**Shared States: Using MVPA to test neural overlap between self-focused emotion imagery and other-focused emotion understanding.**

Suzanne Oosterwijk[1,3*], Lukas Snoek[2*], Mark Rotteveel[1,3], Lisa F. Barrett[4], & H. Steven Scholte[2,3]

[*] These authors contributed equally to this work.


[1] Department of Social Psychology, University of Amsterdam, The Netherlands

[2] Department of Brain and Cognition, University of Amsterdam, The Netherlands

[3] Amsterdam Brain and Cognition Center, Amsterdam, The Netherlands

[4] Affective Science Institute and Department of Psychology, Northeastern University, Department of Psychiatry, Massachusetts General Hospital and Harvard Medical School, and Martinos Center for Biomedical Imaging, Massachusetts General Hospital, Charlestown, MA, U.S.A.


*To whom correspondence should be addressed. E-mail: S.Oosterwijk@uva.nl

**Author Note**

Word count: 6030

**Abstract**

The present study tested whether the neural patterns that support imagining "performing an action", "feeling a bodily sensation" or "being in a situation" are directly involved in understanding *other people's* actions, bodily sensations and situations. Subjects imagined the content of short sentences describing emotional actions, interoceptive sensations and situations (self-focused task), and processed scenes and focused on *how* the target person was expressing an emotion, *what* this person was feeling, and *why* this person was feeling an emotion (other-focused task). Using a linear support vector machine classifier on brain-wide multi-voxel patterns, we accurately decoded each individual class in the self-focused task. When generalizing the classifier from the self-focused task to the other-focused task, we also accurately decoded whether subjects focused on the emotional actions, interoceptive sensations and situations of *others*. These results show that the neural patterns that underlie self-imagined experience are involved in understanding the experience of other people. This supports the theoretical assumption that the basic components of emotion experience and understanding share resources in the brain.

**Introduction**

To navigate the social world successfully it is crucial to understand other people. But how do people generate meaningful representations of other people's actions, sensations, thoughts and emotions? The dominant view assumes that representations of other people's experiences are supported by the same neural systems as those that are involved in generating experience in the self (e.g., Gallese, Keysers & Rizzolatti, 2004; see for an overview Singer, 2012). We tested this principle of self-other neural overlap directly, using multi-voxel pattern analysis, across three different aspects of experience that are central to emotions: actions, sensations from the body, and situational knowledge.

In recent years, evidence has accumulated that suggests a similarity between the neural patterns representing the self and others. For example, a great variety of studies have shown that observing actions and sensations in other people engages similar neural circuits as acting and feeling in the self (see for an overview Bastiaansen, Thioux & Keysers, 2009). Moreover, an extensive research program on pain has demonstrated an overlap between the experience of physical pain and the observation of pain in other people, utilizing both neuroimaging techniques (e.g., Lamm, Decety & Singer, 2011) and analgesic interventions (e.g., Mischkowski, Crocker & Way, 2016; Rütgen, Seidel, Silani, Riečanský, Hummer, Windischberger, et al, 2015). This process of 'vicarious experience' or "simulation" is viewed as an important component of empathy (Carr, Iacoboni, Dubeau, Mazziotta, & Lenzi, 2003; Decety, 2011; Keysers & Gazzola, 2014). In addition, it is argued that mentalizing (e.g., understanding the mental states of other people) involves the same brain networks as those involved in self-generated thoughts (Uddin, Iacoboni, Lange, & Keenan, 2007; Waytz & Mitchell, 2011). Specifying this idea further, a constructionist view on emotion proposes that both emotion experience *and* interpersonal emotion understanding are produced by the same large-scale distributed brain networks that support the processing of sensorimotor, interoceptive and situationally relevant information (Barrett & Satpute, 2013;

Oosterwijk & Barrett, 2014). An implication of these views is that the representation of self- and other-focused emotional actions, interoceptive sensations and situations overlap in the brain.

Although there is experimental and theoretical support for the idea of self-other neural overlap, the present study is the first to directly test this process using multi-voxel pattern analysis (MVPA) across three different aspects of experience (i.e., actions, interoceptive sensations and situational knowledge). Our experimental design consisted of two different tasks aimed at generating self- and other-focused representations with a relatively large weight given to either action information, interoceptive information or situational information.

In the *self-focused* emotion imagery task (SF-task) subjects imagined performing or experiencing actions (e.g., "*pushing someone away*"), interoceptive sensations (e.g., "*increased heart rate*") and situations (e.g., "*alone in a park at night*") associated with emotion. Previous research has demonstrated that processing linguistic descriptions of (emotional) actions and feeling states can result in neural patterns of activation associated with, respectively, the representation and generation of actions and internal states (Pulvermüller & Fadiga, 2010; Oosterwijk, Mackey, Winkielman,Wilson-Mendenhall & Paulus, 2015). Furthermore, imagery-based inductions of emotion have been successfully used in the fMRI scanner before (Oosterwijk, Lindquist, Anderson, Dautoff, Moriguchi & Barrett, 2012; Wilson-Mendenhall, Barrett, Simmons & Barsalou, 2011), and are seen as robust inducers of emotional experience (Lench, Flores & Bench, 2011). In the *other-focused* emotion understanding task (OF-task), subjects viewed images of people in emotional situations and focused on actions (i.e., *How* does this person express his/her emotions?), interoceptive sensations (i.e., *What* does this person feel in his/her body) or the situation (i.e., *Why* does this person feel an emotion?). This task is based on previous research studying the neural basis of emotion oriented mentalizing (Spunt & Lieberman, 2011).

With MVPA, we examined to what extent the SF-task and OF-task evoked similar neural patterns. MVPA allows researchers to assess whether the neural pattern associated with one set of experimental conditions can be used to distinguish between another set of experimental conditions.

This relatively novel technique has been successfully applied to the field of social neuroscience in general (e.g., Brosch, Bar-David & Phelps, 2013; Parkinson, Liu, & Wheatley, 2014), and the field of self-other neural overlap in particular. For example, several MVPA studies recently assessed whether experiencing pain and observing pain in others involved similar neural patterns (Corradi-Dell'Acqua, Tusche, Vuilleumier & Singer, 2016; Krishnan, Woo, Chang, Ruzic, Gu, López-Solà et al., 2016). Although there is an ongoing discussion about the specifics of shared representation in pain based on these MVPA results (see for an overview Zaki, Wager, Singer, Keysers, & Gazzola, 2016), many authors emphasize the importance of this technique in the scientific study of self-other neural overlap (e.g., Corradi-Dell'Acqua et al., 2016; Krishnan et al., 2016).

