University of Glasgow

Arabadzhiyska, Desislava Hristova (2022) *Spatiotemporal neural characterization of social decision making*. PhD thesis.

https://theses.gla.ac.uk/83053/

# Spatiotemporal neural characterization of social decision making

**Desislava Hristova Arabadzhiyska**

MA Psychology, MSc Research Methods of Psychological Sciences

Submitted in fulfilment of the requirements for the Degree of
Doctor of Philosophy

Institute of Neuroscience and Psychology
College of Medical, Veterinary and Life Sciences
University of Glasgow

March 2022

# Author's declaration

This thesis contains work conducted at the University of Glasgow under the supervision of Prof. Marios G. Philiastides, School of Psychology and Neuroscience between October 2018 and March 2022. I hereby declare that except where stated, the work included in this thesis is my own and no part of it has been submitted to any other university or degree.

# Acknowledgments

I would first like to thank my supervisor Prof. Marios Philiastides for the opportunity to be a part his lab and to learn so much over the course of this PhD. I would also like to thank him for his patience and for challenging me to be the best scientist I can be.

I would also like to express my gratitude to Frances Crabbe for her fMRI help and expertise and to Dr. Gemma Learmonth for her tDCS guidance and for her encouragements.

My lab members have also been an invaluable part of my PhD experience. Thank you Kitti, Ralie, Sean, Joana, Filippo, Bharti, Gabby, Belen and Tarryn for the interesting scientific (and also non-scientific) discussions. And I would especially like to thank Andrea - without the opportunity to get involved in his experiment and without his enthusiasm I don't know if I would have done a PhD in social decision making.

I also owe special thanks to Selyem Adam, who has been an endless pool of support - you have been my rock throughout this PhD journey. Thank you for believing in me, especially when I did not believe in myself. And thanks for all the cat videos!

И не на последно място бих искала да благодаря на родителите си. Без вашата подкрепа и разбиране не бих стигнала до тук. Благодаря!

# Contents

# List of Tables

# List of Figures

# List of Publications

**Arabadzhiyska**, D.H., Garrod, O.G.B., Fouragnan, E., De Luca, E., Schyns, P.G. and Philiastides M.G. (submitted) A common neural currency account for social and non-social decisions.

**Arabadzhiyska**, D.H., Learmonth, G. and Philiastides, M.G. (preregistration) Investigating the role of pMFC in social decision making using tDCS. (https://osf.io/y4jsk)

# List of Abbreviations

ACC = anterior cingulate cortex

DDM = drift diffusion model

dlPFC = dorsolateral prefrontal cortex

dmPFC = dorsomedial prefrontal cortex

EA = evidence accumulation

EEG = electroencephalography

FFA = fusiform face area

fMRI = functional magnetic resonance imaging

GLM = genralized linear model

GLMM = genralized linear mixed-effects model

PCC = posterior cingulate cortex

pMFC = posterior medial frontal cortex

PPI = psychophysiological interaction

RT = reaction time

SMA = supplementary motor area

tDCS = transcranial direct current stimulation

TPJ = temporoparietal junction

vmPFC = ventromedial prefrontal cortex

# Abstract

Some of the most important decisions we make over the course of our lives are social in origin. Whether it is to decide to start a business with someone or choosing a retiring home for a family member, these social decisions often have uncertain outcomes. Recent research has begun to elucidate the neurocomputational principles underlying social choices, however many questions about how we process and navigate our uncertain social environments remain. In this thesis, we examine the spatiotemporal neural characteristics of decisions based on social information on an algorithmic (i.e. what mechanistic processes are involved) and an implementational (i.e. which brain structures are involved) level and we assess whether social choices are a part of the same decision-making framework developed to describe non-social decisions.

We present three experiments - a behavioural pilot study, a simultaneous electroencephalography and functional magnetic resonance imaging (EEG-fMRI) experiment and a transcranial direct current stimulation (tDCS) study. We also outline an experimental paradigm based on an economic game, which attempts to ensure the fair comparison between social and non-social choices by manipulating the likelihood of a favourable outcome depending on the decision domain - indirect facial trustworthiness for social decisions and explicit reward probability ranges for non-social choices. As a result, we observed that social decisions display similar behavioural trends to the ones typically seen across the non-social decision-making literature. More importantly, however, we found that a drift diffusion model (DDM), which assumes the accumulation of relevant evidence to an internal boundary (i.e. evidence accumulation process), was able to account for these behavioural patterns and made comparable predictions across both domains, suggesting that social choices may employ similar algorithmic considerations to the ones at the basis of non-social decisions.

To study the implementational specificities of social and non-social choices, we identified neural signatures of evidence accumulation (EA) in our EEG data and used them to discern the neural site that gives rise to these dynamics. This allowed us to implicate the posterior medial-frontal cortex (pMFC) as the potential site for EA for both social and non-social decisions. We also found that the so-

cial and non-social information were initially encoded in distinct regions and that the pMFC clusters co-varied in a task-dependent way with areas of the human valuation system. Taken together, these results suggest that early representations of the two types of uncertainty are encoded in domain-specific areas and then compared in a common human valuation system. Afterwards, the comparison information is accumulated for the decision in the pMFC, thus showing the embodied nature of the choice as this region is adjacent to the relevant motor cortex. We also attempted to examine the mechanistic role of the pMFC in social choices even further in a pre-registered tDCS experiment. Although we were not able to collect our intended sample size due to slow recruitment caused by the COVID-19 pandemic, we present preliminary results to illustrate the types of findings that could be produced by this investigation. Our hierarchical DDM comparison suggested that the pMFC might modulate the rate of EA in addition to determining the amount of evidence necessary for a decision, however this notion was not supported by our formal statistical analysis and we thus conclude that more evidence is needed to establish the details of how social choices are implemented in the pMFC.

Overall, this thesis offers detailed insight into the spatiotemporal neural processes involved in social decision making. It further provides support in favour of a universal decision-making architecture and presents an example for the systematic comparison across decision-making domains.

# 1 Chapter 1, General Introduction

## 1.1 A brief overview of the thesis

Our lives are often defined by social information. We rely on our social circles to survive, we use social cues to understand who might pose a threat to us and we adapt our actions to match our social surroundings (Ruff and Fehr, 2014). However, this social information in not always perfect and how we interpret it may on occasion be erroneous (Lee, 2008). Although social uncertainty plays an important role in our lives, there are still many unanswered questions relating to the underlying neurocomputational mechanisms of how we process and navigate our uncertain social environments. In contrast, the investigation of the basic neural characteristics of how we process non-social sources of information has been the cornerstone of the non-social decision-making literature.

