# The Neural Basis of Shared Preference Learning

Harry Farmer[1,2,*], Uri Hertz[3] and Antonia Hamilton[1]

[1] Institute of Cognitive Neuroscience, University College London

[2] Department of Psychology, University of Bath

[3] Department of Cognitive Sciences, University of Haifa, Israel

*Corresponding Author:

h.farmer@bath.ac.uk

Department of Psychology, 10 West, University of Bath, Bath, BA2 7AY, United Kingdom

Running Title: Shared Preference Learning

# Abstract

During our daily lives, we often learn about the similarity of the traits and preferences of others to our own and use that information during our social interactions. However, it is unclear how the brain represents similarity between the self and others. One possible mechanism is to track similarity to oneself regardless of the identity of the other (Similarity account); an alternative is to track each other person in terms of consistency of their choice similarity with respect to the choices they have made before (consistency account). Our study combined fMRI and computational modelling of reinforcement learning to investigate the neural processes that underlie learning about preference similarity. Participants chose which of two pieces of artwork they preferred and saw the choices of one agent who usually shared their preference and another agent who usually did not. We modelled neural activation with reinforcement learning models based on the similarity and consistency accounts. Our results showed that activity in brain areas linked to reward and social cognition followed the consistency account. Our findings suggest that impressions of other people can be calculated in a person-specific manner which assumes that each individual behaves consistently with their past choices.

# Key Words

fMRI; Reinforcement Learning; Prediction Error; Self; Social Cognition;

# 1. Introduction

The ability to rapidly form and update our impressions about other people is a vital skill in navigating our complex social world. During our daily lives, we frequently learn about the traits and preferences of other people and use that information to inform our social interactions. However, the neural mechanisms which govern our learning of the relationship between our

2

preferences and those of others are currently unclear. The current study investigated these mechanisms by combining fMRI and computational modelling.

Researchers investigating impression formation have sought to determine which brain areas respond when we learn about other people and when our expectations of others are violated. Most have done this by providing participants with some information about a novel person and then presenting either consistent information which confirms the previous impression or inconsistent which requires participants to update their impressions. These studies have shown increased activity in regions like the precuneus/posterior cingulate cortex (PCC), the temporal-parietal junction (TPJ) and the dorsomedial prefrontal cortex (dmPFC) when receiving inconsistent vs. consistent information about another person's moral behaviour (Hughes, Zaki, & Ambady, 2017; Mende-Siedlecki, Baron, & Todorov, 2013; Mende-Siedlecki & Todorov, 2016), competence (Ames & Fiske, 2013; Bhanji & Beer, 2013), traits (Hackel, Doll, & Amodio, 2015; Ma et al., 2012; Van der Cruyssen, Heleven, Ma, Vandekerckhove, & Van Overwalle, 2015), and political beliefs (Cloutier, Gabrieli, Young, & Ambady, 2011). These regions are key nodes in the "mentalising" network which is activated when thinking about the beliefs, preferences and intentions of others (Adolphs, 2009; Frith & Frith, 2012; Schilbach, 2015; Van Overwalle, 2009).

The increased activation to inconsistent information seen in the mentalising network is reminiscent of the prediction error (PE) signal seen in reinforcement learning (RL) models. These signals compute the expectation of a future outcome (or reward) as being a function of the current expectation plus the product of the learning rate and the PE, i.e. the difference between the last expected and actual outcome (Behrens, Hunt, & Rushworth, 2009; Ruff & Fehr, 2014). Reinforcement learning models have been shown to be biologically plausible both at the neuro-chemical level, where the pattern of midbrain dopamine neuron response matches

that of reward PEs (Schultz, 2016), and at the level of whole brain anatomy (Botvinick, Niv, & Barto, 2011). This biological plausibility along with the findings outlined above have led researchers to suggest that regions in the mentalising network may be involved in calculating social prediction errors (Hertz et al., 2017; Mende-Siedlecki, Cai, & Todorov, 2013; Wittmann, Lockwood, & Rushworth, 2018).

Several studies have investigated this possibility directly, using computational modelling to parametrically track prediction error from trial to trial and have found evidence of social prediction error tracking in the dmPFC, the anterior cingulate cortex (ACC), the TJP, the STS, the medial temporal gyrus (MTG), ventrolateral PFC (vlPFC) and the precuneus (Behrens, Hunt, Woolrich, & Rushworth, 2008; Hackel et al., 2015; Lockwood et al., 2018; Stanley, 2016). A recent study by Wittmann et al. (2016) examined the related phenomenon of self-other mergence, in which knowledge about another person's performance reciprocally influences judgements of one's own performance. They found a division between PEs for self-performance, represented in the anterior cingulate cortex, and PEs for other performance, represented in the dmPFC. Interestingly individual variance in the strength of dmPFC activation also predicted how far participants' self PEs were affected by the performance of the others. Such findings have led some researchers (e.g. Bach & Schenke, 2017; Joiner, Piva, Turrin, & Chang, 2017) to argue that predictive processing plays a key role in social cognition.

To date, most studies examining social prediction errors have considered cases where participants learn about other individuals, but do not examine the relationship between those individuals and the self (although see Will, Rutledge, Moutoussis, & Dolan, 2017 for an interesting exception). A distinct literature has examined the role of self-similarity in impression formation (Boer et al., 2011; Montoya & Horton, 2013) and shown that self-similarity can lead to liking and affiliation. Numerous studies have shown that those we

perceive as similar to us in terms of traits (Paunonen & Hong, 2013), attitudes (Montoya & Horton, 2013) and preferences (Boer et al., 2011) tend to be evaluated more favourably than those perceived as different. There is evidence for a ventral-dorsal gradient in the mPFC when processing the similarity of others with similar others being processed in the ventromedial prefrontal cortex (vmPFC) and dissimilar others in the dmPFC (Denny, Kober, Wager, & Ochsner, 2012; Sul et al., 2015).

The current study aims to test how the brain tracks and learns about other people from the self-similarity of their choices. In particular, we distinguish two possible ways in which the brain could track others: the Similarity approach and the Consistency approach. The Similarity approach assumes that, on each trial, we consider 'is this person like me on this trial?' and assign high prediction errors to any trial where an agent makes a different choice to me. The Consistency approach assumes that we model each person we encounter as an individual with a level of overall similarity to me. On each trial, we then consider 'is this person's choice consistent with their overall similarity to me?' and assign high prediction errors to any trial where the agent behaves in a way that is inconsistent with that agent's track record.

To do this we adapted RL models to investigate how the brain tracks the choices of two different agents in terms of how similar they are to the participant's own choices. It is important to note that we are not claiming that the tracking of similarity is necessarily linked to reward based reinforcement in a direct manner. Rather, we use RL models because they can track the accumulation of information and evidence over time. This allows us to look at how the brain represents confirming and disconfirming information about other's similarity to ourselves. For a related approach applied to the learning of others' traits see Zaki, Kallman, Wimmer, Ochsner and Shohamy (2016).

Our task created a context in which participants chose which painting they prefer (an arbitrary aesthetic choice) and then learn the preferences of two agents for the same paintings (see Figure 1). Using fMRI and computational modelling, we can identify which brain areas track agents' preferences relative to self-preferences in a trial-by-trial manner. In each trial our participants saw two paintings and indicated which they preferred. They then saw the preferences of two agents, a similar agent (ASim) who chose the same painting 75% of the time and a different agent (ADiff) who chose the same painting 25% of the time. Using RL models, we are able to calculate the prior probability of the agents' choice and the prediction error of their actual choice separately for each trial and each agent, allowing us to localise brain regions where BOLD signal tracks the model parameters.
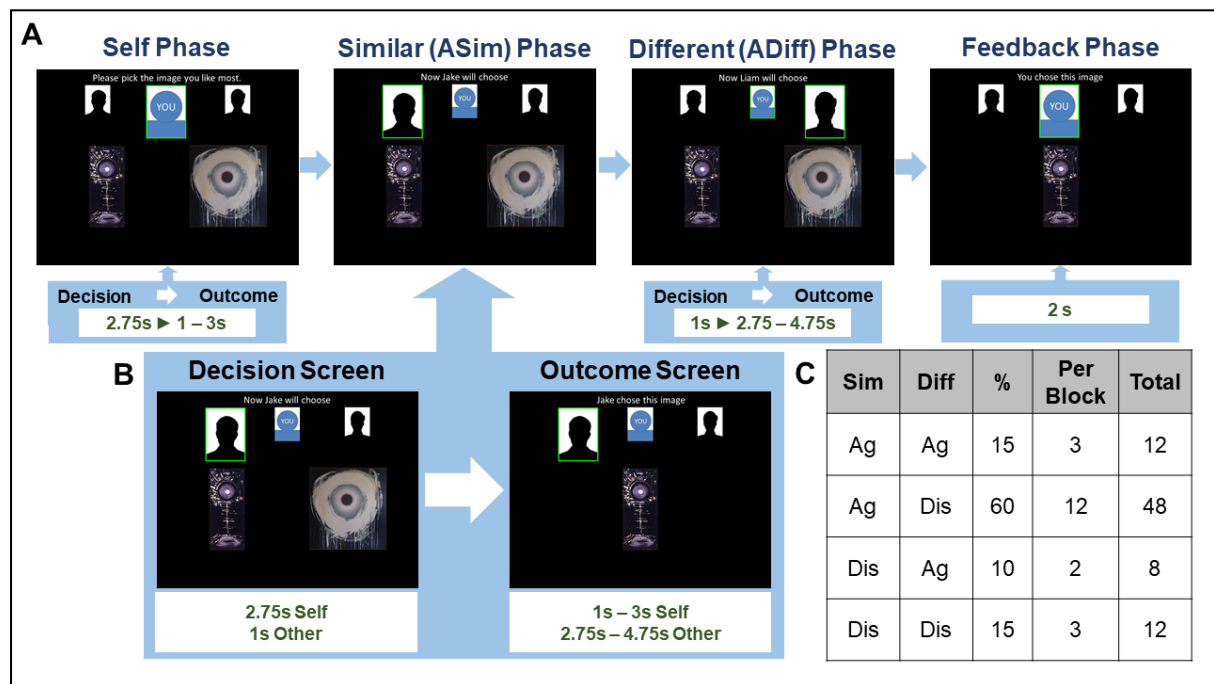


**Figure 1.** *Outline of experimental trial structure and number of trials per condition. A Trial phases & timings. Each trial has four phases (self, similar, different, feedback). On every screen, three icons at the top represent the participant (blue outline in the centre) and the two agents (two photos), with one icon enlarged in a green square to show who is the 'active player' in this phase. In the self-phase, participants chose which of two pictures they prefer. In the*

*ASim phase and ADiff phase, the two agents ASim and ADiff chose pictures and the participant sees the outcome. The order of these two phases was counterbalanced. Finally, in the Feedback phase, the participant sees a reminder of his/her own choice.* **B. Detail of one phase.** *This shows an expanded view of the two different screens within the ASim phase; the same structure was used for the Self phase and ADiff phase. Participant's first see a 'decision screen' with the two pictures used on this trial. During the decision screen participants either chose their own preferred painting (Self phase) or waited to see the choice of the agent (Similar & Different phases). Then they see an 'outcome screen' which shows either the painting they chose (Self phase) or the painting the agent chose (ASim & ADiff phases). The durations of each screen are given at the bottom of the figure, and multiple times separated by a dash represent the jittering in order to effective temporal sampling resolution much finer than one TR.* **C. Number of trials of each type.** *This table shows the breakdown of the four possible combinations of choices made by the two agents, ASim and ADiff. Each agent could agree with the participant's choice (Ag) or disagree (Dis). The columns show the percentage of trials, number of trials by block and total number of trials which had a particular pattern of choices.*