MVPA is an analysis technique that decodes latent categories from fMRI data in terms of multi-voxel patterns of activity (Norman, Polyn, Detre & Haxby, 2006). This technique is particularly suited for our research question for several reasons. First of all, although univariate techniques can demonstrate that tasks activate the same brain regions, only MVPA can *statistically test* for shared representation (see Lamm & Majdandžić, 2015). We will evaluate whether multivariate brain patterns that distinguish between mental events in the SF-task can be used to distinguish, above chance level, between mental events in the OF-task. Second, MVPA analyses are particularly useful in research that is aimed at examining distributed representations (Singer, 2008). Based on our constructionist framework, we indeed hypothesize that the neural patterns that will represent self- and other focused mental events are distributed across large-scale brain networks. To capture these distributed patterns, we used MVPA in combination with data-driven univariate feature selection on whole-brain voxel patterns, instead of limiting our analysis to specific regions-of-interest (Haynes, 2015). And third, in contrast to univariate analyses that aggregate data across subjects, MVPA can be performed within-subjects and is thus able to incorporate individual variation in the representational content of multivariate brain patterns. In that aspect within-subject MVPA is sensitive to individual differences in how people imagine actions, sensations and situations, and how they understand others. In short, for our purpose to explicitly test the

assumption that self and other focused processes share neural resources, MVPA is the designated method.

We tested the following two hypotheses. First, we tested whether we could classify *self-imagined* actions, interoceptive sensations and situations above chance level. Second, we tested whether the multivariate pattern underlying this classification could also be used to classify the how, what and why condition in the *other-focused* task.

## Method

### Subjects

In total, we tested twenty-two Dutch undergraduate students from the University of Amsterdam (14 females; $M_{age}$ = 21.48, $SD_{age}$ = 1.75). Of those twenty-two subjects, thirteen subjects were tested twice in two sessions about one week apart. Half of those sessions were used for the model optimization-procedure. The other half of the sessions, combined with an additional nine subjects (who were tested only once), constituted the model validation set (see Analysis and model optimization procedure section). In total, two subjects were excluded from the model validation dataset: one subject was excluded because there was not enough time to complete the experimental protocol and another subject was excluded due to excessive movement (> 3 mm within data acquisition runs).

All subjects signed informed consent prior to the experiment. The experiment was approved by the University of Amsterdam's ethical review board. Subjects received 22.50 euro per session. Standard exclusion criteria regarding MRI safety were applied and people who were on psychopharmacological medication were excluded a priori.

### Experimental design

**Self-focused emotion imagery task**. The self-focused emotion imagery task (SF-task) was created to preferentially elicit *self-focused* processing of action, interoceptive or situational information associated with emotion. Subjects processed short linguistic cues that described actions (e.g., "*pushing someone away*"; "*making a fist*"), interoceptive sensations (e.g., "*being out of*

*breath*"; "*an increased heart rate*"), or situations (e.g., "*alone in a park at night*"; "*being falsely accused*") and were instructed to imagine performing or experiencing the content. The complete instruction is presented in the Supplementary Materials; all stimuli used in the SF-task are presented in Supplementary Table 1. Linguistic cues were selected from a pilot study performed on an independent sample of subjects (n = 24). Details about this pilot study are available on request. The descriptions generated in this pilot study were used as qualitative input to create short sentences that described actions, sensations or situations that were associated with negative emotions, without including discrete emotion terms. The cues did not differ in number of words, nor in number of characters ($F < 1$).

The SF-task was performed in two runs subsequent to the other-focused task using the software package Presentation (Version 16.4, www.neurobs.com). Each run presented sixty sentences on a black background (20 per condition) in a fully-randomized event-related fashion, with a different randomization for each subject. Note that implementing a separate randomization for each subject prevents inflated false positive pattern correlations between trials of the same condition, which may occur in single-trial designs with short inter-stimulus intervals (Mumford, Davis, & Poldrack, 2014). A fixed inter-trial-interval of two seconds separated trials; twelve null-trials (i.e. a black screen for eight seconds) were mixed with the experimental trials at random positions during each run (see Figure 1).

**Other-focused emotion understanding task**. The other-focused emotion understanding task (OF-task) was created to preferentially elicit *other-focused* processing of action, interoceptive or situational information associated with emotion. Subjects viewed images of people in negative situations (e.g., a woman screaming at a man, a man held at gunpoint). A red rectangle highlighted the face of the person that the subjects should focus on to avoid ambiguity in images depicting more than one person. Image blocks were preceded by a cue indicating the strategy subjects should use in perceiving the emotional state of the people in the images (Spunt & Lieberman, 2011). The cue "*How*" instructed the subjects to identify actions that were informative about the person's emotional

state (i.e., *How* does this person express his/her emotions?). The cue "*What*" instructed subjects to identify interoceptive sensations that the person could experience (i.e., *What* does this person feel in his/her body). The cue "*Why*" instructed subjects to identify reasons or explanations for the person's emotional state (i.e., *Why* does this person feel an emotion?). The complete instruction is presented in the Supplementary Materials

Stimuli for the OF-task were selected from the International Affective Picture System database (IAPS; Lang, Bradley, & Cuthbert, 2008), the image set developed by the Kveraga lab (http://www.kveragalab.org/stimuli.html; Kveraga, Boshyan, Adams, Mote, Betz, Ward, Hadjikhani, Bar & Barret, 2015) and the internet (Google images). We selected images based on a pilot study, performed on an independent sample of subjects (n = 22). Details about this pilot study are available on request.

The OF-task was presented using the software package Presentation. The task presented thirty images on a black background in blocked fashion, with each block starting with a what/why or how cue (see Figure 1). Each image was shown three times, once for each cue type. Images were presented in blocks of six, each lasting six seconds, followed by a fixed inter trial interval of two seconds. Null-trials were inserted at random positions within the blocks. Both the order of the blocks and the specific stimuli within and across blocks were fully randomized, with a different randomization for each subject.

**Procedure**

Each experimental session lasted about two hours. Subjects who underwent two sessions had them on different days within a time-span of one week. On arrival, subjects gave informed consent and received thorough task instructions, including practice trials (see the Supplementary Materials for a translation of the task instructions). The actual time in the scanner was 55 minutes, and included a rough 3D scout image, shimming sequence, three-minute structural T1-weighted scan, one functional run for the OF-task and two functional runs for the SF-task. We deliberately

chose to present the SF-task after the OF-task to exclude the possibility that the SF-task affected the OF-task, thereby influencing the success of the decoding procedure.

After each scanning session, subjects rated their success rate for the SF-task and OF-task (see Supplementary Figure 1). In session 2, subjects filled out three personality questionnaires that will not be further discussed in this paper and were debriefed about the purpose of the study.