One potential reason for this discrepancy lies in the fact that the examination of social decisions has historically focused on the uniquely social factors that might influence these choices. Another possibility comes from the difficulty of defining what social information is, since it may refer to information that is social in origin, information that arises from social interaction or it may even refer to information that requires the understanding and the consideration of social norms (Ruff and Fehr, 2014). Finally, it has also been suggested that neural processes underlying social choices are too difficult to examine since they involve amounts of uncertainty that are far greater than the uncertainty involved in non-social decisions (Lee, 2008). Regardless of the reason for the largely discrepant investigations of social and non-social choices in terms of their underlying neurocomputational processes, recent accounts (Lockwood et al., 2020; Ruff and Fehr, 2014) have showcased the need for a detailed comparison between social and non-social choices to determine whether social decisions take part in a universal, domain-general 'common-value' schema or if they constitute a separate domain-specific process (i.e. the 'social-specific' schema, Ruff and Fehr, 2014). Such a comparison needs to span several levels in order to determine whether social choices are domain-specific. Specifically, it requires evidence demonstrating that it employs different mechanistic computations (i.e. algorithmic level) or that it employs separate brain

regions to carry out these computations (implementational level, Lockwood et al., 2020).

The non-social decision-making literature has examined the neurocomputational principles of non-social choices across several domains (e.g. probabilistic, perceptual, value-based) and across the algorithmic and implementational levels. This has allowed the definition of a common decision-making process model (Rangel et al., 2008; Rangel and Clithero, 2014) involving the early value encoding of the relevant information, followed by the value comparison. Afterwards, the evidence is accumulated until a decision is reached and an action is selected. Finally, the outcomes of the choice are considered and the values for the alternatives are updated. The examination of this decision-making process has benefited from the use of computational modelling, which has provided a mechanistic account of the algorithmic level involved in non-social choices. Drift-diffusion models (DDMs), in which information proportional to the difference in value between the two options is sequentially accumulated to a decision bound, have been widely employed and have been able to account for many phenomena in the decision-making literature (Ratcliff, 1978; Ratcliff and McKoon, 2008). The neural plausibility of such models has been established through the use of electroencephalography (EEG) experiments showing a gradual evidence accumulation (EA) activity arising from a centroparietal cluster in a similar way as suggested by the DDM (e.g. Kelly and O'Connell, 2013; O'Connell et al., 2012). The biological validity of these models further means that they can be used to identify the brain sites and networks responsible for the implementation of these algorithmic processes. Additionally, the functional magnetic resonance imaging (fMRI) literature has been able to implicate a number of regions to the various stages of the decision-making process. Moreover, the combined use of modelling and fMRI has even elucidated the potentially embodied nature of choices, suggesting that decisions are implemented in the same sensorimotor regions responsible for carrying out the outcome selection (e.g. Gold and Shadlen, 2007).

In contrast, the social decision-making literature has identified social cues like facial trustworthiness, which may have a substantial role in determining the outcome of social choices (e.g. Rezlescu et al., 2012). The role of such cues has primarily been investigated in the context of economic games, which allow for

the investigation of the strategic decisions one makes, based on the contingencies between the players' actions (Lee, 2008). In recent years there have been attempts to use DDMs to explain these social phenomena and such models have been successful in accounting for social choices (e.g. Hutcherson, Bushong, et al., 2015). Similarly, recent neuroimaging experiments have been able to implicate a similar network of regions as being involved in social choices as the one responsible for non-social decisions (e.g. Fukuda et al., 2019). There have even been some studies attempting to directly compare social and non-social choices to show that their valuation processes might rely on the same structures in the prefrontal cortex (e.g. Behrens et al., 2008; Janowski et al., 2013). Therefore, the current literature points to the possibility that at least the valuation stages of some social decisions may be a part of a universal decision-making framework.



Figure 1: **Graphical representation of the experiments included in this thesis.** The black panels represent the three experiments included in this thesis - the behavioural pilot (**left**), the main EEG-fMRI experiment (**middle**) and the tDCS study (**right**). The top level corresponds to the sessions involved in each experiment, followed by the tasks involved in each session. Bellow them are the data types produced by each session/task and on the bottom row we present the analyses associated with each data type. The analyses are grouped into three sections and colour-coded based on which experimental chapter they appear in. The analyses in the first group (orange) will be presented in the first experimental chapter - Chapter 2; the second group (green) are presented in the second experimental chapter - Chapter 3 and the final group (blue) are presented in the third experimental chapter - Chapter 4.

However, many questions remain unanswered, specifically regarding whether such universal trends are preserved for social choices, which rely on a social cue and whether these universal trends expand beyond the valuation stages of the

decision-making process. To address some of these questions we aimed to carry out a detailed comparison between social and non-social decisions. The first step in achieving this involves an examination of their algorithmic computations. For such an examination, it is necessary to ensure that the only differences between the social and non-social choices lie in the source of the uncertainty. To this end, in **Chapter 2**, we used a social cue shown to lead to a wide range of behaviours - facial trustworthiness - and we varied it in the context of an economic game. Specifically, we used an indirect trustworthiness measure defined as the perceived likelihood of splitting an investment associated with a range of face displays. To enable the non-social uncertainty to vary along the same 0-1 range, we used reward probability ranges linked to the likelihood of a favourable outcome. In order to test the suitability of our paradigm, we ran a Pilot experiment (Fig. 1, left) and found that the two decision domains were associated with similar behavioural patterns previously seen throughout the non-social decision-making literature (e.g. Bogacz et al., 2006; Drugowitsch et al., 2012; Gold and Ding, 2013; Philiastides and Ratcliff, 2013; Philiastides et al., 2006; Pisauro et al., 2017; Uchida et al., 2006). These results highlighted the suitability of the paradigm to investigate the underlying algorithmic levels of social and non-social choices in our Main EEG-fMRI experiment (Fig. 1, middle). In the Main experiment, we once again found behavioural similarities between the two domains that would enable the comparison of their underlying computations. To this end, we used a DDM, which was able to account for the behavioural patterns and made largely comparable predictions across both decision types, consistent with the decision-making process model defined in the non-social domain. This chapter therefore demonstrated that our paradigm facilitates a direct comparison between social and non-social choices. It also showed the similarity in the underlying algorithmic mechanisms employed by the two decision types by suggesting that they both might employ accumulation-to-bound processes.