We then used reinforcement learning to create signed prediction error models of both the Similarity and Consistency approaches to tracking the agent's choices (see Figure 2). In the Similarity model, agents' are tracked only in relation to the participant's own preferences, on a single dimension of 'distance from me'. This means that the model will tend to have positive prediction errors for ASim and negative prediction errors for ADiff (see Fig2A). In the RL model, each signed prediction error then contributes to an Accumulated Similarity parameter, which will tend to be high for ASim (who is often similar) and low for ADiff (who is often different). To make this model clear, we term the two parameters the 'similarity prediction error' (PE_Sim) and the 'accumulated similarity' (AS).
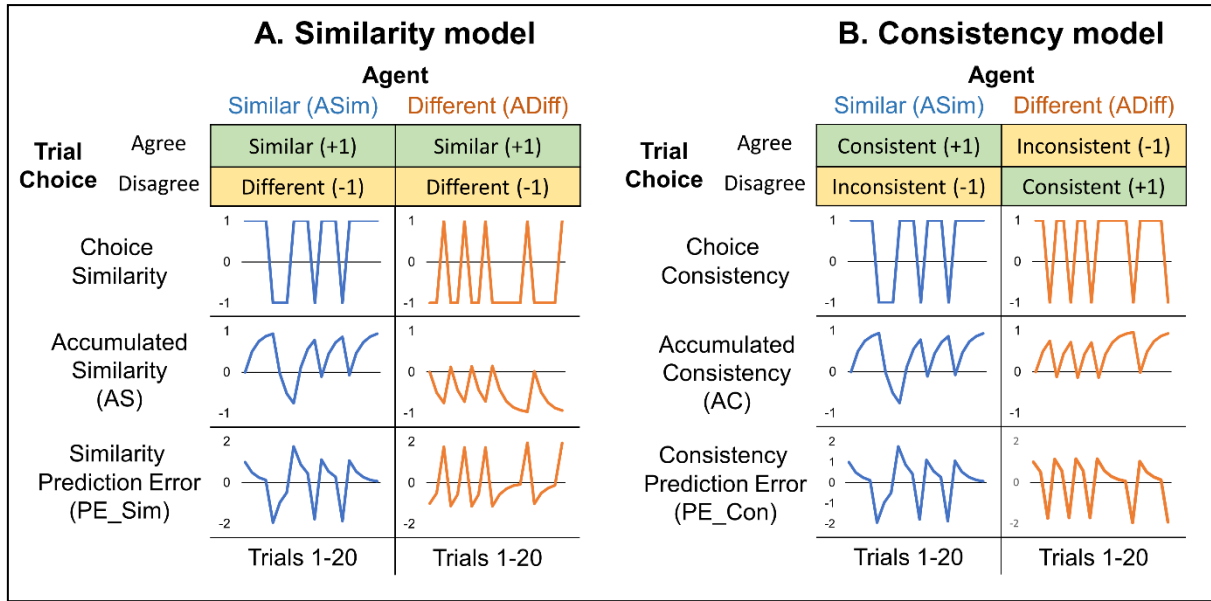
**A. Similarity model**

**Agent**

| | | Similar (ASim) | Different (ADiff) |
|---|---|---|---|
| **Trial Choice** | Agree | Similar (+1) | Similar (+1) |
| | Disagree | Different (-1) | Different (-1) |

Choice Similarity — Similar (ASim): plot (range -1 to 1); Different (ADiff): plot (range -1 to 1)

Accumulated Similarity (AS) — Similar (ASim): plot (range -1 to 1); Different (ADiff): plot (range -1 to 1)

Similarity Prediction Error (PE_Sim) — Similar (ASim): plot (range -2 to 2); Different (ADiff): plot (range -2 to 2)

Trials 1-20 | Trials 1-20

**B. Consistency model**

**Agent**

| | | Similar (ASim) | Different (ADiff) |
|---|---|---|---|
| **Trial Choice** | Agree | Consistent (+1) | Inconsistent (-1) |
| | Disagree | Inconsistent (-1) | Consistent (+1) |

Choice Consistency — Similar (ASim): plot (range -1 to 1); Different (ADiff): plot (range -1 to 1)

Accumulated Consistency (AC) — Similar (ASim): plot (range -1 to 1); Different (ADiff): plot (range -1 to 1)

Consistency Prediction Error (PE_Con) — Similar (ASim): plot (range -2 to 2); Different (ADiff): plot (range -2 to 2)

Trials 1-20 | Trials 1-20

***Figure 2****. Two possible ways that the choices of the two agents, ASim and ADiff, may be tracked in the brain. **A. Similarity approach.** The yellow/green boxes in the top row show how trials are classified as Similar or Different according to whether the agent choose the same picture as the participant or not, and the same classification is used for both agents. Green indicates that a choice is given a positive value and yellow that it has a negative value. This is reflected in the sample sequence of 20 trials, where the 'choice similarity' tends to be high for ASim and low for ADiff. Based on the choice similarity, the Sim_PE and AS parameters are calculated as in Equations 1 and 2. **B. Consistency approach.** Trials are classified as Consistent or Inconsistent according to whether the agent conforms to type. Both agents show high choice consistency most of the time in the sample of 20 trials shown below. Based on the choice consistency the PE_Con and AC parameters are calculated as in Equations 3 and 4.*

The alternative is the consistency model which assumes that participants track agents and choices in terms of whether the agent's choice is consistent with their past level of preference similarity to the participant. Thus, we label each agent's choices as 'consistent' or 'inconsistent' with that agent's past behaviour: agreeing with the participant is *consistent* for ASim but *inconsistent* for ADiff. In this model, a trial will have negative consistency prediction errors when ASim chooses a different picture to the participant, because this is unlike ASim's

typical preference. In the same way a trial will have negative prediction error when ADiff chooses the same picture as the participant (unlike ADiff's typical preference) (see Fig2B). These prediction errors feed into the Accumulated Consistency of each agent, which will be high when that agent conforms to type (i.e. high for both ASim and ADiff most of the time) but will fall if the agent makes atypical choices. To make this model clear, we term the two parameters the 'consistency prediction error' (PE_Con) and the 'accumulated consistency (AC).

Importantly, these two models predict a different pattern of brain activity in our experimental design, as ASim and ADiff's trial-by-trial preferences can have the same sign (both consistent, according to the consistency approach) or opposite sign (as they chose different images, according to the similarity approach, see Figure 2). It is important to note that while our study can test how well each of these models fit activation in different brain areas we are not claiming that they are mutually exclusive competing accounts. Indeed, it is entirely plausible that some brain areas track similarity of choices directly while other track the consistency of choices. Our design allows for us to investigate the neural signature of both models, in two separate GLMs, and thus identify which brain areas (if any) are involved in each of these two ways of processing similarity relationships.

## 2 Methods

### 2.1 Design

In our study participants tracked the choices of two agents on multiple trials, in relation to their own choices. On each trial, the participant and two agents, ASim and ADiff, indicated which of two paintings they preferred. ASim chose the same painting as the participant in 75% of all trials while ADiff only chose the same painting in 25% of trials.

## 2.2 Participants

Twenty-five participants (mean age ± *SD*: 25.1 ± 5.7, 11 male) took part in this study which was approved by The University College London, Institute of Cognitive Neuroscience Research Department's Ethics Committee. All participants gave their informed consent to participate and were paid for their participation. All participants were right-handed and were screened for neurological disorders. Due to technical issues, pre and post ratings data was lost for 7 participants. Therefore, our final sample size for the ratings analysis was n=18. As we did not use this ratings data for model fitting, and data on all 25 participant's choices during the task was collected, this issue did not impact on the fMRI analysis so the full sample n=25 was used for fMRI analysis.

## 2.3 Procedure

### 2.3.1 Experimental Task

The main task in this study was an aesthetic choice task. Participants were told that in each trial they would see a pair of paintings (see Supplementary Materials S1.1) and would have to choose which painting they preferred. They were informed that other participants had previously indicated which of the paintings they preferred and that they would see the choices of two previous participants during the study. Names and faces were assigned to these 'previous participants' but in fact they were computer agents whose choices were determined based on the participant's own choices. Prior to entering the scanner participants completed a training block of the task (see Supplementary Materials S1.2). After the training, participants learnt the names of the agents with whom they would do the experimental task. They also rated their faces for similarity, likeability and attractiveness, using a 10-point scale in order provide us with a manipulation check as to how well the participant's learnt the similarity of the agent to

themselves. Other than being asked to rate their similarity to the agent, participants were not given any information to suggest the relationship between their choices and those of the agents was important to the task.

Each trial was divided into four phases (see Figure 1A). The first three phases were each split into two screens, a *decision screen* and an *outcome screen* (see Figure 1B). In the Self-phase participants were shown a pair of paintings on the *decision* screen and had 2.75 seconds to choose which they preferred using the left and right buttons on a response box. They then saw an *outcome* screen displaying their preferred painting for a jittered interval (1-3 seconds). In the Similar-phase, participants first saw a 1-second *decision* screen which displayed the pair of paintings along with an indicator that ASim was choosing. This was followed by an *outcome* screen which displayed the agent's preferred painting for a jittered interval (2.75 -4.75 seconds). In the Different-phase, participants again saw a *decision* screen with an indicator that ADiff was choosing, followed by a jittered *outcome* screen displaying that agent's preferred painting. The order of the similar & different phases was pseudorandomised across trials. Finally, each trial contained a *feedback* phase in which participants again saw their own choice for an interval of 2-seconds.

Participants completed 4 sessions of 20 trials (see Figure 1C for a breakdown of trial types by block), at the end of each block they rated the similarity, likeability and attractiveness of each agent using a 10-point scale. Using fast event related design, i.e. varying the intervals of the outcome screen in the three choice phases and using many trials, achieved an effective temporal sampling resolution much finer than one TR for each of these periods. The lengths of the intervals were uniformly distributed for each period, ensuring that Evoked Haemodynamic Responses time-locked to the events were sampled evenly across the time period following each choice period.

11

## 2.4 Model-Based fMRI analysis

For full details of image acquisition and fMRI data analysis please see Supplementary Materials S1.3. To examine whether the relationship between the participant preferences and those of the agents was coded in terms of similarity or consistency, two general linear models (GLM) were created, which include different trial types and the parameters of the two RL models. Both GLMs modelled BOLD activation during *outcome* screen for ASim and ADiff separately. Regressors of no interest modelled activity during the *self-choice outcome* screen, the *feedback* phase, the ratings periods, trials where participants failed to make a choice and the residual effects of head motion. In addition, parametric modulators linked to the *outcome* screen regressors allowed us to model the values of our RL parameters on a trial-by-trial basis. Note that we also conducted a more traditional GLM without RL parameters, the details of which can be found in Supplementary materials S2.

In the Similarity GLM we modelled the signed similarity prediction error (PE_Sim) and accumulated similarity (AS) between the agent choice and the participant choice for each agent (n), using the following algorithms:

$$[1] \; PE\_Sim_n(t) = ChoiceSim(t) - AS_n(t)$$

$$[2] \; AS_n(t+1) = AS_n(t) + \lambda * PE\_Sim_n(t)$$

where

$$ChoiceSim(t) = \begin{cases} 1 & agent \; chose \; same \; picture \; as \; participant \\ -1 & agent \; chose \; different \; picture \; from \; participant \end{cases}$$

As we did not fit the model to any response, we set the learning rate ($\lambda$) with a fixed value of 0.5 and initial AS was set to 0. The learning rate of 0.5 was chosen a-priori and fixed for all participants, to indicate the carry-on effect of previous trials to the current trials. This value

was chosen because it is in the middle of the LR range (0-1) and indicates a decaying memory window of about 4 trials. We chose this conservative approach and did not explore learning rates further to avoid double-dipping the data or post-hoc analysis. AS was set at 0 as this represented no a priori expectation of a similarity relationship between the participant and the agents. In total, there were six regressors-of-interest in our Similarity GLM: outcome screens, AS values and PE_Sim values for both ASim and ADiff.