**Image acquisition**

Subjects were tested using a Philips Achieva 3T MRI scanner and a 32-channel SENSE headcoil. A survey scan was made for spatial planning of the subsequent scans. Following the survey scan, a 3-minute structural T1-weighted scan was acquired using 3D fast field echo (TR: 82 ms, TE: 38 ms, flip angle: 8°, FOV: 240 × 188 mm, 220 slices acquired using single-shot ascending slice order and a voxel size of 1.0 × 1.0 × 1.0 mm). After the T1-weighted scan, functional T2* weighted sequences were acquired using single shot gradient echo, echo planar imaging (TR=2000 ms, TE=27.63 ms, flip angle: 76.1°, FOV: 240 × 240 mm, in-plane resolution 64 x 64, 37 slices (with ascending acquisition), slice thickness 3 mm, slice gap 0.3 mm, voxel size 3 × 3 × 3 mm), covering the entire brain. For the SF-task, 301 volumes were acquired; for the OF-task 523 volumes were acquired.

**Analysis and model optimization procedure**

As MVPA is a fairly novel technique, no consistent, optimal analysis pipeline has been established (Etzel, Valchev, & Keysers, 2011). Therefore we adopted a validation strategy in the present study that is advised in the pattern classification field (Kriegeskorte, Simmons, Bellgowan, & Baker, 2009; Kay, Naselaris, Prenger, & Gallant, 2008). We generated an optimization and validation data-set by running the SF-task and OF-task twice, in two identical experimental sessions for a set of thirteen subjects. The sessions were equally split between the optimization and validation set; first and second sessions were counterbalanced between the two sets. Based on a request received during the review process, we added an additional nine subjects to the validation

data-set. Ultimately, the optimization-set held thirteen sessions and the validation-set, after exclusion of two subjects (see Subjects section), held twenty sessions.

In the optimization set, we explored various preprocessing parameters (i.e. smoothing kernel, low-pass filter, and ICA-based denoising strategies) and MVPA hyperparameter values (i.e. univariate feature selection thresholds and train/test size ratio during cross-validation) to find the optimal combination of hyperparameters and preprocessing settings in terms of classification performance. We measured classification performance by using repeated random subsampling cross-validation with 1000 iterations within the optimization-set (see Multi-voxel pattern analysis subsection for details); we determined the combination of optimal preprocessing and model hyperparameters based on mean classification performance (see Supplementary Table 2 and Supplementary Figure 2). The optimal parameters were then applied to the validation-set, in which the findings from the optimization-set were replicated. All findings discussed in the results section follow from the validation-set (see Supplementary Figure 3 for an overview of the findings from the optimization-set).

**Preprocessing and single-trial modeling**

Functional and structural data were preprocessed and analyzed using FSL 5.0 (Jenkinson, Beckmann, Behrens, Woolrich, & Smith, 2012) and MATLAB (2012b; www.mathworks.com/products/matlab), using an in-house developed preprocessing pipeline and the parameters established in the optimization procedure. Functional data was corrected for motion (using FSL MCFLIRT) and slice timing and was spatially smoothed (5 mm isotropic kernel). After preprocessing, individual time series were modeled using a double gamma hemodynamic response function in a single-trial GLM design using FSL's FEAT. Resulting beta-values were converted to $t$-values (Misaki, Kim, Bandettini, & Kriegeskorte, 2010), constituting a whole-brain pattern of $t$-values per trial. Subsequently, the data was indexed by a gray-matter mask (excluding most white-matter, CSF, and brainstem voxels). Thus, the data points for the multi-voxel pattern analysis consist of whole-brain (gray matter) $t$-value patterns per trial. For the optimization analyses, the

data was transformed to standard space (MNI152, 2 mm), but to reduce computation time for the validation data, and in particular its corresponding permutation analysis, analyses on the validation dataset were performed on data in native (functional) space.

**Multi-voxel pattern analysis**

     **MVPA pipeline**. Before fitting the classifier on the train-data in each iteration of the cross-validation scheme (which is described in more detail in the next section), standardization and voxel selection were estimated and applied to the train-set. Standardization ensured that each feature (i.e. voxel) had zero mean and unit variance across trials. After standardization, voxel selection was performed in each iteration on the train-trials by extracting the voxels with the highest average pairwise Euclidian distance across classes, which will be subsequently referred to as a voxel's differentiation score. More specifically, differentiation scores were calculated by subtracting the mean value across trials per class from each other (i.e. action – interoception, action – situation, interoception – situation), normalizing these values across voxels (yielding '$z$-scores'), and taking their absolute value. The three resulting values per voxel were averaged and the most differentiating voxels ($z$-score threshold: 2.3, as determined by the optimization-procedure) were extracted and used as features when fitting the classifier. Importantly, the standardization parameters (voxel mean and variance) and voxel indices (i.e. which voxels had differentiation scores above threshold) were estimated from the train-data only and subsequently applied to the test-data to ensure independence between the train- and test-set. After standardization and voxel selection in each iteration, a support vector classifier was fit on the training data and cross-validated on the test-set, generating a class probability for each trial in the test-set (see Supplementary Figure 4A for a schematic overview of the MVPA pipeline fitting and cross-validation). Our classifier of choice was the support vector classifier (SVC) implementation from the scikit-learn 'svm' module (Pedregosa et al., 2011) with a linear kernel, fixed regularization parameter ('C') of 1.0, one-vs-one multiclass strategy, estimation of class probability output (instead of discrete class prediction), and otherwise default parameters.

**Cross-validation scheme and bagging procedure**. Cross-validation of the classification analysis was implemented using a repeated random subsampling cross-validation paradigm (also known as Monte Carlo cross-validation), meaning that, for each iteration of the analysis, the classification pipeline (i.e. standardization, voxel selection, and SVM fitting) was applied on a random subset of data points (i.e. the train-set) and cross-validated on the remaining data (i.e. the test-set). Each trial belonged to one out of three classes: action, interoception, or situation. Following the results from the parameter optimization process, we selected four trials per class for testing, amounting to twelve test-trials per iteration.

Per iteration, the classification pipeline was fit on the train-trials from the SF-data. Subsequently, this classifier was cross-validated on twelve test SF-trials ('self-analysis') and twelve test OF-trials ('cross-analysis'; see Figure 2). This process was subsequently iterated 100,000 times to generate a set of class distributions for each trial. After all iterations, the final predicted class of each trial was determined by its highest summed class probability across iterations (also known as 'soft voting'; see Supplementary Figure 4B). This strategy of a random sub-sampling cross-validation scheme in combination with majority (soft) voting is more commonly known as 'bagging' (Breiman, 1996). An important advantage of bagging is that it reduces model overfitting by averaging over an ensemble of models, which is especially useful for multi-voxel pattern analyses because fMRI data is known to display high variance (Varoquaux, Raamana, Engemann, Hoyos-Idrobo, Schwartz, & Thirion, preprint).