Having observed algorithmic similarities between social and non-social choices, we turned to their implementational characteristics. To this end, in **Chapter 3** we used simultaneously acquired EEG-fMRI data, which offer not only high temporal information captured by the EEG and spatial specificity offered by the fMRI, but their combination provides additional insight, which cannot be gained by either modality alone (Philiastides et al., 2021). Specifically, this approach allows us

to utilize the latent estimates of internal processing, which vary on a trial-by-trial basis, captured by the EEG data in order to inform the fMRI analysis and thus identify activations, which might not be elucidated through task-dependant or behavioural indices (e.g. Fouragnan et al., 2015). Initially, we performed a univariate analysis of the evokes EEG responses, by utilizing our DDM estimates to create model-predicted EA traces, which we found correlated with the EEG activity of electrodes from the centroparietal cluster previously implicated in EA dynamics. The electrodes associated with the highest correlations further showed difficulty- and speed-related modulations across the two domains, which were further consistent with the domain-general decision-making process (e.g. Philiastides et al., 2006). However, as the validity of the EEG estimates is integral for the identification of the source of the EA dynamics, we ultimately utilized the multidimensional nature of the EEG signal and used a multivariate approach to establish more robust representations of the EA signal. These representations allowed us to implicate the posterior medial frontal cortex (pMFC) as the site for EA for both the social and non-social decisions. These activations were not associated with any of the other task-driven predictors. The task-specific predictors instead highlighted that the two sources of uncertainty were initially encoded in distinct regions and we also observed commonly reported trends in the activations reflecting the difficulty of the trials and the associated value computations (Clithero and Rangel, 2014; Domenech et al., 2018; Grinband et al., 2008; Monosov, 2017; Philiastides and Sajda, 2007). To further investigate the implementational levels, we examined which areas of the brain were co-activated in a task-dependant fashion with the activations in the pMFC. We saw that for both the social and the non-social trials, the pMFC clusters co-varied with areas of the human valuation system (Clithero and Rangel, 2014; Domenech et al., 2018). Taken together, these results are once again consistent with the general decision-making process model, with early representations encoded in domain-specific areas, which are compared in the common human valuation system and then accumulated for decision in the pMFC, thus showing the embodied nature of the choice as this region is adjacent to the relevant motor cortex.

In **Chapter 4** we aimed to use anodal and cathodal transcranial direct current stimulation (tDCS, Fig. 1, right) in order to temporarily change (enhance and diminish, respectively) the excitability of the pMFC during **social** choices in

a pre-registered experiment. Specifically, this stimulation method allowed us to examine whether the pMFC is causally implicated in EA, thus further extending our implementational understanding of the process. It also allowed us to examine the algorithmic role this region has in EA by reconciling two conflicting views - that the pMFC is involved in modulating the amount of evidence necessary for a decision (e.g. Forstmann et al., 2008) and that it reflects the efficiency of the evidence readout, thus affecting the rate of EA (Pisauro et al., 2017). We also considered the possibility that the pMFC is involved in both EA processes. To examine the relative differences produced by the two stimulation types and study the functioning of the pMFC, we used three hierarchical DDMs (hDDMs), which make fewer assumptions about how the parameters are related to each other on a subject-specific and on a group level than traditional DDMs (Wiecki et al., 2013). The first hDDM allowed the parameter modulating the amount of information necessary for a decision to vary between the stimulation conditions; the second varied the rate of EA between the stimulation sessions; the third allowed both parameters to vary between the stimulation conditions. We were not able to collect the necessary sample size due to slow recruitment caused by the pandemic, however we present preliminary results to exemplify the type of results this experiment can produce. We found that the combined model fit the data best, which causally implicated the pMFC as an EA region and suggests that the pMFC has a complex algorithmic role, being responsible for both monitoring how much evidence is needed for a decision, as well as reflecting the rate at which this information is accumulated. These preliminary results correspond to recent accounts demonstrating the embodied nature of decision making involving (pre)motor structures (like the pMFC) used to express the relevant choice (Verdonck et al., 2021). Additionally, by examining the neurocomputational function of this region exclusively in the social domain, we highlight the potential of social decision-making paradigms to be used to examine the underlying basic principles of decision making. Nevertheless, as the formal comparisons for the model parameters across the stimulation conditions were not significant, we conclude that more evidence is needed to establish the role of the pMFC in social decision making.

Lastly, we consider some of the limitations of the studies presented in this thesis and propose how they can be used to guide future research. The limitations include the fact that a large portion of our results were reliant on the EEG sig-

nal, which due to its additive nature, may have contained information arising from areas, not associated with EA and thus may have introduced biases in the data. Nevertheless, our results were consistent with previous findings, however future studies should aim to test the robustness of our results. We also note that in this thesis we focused on one type of social choice - choices based on social cues - and thus it remains to be seen if the same neurocomputational principles would be involved in other social domains. Similarly, it should be noted that the conclusions we present only refer to choices made under time constraint and as such, further research needs to establish if these similarities persist in the absence of such constraints. Furthermore, our tDCS results were limited both by the fact that we did not reach our intended sample size as well as by the limited number of models we used in our model comparison. The models, however, were based on the existing competing theories about the functioning of the pMFC and could be used to reconcile these conflicting views. We also suggest that future research could investigate whether the similarities observed in this thesis extend to the processes involved in motivating social decisions, in learning from these decisions and how confidence affects social choices. We also suggest that our design, which attempts to put social and non-social choices on equal footing may serve as a blueprint for future investigations into the extent to which a universal decision-making process exists by investigating other types of decision making.

## 1.2   The main problem

From trivial decisions such as selecting the best-flavoured bag of crisps, to the more consequential decisions such as choosing who to marry, we rarely have complete information about our potential choices. Uncertainty is further added due to the fact that there might be variability in the association between the decision and the outcome (Lee, 2008). In other words, even if we make the same decision twice, this might lead to different outcomes. This makes the neural systems involved in the computation of uncertainty, as well as the ones that use this information to ultimately arrive at a decision, key players in the decision-making processes (Heekeren et al., 2008; Lee, 2008). However, the sources of uncertainty can be diverse. For instance, some decisions might be based on perceptual information (e.g. de-

termining whether you are seeing a car or a face on a foggy night, Ratcliff et al., 2009), on preference (e.g. deciding which clothing item to buy, Philiastides and Ratcliff, 2013) or on purely probabilistic factors (e.g. deciding whether to place a bet, Zaghloul et al., 2009). In other situations, people might rely on social cues to guide their choices (e.g. scrutinizing someone's trustworthiness when considering their advice, Ennew and Sekhon, 2007; Fouragnan et al., 2013). Even though both social and non-social factors can serve as cues for decision making, so far there have not been many direct comparisons between the two decision domains. Consequently, it is still largely unknown whether all the neurocomputational principles employed by decisions based on social information correspond to the ones guiding non-social choices. This leaves us with two possibilities for social decision making (Ruff and Fehr, 2014). The first one is that it employs different neural networks and processes than non-social decisions (Adolphs, 2010). The second, in accordance with standard economic utility models (Caplin and Glimcher, 2014; Von Neumann and Morgenstern, 2007), suggests that the neurocomputational rules governing both decision-making types should be the same.