In the Consistency GLM we modelled the signed consistency prediction error (PE_Con) and accumulated consistency (AC) between the agent choice and the participant choice for the two agents (n=ASim or ADiff), using the following algorithm.

$$[3] \; PE\_Con_n(t) = ChoiceCon(t) - AC_n(t)$$

$$[4] \; AC_n(t + 1) = AC_n(t) + \lambda * PE\_Con_n(t)$$

where

$$ChoiceCon(t) = \begin{cases} 1 & \textit{Agent's choice was consistent with their overall} \\ & \textit{similarity to the participant's choices} \\ -1 & \textit{Agent's choice was inconsitent with their} \\ & \textit{overall similairty to the participant's choices} \end{cases}$$

Again, the learning rate ($\lambda$) was set to 0.5 and initial AC was set to 0 (see Figure 2C for examples of how AS and PE varied across 20 trials). In total, there were six regressors-of-interest in our Consistency GLM: outcome screens; AC values and PE_Con values for both ASim and ADiff.

# 3 Results

## 3.1 Behavioural Results

To examine whether learning about the preferences of the agents changed participants feelings of affiliation towards them, we collected ratings of similarity, likeability and trustworthiness at the start of the study and after every 20 trials. This meant that each participant contributed five ratings of each of the three attributes across the study. These ratings were then z-scored within participant to remove baseline differences between participants, before the next analysis. Three separate 2 (agent: similar/different) x 5 (session number: pre/S1/S2/S3/S4) repeated measures ANOVAs were carried out on the z-scored ratings of similarity, liking and trust (See Figure 3). Due to problems with data recording, the ratings from 7 participants were incomplete and were excluded from the behavioural analysis leaving a remaining sample of 18 participants.

The ANOVA on similarity ratings found a significant main effect of agent, $F(1,17) = 23.52$, $p < .001$, $\eta^2_p = .58$. Overall participants rated ASim as being more similar ($M = 0.33$, $MSE = 0.15$) to them than ADiff ($M = -0.68$, $MSE = 0.12$). There was also a significant interaction between agent and session $F(1,17) = 5.65$, $p = .001$, $\eta^2_p = .25$. To examine this interaction further, ratings for ADiff were subtracted from the ratings of ASim for each session to create a difference score. Pairwise comparisons (Bonferroni corrected) showed that the difference score for the pre-session ($M = -0.16$, $MSE = 0.36$) significantly differed from the scores after sessions S1 ($M = 1.49$, $MSE = 0.35$), $p < .05$, S3 ($M = 1.26$, $MSE = 0.29$), $p < .05$, and S4 ($M = 1.43$, $MSE = 0.25$), $p < .01$. No other pairwise comparisons were significant.

The ANOVA on liking ratings found a significant main effect of agent, $F(1,17) = 23.8$, $p < .001$, $\eta^2_p = .58$. Overall participants rated ASim as being more likeable ($M = 0.55$, $MSE = 0.07$) than ADiff ($M = -0.2$, $MSE = 0.12$). There was no significant effect of session and no interaction between session and agent. The ANOVA on trust ratings found a significant

main effect of agent, $F(1,17) = 7.67$, $p < .05$, $\eta^2_p = .31$. Overall participants rated ASim as being more trustworthy ($M = 0.23$, $MSE = 0.11$ than ADiff ($M = -0.24$, $MSE = 0.01$). There was no significant main effect of session and no interaction between session and agent.
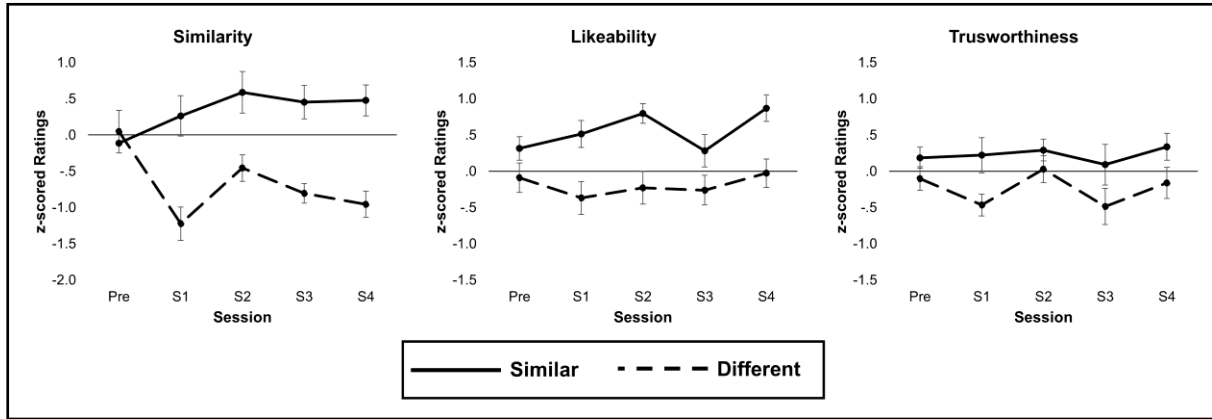


*Figure 3. Z-scored ratings of liking similarity and trustworthiness for the similar and different agents across rating sessions.*

## 3.2 fMRI Results

### 3.2.1 Main Effect of Agent Preference Similarity

Two contrasts investigated the main effect of agent identity (ASim/ADiff) on BOLD response. The regressors which contribute to these contrasts were identical in the Similarity GLM and the Consistency GLM, so the results here are the same for both. The ADiff > ASim contrast revealed that observing the choice of ADiff compared to ASim led to greater activation in the right inferior frontal sulcus (rIFS) and in a cluster centred on the right fusiform gyrus (rFG) (see Table 1 and Figure 4A). No significant activations were found in the ASim > ADiff contrast.

**Table 1.** Peak voxel coordinates in MNI space, z-values and cluster sizes for analyses of the outcome screen showing significant effects after cluster correction for main effect of similarity.

Same shading indicates local maxima in distinct anatomical regions within the same cluster, BA indicates Brodman Area, k indicates the cluster size threshold for whole brain significance of $p < 0.05$.

| Region | Hem. | X | Y | Z | Z-Score | Cluster Size |
|---|---|---|---|---|---|---|
| **Different > Similar (k = 33)** | | | | | | |
| Inferior Frontal Sulcus (BA 44) | R | 38 | 10 | 34 | 3.86 | 57 |
| Fusiform Gyrus (BA 18) | R | 14 | -82 | -10 | 3.51 | 72 |
| Lateral Occipital Gyrus (BA 19) | R | 30 | -82 | -14 | 3.40 | |

### 3.2.2 Parametric Analysis of the Similarity GLM

To identify brain regions which tracked Accumulated Similarity (AS) across both agents, we calculated a conjunction of the RL parameters for each of the agents, that is: $AS_{ASim} \cap AS_{ADiff}$. This did not reveal any significant clusters in either a positive or negative direction, suggesting that no brain areas directly tracked preference similarity between agents and participant. Similarly there were no significant clusters that tracked the positive conjunction of similarity prediction error for both agents, that is, $PE\_Sim_{ASim} \cap PE\_Sim_{ADiff}$. This means that no areas showed increased activation when both agents preference were unexpectedly similar to that of the participant. However, the negative PE_Sim conjunction analysis revealed that unexpected dissimilarity between either agent choice and participant choice correlated with activation in a number of clusters within the occipital cortex including the bilateral lateral occipital cortex (LOC) and the lingual gurus (see Table 2 and Figure 4B).

### 3.2.3 Parametric Analysis of the Consistency GLM

To identify brain regions tracking the consistency of agents' choices across both agents, we first examined the conjunction of areas tracking accumulated consistency (AC), that is $AC_{ASim} \cap AC_{ADiff}$. The positive conjunction showed a significant activation in a cluster-corrected region centred on the superior medial frontal gyrus (smFG) (see Table 3 and Figure 5A). This region showed greater activation as evidence for the consistency of the agents' choice similarity to the self increased, and lower activation during inconsistence periods. No

significant activations were found in the conjunction analysis testing for areas negatively correlated with AC.
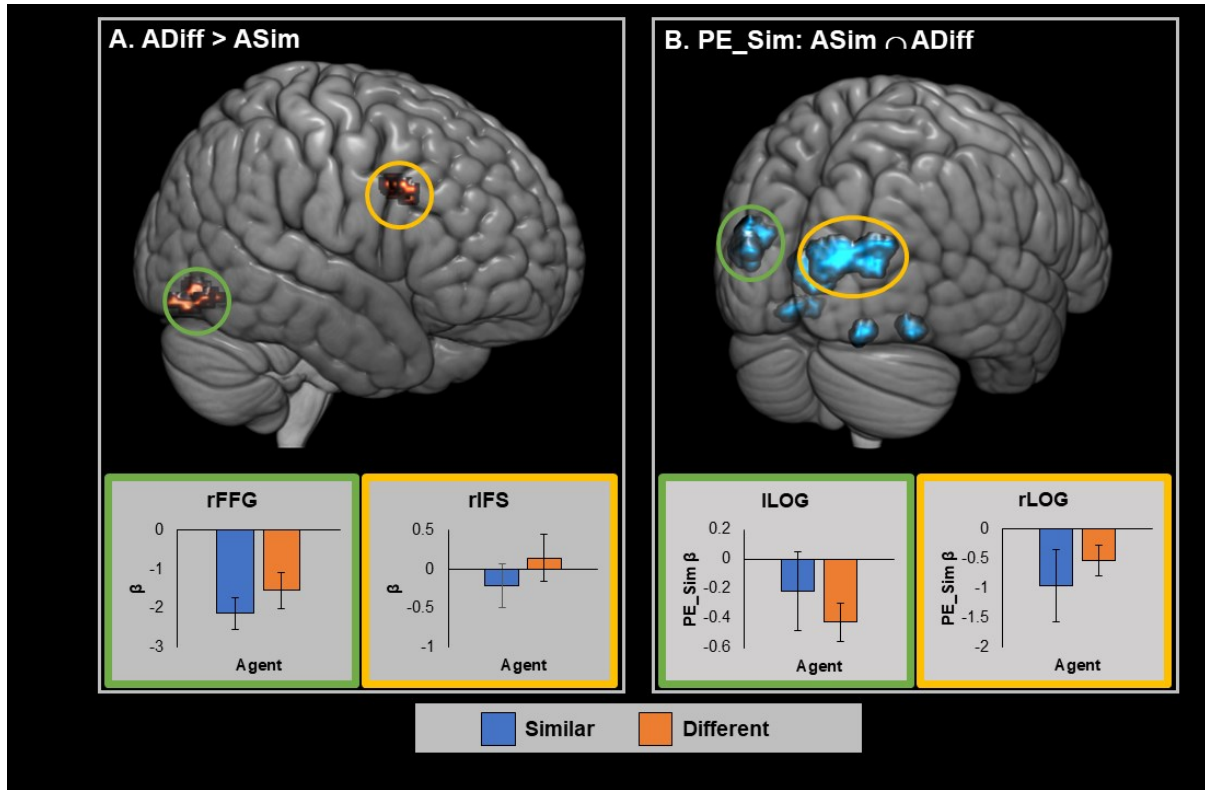


*Figure 4. A) Brain areas showing significant cluster corrected results in the ADiff > ASim contrast for the Outcome screen. B) Brain areas tracking the PE_Sim parameter (Similarity prediction error) for the Outcome screen across both agents, cluster corrected. Parameter estimates in the lower panel are averaged across the whole cluster. Error bars represent SEM. Graph border colours indicate matching circled area. Red/yellow represents positive activations and blue/green represents negative activations.*

The conjunction analysis testing for areas tracking prediction error in consistency (PE_Con$_{ASim}$ ∩ PE_Con$_{ADiff}$) identified significant cluster-corrected activations bilaterally in a dorsal region of the caudate nucleus as well as in a more ventral midbrain region of the left hemisphere (see Table 3 and Figure 5B). These areas showed increased BOLD response when the agents' choices were unexpectedly consistent with their overall preference, and decreased

activation when agents' choices were unexpectedly inconsistent. Note that while the peak activation in the more dorsal left hemisphere cluster is in fact found in the neighbouring corpus callosum both dorsal clusters showed considerable overlap with the caudate nucleus. The conjunction analysis testing for areas tracking PE_Con in a negative direction identified significant clusters in several right hemisphere regions, namely the angular gyrus (rAG), the superior frontal sulcus (rSFS), the superior temporal sulcus (rSTS), the medial temporal gyrus (rMTG) and the Precuneus (see Table 3 and Figure 5C). These areas showed increased BOLD response when the agents' choices were unexpectedly inconsistent with their overall preference, and reduced activity when the agents' choices were highly predictable.