After generating a final prediction for all trials using this soft voting method, we constructed confusion matrices for both the self- and cross-analysis. In each raw confusion matrix with prediction counts per class, cells were normalized by dividing prediction counts by the sum over rows (i.e. the total amount of predictions per class), yielding precision-scores (also known as Positive Predictive Value). In other words, this metric represents the ratio of true positives to the sum of true positives and false positives (see Supplementary Figure 5 for a description of the results expressed as *recall* estimates, or the ratio of true positives to the total number of samples in that

class). This classification pipeline generated subject-specific confusion matrices that were subsequently averaged to generate the final classification scores.

**Statistical evaluation.** To evaluate the statistical significance of the observed average precision-scores in the confusion matrices, we permuted the original self- and cross-analysis 1300 times per subject with randomly shuffled class labels, yielding 1300 confusion matrices (with precision-scores). We then averaged the confusion matrices across subjects, yielding 1300 permuted confusion matrices reflecting the null-distribution of each cell of the matrix (which is centered around chance level classification, i.e. 33%). For each cell in the diagonal of the observed confusion matrix, *p*-values were calculated as the proportion of instances of values in the *permuted* matrix which were higher than the values in the *observed* matrix (Nichols & Holmes, 2001). To correct for multiple comparisons, *p*-values were tested against a Bonferroni-corrected threshold. The distribution of precision-scores and the relationship between precision-scores in the self- and cross-analysis is reported in Supplementary Figure 6.

**Spatial representation.** To visualize the classifier feature weights, we plotted the absolute feature weights averaged over iterations, subjects, and pairwise classifiers (action vs. interoception, action vs. situation, interoception vs. situation) that underlie our multiclass classification analysis. We chose to visualize the spatial representation of our model by plotting the average absolute feature weights, because the absolute value of feature weights in linear SVMs can be interpreted as how important the weights are in constructing the model's decision hyperplane (Guyon, Weston, Barnhill, & Vapnik, 2002; Ethofer, Van De Ville, Scherer, & Vuilleumier, 2009; Stelzer, Buschmann, Lohmann, Margulies, Trampel, & Turner, 2014). To correct for a positive bias in plotting *absolute* weights, we ran the main classification analysis again with permuted labels to extract the average absolute feature weights that one would expect by chance. Subsequently, a voxel-wise independent *t*-test was performed for all feature weights across subjects, using the average permuted feature weights as the null-hypothesis, yielding an interpretable *t*-value map (see the supplementary code notebook on our Github repository for computational details).

**Additional analyses**

In addition to the self-analysis and the self-to-other cross-analysis presented in the main text, we also performed a within-subjects other-to-self cross-analysis (see for a similar approach Corradi-Dell'Acqua et al., 2016) and a between-subjects self-analysis and self-to-other cross-analysis. These analyses forward largely similar results as the analyses presented in the main text. Due to space constraints, we present these two additional analyses in the Supplementary Materials. Supplementary Figure 7 represents confusion matrices with precision and recall estimates for the other-to-self cross-analysis. Supplementary Figure 8 presents the results of MVPA analyses using condition-average voxel patterns across subjects instead of single-trial patterns within subjects.

**Univariate analysis**

To be complete, we also report a set of univariate analyses performed on the SF-task and the OF-task data. The univariate analyses were performed on the validation data-set, and were subject to the same preprocessing steps as the MVPA analysis, except that we did not model each trial, but each condition as a separate regressor. The group-level analysis was performed with FSL's FLAME1 option. To examine differences in neural activity between conditions, we calculated contrasts between the three classes in the SF-task (self-action vs. self-interoception; self-action vs. self-situation and self-interoception vs. self-situation) and the three classes in the OF-task (other-action vs. other-interoception; other-action vs. other-situation and other-interoception vs. other-situation). We report clusters that were corrected using cluster-correction with a voxel-wise threshold of .005 (z = 2.7) and a cluster-wise threshold of .05.

**Code availability**.

The MVPA-analysis and subsequent (statistical) analyses were implemented using custom Python scripts, which depend heavily on the *skbold* package, a set of tools for machine learning analyses of fMRI data developed in-house (see https://github.com/lukassnoek/skbold). The original scripts were documented and are hosted at the following Github repository:

https://github.com/lukassnoek/SharedStates.

**Results**

**Multi-voxel pattern analysis**

The analyses of the SF-task demonstrated that voxel patterns reflecting imagined self-focused actions, interoceptive sensations and situations associated with emotion could be decoded accurately for each individual class (all *p* < 0.001, see Figure 3). Furthermore, when we generalized the classifier based on the SF-task to the data from the OF-task (i.e., cross-analysis), we found that neural representations of emotional actions, interoceptive sensations, and situations of *others* could also be reliably decoded above chance (all *p* < 0.001; see Figure 3). Supplementary Table 3 presents mean precision-scores across classes for each subject separately. As predicted, our findings demonstrate that *self-imagined* actions, interoceptive sensations and situations are associated with distinct neural patterns. Furthermore, and as predicted, our findings demonstrate that the patterns associated with *self-imagined* actions, sensations and situations can be used to decode *other-focused* actions, interoceptive sensations and situations (see Supplementary Figure 7 for the complementary other-to-self cross-analysis).

To visualize which neural regions were involved in the successful decoding of the three classes in the OF-task and SF-task, we display in Figure 4 the averaged absolute values of the SVM feature weights. Note that Figure 4 only displays one feature map, as both the self and cross-analysis depend on the same model. Regions displaying high and consistent feature weights across subjects were frontal pole (including parts of the dorsomedial prefrontal cortex and ventromedial prefrontal cortex), orbitofrontal cortex (OFC), inferior frontal gyrus (IFG), superior frontal gyrus (SFG), middle frontal gyrus (MFG), insular cortex, precentral gyrus, postcentral gyrus, posterior cingulate gyrus/precuneus, superior parietal lobule (SPL), supramarginal gyrus (SMG), angular gyrus (AG), medial temporal gyrus (MTG), temporal pole (TP), lateral occipital cortex (LOC) and occipital pole (see Supplementary Table 4  for an overview of all involved regions).

**Univariate analyses**

Figure 5 displays the pattern of neural activity revealed by univariate contrasts between the three different classes in the SF-task and the OF-task. For the sake of brevity, we summarize the most relevant univariate results here. Please see the Supplementary Table 5 and 6 for an overview of all clusters.