## 1.3 Difficulty in comparing social and non-social choices

One of the main reasons why the degree of overlap between the neurocomputational principles of social decisions and non-social choices is unknown, comes from the fact that they have largely been studied in isolation, with the two types of decision making historically focusing of different objectives (Ruff and Fehr, 2014). Specifically, the non-social decision-making literature has primarily concerned itself with uncovering the processes behind the evaluation of the evidence available for certain options and choosing the best alternative (Rilling and Sanfey, 2011). The social domain, however, has prioritized the investigation of the neural underpinnings of uniquely social phenomena like altruism, cooperation and goods distributions between self and others (Ruff and Fehr, 2014). While in recent years research has begun to compare social and non-social decisions in terms of their algorithmic characteristics (e.g. Krajbich et al., 2015), their learning principles (Behrens et al., 2008; Lockwood and Klein-Flügge, 2021; Tarantola et al., 2017), as well as the brain areas they employ (Harris et al., 2007; Janowski et al., 2013;

**Figure 2:** **Graphical representation of the two main theoretical schemas.** The 'common currency' schema **a)** and the 'social valuation' schema **b)**. The 'common currency' schema posits the existence of a shared neural framework (shown in purple) responsible for the integration of all domain-specific information into a 'common currency' value. The activity in this shared circuit should therefore be the same for social (red) and non-social (blue) value computation, however the functional connectivity with other regions responsible for domain-specific processing should differ. The 'social valuation' schema, instead predicts the existence of domain-specific architectures. The valuation process involved in social value computation may follow the same principle as the one involved in non-social choices, however, they are implemented in separate brain structures. Adapted from Ruff and Fehr (2014).

Konovalov et al., 2021), many questions remain, such as whether any commonalities extend to choices where the social uncertainty is determined by a social cue as well as across all stages of the decision-making process. As a consequence, there are two options for the relationship between how we make decisions based on social and non-social uncertainty. The first option is described by the 'common currency' schema (Fig. 2, a), which suggests domain-specific early value encoding, followed by the conversion to a 'common currency' within the same domain-universal structures. The second option - the 'social valuation' schema - instead suggests that all computations involved in social decision making are made in domains-specific regions, even though they may rely on the same mechanistic principles (Fig. 2, b).

Another difficulty in establishing whether social choices are processed in the same manner as non-social ones stems from the fact there is no clear definition of what social decision making is. According to one view (Ruff and Fehr, 2014), there are three broad situations that qualify as social decision making. The first refers to cases where the choice is based on social cues, which determine how we value

others, their actions and how their actions affect us (e.g. deciding whether to ask someone on a date). In the second set of social decisions, the outcome of the decision affects someone other than ourselves (e.g. choosing which desert to order for your sibling). And the final decision sub-type refers to situations requiring the consideration of one's actions in the context of social norms (e.g. deciding how to respond to an unfair offer). Given this wide variety of phenomena that qualify under the umbrella term of social decision making it is often hard to draw a direct comparison between social and non-social choices due to some social scenarios not generalizing to non-social settings (e.g. learning from advice or from observing others).

Some have even argued that the role of uncertainty in social decision making is more variable since the behaviour of other individuals might change frequently and thus is not comparable to the one encountered in non-social decisions (Lee, 2008; Rilling et al., 2008; Sanfey, 2007; Suzuki and O'Doherty, 2020). Specifically, since these behaviours are based on multiple hidden states (e.g. the other's intentions, preferences, state of mind), they come with their own added uncertainty and the optimal solutions are determined with higher difficulty. This also implies that ultimately, the integrated value computation involved in social choices may be considerably more complex than the one involved in non-social choices. To this end it has been proposed (FeldmanHall and Shenhav, 2019) that social situations are prone to increased levels of uncertainty as these multiple uncertain factors can even build on each other. FeldmanHall and Shenhav (2019) further argue that unlike in non-social situations, this uncertainty can be interpreted as a beneficial feature of certain choices such as when exploring aspects of a social environment. However, they also argue that in some situations, resolving social uncertainty may also have negative consequences (e.g. by not considering the range of potential outcomes, people may start to develop stereotypical beliefs about different populations). While these points suggest that certain aspects of social decision making may be associated with domain-specific idiosyncrasies, it is also worthwhile to highlight that a subset of non-social choices - value-based choices, which involve internal assessments of items - may be associated with similarly complex considerations (Polanía et al., 2019). Specifically, it has been proposed that such choices are a product of a complex inference process, which takes into consideration how the values are structured in the environment. This process also attempts to maxi-

mize information in value representation to meet the demands imposed by limited coding resources (we receive a lot more information than we can handle). It is also conceivable that like in the social domain, uninformed reduction of uncertainty in value-based choices could also lead to negative outcomes. For instance, not considering the potential benefits or enjoyment of certain food groups may lead a person to only consume their preferred food items and thus not get proper nutrition. Examples like these therefore highlight that determining the degree of the potential specificity of social decision-making is not a trivial endeavour, however it is still noteworthy that some of the seemingly different aspects of social and non-social decisions may not be as dissimilar as previously thought, especially when we consider the appropriate domain-specific subtypes.

Nonetheless, the challenges highlighted above may have contributed to the fact that until recently not many studies have attempted to investigate the degree to which social decisions share their underlying neural computations with non-social choices. Nevertheless, frameworks have begun to be put forward to address this current knowledge gap. For instance, it has been proposed that social behaviours may be explained through the same motivation framework developed from the study of non-social behaviours (Contreras-Huerta et al., 2020). Similarly, it has been suggested that Marr's framework (Marr and Poggio, 1979) may offer an exhaustive approach towards determining whether the mechanisms involved in social choices are domain-specific (Lockwood et al., 2020). This framework suggests that the comparison could span three levels - the computational level responsible for the goal of the behaviour, the algorithmic level governing the mechanistic principles of the actions and the implementational level where social and non-social choices are compared based on the specific brain regions employed for behaviour. Therefore, in order to determine whether a process is domain-specific, it has to be different at either the implementational and/or algorithmic level i.e. it has to employ a different mechanism or it should involve separate brain sites. These frameworks also highlight a key aspect of such a comparison, which is that in order to robustly investigate the potential commonalities between social and non-social choices, proper controls need to be used to ensure that the two decision-making domains share all the same attributes apart from the source of the uncertainty. It is also highlighted that to adequately investigate the implementational level of social and non-social decisions, correlational methods demonstrating simple associations

are not enough and that stimulation methods are needed to examine whether the involvement of a certain region is necessary for the identified processes. Similarly careful considerations are needed to establish a potentially shared algorithmic level and it has been proposed that the most robust candidate for this assessment is the use of computational modelling (Lockwood et al., 2020).

