior of SN regions (e.g. the IFG and ATL) than DN regions (e.g. vmPFC; Chiou et al. 2020). However, using multivariate analysis, in the present study we found that representational content of the dmPFC was more similar to those of DN regions relative to SN regions, making it more affiliated to the DN. Together, our data are consistent with the previous literature concerning this region's somewhat inconclusive affiliation with different networks; it also highlights the flexibility of a “connector” region and the peril of imposing a rigid demarcation on the perimeter between networks.

The representational structure of self- versus other-concept

Our findings add to the literature that deciphering self-/other-referential thoughts is possible based on the patterns of DN and SN regions (Hassabis et al. 2013; Thornton and Mitchell 2017; Courtney and Meyer 2020; Peer et al. 2021). Moreover, our cross-classification analyses have important implications for the cognitive theories of self-processing. The ability to successfully cross-classify has been considered as a benchmark of testing whether there is genuine abstraction of neural coding across domains (Kaplan et al. 2015). We found that the brain uses robustly contextually generalizable neural patterns to encode selfness and otherness such that this cardinal representational code applies across near and far interpersonal distances and across positive- and negative-valence depictions. On top of the separation between self and other, social distance was another factor that sculpted the representational landscape of neural codes. As discussed earlier, “self vs. other” serves as the primary dimension that delineates the most salient difference between classes, while other auxiliary dimensions (e.g. social distance) are couched within the primary segregation. This is reminiscent of representational dimensions of visual perception (Konkle and Caramazza 2013): Animacy (living vs. nonliving) is the primary dimension that separates animate entities from inanimate items, while object size serves as the secondary dimension that assorting subgroups within the inanimate class (Long et al. 2018). Interestingly, participants' behavioral rating on personality was not correlated with the neural pattern of any region. A possible reason of this result might be that while the local pattern of an ROI is sufficient to represent (i)

the binary/broad-stroke difference of self versus other and (ii) the extent of social distance that intersects with the “self/other” dimension, it was insufficient to capture the more sophisticated pattern inherent in the thoughts behind the behavioral rating. As the matrix of behavioral rating illustrates (Fig. 9), our participants rated their own mother as more similar to themselves and as highly different from the *Queen*; this led to a clear cluster that lumped the 3 personally familiar targets together (*Present Self*, *Past Self*, and *Mother*), separate from the *Queen*. The majority of participants rated that the similarity of *Present Self* and *Mother* was even higher than that of *Present Self* and *Past Self*, which reflects complex, metacognitive thoughts that apparently transcend the binary “self vs. other” boundary. This implies elaborated and multifaceted considerations behind the ratings might require additional neurocomputation beyond the information within a single DN/SN region, whose local pattern cares primarily about “self vs. other” plus social distance.

Conclusion

In the present study, we reported evidence of multivoxel decoding that manifested a robust bipartite split within the brain’s social system into the DN and SN. Our findings showed that the 2 subnetworks within the social system contribute differentially to the representations of social concepts, highlighting the peril that the DN and SN should not be conflated as a single, homogenous system. These findings inform the burgeoning field of human-connectome research and its relationship with social neuroscience; our results also shed light on the decade-long investigation into how the brain implements the conceptual distinction of “self vs. other.”

Supplementary material

Supplementary material is available at *Cerebral Cortex Journal* online.

Funding

This research was funded by a Sir Henry Wellcome Fellowship (201381/Z/16/Z) to RC, and an MRC programme grant and intramural funding to MALR (MR/R023883/1; MC_UU_00005/18).

Conflict of interest statement: The authors declare no competing financial interests.