In the SF-task, action was associated with increased involvement of the MFG, SFG, AG, SMG, LOC and medial temporal gyrus (temporo-occipital) as compared to interoception, and increased involvement of the IFG, MFG, SFG, anterior cingulate cortex (ACC), supplementary motor area (SMA), precentral gyrus, postcentral gyrus, insular cortex, SMG, SPL, LOC and medial temporal gyrus (temporo-occipital) as compared to situation. Interoception was associated with increased involvement of the insular cortex, precentral gyrus, postcentral gyrus and central operculum as compared to action, and increased involvement of the insular cortex, central operculum, parietal operculum, IFG, frontal pole, ACC, SMA, precentral gyrus, postcentral gyrus, SMG and SPL as compared to situation. The situation vs. action contrast and the situation vs. interoception contrast forwarded clusters in similar regions, including the temporal pole, superior/middle temporal gyrus, IFG, frontal pole, medial prefrontal cortex (mPFC), OFC, precuneus, PCC, LOC, fusiform gyrus, hippocampus and lingual gyrus.

In the OF-task, action was associated with increased involvement of the IFG, MFG, SFG , precentral gyrus, postcentral gyrus, SMG, SPL, middle/inferior temporal gyrus (temporo-occipital), lOC and fusiform gyrus, as compared to interoception, and increased involvement of the IFG, MFG, SFG, frontal pole, precentral gyrus, postcentral gyrus, SMG, SPL, middle/inferior temporal gyrus (temporo-occipital) and LOC, as compared to situation. Interoception was associated with increased involvement of the left frontal pole as compared to action, and increased involvement of the SMG, SPL, precentral gyrus, postcentral gyrus, PCC, IFG and frontal pole, as compared to situation. The situation vs. action contrast and the situation vs. interoception contrast forwarded clusters in similar regions, including the temporal pole, superior/middle temporal gyrus (STG/MTG), frontal pole, mPFC, PCC, precuneus, AG, LOC, occipital pole, fusiform gyrus and lingual gyrus.

**Discussion**

In this study, we investigated the neural overlap between self-focused emotion imagery and other-focused emotion understanding using a decoding approach. The results confirmed our hypothesis that other-focused representations of emotion-related actions, bodily sensations and situations can be decoded from neural patterns associated with accessing similar sources of information in a self-focused task. This cross-classification was successful even though the tasks employed different stimulus materials and instructions. Thus, the observed neural overlap between the underlying processes in the SF-task and OF-task cannot be attributed to similarities in stimulus dimensions or task instructions. Rather, we conclude from our findings that emotion experience and emotion understanding have basic psychological processes in common.

Although we could successfully classify the interoception class in the SF-task and in the OF-task in the validation dataset, we were not able to do this in the optimization dataset. Furthermore, although precision and recall metrics demonstrated similar results for the action and situation cross-classification, these metrics demonstrated different results for the classification of the interoception class (see Supplementary Figure 5). This difference was partly driven by the fact that trials were very infrequently classified as interoception in the cross-classification analysis. The finding that subjects reported lower success rates for the WHAT trials in which they were asked to identify interoceptive sensations in other people than for the HOW (action) and WHY (situation) trials may point to a possible explanation for the inconsistent findings regarding interoception. Although speculative, it may be relatively easy to recognize (and represent) interoceptive sensations when they are described in words (as in the SF-task), but relatively hard to deduce these sensations when only diffuse cues about someone's internal state are available (e.g., posture, frowning facial expression, as in the OF-task).

An exploration of the spatial characteristics of the distributed neural pattern associated with successful decoding of the SF-task and OF-task revealed regions that are commonly active during self- and other-focused processing. First, we found that successful classification was associated with

voxels in the precentral gyrus, IFG, SMA and SPL. These same regions were also revealed by the univariate analyses, in particular for the action and interoception classes. These regions are part of the so-called 'mirror' network, which is argued to support both action planning and action understanding (Gallese, et al., 2004; Bastiaansen, et al., 2009; Van Overwalle & Baetens, 2009; Spunt & Lieberman, 2013). Furthermore, we found that successful classification was associated with voxels in the lateral occipital cortex and fusiform gyrus, which have been linked in the literature to the processing of both concrete and abstract action (Wurm, Ariani, Greenlee, & Lingnau, 2015) and the (visual) processing of emotional scenes, faces and bodies (Sabatinelli, Fortune, Li, Siddiqui, Krafft, Oliver, Beck and Jeffries, 2011; De Gelder, Van den Stock, Meeren, Sinke, Kret, & Tamietto, 2010). The univariate analyses demonstrated activity in the LOC and the fusiform gyrus in particular for the situation class, both when subjects viewed images of other people in emotional situations, and when subjects imagined being in an emotional situation themselves.

Second, we found that successful classification was associated with voxels in regions associated with somatosensory processing (postcentral gyrus) and the representation of interoceptive sensations (insular cortex, see Medford & Critchley, 2010; Craig, 2009). Univariate analyses of the SF-task also demonstrated involvement of these regions for both the action and interoception classes. This pattern of activation is consistent with embodied cognition views that propose that thinking about or imagining bodily states is grounded in simulations of somatosensory and interoceptive sensations (Barsalou, 2009). In contrast to previous work on interoceptive simulation when observing pain or disgust in other people (cf. Bastiaansen, Thioux & Keysers, 2009; Lamm, Decety & Singer, 2011), the univariate analyses of the OF-task did not demonstrate insular cortex activation for the interoception class.

And third, we found that successful classification was associated with voxels in the middle temporal gyrus (including the temporal pole), PCC/precuneus, dmPFC and vmPFC. These regions are part of the so-called 'mentalizing' network (or 'default' network). This same network was also

revealed by the univariate analyses, in particular for the situation class. Meta-analyses have demonstrated that the mentalizing network is commonly active during tasks involving emotion experience and perception (Lindquist, Wager, Kober, Bliss-Moreau & Barrett, 2012), mentalizing/theory of mind (Van Overwalle & Baetens, 2009; Spreng, Mar & Kim, 2008), judgments about the self and others (Denny, Kober, Wager & Ochsner, 2012) and semantic/conceptual processing in general (Binder, Desai, Graves & Conant, 2009). Moreover, this network contributes to the representation of emotion knowledge (Peelen, Atkinson & Vuilleumier, 2010) and is involved in both empathy (Keysers & Gazzola, 2014; Zaki & Ochsner, 2012) and self-generated thought (Andrews-Hanna, Smallwood, & Spreng, 2014). We propose that this network supports the implementation of situated knowledge and personal experience that is necessary to generate rich mental models of emotional situations, both when experienced individually, and when understood in someone else (cf. Barrett & Satpute, 2013; Oosterwijk & Barrett, 2014).