## 1.4 Non-social decision making - approaches to studying non-social choices, computational models, neural correlates

As pointed out in the previous section, a crucial part of the investigation of the potential universality of the decision-making process is the selection of appropriate paradigms[1]. Some of the most commonly used paradigms to investigate the perceptual decision-making (PDM) field, where the uncertainty for the decision comes from external sources, include discriminating between noisy perceptual stimuli such as distorted images of cars and faces (e.g. Philiastides and Sajda, 2006; Ratcliff et al., 2009). The use of these paradigms has allowed researchers to showcase that neural responses scale with the amount of available evidence for either alternative. Other examples in this field include discriminating the direction of a moving cloud of dots and the difficulty of the trials is established by the percentage of dots moving in the same direction (e.g. Gherman and Philiastides, 2018). For value-based choices (VBDM), where the uncertainty is guided by internal assessments, a typical paradigm involves the choice between different food items, which have previously been assessed for their subject-specific preference (e.g. Pisauro et al., 2017). Examples for probabilistic decisions include displaying two card decks with different reward probabilities that are unknown to the participant and the participants' goal is to identify the more favourable one by drawing cards from the decks (Zaghloul et al., 2009). These paradigms have been invaluable in examining the algorithmic and implementational levels of Marr's framework in non-social decision making due to their use of simple stimuli that allow to ex-

---

[1]Note that it is beyond the scope of this section to exhaustively summarize all existing literature in the non-social decision-making domain. Instead, this section aims to provide a general understanding of some of the relevant approaches and findings related to the neurocomputational principles involved in decision making and to thus provide a benchmark for the discussion of social choices later on.

perimentally manipulate the available information to parametrically modulate the difficulty of the trials and to vary the uncertainty across clearly defined ranges.

These paradigms have also facilitated the examination of the mechanism underlying non-social decision making and have led to the development of a popular understanding of the three main stages involved in these processes (Fig. 3, Rangel et al., 2008; Rangel and Clithero, 2014) : 1) the assignment of value to the choice alternatives, 2) the option comparison, EA and choice selection and 3) the value updating following the revealed outcome. To illustrate these stages, we will consider examples from the PDM and VBDM domains, which make use of two different sources of uncertainty - externally and internally-generated, respectively - to highlight the similarities across the two domains. These two examples are presented in the bottom two rows of Fig: 3 and showcase that regardless of the source of uncertainty i.e. whether the person has to properly identify if a distorted image shows a house or a face, or if they need to select whether they would enjoy an apple more than an orange, value is assigned to the alternative to reflect the evidence in favour of the given options. Following the initial valuation stage, in both examples, the alternatives are compared and the comparison information is accumulated until an internal boundary (corresponding to one of the two choice options) is reached and a decision and the associated action are made. Finally, once the outcome is known, the values associated with the two options are reevaluated, so as to be used in future choices. These examples therefore highlight that the main differences between these two domains (and it has been suggested that this is the case with other non-social domains, e.g. van Vugt et al., 2019) is the initial valuation stage and the information is subsequently used in the same way as described above. These processes and the details of their algorithmic characteristics have been elucidated primarily through the use of computational models, which have given rise to mechanistic accounts of non-social decision making at the different levels of processing. These models allow the prediction of outcomes from a set of options through algebraic calculations, which transform the option characteristics and produce a ranking of the possible choices based on their attractiveness (Johnson and Ratcliff, 2014).

Some of the most widely-used computational models have been process models, which have been useful in providing an explanation for the underlying neurocompu-

**Figure 3:** **Decision-making process model and examples in the perceptual (PDM) and value-based (VBDM) domains.** Decision making broadly involves processes spanning three levels (Rangel et al., 2008). 1) First, value is assigned to the amount of information favouring an alternative. In the PDM example, where uncertainty is externally-generated, that corresponds to the evidence that a distorted image represents a house (blue bar) or a face (orange bar). In the VBDM example, where uncertainty is internally-generated, value is assigned to the apple (blue bar) and the orange (orange bar) in correspondence with subjective preference. 2) Then, these values are integrated and compared. Over time the comparison information is accumulated (green line) until a decision is reached after crossing one of the boundaries mapped to the two alternatives and an action is selected. In the PDM example, this comparison determines that it is more likely that the distorted image is a house than a face and the house is selected. In the VBDM example the person values the apple more highly than the orange and therefore chooses the apple. 3) Finally, the values for the two options are updated based on the outcome. These values later inform future value assignment stages. In the PDM example, the image was correctly identified as a house and the value of the house increases (purple bar), whereas the one for the face decreases (yellow bar). Similarly, if the person enjoyed their apple, they would update positively how they value apples (purple bar) and decrease their preference for oranges (yellow bar).

tational characteristics of two-alternative forced-choice tasks by taking into consideration reaction times (RTs) and selected choices, thus offering detailed mechanistic accounts of the processes involved in decision making, described above (Fig. 3). They are especially popular since they make predictions about the underlying neural computations and can be applied in both human neuroimaging studies and animal work. These include drift diffusion models (DDMs), where information is sequentially accumulated to a decision bound and the time taken for the accumu-

lation is reflected in the RTs (Ratcliff, 1978; Ratcliff and McKoon, 2008). At each time point an accumulator moves towards one of two alternatives, mapped to the two possible outcomes. The magnitude of this shift is modulated by the amount of evidence processed at that moment in favour of either alternative. Once the accumulator crosses a criterion bound, a decision is made (Bogacz et al., 2006). Even though the most optimal strategy would be to always choose an option if the evidence in its favour is higher than for the alternative, the process is accompanied by noise and stochasticity leading to occasional errors (Rangel and Clithero, 2014). This is often accounted for in the model with the inclusion of noise parameters. Another feature of these models is the inclusion of non-decision components to account for information encoding, memory access and response-related movements. One benefit of DDMs is that they can also reflect certain biases. For instance, when participants are told to be either fast or accurate, small adjustments to the model parameters such as the boundary can capture the relevant differences (Bogacz, Hu, et al., 2010). Furthermore, the parameters estimated through DDMs have been able to account for how certain cognitive impairments occur as well as how they change over time with age (Brosnan et al., 2020).