**Table 2.** Peak voxel coordinates in MNI space, z-values and cluster sizes for analyses of the outcome screen in the Similarity GLM showing significant effects after cluster correction for conjunction analyses of the AS and PE parametric modulators. Same shading indicates local maxima in distinct anatomical regions within the same cluster, BA indicates Brodman Area, k indicates the cluster size threshold for whole brain significance of $p < 0.05$.

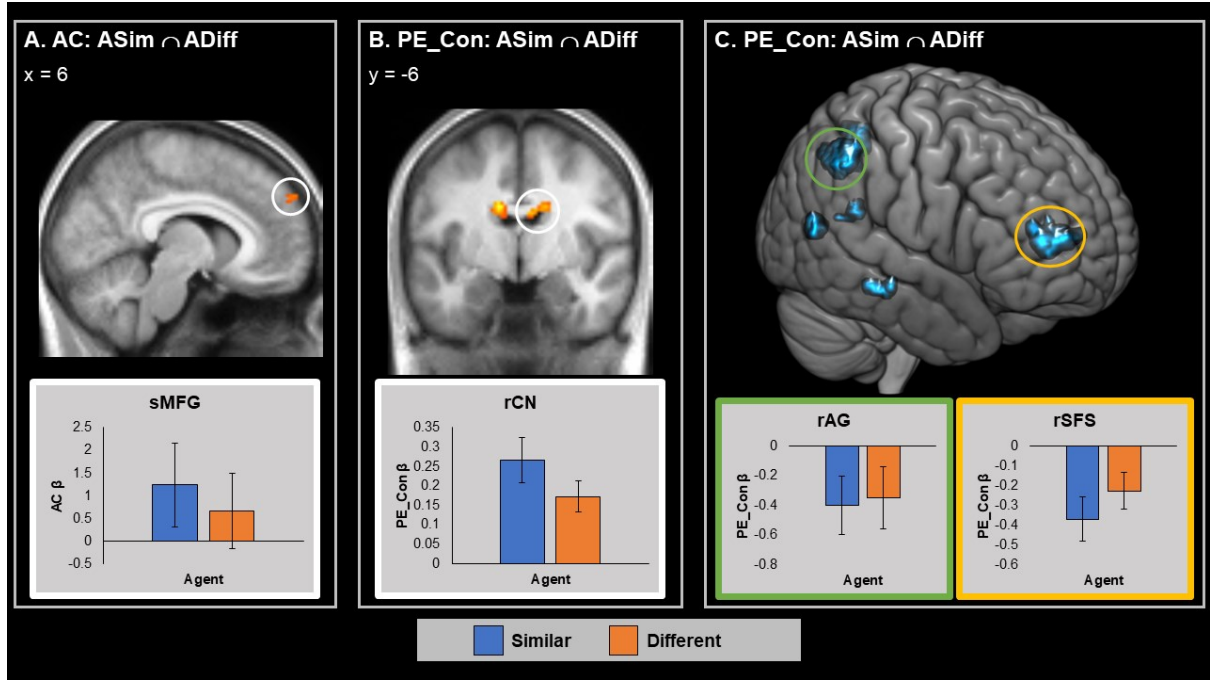| Region | Hem. | X | Y | Z | Z-Score | Cluster Size |
|---|---|---|---|---|---|---|
| **Negative PE_Sim Similar ∩ Different (k = 42)** | | | | | | |
| Lateral Occipital Gyrus (18) | L | -28 | -94 | 16 | 4.06 | 324 |
| Lateral Occipital Gyrus (37) | R | 32 | -54 | -16 | 3.81 | 86 |
| Lateral Occipital Gyrus (18) | R | 24 | -90 | 18 | 3.80 | 457 |
| Middle Occipital Gyrus (19) | R | 36 | -80 | 22 | 3.74 | |
| Lingual Gyrus (17) | L | -6 | -78 | 8 | 3.79 | 249 |
| Lateral Occipital Gyrus (19) | R | 28 | -82 | -16 | 3.67 | 64 |
| Lateral Occipital Gyrus (37) | L | -28 | -60 | -16 | 3.54 | 100 |
| Fusiform Gyrus (37) | L | -26 | -48 | -14 | 3.39 | |

*Figure 5.* Brain areas showing significant cluster corrected tracking of AC and PE_Con for the Outcome screen. A) Areas significantly tracking AC in the positive ASim ∩ ADiff conjunction. B) Areas significantly tracking PE_Con in the positive ASim ∩ ADiff conjunction. C) Areas significantly tracking PE_Con in the negative ASim ∩ ADiff conjunction. Parameter estimates averaged across whole cluster. Error bars represent SEM. Graph border colours indicate matching circled area. Red/yellow represents positive activations and blue/green represents negative activations. sMFG = superior Medial Frontal Gyrus, rCN = right Caudate Nucleus, rAG = right Angular Gyrus, rSFS = right Superior Frontal Sulcus.

# 4 Discussion

Our study examined the neural basis of learning about preference similarity between self and others and its role in promoting affiliation. We created a context where participants could express a preference for a painting and learn about the preferences of two agents for the same paintings. Our behavioural data shows that similar preferences lead to higher ratings of liking, trustworthiness and similarity indicating that participants tracked the agents'

preferences in relation to their own preferences.

**Table 3.** Peak voxel coordinates in MNI space, z-values and cluster sizes for analyses of the outcome screen in the Consistency GLM showing significant effects after cluster correction for conjunction analyses of the AS and PE parametric modulators. Same shading indicates local maxima in distinct anatomical regions within the same cluster, BA indicates Brodman Area, k indicates the cluster size threshold for whole brain significance of $p < 0.05$.

| Region | Hem. | X | Y | Z | Z-Score | Cluster Size |
|---|---|---|---|---|---|---|
| **Positive AC ASim ∩ ADiff (k = 43)** | | | | | | |
| Superior Medial Frontal Gyrus (9) | R | 8 | 56 | 34 | 3.37 | 76 |
| Superior Medial Frontal Gyrus (10) | L | -2 | 54 | 24 | 3.25 | |
| Superior Medial Frontal Gyrus (10) | R | 6 | 56 | 22 | 3.17 | |
| **Positive PE_Con ASim ∩ ADiff (k = 42)** | | | | | | |
| Corpus Callosum | L | -12 | -6 | 28 | 4.54 | 52 |
| Caudate Nucleus | R | 16 | -6 | 28 | 3.90 | 71 |
| Corpus Callosum | L | -4 | 14 | 12 | 3.64 | 56 |
| **Negative PE_Con ASim ∩ ADiff (k = 42)** | | | | | | |
| Angular Gyrus (40) | R | 56 | -46 | 50 | 4.22 | 341 |
| Interparietal Sulcus (40) | R | 32 | -50 | 40 | 3.70 | |
| Superior Frontal Sulcus (10) | R | 34 | 50 | 10 | 4.20 | 270 |
| Superior Temporal Sulcus (37) | R | 60 | -58 | 16 | 3.85 | 76 |
| Superior Temporal Sulcus (41) | R | 44 | -42 | 20 | 3.72 | 43 |
| Superior Temporal Sulcus (39) | R | 42 | -54 | 16 | 3.19 | |
| Precuneus (39) | R | 10 | -56 | 48 | 3.72 | 106 |
| Middle Temporal Gyrus (21) | R | 60 | -20 | -16 | 3.46 | 57 |
| Superior Temporal Sulcus (21) | R | 62 | -28 | -10 | 3.43 | |

Our introduction outlined two possible, non-mutually exclusive, ways in which preference similarity might be tracked in the brain: either by a general mechanism which tracks an agent's choice in relation to one's own, i.e. how similar or dissimilar they are from the self, or via a model of consistency, which tracks agent's choices in terms of their consistency to that

agent's previous choice, i.e. how *consistently* similar or dissimilar they are from the self. To examine the evidence for each of these two mechanisms, we created two reinforcement learning (RL) models which tracked the agents' choices based on similarity and consistency respectively. Our results from the similarity model indicated that regions of the visual cortex negatively tracked similarity prediction error (PE_Sim). Results from the consistency model showed a number of brain areas tracking different variables associated with the consistency model; the dorsomedial pre-frontal cortex (dmPFC) tracking Accumulated Consistency (AC), and the caudate nucleus, angular gyrus and precuneus tracked consistency prediction error (PE_Con). The caudate is involved in value updating (Bhanji & Delgado, 2014; O'Doherty et al., 2004), while the angular gyrus and precuneus are associated with social cognition (Murray, Debbané, Fox, Bzdok, & Eickhoff, 2015; Spreng, Mar, & Kim, 2009). Below we elaborate on the results of the AC conjunction before moving on to discuss the findings on PE_Con and PE_Sim.

## 4.1 dmPFC Tracks Accumulated Consistency

The AC parameter represents a trial-by-trial estimate of the probability that a person makes choices in line with his previous choices, this is, that the similar agent (ASim) should choose the same painting as the participant while the different agent (ADiff) should choose differently. The only area we found tracking AC was a cluster in the bilateral superior medial frontal gyrus (smFG) corresponding to the anterior region of the dmPFC. The dmPFC is known to be a key area for the processing of information about both self and other (Amodio & Frith, 2006; Eickhoff, Laird, Fox, Bzdok, & Hensel, 2014; Mitchell, Banaji, & Macrae, 2005). See supplementary materials S3 for a more detailed survey of previous results.

The dmPFC's involvement in coding prior knowledge of other people is supported by previous research suggesting that the dmPFC encodes reputational priors of one's partners

during economic games (Fouragnan et al., 2013; Hampton, Bossaerts, & O'Doherty, 2008). Our results build on these findings by suggesting that dmPFC prediction errors track the *consistency* of the agent's similarity to the self rather than simply tracking preference similarity.

## 4.2 Consistency Prediction Errors are tracked by Regions Involved in Reward and Social Cognition

PE_Con reflects the difference between the agent's choice and the participant's expectation of what choice the agent will make. For example, the model assigns a positive update signal when ADiff picked the painting not chosen by the participant, and a negative signal when ADiff picked the same painting (See Figure 2). Areas that tracked PE_Con revealed two distinct patterns of activation. Clusters in the bilateral caudate nucleus (Figure 4B) showed increased activity when the agents chose consistently with their type. Meanwhile clusters in regions associated with social cognition including the superior temporal sulcus (STS), the angular gyrus (AG), Precuneus and SFS (superior frontal sulcus; Figure 4C) showed increased activations when the agent's choice was inconsistent with their type. Overall, this pattern shows that prediction error (PE) tracking in these regions is not a 'generic' signal of how similar a person is to me, but rather reflects how much each person's choice conforms to their typical pattern of similarity to me.