References

- Andrews-Hanna JR, Reidler JS, Sepulcre J, Poulin R, Buckner RL. Functional-anatomic fractionation of the brain’s default network. *Neuron*. 2010;65:550–562.
- Ashburner J. A fast diffeomorphic image registration algorithm. *NeuroImage*. 2007;38:95–113.
- Binney RJ, Hoffman P, Lambon Ralph MA. Mapping the multiple graded contributions of the anterior temporal lobe representational hub to abstract and social concepts: evidence from distortion-corrected fMRI. *Cereb Cortex*. 2016;26:4227–4241.
- Braga RM, Buckner RL. Parallel interdigitated distributed networks within the individual estimated by intrinsic functional connectivity. *Neuron*. 2017;95:457–471.e5.
- Braga RM, van Dijk KRA, Polimeni JR, Eldaief MC, Buckner RL. Parallel distributed networks resolved at high resolution reveal close juxtaposition of distinct regions. *J Neurophysiol*. 2019;121:1513–1534.
- Button KS. Double-dipping revisited. *Nat Neurosci*. 2019;22:688–690.
- Bzdok D, Langner R, Schilbach L, Engemann DA, Laird AR, Fox PT, Eickhoff S. Segregation of the human medial prefrontal cortex in social cognition. *Front Hum Neurosci*. 2013;7:232.
- Bzdok D, Heeger A, Langner R, Laird AR, Fox PT, Palomero-Gallagher N, Vogt BA, Zilles K, Eickhoff SB. Subspecialization in the human posterior medial cortex. *NeuroImage*. 2015;106:55–71.
- Carhart-Harris RL, Friston KJ. The default-mode, ego-functions and free-energy: a neurobiological account of Freudian ideas. *Brain*. 2010;133:1265–1283.
- Chang C-C, Lin C-J. LIBSVM: a library for support vector machines. *ACM Trans Intell Syst Technol*. 2011;2:1–27.
- Chavez RS, Heatherton TF, Wagner DD. Neural population decoding reveals the intrinsic positivity of the self. *Cereb Cortex*. 2016;27:5222–5229.
- Chiou R, Lambon Ralph MA. Unveiling the dynamic interplay between the hub-and-spoke-components of the brain’s semantic system and its impact on human behaviour. *NeuroImage*. 2019;199:114–126.
- Chiou R, Humphreys GF, Jung J, Lambon Ralph MA. Controlled semantic cognition relies upon dynamic and flexible interactions between the executive ‘semantic control’ and hub-and-spoke ‘semantic representation’ systems. *Cortex*. 2018;103:100–116.
- Chiou R, Humphreys GF, Lambon Ralph MA. Bipartite functional fractionation within the default network supports disparate forms of internally oriented cognition. *Cereb Cortex*. 2020;30:5484–5501.
- Clithero JA, Carter RM, Huettel SA. Local pattern classification differentiates processes of economic valuation. *NeuroImage*. 2009;45:1329–1338.
- Courtney AL, Meyer ML. Self-other representation in the social brain reflects social connection. *J Neurosci*. 2020;40:5616–5627.
- Deković M, Meeus W. Peer relations in adolescence: effects of parenting and adolescents’ self-concept. *J Adolesc*. 1997;20:163–176.
- Dimsdale-Zucker HR, & Ranganath C. Representational similarity analyses: a practical guide for functional MRI applications. In: *Handbook of behavioral neuroscience*. Elsevier; 2018;28:509–525.
- DiNicola LM, Braga RM, Buckner RL. Parallel distributed networks dissociate episodic and social functions within the individual. *J Neurophysiol*. 2020;123:1144–1179.
- Dodell-Feder D, Koster-Hale J, Bedny M, Saxe R. fMRI item analysis in a theory of mind task. *NeuroImage*. 2011;55:705–712.
- Doré BP, Zerubavel N, Ochsner KN. Social cognitive neuroscience: a review of core systems. In: *APA handbook of personality and social psychology*. Vol. 1: Attitudes and social cognition; 2014. pp. 693–720.
- Eickhoff SB, Laird AR, Fox PT, Bzdok D, Hensel L. Functional segregation of the human dorsomedial prefrontal cortex. *Cereb Cortex*. 2016;26:304–321.
- Halai AD, Wellbourne SR, Embleton KV, Parkes LM. A comparison of dual gradient-echo and spin-echo fMRI of the inferior temporal lobe. *Hum Brain Mapp*. 2014;35:4118–4128.