Nevertheless, although they have strong mechanistic explanatory power, DDMs do not directly inform about how the proposed computational processes might be implemented in the brain. And such a comparison with neural systems is vital for the validation and creation of more robust models. The assessment of the biological validity of DDMs has necessitated the use of various tools to probe the spatiotemporal implementation of the decision-making processes. Specifically, these include tools with good temporal and spatial resolution such as electroencephalogram (EEG) and functional magnetic resonance imaging (fMRI), respectively, as well as stimulation methods such as transcranial magnetic stimulation and transcranial electrical stimulation, which can be further divided into transcranial direct current stimulation (tDCS) and transcranial alternating current stimulation. For instance, the EEG findings in this area, which benefit from high temporal resolution on a millisecond scale, have further elucidated the algorithmic levels of non-social decision-making temporal processing by suggesting that the amount of evidence available on a given trial guides the accumulation rate, with low uncertainty leading to faster accumulation (e.g. Philiastides and Sajda, 2006; Ratcliff et al., 2009). EEG experiments have also shown that at response, the neural activity,

regardless of the domain of the decision, reaches a common boundary (Kelly and O'Connell, 2013; O'Connell et al., 2012; Polanía et al., 2014). Thus, these studies provide evidence in favour of the neural plausibility of the DDMs, which describe a similar accumulation-to bound mechanism to the one observed across these EEG studies. This accumulating activity has also been consistently shown to arise from a centroparietal cluster in a number of non-social decision-making subtypes: in value-based (Pisauro et al., 2017; Polanía et al., 2014), perceptual (Gherman and Philiastides, 2015; Kelly and O'Connell, 2013) and even memory-based decisions (van Vugt et al., 2019).

Conversely, the synergy between modelling and fMRI data has been even more apparent since the predictions from a model can be subsequently correlated with the fMRI data in order to discover the neural network involved in the identified computations (Busemeyer and Diederich, 2014). fMRI studies have highlighted the consistency across the various non-social decision domains in the early evidence representation stages of the decision-making process by demonstrating that the activation in the relevant early visual processing areas (e.g. the parahippocampal place area for houses) correlates with the amount of available evidence in favor of either alternative (e.g. Heekeren et al., 2004; Philiastides et al., 2010; Philiastides and Sajda, 2007). Similarly universal trends have been found during the other stages of the decision-making process within a wide network of regions, some of which we will briefly presented here.

For instance, the activity in the dorsolateral prefrontal cortex (dlPFC) has been found to correlate with a comparator operation that reflects the difference between the evidence for the choice alternatives (Heekeren et al., 2004; Wan et al., 2015). The medial prefrontal cortext (mPFC) has been associated with many functions such as the processing of risk (Levy et al., 2010) and strategic choices (Coricelli and Nagel, 2009). However, most consistently, the activity in the ventromedial prefrontal cortex (vmPFC) has been shown to correlate with the probability of receiving a reward, but also with the absolute difference between the activity in the brain areas responsible for the processing of the two alternatives (e.g. Chib et al., 2009; Lim et al., 2011). This therefore, implicates the vmPFC as actively involved in the computation of value signals and shows that it does not simply reflect the consequences of the choice. This has been shown to be especially true

when that value is relevant to the decision (Grueschow et al., 2015). The striatum has been linked to the computations of risk (Levy et al., 2010), losses and gains (Tom et al., 2007), reward probability and reward magnitude (Yacubian et al., 2007) as well as in modulating learning (Fouragnan et al., 2015). The insula has been implicated in the assessment of risk (Levy et al., 2010; Mohr et al., 2010) as well as in the processing of feelings, empathy and uncertainty (Singer et al., 2009). Specifically, it has been argued that this area is responsible for integrating information into a general 'subjective feeling state' related to subject-specific levels of risk aversion and situational assessment (Singer et al., 2009). This region has further been associated with uncertainty across different contexts, including situations associated with threat and reward as well as choice selection and associative learning under uncertainty (Morriss et al., 2019). The anterior cingulate cortex (ACC) has been linked to the processing of value (Wan et al., 2015) and uncertainty (Monosov, 2017; Paulus and Frank, 2006) and it has been argued that it controls information-seeking strategies under uncertainty (Monosov, 2020). Another line of research has implicated the region in the processing of difficulty in foraging tasks (Shenhav et al., 2016; Shenhav et al., 2014). The parietal cortex has been linked to ambiguity (Bach et al., 2011; Huettel et al., 2006) and risk processing (Mohr et al., 2010) and has been shown to encode value (Wan et al., 2015). The posterior cingulate cortex (PCC), specifically, has also been implicated in being part of a common framework for value processing after the choice alternatives are initially valued (Levy et al., 2010).

The non-social decision-making literature has also investigated the potential embodied nature of choices, which puts forward the notion that the sensorimotor areas responsible for carrying out the response are also where the decisions are implemented (Gold and Shadlen, 2007). For instance, it has been proposed that when the response modality is unknown, evidence integration occurs in frontal areas, but once the link between a response and the decision outcome is established, the relevant response areas receive the integrated information (Filimon et al., 2013). Knowing this pairing in advance, however does not improve the integration strategies during the course of the decision-making process (Tsetsos et al., 2015). This may also be why frontal areas have been implicated in computing choice value and comparison signals, whereas motor areas have been found to reflect the value of the action and provide EA readout (Klein-Flügge and Bestmann, 2012; Wunderlich

et al., 2009). However, in most real-life situation the association between the decision and the outcome are known and it has been suggested that the involvement of the relevant motor system is engaged prior to the accumulation region reaching its internal threshold (McBride et al., 2018). This is why it has been proposed that there may be a second motor-preparation accumulation process, which receives continuous input from the original accumulator until it reaches its own threshold (Verdonck et al., 2021). Furthermore, the use of models that employ the assumption that the action preparation is part of the decision-making process, leads to better fits and tends to explain the speed-accuracy trade-off better than models, which do not make these assumptions (Lepora and Pezzulo, 2015). Recent evidence has also identified (pre)motor structures, such as the posterior medial frontal cortex (pMFC) as the site for EA for value-based decision (Pisauro et al., 2017), by associating its activity to the neural signal arising from a centroparietal cluster, implicated in such accumulation-to-bound processing (e.g. O'Connell et al., 2012). Similarly, there has been evidence linking the centroparietal activity with a dorsal frontoparietal network (dFPN, which includes the pMFC) in perceptual decision making (Brosnan et al., 2020) through the use of DDMs, EEG and diffusion and resting-state fMRI. This experiment also showed that the build-up rate (or drift rate) is associated with the connectivity between the structures of the dFPN and with the speed of decisions. The higher the build-up, the stronger the connection and the faster the decisions. Finally, the animal literature has also corroborated this view by showing that there are brain structures in the midbrain (Ding and Gold, 2010), premotor regions (Cisek and Kalaska, 2005) and within the dFPN (Ding and Gold, 2012) that display this type of EA activity. Changing the rules of the decision-making task leads to altered activity of several of those regions. The timescales of their activation also suggests the necessity for communication among these regions for successful choice selection and thus further demonstrates the involvement of such (pre)motor areas in decision making.