The most important contribution of our study is that it provides direct evidence for the idea of shared neural resources between self-and other focused processes. It is important, however, to specify what we think this "sharedness" entails. In research on pain, there is an ongoing discussion about whether experiencing pain and observing pain in others are distinct processes (Krishnan et al., 2016), or whether experiencing and observing pain involve a shared domain-specific representation (e.g., a discrete pain-specific brain state; Corradi-Dell'Acqua et al., 2016) and/or the sharing of domain-general processes (e.g., general negative affect; Zaki et al., 2016). Connecting to this discussion, we think that it is unlikely that our decoding success reflects the sharing of discrete experiential states between the SF-task and OF-task. After all, unlike in studies on pain, the stimuli in our tasks referred to a large variety of different actions, sensations and situations. Instead, decoding success in our study is most likely due to shared brain state configurations, reflecting the similar engagement of domain-general processes evoked by self- and other-focused instances of action (or interoceptive sensation or situation). This interpretation is consistent with views that suggests that global processes are shared between pain experience and pain observation (Lamm, et

al., 2011; Zaki et al., 2016) or between self- and other-focused tasks in general (e.g., Legrand & Ruby, 2009). Moreover, this interpretation is consistent with the suggestion that neural re-use is a general principle of brain functioning (e.g., Anderson, 2016).

In our constructionist view, we posit that emotion imagery and understanding share basic psychological processes (cf. Oosterwijk & Barrett, 2014). More specifically, both emotion imagery and understanding are "conceptual acts" in which the brain generates predictions based on concept knowledge (including sensorimotor and interoceptive predictions) that are meaningful within a particular situational context (see Barrett, 2012; Barrett & Simmons, 2015). Based on accumulating evidence, we propose that these predictions are implemented in domain-general brain networks (cf. Oosterwijk et al., 2012; Barrett & Satpute, 2013). The relative contribution of these networks depends on the demands of the situational context. Specifically, in contexts where people are focused on actions and expressions (their own, or someone else's) a network that supports the representation of sensorimotor states (i.e., the mirror system) may contribute relatively heavily; in contexts where people are focused on bodily states (their own, or someone else's) a network that supports the representation of interoceptive states (i.e., the salience network) may contribute relatively heavily; and in contexts where people are focused on interpreting a situation (their own, or someone else's) a network that supports a general inferential meaning-making function (i.e., the mentalizing network) may contribute relatively heavily (see also Oosterwijk et al., 2015). We believe that it is likely that our ability to successfully distinguish between classes in the self-task relies on the relatively *different* patterns of activity across these networks for actions, interoceptive sensations and situations. Regarding our ability to successfully generalize from the self- to the other-focused task, we believe that this relies on the relatively *similar* pattern of activity across these networks when people generate self-focused or other-focused instances of action (or interoceptive sensation or situation).

Our explicit manipulation of the weight of action, interoceptive and situational information in the SF-task and the OF-task tests the possibility of shared representation in a novel way.

Although this procedure may seem artificial, social neuroscience studies support the notion that

there is contextual variety in the contribution of action, interoceptive, and situation information

when understanding other people (Oosterwijk et al., 2015; Van Overwalle & Baetens, 2009).

Moreover, this weighting may mimic the variability with which these sources of information

contribute to different instances of subjective emotional experience in reality (Barrett, 2012). In

future directions, it may be relevant to apply the current paradigm to the study of individuals in

which access to these sources of information is disturbed (e.g., individuals with different types of

psychopathology) or facilitated (e.g., individuals with high interoceptive sensitivity).

In short, the present study demonstrates that the neural patterns that support imagining

"performing an action", "feeling a bodily sensation", or "being in a situation" are directly involved

in understanding *other people's* actions, sensations and situations. This supports our prediction that

self- and other-focused emotion processes share resources in the brain.

**References**

Andrews-Hanna, J. R., Smallwood, J., & Spreng, R. N. (2014). The default network and self-generated thought: component processes, dynamic control, and clinical relevance. *Annals of the New York Academy of Sciences*,*1316*(1), 29-52.

Barsalou, L. W. (2009). Simulation, situated conceptualization, and prediction. *Philosophical Transactions of the Royal Society of London: Biological Sciences*, 364, 1281-1289.

Barrett, L. F. (2012). Emotions are real. *Emotion, 12*, 413-429.

Barrett, L.F., & Satpute, A.B. (2013). Large-scale brain networks in affective and social neuroscience: towards and integrative functional architecture of the brain. *Current Opinion in Neurobiology, 23*, 1-12.

Bastiaansen, J.A.C.J., Thouix, M., & Keysers, C. (2009). Evidence for mirror systems in emotions. *Philosophical Transactions of the Royal Society B: Biological Science, 364*, 2391-2404.

Binder, J. R., Desai, R. H., Graves, W. W., & Conant, L. L. (2009). Where is the semantic system? A critical review and meta-analysis of 120 functional neuroimaging studies. *Cerebral Cortex*, *19*, 2767-2796.

Breiman, L. (1996). Bagging predictors. *Machine learning, 24*(2), 123-140.

Brosch, T., Bar-David, E., Phelps, E.A. (2013). Implicit race bias decreases the similarity of neural representations of black and white faces. *Psychological Science, 24*, 160-166.

Carr, L., Iacoboni, M., Dubeau, M., Mazziotta, J.C., & Lenzi, G.L. (2003). Neural mechanisms of empathy in humans: a relay from neural systems for imitation to limbic areas. *Proceedings of the National Academy of Sciences*,*100*, 5497-5502.

Corradi-Dell'Acqua, C., Tusche, A., Vuilleumier, P., & Singer, T. (2016). Cross-modal representations of first-hand and vicarious pain, disgust and fairness in insular and cingulate cortex. *Nature communications*, *7*.

Craig, A. D. (2009). How do you feel--now? The anterior insula and human awareness. *Nature Reviews Neuroscience, 10*, 59-70.

Decety, J. (2011). Dissecting the neural mechanisms mediating empathy. *Emotion Review, 3*, 92-108.

de Gelder, B., Van den Stock, J., Meeren, H. K., Sinke, C. B., Kret, M. E., & Tamietto, M. (2010). Standing up for the body. Recent progress in uncovering the networks involved in the perception of bodies and bodily expressions. *Neuroscience & Biobehavioral Reviews, 34*, 513-527.

Denny, B.T., Kober, H., Wager, T.D. & Ochsner, K.N. (2012). A meta-analysis of functional neuroimaging studies of self- and other judgments reveals a spatial gradient for mentalizing in medial prefrontal cortex. *Journal of Cognitive Neuroscience, 24*, 1742-1752.