- Hassabis D, Spreng RN, Rusu AA, Robbins CA, Mar RA, Schacter DL. Imagine all the people: how the brain creates and uses personality models to predict behavior. *Cereb Cortex*. 2013;24:1979–1987.
- Haynes J-D. A primer on pattern-based approaches to fMRI: principles, pitfalls, and perspectives. *Neuron*. 2015;87:257–270.
- Heatherton TF, Wyland CL, Macrae CN, Demos KE, Denny BT, Kelley WM. Medial prefrontal activity differentiates self from close others. *Soc Cogn Affect Neurosci*. 2006;1:18–25.
- Hebart MN, Görden K, Haynes J-D. The Decoding Toolbox (TDT): a versatile software package for multivariate analyses of functional imaging data. *Front Neuroinform*. 2015;8:88.
- Hiser J, Koenigs M. The multifaceted role of the ventromedial prefrontal cortex in emotion, decision making, social cognition, and psychopathology. *Biol Psychiatry*. 2018;83:638–647.
- Humphreys GF, Lambon Ralph MA. Mapping domain-selective and counterpointed domain-general higher cognitive functions in the lateral parietal cortex: evidence from fMRI comparisons of difficulty-varying semantic versus visuo-spatial tasks, and functional connectivity analyses. *Cereb Cortex*. 2017;27:4199–4212.
- Humphreys GF, Hoffman P, Visser M, Binney RJ, Lambon Ralph MA. Establishing task-and modality-dependent dissociations between the semantic and default mode networks. *Proc Natl Acad Sci*. 2015;112:7857–7862.
- Huntenburg JM, Bazin P-L, Goulas A, Tardif CL, Villringer A, Margulies DS. A systematic relationship between functional connectivity and intracortical myelin in the human cerebral cortex. *Cereb Cortex*. 2017;27:981–997.
- Jackson RL, Hoffman P, Pobric G, Lambon Ralph MA. The nature and neural correlates of semantic association versus conceptual similarity. *Cereb Cortex*. 2015;25:4319–4333.
- Jackson RL, Hoffman P, Pobric G, Lambon Ralph MA. The semantic network at work and rest: differential connectivity of anterior temporal lobe subregions. *J Neurosci*. 2016;36:1490–1501.
- Jackson RL, Cloutman LL, Ralph MAL. Exploring distinct default mode and semantic networks using a systematic ICA approach. *Cortex*. 2019;113:279–297.
- Jackson RL, Bajada CJ, Lambon Ralph MA, Cloutman LL. The graded change in connectivity across the ventromedial prefrontal cortex reveals distinct subregions. *Cereb Cortex*. 2020;30:165–180.
- Kaplan JT, Man K, Greening SG. Multivariate cross-classification: applying machine learning techniques to characterize abstraction in neural representations. *Front Hum Neurosci*. 2015;9:151.
- Kelley WM, Macrae CN, Wyland CL, Caglar S, Inati S, Heatherton TF. Finding the self? An event-related fMRI study. *J Cogn Neurosci*. 2002;14:785–794.
- Klein SB, Loftus J, Kihlstrom JF. Memory and temporal experience: the effects of episodic memory loss on an amnesic patient's ability to remember the past and imagine the future. *Soc Cogn*. 2002;20:353–379.
- Klein A, Andersson J, Ardekani BA, Ashburner J, Avants B, Chiang M-C, Christensen GE, Collins DL, Gee J, Hellier P. Evaluation of 14 nonlinear deformation algorithms applied to human brain MRI registration. *NeuroImage*. 2009;46:786–802.
- Koban L, Gianaros PJ, Kober H, Wager TD. The self in context: brain systems linking mental and physical health. *Nat Rev Neurosci*. 2021;22:309–322.
- Konkle T, Caramazza A. Tripartite organization of the ventral stream by animacy and object size. *J Neurosci*. 2013;33:10235–10242.
- Koster-Hale J, Saxe R, Dungan J, Young LL. Decoding moral judgments from neural representations of intentions. *Proc Natl Acad Sci*. 2013;110:5648–5653.