As suggested by the previous paragraphs, the non-social decision-making literature has not only defined a process model describing how decisions are made, but has also managed to highlight the potential universality of the process across the non-social sutypes. The strongest evidence in favour of this notion comes from studies, which have directly compared how decisions are made across the non-social domains. For instance, it has been shown that perceptual decisions based

on somatosensory and visual information follow the same EA patterns (von Lautz et al., 2019). Further similarities have been drawn between PDM and VBDM, for which the literature has argued that they both necessitate stimulus disambiguation as well as item valuation across a common decision framework (Summerfield and Tsetsos, 2012). The direct attempts to compare PDM and VBDM have shown that the two decisions follow a similar accumulation-to-bound process (Pisauro et al., 2017; Polanía et al., 2014). Even though the VBDM has also been associated with accumulation patterns in the frontal region as well as with a fronto-parietal synchronization (absent in the PDM, Polanía et al., 2014), the parietal gamma frequency oscillations displayed comparable evidence of accumulation for both decision-making types. This suggests that even though there might be some differences with VBDM involving the activation of additional brain regions, VBDM and PDM might share the same underlying mechanism (Polanía et al., 2014). Taken together the above-mentioned examples suggest that the non-social decision-making literature has been able to find similarities between various subtypes of non-social decision making. Consequently, this has helped identify a common mechanism involved in neural computation of non-social evidence which relies on an accumulation-to-bound processes.

In summary, the non-social decision-making literature has thoroughly investigated both the algorithmic and the implementational levels of the processes involved in non-social choice selection. This has been made possible through the use of carefully selected paradigms, which allow to manipulate crucial features related to the decision-making process, while avoiding potential confounds. This field has largely benefited from the use of a variant of sequential sampling computational models - DDMs, which describe an accumulation-to-bound process as the basis for decision making. These models have demonstrated that they are consistent with observed neural signatures found in EEG data as well as with activations seen in fMRI data. The parameters produced by these models have also been instrumental in implicating a large network of cortical and subcortical regions in the process of decision making and have allowed to demonstrate the embodied nature of this process. Finally, these observations have been commonly found across non-social decision-making domains pointing to the potential universality of the underlying decision-making mechanism.

## 1.5 Social decision making - approaches to studying social choices and their neural correlates

Historically, the literature investigating social decision making has primarily focused on specific processes leading to social behaviour and less on the basic decision-making processes (Ruff and Fehr, 2014). Consequently, the decision-making aspect has often been used as a tool to investigate social behaviour. Nevertheless, the social decision-making literature has managed to highlight factors such as trustworthiness that can influence social decisions (e.g. Rezlescu et al., 2012; Wilson and Rule, 2016). Such experiments have uncovered that social choices often depend on unmalleable features such as facial structure and that similar biases are present in cross-cultural settings (e.g. Xu et al., 2012) and clinical populations (Hooper et al., 2019). It has also been suggested that these trustworthiness associations are not coincidental and have been shown to correlate with hereditary facial features (Lee et al., 2017). Furthermore, facial features (Stirrat and Perrett, 2010) and expressions (Reed et al., 2012) have been shown to predict altruistic behaviours.

A common approach to studying social decisions has been to use game theory, which examines the effects of incentives on decisions in strategic environments (Houser and McCabe, 2014; Lee, 2008; Sanfey, 2007). The strategic element refers to the observation that one player's choice influences the other player's opportunities and payoffs, and a crucial part of these games is the fact that both parties need to be aware of these dependencies. This allows to study the relationship between individual decisions and group-level outcomes, which are linked through clearly defined mechanisms. These mechanisms refer to the notion of a 'game tree', which specifies who moves when, which movements are allowed, what information is available to them at different stages and how other's moves influence the joint outcome.

Furthermore, game theory offers simple yet sophisticated tasks since they are easy to implement and for participants to understand, but offer the investigation of multiple social scenarios and social interaction (Rilling and Sanfey, 2011). They have been useful in the investigation of many aspects of social decision making such as reciprocal coordination, fairness, altruism/punishment and reciprocal ex-

**a**
**TRUST GAME**

Player 1 (Investor)

Player 1 (Investor)

1 point

Keep

Play

1 point

Player 2 (Trustee)

4 points

Keep

Split

Player 1 (Investor)

0 points

Player 2 (Trustee)

4 points

Player 1 (Investor)

2 points

Player 2 (Trustee)

2 points

**b**
**PRISONER'S DILEMMA**

Player 1 cooperates
Player 2 cooperates

Player 1   Player 2

1 year     1 year

Player 1 defects
Player 2 cooperates

Player 1   Player 2

0 years    5 years

Player 1 cooperates
Player 2 defects

Player 1   Player 2

5 years    0 years

Player 1 defects
Player 2 defects

Player 1   Player 2

3 years    3 years

**c**
**ULTIMATUM GAME**

Player 1 (Proposer)

2 points
Fair offer

4 points

1 point
Unfair offer

Player 2 (Responder)

Accept      Reject

Player 1 (Proposer)

2 points

Player 2 (Responder)

2 points

Player 1 (Proposer)

0 points

Player 2 (Responder)

0 points

Player 2 (Responder)

Accept      Reject

Player 1 (Proposer)

3 points

Player 2 (Responder)

1 point

Player 1 (Proposer)

0 points

Player 2 (Responder)

0 points

Figure 4: **Common economic games. a)** Trust game: An Investor (Player 1) is allocated a certain sum (e.g. 1 point), which they can 'Keep' or pass on to the Trustee (Player 2) for the chance to double their reward ('Play'). If the Investor chooses to 'Keep', the game ends. If they decided to 'Play', the sum is transferred to the Trustee and multiplied (e.g. by 4). Then the Trustee decides whether to 'Keep' the augmented sum or to 'Split' it equally between the two players. **b)** Prisoner's Dilemma: Two prisoners (players) are offered a bargain. They can either stay silent (i.e cooperate with the co-player) or they can betray the co-player (defect). If both cooperate (leftmost panel), they will both have to serve 1 year in prison. If one of them cooperates, but the other defects (middle two panels), the defector does not have to serve any jail time, but the other has to serve 5 years in prison. Finally, if both defect (rightmost panel), they will each have to serve 3 years in prison. **c)** Ultimatim game: A Proposer (Player 1) is allocated a sum (e.g. 4 points), which they have to distribute between themselves and the Responder (Player 2). The proposer can make a 'Fair' (each receives 2 points) or 'Unfair' offer (they keep 3 points and the Responder receives 1 point). The Responder can then 'Accept' the offer and the two players receive the amounts suggested by the Proposer, or they can 'Reject' the offer and neither player receives anything.