Ethofer, T., Van De Ville, D., Scherer, K., & Vuilleumier, P. (2009). Decoding of emotional information in voice-sensitive cortices. *Current Biology, 19*, 1028-1033.

Etzel, J. A., Valchev, N., & Keysers, C. (2011). The impact of certain methodological choices on multivariate analysis of fMRI data with support vector machines. *Neuroimage, 54*, 1159-1167.

Gallese, V., Keysers, C., & Rizzolatti, G. (2004). A unifying view on the basis of social cognition. *Trends in Cognitive Sciences, 8*, 396-403.

Guyon, I., Weston, J., Barnhill, S., & Vapnik, V. (2002). Gene selection for cancer classification using support vector machines. *Machine learning, 46*(1-3), 389-422.

Haynes, J.D. (2015). A primer on pattern-based approaches to fMRI: principles, pitfalls, and perspectives. *Neuron, 87*(2), 257-270.

Jenkinson, M., Beckmann, C.F., Behrens, T.E., Woolrich, M.W., and Smith, S.M. (2012). FSL. *NeuroImage, 62*, 782-90.

Kay, K.N., Naselaris, T., Prenger, R.J., & Gallant, J.L. (2008). Identifying natural images from human brain activity. *Nature, 452*, 352-355.

Keysers, C. & Gazzola, V. (2014). Dissociating the ability and propensity for empathy. *Trends in Cognitive Sciences, 18*, 163-166.

Kriegeskorte, N., Simmons, W.K., Bellgowan, P.S., & Baker, C.I. (2009). Circular analysis in systems neuroscience: the dangers of double dipping. *Nature Neuroscience*, *12*, 535-540.

Krishnan, A., Woo, C. W., Chang, L. J., Ruzic, L., Gu, X., López-Solà, M., ... & Wager, T. D. (2016). Somatic and vicarious pain are represented by dissociable multivariate brain patterns. *Elife*, *5*, e15166.

Lamm, C., Decety, J., & Singer, T. (2011). Meta-analytic evidence for common and distinct neural networks associated with directly experienced pain and empathy for pain. *Neuroimage*, *54*, 2492-2502.

Lamm, C., & Majdandžić, J. (2015). The role of shared neural activations, mirror neurons, and morality in empathy–A critical comment. *Neuroscience Research*, *90*, 15-24.

Lang, P.J., Bradley, M.M., & Cuthbert, B.N. (2008). *International affective picture system (IAPS): Affective ratings of pictures and instruction manual. Technical Report A-8.* University of Florida, Gainesville, FL.

Legrand, D., & Ruby, P. (2009). What is self-specific? Theoretical investigation and critical review of neuroimaging results. *Psychological Review, 116*, 252-282.

Lench, H.C., Flores, S.A., Bench, S.W. (2011). Discrete emotions predict changes in cognition, judgment, experience, behavior, and physiology: A meta-analysis of experimental emotion elicitations. *Psychological Bulletin*, 137, 834-855.

Lindquist, K. A., Wager, T. D., Kober, H., Bliss-Moreau, E. & Barrett, L. F. (2012). The brain basis of emotion: A meta-analytic review. *Behavavioral Brain Sciences. 35*, 121-143.

Medford, N., Critchley, H.D., 2010. Conjoint activity of anterior insula and anterior cingulate cortex: awareness and response. *Brain Structure Function, 214*, 535–549.

Misaki, M., Kim, Y., Bandettini, P. A., & Kriegeskorte, N. (2010). Comparison of multivariate classifiers and response normalizations for pattern-information fMRI. *Neuroimage*, 53(1), 103-118.

Mischkowski, D., Crocker, J., & Way, B. M. (2016). From painkiller to empathy killer: acetaminophen (paracetamol) reduces empathy for pain. *Social cognitive and affective neuroscience*, nsw057.

Nichols, T.E., & Holmes, A.P. (2002). Nonparametric permutation tests for functional neuroimaging: A primer with examples. *Human Brain Mapping, 15*, 1-25.

Norman, K. A., Polyn, S. M., Detre, G. J., & Haxby, J. V. (2006). Beyond mind-reading: multi-voxel pattern analysis of fmri data. *Trends in cognitive sciences, 10*, 424–430.

Oosterwijk, S., & Barrett, L.F. (2014). Embodiment in the construction of emotion experience and emotion understanding. In: L. Shapiro (Ed.), *Routledge Handbook of Embodied Cognition* (pp. 250-260). New York: Routledge.

Oosterwijk, S., Lindquist, K. A., Anderson, E., Dautoff, R., Moriguchi, Y., & Barrett, L. F. (2012). Emotions, body feelings, and thoughts share distributed neural networks. *Neuroimage, 62*, 2110-2128.

Oosterwijk, S., Mackey, S., Winkielman, P., Wilson-Mendenhall, C., & Paulus, M.P. (2015). Concepts in context: Processing mental state concepts with internal or external focus involves different neural systems. *Social Neuroscience, 10*, 294-307.

Parkinson, C., Liu, S., & Wheatley, T. (2014). A common cortical metric for spatial, temporal and social distance. *Journal of Neuroscience, 34*, 1979-1987.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O. …, & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. *The Journal of Machine Learning Research, 12*, 2825-2830.

Peelen, M.V., Atkinson, A.P., Vuilleumier, P. (2010). Supramodal representations of perceived emotions in the human brain. *Journal of Neuroscience, 30*, 10127-134.

Pulvermüller, F., & Fadiga, L. (2010). Active perception: Sensorimotor circuits as a cortical basis for language. *Nature Reviews Neuroscience, 11*, 351–360.

Rütgen, M., Seidel, E. M., Silani, G., Riečanský, I., Hummer, A., Windischberger, C., ... & Lamm, C. (2015). Placebo analgesia and its opioidergic regulation suggest that empathy for pain is grounded in self pain. *Proceedings of the National Academy of Sciences, 112*, E5638-E5646.

Sabatinelli, D., Fortune, E. E., Li, Q., Siddiqui, A., Krafft, C., Oliver, W. T., ... & Jeffries, J. (2011). Emotional perception: meta-analyses of face and natural scene processing. *Neuroimage, 54*, 2524-2533.

Singer, T. (2012). The past, present and future of social neuroscience: a European perspective. *Neuroimage, 61*, 437-449.

Spreng, R. N., & Grady, C. L. (2010). Patterns of brain activity supporting autobiographical memory, prospection, and theory of mind, and their relationship to the default mode network. *Journal of Cognitive Neuroscience, 22*, 1112-1123.

Spunt, R.P., & Lieberman, M.D. (2011). An integrative model of the neural systems supporting the comprehension of observed emotional behavior. *Neuroimage, 59*, 3050-3059.