The caudate nucleus, along with other parts of the striatum, has been heavily implicated in the generation of prediction errors during RL of rewards for self (Balleine, Delgado, & Hikosaka, 2007; O'Doherty et al., 2004; Schultz, 2015) and others (Báez-Mendoza & Schultz, 2013; Bhanji & Delgado, 2014; Ruff & Fehr, 2014). Previous studies have shown that the caudate nucleus is also involved in signalling PEs when learning the characteristics of others. King-Casas et al. (2005) found that the caudate nucleus activity tracked PEs regarding the

trustworthiness of other during an economic game. Subsequent studies have found similar results for trustworthiness (Fareri, Chang, & Delgado, 2012; Fett, Gromann, Giampietro, Shergill, & Krabbendam, 2014; Fouragnan et al., 2013), generosity (Fareri et al., 2012), reliability in advice giving (Diaconescu et al., 2017) and general behavioural traits (Mende-Siedlecki & Todorov, 2016). Our findings add to this literature by showing that caudate nucleus activity also track prediction error when learning about the similarity of others' preferences to one's own.

The regions showing greater activations when PE_Con was negative, i.e. when the agents' choice was inconsistent with their typical choices, are key nodes of the mentalising network involved in processing information about self and others (Barrett & Satpute, 2013; Murray et al., 2015; Spreng et al., 2009; Van Overwalle, 2009). These areas have been implicated in the formation of impressions about other peoples' traits (Gilron & Gutchess, 2012; Hackel et al., 2015; Hughes et al., 2017; Ma et al., 2012; Mende-Siedlecki, Cai, et al., 2013), beliefs (Cloutier et al., 2011) and abilities (Bhanji & Beer, 2013; Mende-Siedlecki, Baron, et al., 2013). Of particular note are two studies which directly modelled PEs for learning about the traits of other. Hackel et al. (2015) found that the precuneus and STS tracked PEs for other generosity during an economic game while Stanley (2016) found that only the precuneus showed greater tracking of PEs in a social verses non-social setting. The current study shows that these regions also tracks PEs regarding the similarity relationship between self and others, underlining the role of PEs in social learning (Joiner et al., 2017).

It is also notable that while previous studies on social impression formation have tended to show bilateral activations of the mentalising network, in the current studies activity was limited to the right hemisphere. This is consistent with previous research demonstrating right lateralisation for tasks involving self and other differentiation (Decety, 2003; Hu et al., 2016;

Kaplan, Aziz-Zadeh, Uddin, & Iacoboni, 2008; Uddin, Kaplan, Molnar-Szakacs, Zaidel, & Iacoboni, 2005).

## 4.4 Similarity Related Responses in Regions Involved in Visual Attention

In addition to modelling the RL parameters, we also directly contrasted the outcome screen where participants see the choices of ASim with the outcome screen for ADiff. This contrast shows greater activation for ADiff in two clusters; one centred on the rIFS and the other on the rFG. The IFS has been implicated in attentional processing and in particular in the control of attentional shifts by both internal goals and by salient external stimuli (Aron, Robbins, & Poldrack, 2004, 2014; Asplund, Todd, Snyder, & Marois, 2010; Filimon, Philiastides, Nelson, Kloosterman, & Heekeren, 2013; Levy & Wagner, 2012), while the FG is known to play a key role in the visual perception of faces (Contreras, Banaji, & Mitchell, 2013; Kanwisher & Yovel, 2006; Rotshtein, Henson, Treves, Driver, & Dolan, 2005). Interestingly a previous study found greater FG activation when participant observed faces of individuals judged to have different traits to themselves (Leshikar, Cassidy, & Gutchess, 2016). These findings were also consistent with our conjunction analysis of regions that showed a negative relationship to the value of PE_Sim. This analysis revealed that when an agent made an unexpectedly dissimilar choice to that of the participant it led to increased activation across a series of visual areas including regions in the bilateral LOC and in the left FG.

The activation of these areas suggests that participants may have found the choices of ADiff to be more attention grabbing than those of ASim in a comparable way to studies that have demonstrated an attentional bias towards untrustworthy as opposed to trustworthy agents (Dzhelyova, Perrett, & Jentzsch, 2012; Farmer, Apps, & Tsakiris, 2016; Vanneste, Verplaetse, Van Hiel, & Braeckman, 2007).

## 4.5 Comparison with non-Reinforcement Learning GLM

In addition to running our main RL analysis we also conducted a more traditional GLM which divided our trails using a 2 x 2 design with confederate/agent identity (Similar vs Different) as one factor and choice decision (Agree vs Disagree) as the other factor, the interaction between them (i.e. Similar Agree and Different Disagree vs Similar Disagree and Different Agree) was equivalent to our consistency model. This allowed us to compare the results of our RL model to more traditional non-parametric approaches (see Supplementary Materials S2 for full details and results). When comparing the RL models and the traditional GLM the results the activations for the choice main effects and the consistency (interaction effects) are largely similar with the Disagree > Agree contrast showing activations equivalent to the clusters shown for areas that negatively tracked similarity prediction errors, the Consistent > Inconsistent contrast showing activations for two of the three clusters we identified that positively tracked consistency PE and the results for the Inconsistent > Consistent contrast showing results largely consistent with areas negatively tracking consistency PE.

Despite these similarities, our model has two advantages over the non RL GLM. First it is more sensitive to the temporal order of observations, as it takes history into account. For example, it treats differently two consecutive inconsistencies as the first one is more surprising than the second one, while the standard GLM treats them in the same way. This makes our approach more sensitive, more powerful (statistically) and more relevant to our research question. The second advantage is that we can estimate the hidden variables of accumulated consistency/similarity which the standard GLM cannot. This allowed our model to identify the dMPFC area which is involved in the tracking of accumulated consistency.

## 4.6 Limitations

One key limitation of the current study is that our task did not allow us to collect trial-by-trial behavioural data showing what participants had learnt about the agents. This is because we wanted participants to learn implicitly, rather than making explicit predictions of the agent's choice on each trial. Because of this we approximated a learning rate (0.5) and used it in our RL models to track changes in preference tracking according to the actual choices made by the

agents. This raises the possibility that there may only be a weak fit between the learning rate used in our model and the actual learning rate of our participants. However, our main predictions related to the direction of the tracked prediction errors and accumulated preferences, and not with the specific magnitude of these variables, which are less likely to be affected by our approximation. This is in line with a recent theoretical paper (Wilson & Niv, 2015) that demonstrated that model based fMRI results are, under some conditions, insensitive to changes in individual learning rates. While it is possible that our approximation may be lead to lower power at detecting brain responses to prediction errors, we feel that the main hypothesis concerning the direction of the effects (Similarity approach vs Consistency approach) is supported by our analysis.

## 4.6 Conclusions

In this study, we combined computational modelling and fMRI to investigate the neural processes that underlie learning about the similarity of other people's preferences to one's own. We found that more regions of the brain encode information about the similarity of others' choices in a consistency driven manner than encode that information purely based on each particular preference's similarity to one's own. This was particularly the case for the accumulated information about the other's similarity with no areas showing sensitivity to purely accumulated similarity while a region of the dmPFC showed significant tracking of accumulated consistency.

These findings suggest that higher level neural representations of similarity to the self are coded in a person specific manner which reflects how consistent are that person's preference related to the self, i.e. do we usually agree or disagree in our preferences. As such our study highlights the role of context dependent predictive processing in the learning of preference similarity between self and others and, by extension, in the formation of social

26

impressions more generally. Further research in this area could build on our results by examining whether the neural correlates of similarity learning are modulated by having pre-existing cues about how similar that person is to oneself. In addition, it is possible that this consistency approach also applies to learning about other domains including people's traits, attitudes and competence.

## Word Count

5000

## References

Adolphs, R. (2009). The social brain: Neural basis of social knowledge. *Annual Review of Psychology*, *60*, 693–716.

Ames, D. L., & Fiske, S. T. (2013). Outcome dependency alters the neural substrates of impression formation. *NeuroImage*, *83*, 599–608.

Aron, A. R., Robbins, T. W., & Poldrack, R. A. (2004). Inhibition and the right inferior frontal

cortex. *Trends in Cognitive Sciences*, *8*(4), 170–177.

Aron, A. R., Robbins, T. W., & Poldrack, R. A. (2014). Inhibition and the right inferior frontal cortex: One decade on. *Trends in Cognitive Sciences*, *18*(4), 177–185.

Asplund, C. L., Todd, J. J., Snyder, A. P., & Marois, R. (2010). A central role for the lateral prefrontal cortex in goal-directed and stimulus-driven attention. *Nature Neuroscience*, *13*(4), 507–512.

Bach, P., & Schenke, K. C. (2017). Predictive social perception: Towards a unifying framework from action observation to person knowledge. *Social and Personality Psychology Compass*, *11*(7), 1–17.

Báez-Mendoza, R., & Schultz, W. (2013). The role of the striatum in social behavior. *Frontiers in Neuroscience*, *7*, 233.

Balleine, B. W., Delgado, M. R., & Hikosaka, O. (2007). The role of the dorsal striatum in reward and decision-making. *Journal of Neuroscience*, *27*(31), 8161–8165.

Barrett, L. F., & Satpute, A. B. (2013). Large-scale brain networks in affective and social neuroscience: Towards an integrative functional architecture of the brain. *Current Opinion in Neurobiology*, *23*(3), 361–372.

Behrens, T. E. J., Hunt, L. T., & Rushworth, M. F. S. (2009). The computation of social behavior. *Science*, *324*, 1160–1164.

Behrens, T. E. J., Hunt, L. T., Woolrich, M. W., & Rushworth, M. F. S. (2008). Associative learning of social value. *Nature*, *456*(7219), 245–249.

Bhanji, J. P., & Beer, J. S. (2013). Dissociable neural modulation underlying lasting first

impressions, changing your mind for the better, and changing it for the worse. *Journal of Neuroscience*, *33*(22), 9337–9344.

Bhanji, J. P., & Delgado, M. R. (2014). The social brain and reward: Social information processing in the human striatum. *Wiley Interdisciplinary Review of Cognitive Science*, *5*(1), 61–73.

Boer, D., Fischer, R., Strack, M., Bond, M. H., Lo, E., & Lam, J. (2011). How shared preferences in music create bonds between people: Values as the missing link. *Personality and Social Psychology Bulletin*, *37*(1), 1159–1171.

Botvinick, M. M., Niv, Y., & Barto, A. G. (2011). Hierarchically organised behaviour and its neural foundations: A reinforcement-learning perspective. *Modelling Natural Action Selection*, *113*(3), 264–299.

Cloutier, J., Gabrieli, J. D. E., Young, D. O., & Ambady, N. (2011). An fMRI study of violations of social expectations: When people are not who we expect them to be. *NeuroImage*, *57*(2), 583–588.

Contreras, J. M., Banaji, M. R., & Mitchell, J. P. (2013). Multivoxel patterns in fusiform face area differentiate faces by sex and race. *PLoS One*, *8*(7), e69684.

Decety, J. (2003). When the self represents the other: A new cognitive neuroscience view on psychological identification. *Consciousness and Cognition*, *12*(4), 577–596.

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*(8), 1742–52.

Diaconescu, A. O., Mathys, C., Weber, L. A. E., Kasper, L., Mauer, J., & Stephan, K. E. (2017). Hierarchical prediction errors in midbrain and septum during social learning. *Social Cognitive and Affective Neuroscience*, (2016), 618–634.

Dzhelyova, M., Perrett, D. I., & Jentzsch, I. (2012). Temporal dynamics of trustworthiness perception. *Brain Research*, *1435*, 81–90.

Fareri, D. S., Chang, L. J., & Delgado, M. R. (2012). Effects of direct social experience on trust decisions and neural reward circuitry. *Frontiers in Neuroscience*, *6*, 148.

Farmer, H., Apps, M. A. J., & Tsakiris, M. (2016). Reputation in an economic game modulates premotor cortex activity during action observation. *European Journal of Neuroscience*, *44*, 2191–2201.

Fett, A.-K. J., Gromann, P. M., Giampietro, V. P., Shergill, S. S., & Krabbendam, L. (2014). Default distrust? An fMRI investigation of the neural development of trust and cooperation. *Social Cognitive and Affective Neuroscience*, *9*(4), 395–402.