- Kriegeskorte N, Goebel R, Bandettini P. Information-based functional brain mapping. *Proc Natl Acad Sci*. 2006;103:3863–3868.
- Kriegeskorte N, Mur M, Bandettini PA. Representational similarity analysis-connecting the branches of systems neuroscience. *Front Syst Neurosci*. 2008;2:4.
- Krienen FM, Tu P-C, Buckner RL. Clan mentality: evidence that the medial prefrontal cortex responds to close others. *J Neurosci*. 2010;30:13906–13915.
- Lambon Ralph MA, Jefferies E, Patterson K, Rogers TT. The neural and computational bases of semantic cognition. *Nat Rev Neurosci*. 2017;18:42–55.
- Leech R, Sharp DJ. The role of the posterior cingulate cortex in cognition and disease. *Brain*. 2014;137:12–32.
- Leech R, Braga R, Sharp DJ. Echoes of the brain within the posterior cingulate cortex. *J Neurosci*. 2012;32:215–222.
- Lieberman MD, Straccia MA, Meyer ML, Du M, Tan KM. Social, self (situational), and affective processes in medial prefrontal cortex (MPFC): causal, multivariate, and reverse inference evidence. *Neurosci Biobehav Rev*. 2019;99:311–328.
- Long B, Yu C-P, Konkle T. Mid-level visual features underlie the high-level categorical organization of the ventral stream. *Proc Nat Acad Sci U S A*. 2018;115:E9015–E9024.
- Margulies DS, Ghosh SS, Goulas A, Falkiewicz M, Huntenburg JM, Langs G, Bezgin G, Eickhoff SB, Castellanos FX, Petrides M. Situating the default-mode network along a principal gradient of macroscale cortical organization. *Proc Natl Acad Sci*. 2016;113:12574–12579.
- Meyer ML, Lieberman MD. Why people are always thinking about themselves: medial prefrontal cortex activity during rest primes self-referential processing. *J Cogn Neurosci*. 2018;30:714–721.
- Mitchell JP, Macrae CN, Banaji MR. Dissociable medial prefrontal contributions to judgments of similar and dissimilar others. *Neuron*. 2006;50:655–663.
- Moran JM, Macrae CN, Heatherton TF, Wyland CL, Kelley WM. Neuroanatomical evidence for distinct cognitive and affective components of self. *J Cogn Neurosci*. 2006;18:1586–1594.
- Murphy C, Poerio G, Sormaz M, Wang H-T, Vatansever D, Allen M, Margulies DS, Jefferies E, Smallwood J. Hello, is that me you are looking for? A re-examination of the role of the DMN in social and self-relevant aspects of off-task thought. *PLoS One*. 2019a;14:e0216182.
- Murphy C, Wang H-T, Konu D, Lowndes R, Margulies DS, Jefferies E, Smallwood J. Modes of operation: a topographic neural gradient supporting stimulus dependent and independent cognition. *NeuroImage*. 2019b;186:487–496.
- Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N. A toolbox for representational similarity analysis. *PLoS Comput Biol*. 2014;10:e1003553.
- Northoff G. Personal identity and cortical midline structure (CMS): do temporal features of CMS neural activity transform into “self-continuity”? *Psychol Inq*. 2017;28:122–131.
- Olson IR, McCoy D, Klobusicky E, Ross LA. Social cognition and the anterior temporal lobes: a review and theoretical framework. *Soc Cogn Affect Neurosci*. 2013;8:123–133.
- Oyserman D, Elmore K, Smith G. Self, self-concept, and identity. In: *Handbook of self and identity*. New York, NY: The Guilford Press; 2012. pp. 69–104
- Peer M, Hayman M, Tamir B, Arzy S. Brain coding of social network structure. *J Neurosci*. 2021;41:4897–4909.
- Pereira F, Mitchell T, Botvinick M. Machine learning classifiers and fMRI: a tutorial overview. *NeuroImage*. 2009;45:S199–S209.