Spunt, R.P., & Lieberman, M.D. (2013). The busy social brain: Evidence for automacity and control in the neural systems supporting social cognition and action understanding. *Psychological Science, 24*, 80-86.

Stelzer, J., Buschmann, T., Lohmann, G., Margulies, D. S., Trampel, R., & Turner, R. (2014). Prioritizing spatial accuracy in high-resolution fMRI data using multivariate feature weight mapping. Frontiers in neuroscience, 8, 66.

Uddin, L.Q., Iacoboni, M., Lange, C. & Keenan, J.P. (2007). The self and social cognition: The role of cortical midline structures and mirror neurons. *Trends in Cognitive Sciences, 11*, 153-157.

Van Overwalle, F., & Baetens, K. (2009). Understanding others' actions and goals by mirror and mentalizing systems: a meta-analysis. *Neuroimage, 48*, 564-584.

Varoquaux, G., Raamana, P.R., Engemann, D.A., Hoyos-Idobro, A., Schwartz, Y., & Thirion, B. (preprint). Assessing and tuning brain decoders: cross-validation, caveats, and guidelines. Retrieved from https://arxiv.org/abs/1606.05201.

Waytz, A. & Mitchell, J.P (2011). The mechanisms for simulating other minds: Dissociations

    between mirroring and self-projection. *Current Directions in Psychological Science, 20,*

    197-200.

Wilson-Mendenhall, C. D., Barrett, L. F., Simmons, W. K., & Barsalou, L. (2011). Grounding

    emotion in situated conceptualization. *Neuropsychologia, 49,* 1105-1127.

Zaki, J. & Ochsner, K.N. (2012). The neuroscience of empathy: progress, pitfalls and promise.

    *Nature Neuroscience, 15,* 675-680.

Zaki, J., Wager, T. D., Singer, T., Keysers, C., & Gazzola, V. (2016). The anatomy of suffering:

    Understanding the relationship between nociceptive and empathic pain. *Trends in Cognitive*
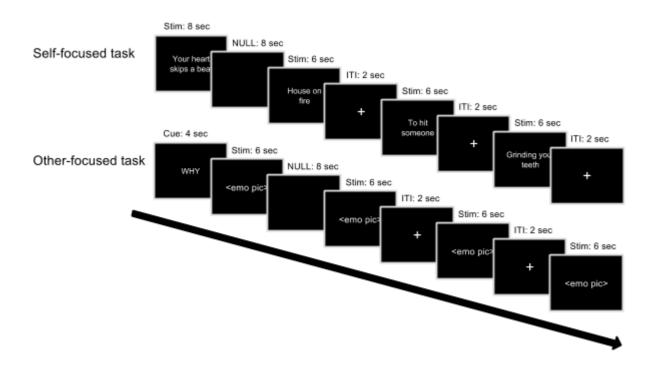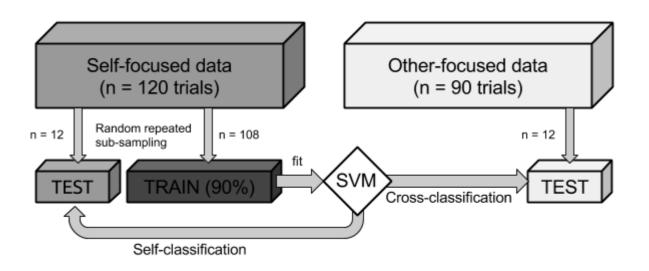
    *Sciences, 20*(4), 249-259.

**Figure Captions**

**Figure 1.** Overview of the self-focused and other-focused task.

**Figure 2**. Schematic overview of data partitioning and classifier training- and cross-validation. Train/test partitioning was done using a repeated random subsampling procedure, also known as Monte Carlo cross-validation. The test-set always consisted of twelve trials (i.e. 4 per class) for both the SF-data and the OF-data.
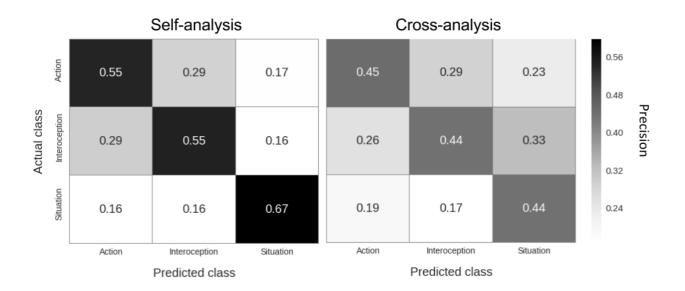
**Figure 3.** Confusion matrices for the self- (left diagram) and cross-analysis (right diagram). Values indicate precision-scores, representing the proportion of true positives given all predictions for a certain class. Note that action and interoception columns in the cross-analysis confusion matrix do not add up to 1, which is caused by the fact that, for some subjects, no trials were predicted as action or interoception, rendering the calculation of precision ill-defined (i.e. zero divided by zero). In this case, precision scores were manually set to zero.
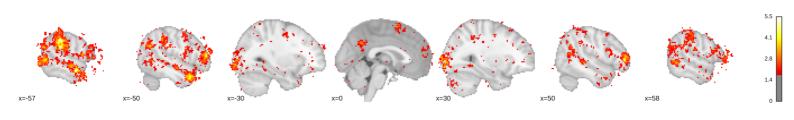
**Figure 4.** Uncorrected $t$-value map of average feature weights across subjects; $t$-values were calculated by dividing the average absolute feature weights, which was corrected for positive bias by subtracting the mean permuted weight across all iterations, by the standard error across subjects.

**Figure 5.** Univariate contrasts for the self-focused and other-focused task.

**Self-focused task**

Stim: 8 sec — Your heart skips a bea[t]

NULL: 8 sec

Stim: 6 sec — House on fire

ITI: 2 sec — +

Stim: 6 sec — To hit someone

ITI: 2 sec — +

Stim: 6 sec — Grinding your teeth

ITI: 2 sec — +

**Other-focused task**

Cue: 4 sec — WHY

Stim: 6 sec — <emo pic>

NULL: 8 sec

Stim: 6 sec — <emo pic>

ITI: 2 sec — +

Stim: 6 sec — <emo pic>

ITI: 2 sec — +

Stim: 6 sec — <emo pic>

x=-57    x=-50    x=-30    x=0    x=30    x=50    x=58

**Self-focused task**

z = 45    y = 65    x = 71

*Contrast*

action > interoception

action > situation

z = 55    y = 59    x = 27

interoception > situation

interoception > action

z = 52    x = 46    x = 68

situation > action

situation > interoception

**Other-focused task**

z = 45    y = 65    x = 71

z = 55    y = 59    x = 27

z = 52    x = 46    x = 68

2.6    *z-value*    5