Filimon, F., Philiastides, M. G., Nelson, J. D., Kloosterman, N. A., & Heekeren, H. R. (2013). How embodied is perceptual decision making? Evidence for separate processing of perceptual and motor decisions. *Journal of Neuroscience*, *33*(5), 2121–2136.

Fouragnan, E., Chierchia, G., Greiner, S., Neveu, R., Avesani, P., & Coricelli, G. (2013). Reputational priors magnify striatal responses to violations of trust. *Journal of Neuroscience*, *33*(8), 3602–3611.

Frith, C. D., & Frith, U. (2012). Mechanisms of Social Cognition. *Annual Review of Psychology*, *63*, 287–313.

Gilron, R., & Gutchess, A. H. (2012). Remembering first impressions: Effects of intentionality

and diagnosticity on subsequent memory. *Cognitive, Affective and Behavioral Neuroscience*, *12*(1), 85–98.

Hackel, L. M., Doll, B. B., & Amodio, D. M. (2015). Instrumental learning of traits versus rewards: Dissociable neural correlates and effects on choice. *Nature Neuroscience*, *18*(9), 1233–1235.

Hampton, A. N., Bossaerts, P., & O'Doherty, J. P. (2008). Neural correlates of mentalizing-related computations during strategic interactions in humans. *Proceedings of the National Academy of Sciences*, *105*(18), 6741–6746.

Hertz, U., Palminteri, S., Brunetti, S., Olesen, C., Frith, C. D., & Bahrami, B. (2017). Neural computations underpinning the strategic management of influence in advice giving. *Nature Communications*, *8*(2191), 1–12.

Hu, C., Di, X., Eickhoff, S. B., Zhang, M., Peng, K., Guo, H., & Sui, J. (2016). Distinct and common aspects of physical and psychological self-representation in the brain: A meta-analysis of self-bias in facial and self-referential judgements. *Neuroscience and Biobehavioral Reviews*, *61*, 197–207.

Hughes, B. L., Zaki, J., & Ambady, N. (2017). Motivation alters impression formation and related neural systems. *Social Cognitive and Affective Neuroscience*, *12*(1), 49–60.

Joiner, J., Piva, M., Turrin, C., & Chang, S. W. C. (2017). Social learning through prediction error in the brain. *Npj Science of Learning*, *2*(1), 8.

Kanwisher, N., & Yovel, G. (2006). The fusiform face area: a cortical region specialized for the perception of faces. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *361*(1476), 2109–2128.

Kaplan, J. T., Aziz-Zadeh, L., Uddin, L. Q., & Iacoboni, M. (2008). The self across the senses: An fMRI study of self-face and self-voice recognition. *Social Cognitive and Affective Neuroscience*, *3*(3), 218–223.

King-Casas, B., Tomlin, D., Anen, C., Camerer, C. F., Quartz, S. R., & Montague, P. R. (2005). Getting to know you: Reputation and trust in a two-person economic exchange. *Science*, *308*(78), 78–83.

Leshikar, E. D., Cassidy, B. S., & Gutchess, A. H. (2016). Similarity to the self influences cortical recruitment during impression formation. *Cognitive, Affective, & Behavioral Neuroscience*, *16*(2), 302–314.

Levy, B. J., & Wagner, A. D. (2012). Cognitive control and right ventrolateral prefrontal cortex: Reflexive reorienting, motor inhibition, and action updating. *Annals of the New York Academy of Sciences*, *1224*(1), 40–62.

Lockwood, P. L., Wittmann, M. K., Apps, M. A. J., Klein-Flügge, M. C., Crockett, M. J., Humphreys, G. W., & Rushworth, M. F. S. (2018). Neural mechanisms for learning self and other ownership. *Nature Communications*, *9*, 4747.

Ma, N., Vandekerckhove, M., Baetens, K., Overwalle, F. Van, Seurinck, R., & Fias, W. (2012). Inconsistencies in spontaneous and intentional trait inferences. *Social Cognitive and Affective Neuroscience*, *7*(8), 937–950.

Mende-Siedlecki, P., Baron, S. G., & Todorov, A. (2013). Diagnostic value underlies asymmetric updating of impressions in the morality and ability domains. *Journal of Neuroscience*, *33*(50), 19406–19415.

Mende-Siedlecki, P., Cai, Y., & Todorov, A. (2013). The neural dynamics of updating person

impressions. *Social Cognitive and Affective Neuroscience*, *8*(6), 623–631.

Mende-Siedlecki, P., & Todorov, A. (2016). Neural dissociations between meaningful and mere inconsistency in impression updating. *Social Cognitive and Affective Neuroscience*, *11*(9), 1489–1500.

Montoya, R. M., & Horton, R. S. (2013). A meta-analytic investigation of the processes underlying the similarity-attraction effect. *Journal of Social and Personal Relationships*, *30*(1), 64–94.

Murray, R. J., Debbané, M., Fox, P. T., Bzdok, D., & Eickhoff, S. B. (2015). Functional connectivity mapping of regions associated with self- and other-processing. *Human Brain Mapping*, *36*(4), 1304–1324.

O'Doherty, J. P., Dayan, P., Schultz, J., Deichmann, R., Friston, K. J., & Dolan, R. J. (2004). Dissociable roles of ventral and dorsal striatum in instrumental conditioning. *Science*, *304*(5669), 452–454.

Paunonen, S. V, & Hong, R. Y. (2013). The many faces of assumed similarity in perceptions of personality. *Journal of Research in Personality*, *47*(6), 800–815.

Rotshtein, P., Henson, R. N. A., Treves, A., Driver, J., & Dolan, R. J. (2005). Morphing Marilyn into Maggie dissociates physical and identity face representations in the brain. *Nature Neuroscience*, *8*(1), 107–113.

Ruff, C. C., & Fehr, E. (2014). The neurobiology of rewards and values in social decision making. *Nature Reviews Neuroscience*, *15*(8), 549–562.

Schilbach, L. (2015). The neural correlates of social cognition and social interaction. *Brain Mapping*, *3*, 159–164.

Schultz, W. (2015). Neuronal reward and decision signals: From theories to data. *Physiological Reviews*, *95*(3), 853–951.

Schultz, W. (2016). Dopamine reward prediction error coding. *Dialogues in Clinical Neuroscience*, *18*(1), 23–32.

Spreng, R. N., Mar, R. A., & Kim, A. S. N. (2009). The common neural basis of autobiographical memory, prospection, navigation, theory of mind, and the default mode: A quantitative meta-analysis. *Journal of Cognitive Neuroscience*, *21*(3), 489–510.

Stanley, D. A. (2016). Getting to know you: General and specific neural computations for learning about people. *Social Cognitive and Affective Neuroscience*, *11*(4), 525–536.

Sul, S., Tobler, P. N., Hein, G., Leiberg, S., Jung, D., Fehr, E., & Kim, H. (2015). Spatial gradient in value representation along the medial prefrontal cortex reflects individual differences in prosociality. *Proceedings of the National Academy of Sciences*, *112*(25), 201423895.

Uddin, L. Q., Kaplan, J. T., Molnar-Szakacs, I., Zaidel, E., & Iacoboni, M. (2005). Self-face recognition activates a frontoparietal "mirror" network in the right hemisphere: An event-related fMRI study. *NeuroImage*, *25*(3), 926–935.

Van der Cruyssen, L., Heleven, E., Ma, N., Vandekerckhove, M., & Van Overwalle, F. (2015). Distinct neural correlates of social categories and personality traits. *NeuroImage*, *104*, 336–346.

Van Overwalle, F. (2009). Social cognition and the brain: A meta-analysis. *Human Brain Mapping*, *30*(3), 829–858.

Vanneste, S., Verplaetse, J., Van Hiel, A., & Braeckman, J. (2007). Attention bias toward

noncooperative people. A dot probe classification study in cheating detection. *Evolution and Human Behavior*, *28*(4), 272–276.

Will, G. J., Rutledge, R. B., Moutoussis, M., & Dolan, R. J. (2017). Neural and computational processes underlying dynamic changes in self-esteem. *ELife*, *6*, 1–21.

Wilson, R. C., & Niv, Y. (2015). Is model fitting necessary for model-based fMRI? *PLoS Computational Biology*, *11*(6), e1004237.

Wittmann, M. K., Kolling, N., Faber, N. S., Scholl, J., Nelissen, N., Rushworth, M. F. S., … Nelissen, N. (2016). Self-other mergence in the frontal cortex during cooperation and competition. *Neuron*, *91*(2), 482–493.

Wittmann, M. K., Lockwood, P. L., & Rushworth, M. F. S. (2018). Neural mechanisms of social cognition in primates. *Annual Review of Neuroscience*, *41*(1), 99–118.

Zaki, J., Kallman, S., Wimmer, G. E., Ochsner, K. N., & Shohamy, D. (2016). Social cognition as reinforcement learning: Feedback modulates emotion inference. *Journal of Cognitive Neuroscience*, *28*(9), 1270–1282.

# The Neural Basis of Shared Preference Learning – Supplementary Materials

## S1. Supplementary Methods

### S1.1. Materials

The picture stimuli used in this study comprised of 40 pairs of abstract paintings and 40 pairs of landscape paintings that were matched as closely as possible in terms of their visual and aesthetic properties. This was done to ensure that different pairs of paintings would not differ wildly in their characteristics so that the choices of the different agent would not suggest an abnormal set of preferences. To construct these sets 120 abstract and 120 landscape images were downloaded from the internet and resized to 390x390 JPEG images with any remaining space on either dimension filled in with black. These images were then rated in a pre-study by a group of 20 participants on their complexity, concreteness, attractiveness, valence, affectivity and interest using a 7-point scale. In addition each images' luminance and contrast were calculated using MATLAB (Mathworks 2015).

The mean ratings and luminance and contrast measures for each image were standardised across all images. Similarity scores were created for each measure by subtracting each images score from the score of every other image of its group (abstract or landscape). The similarity scores for all measures were then combined into one value using the following algorithm:

$$(Attractiveness * 2) + (Interest * 2) + Complexity + Valance + Affectivity$$
$$+ \left(\frac{Concreteness}{2}\right) + \left(\frac{Luminance}{2}\right) + \left(\frac{Contast}{2}\right)$$

Each of the images was then paired with its closest neighbour and each pair was then removed from the array. The 40 closest pairs in each group were used in the fMRI experiment and the next 5 closest pairs were used in the training block.

## S1.2 Pre Scanning Training

Prior to entering the scanner participants completed a training block of the task consisting of 10 trials, to ensure that they understood the task. In the training block, they saw two agents of the opposite gender to themselves, and choices were made in the order: Agent 1, Participant, Agent 2. This was to create the belief that the order of making choices was random and that each agent made choices independently. In fact, in the experimental trials, the choices of the agents were determined from the participant choices to create appropriate levels of similarity. After the training, participants learnt the names of the agents with whom they would do the experimental task and rated their faces for similarity, likeability and attractiveness, using a 10-point scale. In the experimental blocks, participants saw these two agents of the same gender as them and the participants' choice always came before the choices of the two agents. Images of the agents were taken from the Karolinska Directed Emotional Face database (Lundqvist, Flykt, & Öhman, 1998) but participants were informed that these photos were stand in images for the actual other participants.