- Pessoa L. Understanding brain networks and brain organization. *Phys Life Rev.* 2014;11:400–435.
- Petersen SE, Sporns O. Brain networks and cognitive architectures. *Neuron.* 2015;88:207–219.
- Philippi CL, Duff MC, Denburg NL, Tranel D, Rudrauf D. Medial PFC damage abolishes the self-reference effect. *J Cogn Neurosci.* 2012;24:475–481.
- Poser BA, Versluis MJ, Hoogduin JM, Norris DG. BOLD contrast sensitivity enhancement and artifact reduction with multiecho EPI: parallel-acquired inhomogeneity-desensitized fMRI. *Magn Reson Med.* 2006;55:1227–1235.
- Raichle ME, MacLeod AM, Snyder AZ, Powers WJ, Gusnard DA, Shulman GL. A default mode of brain function. *Proc Natl Acad Sci.* 2001;98:676–682.
- Saxe R, Kanwisher N. People thinking about thinking people: the role of the temporo-parietal junction in “theory of mind”. *NeuroImage.* 2003;19:1835–1842.
- Saxe R, Wexler A. Making sense of another mind: the role of the right temporo-parietal junction. *Neuropsychologia.* 2005;43:1391–1399.
- Skerry AE, Saxe R. A common neural code for perceived and inferred emotion. *J Neurosci.* 2014;34:15997–16008.
- Skerry AE, Saxe R. Neural representations of emotion are organized around abstract event features. *Curr Biol.* 2015;25:1945–1954.
- Smith DV, Clithero JA, Rorden C, Karnath H-O. Decoding the anatomical network of spatial attention. *Proc Natl Acad Sci.* 2013;110:1518–1523.
- Spreng RN, Grady CL. Patterns of brain activity supporting autobiographical memory, prospection, and theory of mind, and their relationship to the default mode network. *J Cogn Neurosci.* 2010;22:1112–1123.
- Spreng RN, Stevens WD, Chamberlain JP, Gilmore AW, Schacter DL. Default network activity, coupled with the frontoparietal control network, supports goal-directed cognition. *NeuroImage.* 2010;53:303–317.
- Spreng RN, Sepulcre J, Turner GR, Stevens WD, Schacter DL. Intrinsic architecture underlying the relations among the default, dorsal attention, and frontoparietal control networks of the human brain. *J Cogn Neurosci.* 2013;25:74–86.
- Steinbeis N. The role of self–other distinction in understanding others’ mental and emotional states: neurocognitive mechanisms in children and adults. *Philos Trans R Soc B Biol Sci.* 2016;371:20150074.
- Tamir DI, Thornton MA, Contreras JM, Mitchell JP. Neural evidence that three dimensions organize mental state representation: rationality, social impact, and valence. *Proc Natl Acad Sci U S A.* 2016;113:194–199.
- Thornton MA, Mitchell JP. Consistent neural activity patterns represent personally familiar people. *J Cogn Neurosci.* 2017;29:1583–1594.
- Thornton MA, Mitchell JP. Theories of person perception predict patterns of neural activity during mentalizing. *Cereb Cortex.* 2018;28:3505–3520.
- Uddin LQ, Yeo BT, Spreng RN. Towards a universal taxonomy of macro-scale functional human brain networks. *Brain Topogr.* 2019;32:926–942.
- Vatanserver D, Menon DK, Stamatakis EA. Default mode contributions to automated information processing. *Proc Natl Acad Sci U S A.* 2017;114:12821–12826.
- Verschuere K, Doumen S, Buyse E. Relationships with mother, teacher, and peers: unique and joint effects on young children’s self-concept. *Attach Hum Dev.* 2012;14:233–248.
- Visser M, Jefferies E, Lambon Ralph MA. Semantic processing in the anterior temporal lobes: a meta-analysis of the functional neuroimaging literature. *J Cogn Neurosci.* 2010;22:1083–1094.
- Wagner DD, Chavez RS, Broom TW. Decoding the neural representation of self and person knowledge with multivariate pattern analysis and data-driven approaches. *Wiley Interdiscip Rev Cogn Sci.* 2019;10:e1482.
- Wang Y, Collins JA, Koski J, Nugiel T, Metoki A, Olson IR. Dynamic neural architecture for social knowledge retrieval. *Proc Natl Acad Sci U S A.* 2017;114:E3305–E3314.
- Yankouskaya A, Humphreys G, Stolte M, Stokes M, Moradi Z, Sui J. An anterior–posterior axis within the ventromedial prefrontal cortex separates self and reward. *Soc Cogn Affect Neurosci.* 2017;12:1859–1868.
- Yeo BTT, Krienen FM, Sepulcre J, Sabuncu MR, Lashkari D, Hollinshead M, Roffman JL, Smoller JW, Zöllei L, Polimeni JR, et al. The organization of the human cerebral cortex estimated by intrinsic functional connectivity. *J Neurophysiol.* 2011;106:1125–1165.
- Yeshurun Y, Nguyen M, Hasson U. The default mode network: where the idiosyncratic self meets the shared social world. *Nat Rev Neurosci.* 2021;22:181–192.