## S1.3 Image Acquisition and Data Analysis

A 1.5 T Siemens TIM Avanto scanner with a 32-channel head coil was used to acquire both T1-weighted structural images and T2*-weighted echoplanar images using the multiband method ($64\times64$ pixels; $3.2\times3.2$ mm; echo time, 55 ms, multiband factor=2) with blood oxygen level-dependent (BOLD) contrast. Each volume comprised 40 axial slices (3.2 mm thick, oriented approximately to the anterior commissure–posterior commissure plane), covering most of the brain but omitting inferior portions of the cerebellum. Functional scans were acquired in four sessions, each comprising 222 volumes (~7.4 min). Volumes were acquired continuously with an effective repetition time of 2s per volume. The first four volumes in each session were discarded to allow for T1 equilibration effects. Prior to functional scanning, a 6 min T1-weighted MPRAGE structural scan was collected at a resolution of $1\times1\times1$ mm. Stimuli were projected onto a screen behind the participant and viewed in a mirror. Participants responded using a 4-button response box. All stimuli were presented with Cogent running under Matlab2014, permitting synchronisation with the scanner and accurate timing of stimuli presentation.

Data were processed and analysed using SPM12 (www.fil.ion.ucl.ac.uk/spm). The EPI images from all four sessions of each participant were realigned to a mean EPI image for that

participant. Images in which the participant moved more than 1.5mm or had rotation of more than 1 degree were visually examined and if seen to contain artefacts were removed from the analysis and replaced with volumes interpolated from the preceding and subsequent images. No participant had artefacts in more than 5% of images. Each participant's structural image was processed using a unified segmentation procedure combining segmentation, bias correction, and spatial normalization to the MNI template (Ashburner & Friston, 2005). The same normalization parameters were then used to normalize the EPI images. Finally, the images were spatially smoothed to conform to the assumptions of the GLM implemented in SPM12 by applying a Gaussian kernel of 8 mm FWHM.

For each of our two GLMs SPM12 was used to compute first level parameter estimates (beta) and t-contrast images (containing weighted parameter estimates) for each contrast at each voxel. To examine regions showing a main effect of agent similarity, two contrasts were carried out between the *outcome screen* regressors (ASim > ADiff, ADiff > ASim). In addition, to examine regions that tracked the RL model parameters in each model, conjunction images were calculated for each RL parameter ($AS_{ASim} / AC_{ASim} \cap AS_{ADiff}/AC_{ADiff}$) and ($PE\_Sim_{ASim} /PE\_Con_{ASim} \cap PE\_Sim_{ADiff}/PE\_Con_{ADiff}$).

For the group-level analysis, the first level images from all participants were subjected to two one-sample t-tests, one in the positive direction and the other in the negative direction. Images derived from these second level analyses were thresholded at $p < 0.001$, uncorrected. For each analysis, a separate Monte Carlo simulation implemented in 3dClustSim (Forman et al., 1995) was used to determine the correct cluster extent threshold needed for a whole brain cluster-wise significance level of $p < 0.05$. Anatomical Regions were determined using the AICHA atlas (Joliot et al., 2015) to for gray matter and the Tractography based Atlas of human brain connections Projection Network (Natbrainlab, Neuroanatomy and Tractography Laboratory) (Catani & de Schotten, 2012; de Schotten et al., 2011) for the white matter.

## S2. Alternative Non Parametric Data Analysis

In addition to our two parametric GLMs we conducted an additional analysis to examine to what extent the use of our parametric modulators shed additional light on our findings when compared with a more traditional factorial GLM.

### S2.1 GLM Design and Data Analysis

This GLM modelled BOLD activation during agent *outcome* screens categorised across the factors of agent (ASim, ADiff) and choice (agree (Ag), disagree (Dis)). Regressors of no interest modelled activity during the *self-choice outcome* screen, the *feedback* phase, the ratings periods, trials where participants failed to make a choice and the residual effects of head motion.

SPM12 was used to compute first level parameter estimates (beta) and t-contrast images (containing weighted parameter estimates) for each contrast at each voxel. To examine regions showing a main effect of agent similarity two t-contrasts were carried out between the *outcome screen* regressors (ASimAg + ASimDis > ADiffAg + ADiffDis , ADiffAg + ADiffDis > ASimAg + ASimDis). To examine regions showing a main effect of choice similarity two t-contrasts were carried out between the *outcome screen* regressors (ASimAg + ADiffAg > ASimDis + ADiffDis, ASimDis + ADiffDis > ASimAg + ADiffAg). Finally, to examine regions showing a main effect of choice consistency two t-contrasts were carried out between the *outcome screen* regressors (ASimAg + ADiffDis > ASimDis + ADiffAg, ASimDis + ADiffAg > ASimAg + ADiffDis).

For the group-level analysis, the first level images from all participants were subjected to one-sample t-tests. Images derived from these second level analyses were thresholded at $p < 0.001$, uncorrected. For each analysis, a separate Monte Carlo simulation implemented in 3dClustSim (Forman et al., 1995) was used to determine the correct cluster extent threshold needed for a whole brain cluster-wise significance level of $p < 0.05$. Anatomical Regions were determined using the AICHA atlas (Joliot et al., 2015) for gray matter and the Tractography based Atlas of human brain connections Projection Network (Natbrainlab, Neuroanatomy and Tractography Laboratory) (Catani & de Schotten, 2012; de Schotten et al., 2011) for white matter.

## S2.2 Factorial GLM Results

Full results of this GLM analysis can be seen below in table S1. As can be seen the results the activations for the choice main effects and the consistency are largely in with those of our parametric modulator GLMs. The Disagree > Agree contrast results showed activations equivalent to the clusters shown for areas that negatively tracked similarity prediction errors, the results for the Consistent > Inconsistent contrast showed significant activations for two of the three clusters we identified that positively tracked consistency PE and the results for the Inconsistent > Consistent contrast showed results largely consistent with areas negatively

tracking consistency PE.

**Table S1.** Peak voxel coordinates in MNI space, z-values and cluster sizes for analyses of the outcome period in the Consistency GLM showing significant effects after cluster correction for conjunction analyses of the AS and PE parametric modulators. Same shading indicates local maxima in distinct anatomical regions within the same cluster, BA indicates Brodman Area, k indicates the cluster size threshold for whole brain significance of $p < 0.05$.

| Region | Hem. | X | Y | Z | Z-Score | Cluster Size |
|---|---|---|---|---|---|---|
| **Disagree > Agree (k = 35)** | | | | | | |
| Lateral Occipital Gyrus (18) | L | -30 | -92 | 22 | 4.29 | 395 |
| Cuneus (18) | L | -12 | -88 | 16 | 3.30 | |
| Lateral Occipital Gyrus (37) | R | 32 | -54 | -16 | 3.94 | 129 |
| Fusiform Gyrus (19) | R | 30 | -64 | -14 | 3.37 | |
| Fusiform Gyrus (37) | R | 28 | -46 | -14 | 3.18 | |
| Lingual Gyrus (17) | L | -6 | -78 | 8 | 3.85 | 233 |
| Lingual Gyrus (18) | L | -8 | -70 | -2 | 3.34 | |
| Occipital Superior Gyrus (18) | R | 24 | -92 | 16 | 3.80 | 337 |
| Middle Occipital Gyrus (19) | R | 36 | -80 | 22 | 3.46 | |
| Lateral Occipital Gyrus (19) | R | 28 | -82 | -16 | 3.78 | 70 |
| Lateral Occipital Gyrus (37) | L | -28 | -60 | -16 | 3.64 | 126 |
| Fusiform Gyrus (37) | L | -26 | -48 | -14 | 3.46 | |
| **Consistent > Inconsistent (35)** | | | | | | |
| Corpus Callosum | L | -2 | 14 | 10 | 3.70 | 37 |
| Corpus Callosum | R | 16 | -6 | 28 | 3.69 | 35 |
| **Inconsistent > Consistent (35)** | | | | | | |
| Superior Temporal Sulcus (37) | R | 62 | -58 | 12 | 4.79 | 108 |
| Supramarginal Gyrus (40) | R | 58 | -42 | 46 | 3.92 | 192 |
| Interparietal Sulcus (7) | R | 28 | -50 | 42 | 3.82 | 65 |
| Superior Frontal Sulcus (10) | R | 34 | 50 | 10 | 3.61 | 62 |
| Precuneus | R | 8 | -58 | 48 | 3.56 | 99 |
| Middle Temporal Gyrus (20) | R | 62 | -24 | -14 | 3.49 | 61 |
| Middle Temporal Gyrus (21) | R | 64 | -18 | -8 | 3.28 | |

However, there were two important differences which point to the increased value of our modelling analysis. First, the simple contrast GLM did not identify any areas that could be explained by agent identity alone as opposed to choice similarity, presumably because the close relationship between these two factors meant that the variance in activation was captured by the choice analysis instead. Second, this analysis does not capture our finding of the dMPFC area which is involved in the tracking of accumulated consistency. Thus while the prediction error aspects of our model based analysis can be captured by a more traditional analysis, our use of computational modelling has  have additional value in giving greater insight into the mechanism behind such activations and in allowing us to additionally understand which brain areas are involved in the representation of the current predictions regarding the others consistency in similarity.

## S3. dmPFC literature comparison

To place our findings of accumulated consistency being tracked within the dmPFC in context, we identified 15 studies that had found activation in the dMPFC and categorised the contrasts those activations came from into the following four categories: 1) Diagnostic > Non-Diagnostic cases: contrasting information relevant to a trait judgement about an individual with irrelevant information; 2) Inconsistent > Consistent: cases contrasting novel information that was inconsistent with previous knowledge about an individual with novel information that was consistent with previous knowledge; 3) Other Impression Formation: other contrasts relevant to impression formation, often linking photographs of individuals with information about their traits; 4) Self Relevant: contrasts in which participants judged whether information was self-relevant or not (see Table S2). We then collapsed the clusters found in these studies across the x-axis to create Figure S1. Our result falls in the middle of the region activated by previous studies, with most contrasts investigating Inconsistent > Consistent falling more dorsally and most contrasts investigating Diagnostic > Non-Diagnostic falling more ventrally suggesting that the activation found in our study is in agreement with the previous literature.
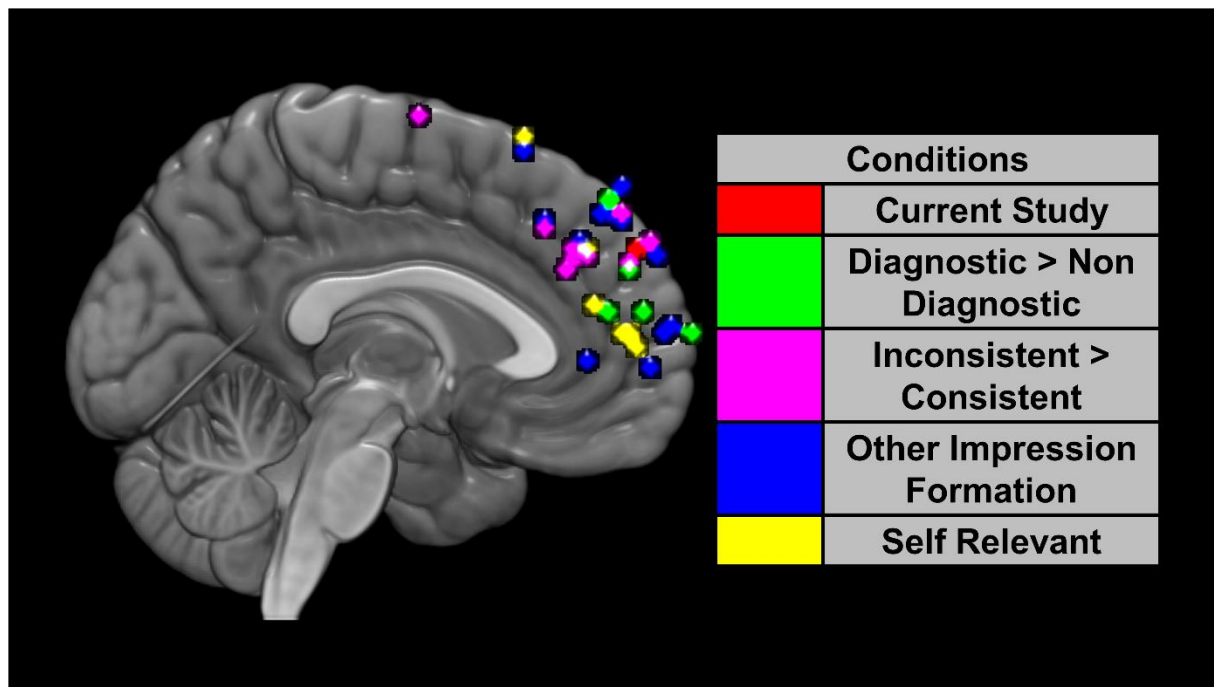
*Figure S1. dmPFC activations across 15 studies investigating impression formation and the self together with the result from the current study. For ease of presentation we have collapsed the studies across the x-axis (max x = 14, min x = -13).*

**Table S2.** Details of studies used in the comparison of dmPFC activation across impression formation studies (see Figure S1). Coordinates are reported in MNI space and coordinates from studies using Talairach space were converted using the WFU PickAtlas version 2.4 (Maldjian, Laurienti, and Burdette 2004; Maldjian et al. 2003).

| Study | Information presented | Impression studied | X | Y | Z |
|---|---|---|---|---|---|
| **Current Study (Red)** | Choice of preferred painting | Similarity to Self | 8 | 5 | 34 |
| **Diagnostic > Non Diagnostic (Green)** | | | | | |
| Gilron & Gutchess, 2012 | Moral and Neutral Behaviours | Moral Impression vs Location of Behaviour | 0 | 57 | 15 |
| | Moral and Neutral Behaviours | Moral Impression vs Location of Behaviour | 3 | 48 | 48 |
| | Moral and Neutral Behaviours | Moral Impression vs Location of Behaviour | 9 | 72 | 9 |

| | | | | | |
|---|---|---|---|---|---|
| Ma, Vandekerckhove, Van Overwalle, Seurinck, & Fias, 2010 | Moral and Neutral Behaviours | Personality Traits | 4 | 54 | 28 |
| | Moral and Neutral Behaviours | Personality Traits | 14 | 48 | 16 |

**Inconsistent > Consistent (Pink)**

| | | | | | |
|---|---|---|---|---|---|
| Cloutier, Gabrieli, Young, & Ambady, 2011 | Political Views | Political Affiliation | 2 | 54 | 30 |
| | Political Views | Political Affiliation | 6 | 52 | 44 |
| Ma et al., 2012 | Moral vs Neutral Behaviours | Moral Traits | 4 | 42 | 32 |
| | Moral vs Neutral Behaviours | Moral Traits | 4 | 38 | 34 |
| | Moral vs Neutral Behaviours | Moral Traits | 4 | 35 | 28 |
| Mende-Siedlecki, Cai, & Todorov, 2013 | Moral Behaviours | Moral Traits | 2 | 31 | 40 |
| Mende-Siedlecki & Todorov, 2016 | Moral and Neutral Behaviours | Trustworthiness and Surprise | -5 | 59 | 36 |
| | Moral and Neutral Behaviours | Trustworthiness and Surprise | -2 | -5 | 72 |

**Other Impression Formation (Blue)**

| | | | | | |
|---|---|---|---|---|---|
| Ames & Fiske, 2013 | Assessment of Teaching Ability | Expertise | -4 | 45 | 45 |
| Baron, Gobbini, Engell, & Todorov, 2011 | Moral Behaviours | Trustworthiness | -4 | 41 | 36 |
| Freeman, Schiller, Rule, & Ambady, 2010 | Individuated vs Superficial for Racial In-group vs Out-group | Personality Traits | -13 | 43 | 4 |
| Fouragnan et al., 2013 | Choices in Trust Games and Prior information about Trustworthiness | Trustworthiness | -2 | 64 | 10 |
| | Choices in Trust Games and Prior information about Trustworthiness | Trustworthiness | 0 | 62 | 31 |

| | | | | | |
|---|---|---|---|---|---|
| Gilron & Gutchess, 2012 | Moral Behaviours vs Neutral Behaviours | Moral Impession vs Location of Behaviour | 3 | 30 | 42 |
| Mende-Siedlecki, Baron, & Todorov, 2013 | Moral and Ability Behaviours | Competence and Trusworthiness | -5 | 66 | 13 |
| | Moral and Ability Behaviours | Competence and Trusworthiness | 32 | 60 | -1 |
| Mende-Siedlecki, Cai, et al., 2013 | Moral Behaviours | Moral Impression | 5 | 52 | 51 |
| Schiller, Freeman, Mitchell, Uleman, & Phelps, 2009 | Moral Behaviours | Moral Impression | -9 | 24 | 61 |
| Schiller et al., 2009 | Moral Behaviours | Moral Impression | -7 | 52 | 43 |
| **Self (Yellow)** | | | | | |
| Martinelli et al. 2013 | Memory Meta-Analysis | Semantic Autobiographic Memory | -10 | 45 | 18 |
| | Memory Meta-Analysis | Episodic Autobiographic Memory | -6 | 51 | 9 |
| | Memory Meta-Analysis | Conceptual Self | 6 | 55 | 5 |
| Moran, J.M. et al., 2006 | Personality Traits | Self-Relevance | -6 | 53 | 6 |
| Phan, K.L. et al., 2004 | Emotional Pictures | Self-Relatedness | 0 | 42 | 33 |
| Schneider et al. 2008 | Emotional or Neutral Pictures | Self-Relatedness | -3 | 24 | 66 |

# **Supplementary References**

Ashburner, J., & Friston, K. J. (2005). Unified segmentation. *NeuroImage*, *26*(3), 839–851.

Ames, Daniel L, and Susan T Fiske. 2013. "Outcome Dependency Alters the Neural Substrates of Impression Formation." *NeuroImage* 83: 599–608.

Baron, Sean G, Maria Ida Gobbini, Andrew D. Engell, and Alexander Todorov. 2011. "Amygdala and Dorsomedial Prefrontal Cortex Responses to Appearance-Based and Behavior-Based Person Impressions." *Social Cognitive and Affective Neuroscience* 6 (5): 572–81.

Brett, Matthew, Jean-Luc L Anton, Romain Valabregue, and Jean-Baptiste Poline. 2002. "Region of Interest Analysis Using an SPM Toolbox." *NeuroImage* 16 (2.1): 1140.

Catani, M., & de Schotten, M. T. (2012). *Atlas of human brain connections*. Oxford: Oxford University Press.

Cloutier, Jasmin, J D E Gabrieli, D O Young, and Nalini Ambady. 2011. "An fMRI Study of Violations of Social Expectations: When People Are Not Who We Expect Them to Be." *NeuroImage* 57 (2): 583–88.

de Schotten, M. T., Ffytche, D. H., Bizzi, A., Dell'Acqua, F., Allin, M., Walshe, M., … Catani, M. (2011). Atlasing location, asymmetry and inter-subject variability of white matter tracts in the human brain with MR diffusion tractography. *NeuroImage*, *54*(1), 49–59.

Forman, S. D., Cohen, J. D., Fitzgerald, M., Eddy, W. F., Mintun, M. A., & Noll, D. C. (1995). Improved assessment of significant activation in functional magnetic resonance imaging (fMRI): Use of a cluster-size threshold. *Magnetic Resonance in Medicine*, *33*(5), 636–647.

Fouragnan, Elsa, Gabriele Chierchia, Susanne Greiner, Remi Neveu, Paolo Avesani, and Giorgio Coricelli. 2013. "Reputational Priors Magnify Striatal Responses to Violations of Trust." *Journal of Neuroscience* 33 (8): 3602–11.

Freeman, Jonathan B, Daniela Schiller, Nicholas O Rule, and Nalini Ambady. 2010. "The Neural Origins of Superficial and Individuated Judgments about Ingroup and Outgroup Members." *Human Brain Mapping* 31 (1): 150–59.

Gilron, Roee, and Angela H Gutchess. 2012. "Remembering First Impressions: Effects of Intentionality and Diagnosticity on Subsequent Memory." *Cognitive, Affective and Behavioral Neuroscience* 12 (1): 85–98.

Joliot, M., Jobard, G., Naveau, M., Delcroix, N., Petit, L., Zago, L., … Tzourio-Mazoyer, N. (2015). AICHA: An atlas of intrinsic connectivity of homotopic areas. *Journal of Neuroscience Methods*, *254*, 46–59.

Lundqvist, D., Flykt, A., & Öhman, A. (1998). The Karolinska directed emotional faces. Stockholm: Department of Clinical Neuroscience, Psychology section, Karolinska Institutet.

Ma, Ning, Marie Vandekerckhove, Kris Baetens, Frank Van Overwalle, Ruth Seurinck, and Wim Fias. 2012. "Inconsistencies in Spontaneous and Intentional Trait Inferences." *Social*

*Cognitive and Affective Neuroscience* 7 (8): 937–50.

Ma, Ning, Marie Vandekerckhove, Frank Van Overwalle, Ruth Seurinck, and Wim Fias. 2010. "Spontaneous and Intentional Trait Inferences Recruit a Common Mentalizing Network to a Different Degree: Spontaneous Inferences Activate Only Its Core Areas." *Social Neuroscience* 6 (2): 123–38.

Maldjian, Joseph A, Paul J Laurienti, and Jonathan H Burdette. 2004. "Precentral Gyrus Discrepancy in Electronic Versions of the Talairach Atlas." *NeuroImage* 21 (1): 450–55.

Maldjian, Joseph A, Paul J Laurienti, Robert A Kraft, and Jonathan H Burdette. 2003. "An Automated Method for Neuroanatomic and Cytoarchitectonic Atlas-Based Interrogation of fMRI Data Sets." *NeuroImage* 19 (3): 1233–39.

Martinelli, Pénélope, Marco Sperduti, and Pascale Piolino. 2013. "Neural Substrates of the Self-Memory System: New Insights from a Meta-Analysis." *Human Brain Mapping* 34 (7): 1515–29.

Mathworks. 2015. "Matlab R2015b." Natick, MA: The Mathworks Inc.

Mende-Siedlecki, Peter, Sean G Baron, and Alexander Todorov. 2013. "Diagnostic Value Underlies Asymmetric Updating of Impressions in the Morality and Ability Domains." *Journal of Neuroscience* 33 (50): 19406–15.

Mende-Siedlecki, Peter, Yang Cai, and Alexander Todorov. 2013. "The Neural Dynamics of Updating Person Impressions." *Social Cognitive and Affective Neuroscience* 8 (6): 623–31.

Mende-Siedlecki, Peter, and Alexander Todorov. 2016. "Neural Dissociations between Meaningful and Mere Inconsistency in Impression Updating." *Social Cognitive and Affective Neuroscience* 11 (9): 1489–1500.

Moran, Joseph M, C Neil Macrae, Todd F Heatherton, C L Wyland, and W M Kelley. 2006. "Neuroanatomical Evidence for Distinct Cognitive and Affective Components of Self." *Journal of Cognitive Neuroscience* 18 (9): 1586–94.

Phan, K Luan, Stephan F. Taylor, Robert C. Welsh, Shao Hsuan Ho, Jennifer C. Britton, and Israel Liberzon. 2004. "Neural Correlates of Individual Ratings of Emotional Salience: A Trial-Related fMRI Study." *NeuroImage* 21 (2): 768–80.

Schiller, Daniela, Jonathan B Freeman, Jason P Mitchell, James S Uleman, and Elizabeth A Phelps. 2009. "A Neural Mechanism of First Impressions." *Nature Neuroscience* 12 (4): 508–14.

Schneider, Felix, F. Bermpohl, A. Heinzel, M. Rotte, M. Walter, C. Tempelmann, C. Wiebking, H. Dobrowolny, H. J. Heinze, and G. Northoff. 2008. "The Resting Brain and Our Self: Self-Relatedness Modulates Resting State Neural Activity in Cortical Midline Structures." *Neuroscience* 157 (1): 120–31.