change, which have been investigated with Trust games and Prisoner's Dilemma tasks (Fig. 4, a and b). A typical Trust game (Fig. 4, a, Berg et al., 1995) involves two individuals – an Investor (Player 1) and a Trustee (Player 2). The Investor is allocated a certain endowment (e.g. 1 point), which they can either keep or offer to the Trustee ('Keep' and 'Play' options, respectively). If the endowment is offered to the Trustee, it is augmented (e.g. quadrupled). The Trustee then has the option to keep the augmented endowment (e.g. 4 points, 'Keep' option) or to split it with the Investor (e.g. each receiving 2 points, 'Split' option). A Prisoner's Dilemma task (Fig. 4, b, Rapoport et al., 1965) involves two individuals, each presented with the same two options. They can either collaborate or defect. The collaboration of both players leads to the overall best outcome for the pair, whereas the defection of both leads to the worst outcome for the pair. However, if one defects and the other collaborates, the one who defects receives the best individual outcome, whereas the one who collaborated, receives the worst individual outcome. Finally, another commonly used economic game is the Ultimatum game (Fig. 4, c, Güth et al., 1982), which also involves a pair of players - a Proposer (Player 1) and a Responder (Player 2). The proposer is allocated an endowment and has to choose how to split the endowment between the two players. They can make two types of offers - fair offers and unfair offers. Once the offer has been made, the Responder can accept the offer, in which case the two players receive the suggested amounts proposed by the Proposer. However, if the Responder rejects the proposal, neither player receives anything.

These economic games can be divided based on the information available to the player. For instance, in perfect information games each person knows everything about the payoff structure (Kuhn, 1950), whereas in incomplete information games, like in a Trust game (Harsanyi, 1967, 1968), players do not have complete information about their opponents/co-players. The optimal strategies in economic games can usually be defined by a form of Nash equilibrium where each player is getting as much as they can, given the actions of the other players (Houser and McCabe, 2014). In Trust games, the Nash equilibrium is for Player 1 to always 'Keep', since if they do not, they allow Player 2 to 'Keep'. However, this optimal behaviour is often not observed and both players often choose the other possibilities (Berg et al., 1995; Camerer, 2003; Sanfey, 2007). Such tendencies towards cooperation may arise because of evolutionary pressures of guilt as the decision whether or not

to 'Keep' as Player 2 often depends on how guilty they feel (Houser and McCabe, 2014). However, the guilt-state of Player 2 is unknown to Player 1. Consequently, there is only one subgame perfect equilibrium (which treats each subnode as the start of a separate game, each associated with its own Nash equilibrium), which is for Player 1 to always 'Play' if they believe that at least half of the population are high-guilt individuals. Therefore, Trust games are not perfect information games since there are 2 possibilities from the nature of Player 2.

More recently, such economic games have been combined with some of the methodologies used by non-social decision-making investigations, including the use of computational models (developed in non-social contexts), which allow the investigation of trial-by-trial changes as well as predicting future behaviour while investigating the algorithmic underpinnings of social choices (e.g. Hackel and Amodio, 2018). This has led to the suggestion that a DDM can account for altruistic choices as the benefits for one's self and those for others are weighed and compared to guide a decision (Hutcherson, Bushong, et al., 2015; Tusche and Bas, 2021). Research has also shown that DDMs can account for biases in speeded and time-delayed social decisions (Chen and Krajbich, 2018). Specifically, prosocial participants were more likely to make pro-social decisions under time pressures and would be more selfish when allowed to deliberate, while the opposite trend was found for more competitive participants. Similar models have also been used for decisions concerning others, such as when participants had to decide if they wanted to delegate a choice between risky and safe options to themselves or a group (Edelson et al., 2018). Specifically, it was shown that participants tended to avoid responsibility and this tendency was attributed to a second-order process, whereupon people wanted to have high certainty when the outcome could affect many other individuals. However, in an even more notable example it was shown that a DDM using estimates from a food-choice game could predict social decisions made during economic games (Krajbich et al., 2015). This therefore shows the potential for a shared underlying mechanism between social and non-social value computations. However, because the experiment did not use the same paradigms for the social and the non-social condition, it is still unclear whether these observations would be preserved following a direct comparison of the two sources of uncertainty.

Interestingly, there has also been evidence to suggest that on the implementational level, the underlying principles driving social choices may be similar to the ones at the basis of non-social ones (Lee, 2008). For instance, the same networks often associated with reward processing - striatum, insula, orbitofrontal cortex are also activated during social judgments. Specifically, the striatum, which has been linked to reward prediction error changes, correlates positively with cooperation and negatively with the lack of reciprocation. Additionally, monetary and social rewards have been found to produce the same type of activation within the striatum (Izuma et al., 2008; Williams et al., 2020). Furthermore, the striatum's activity has been shown to track rewards relative to the winnings of the opponent. Therefore its role depends on whether the behaviour is thought to be competitive or cooperative. The anterior insula on the other hand, which has been linked to risk and uncertainty, has been found to respond to unfair offers during an Ultimatum game (Harlé et al., 2012). Similarly, there is evidence to suggest that the same network composed of the insula, dmPFC and the PCC, which is involved in the processing of risk with uncertain probabilities, is also active during strategic social choices (Nagel et al., 2018).

There has also been evidence to suggest that mPFC is related to social processing. For instance, higher activations in the vmPFC, which have been linked to value processing, have also been associated with trusting the co-player in a Trust game and has similarly been shown to track monetary rewards gained on the basis of trust (Rilling and Sanfey, 2011). This region has also been shown to reflect unfair offers proposed during an Ultimatum game (Grecucci et al., 2013) as part of an intuitive, fast-processing network responsible for norm violation monitoring, whereas more deliberate, long-term decisions about personal gains are processed in the dorsal ACC (Feng et al., 2015). The vmPFC has also been linked to social and affective processing as well as situational value (Lieberman et al., 2019) and self-referential processing (Gusnard et al., 2001). Perhaps even more importantly, the literature on decisions made for one's self vs the ones made for others has argued that similar value computations are made in the vmPFC for non-social choices (i.e. self choices) and for social decisions (i.e. other decisions, e.g. Janowski et al., 2013). Similar observations have been made in the context of learning on one's own vs learning based on observing others (Behrens et al., 2008) implicating a conversion of domain-specific information in the vmPFC. Furthermore, the poste-

rior subsections of the mPFC have also been related to personal traits, physical characteristics and feelings (Jenkins and Mitchell, 2011). The dorsomedial prefrontal cortex (dmPFC) has similarly been associated with social (Lieberman et al., 2019) and self-referential information (Gusnard et al., 2001) and the activity in the anteromedial prefrontal cortex has been shown to relate to self and affective processes (Lieberman et al., 2019).

The potential commonalities between social and non-social choices can further be highlighted by the literature comparing neurotypical populations with cohorts who face challenges in processing social stimuli. For instance, it has been shown that in individuals with Autism, both social and non-social stimuli are associated with decreased activity in ventral prefrontal cortex (Shafritz et al., 2015). However, others have reported opposing EEG patterns arising between individuals with autism and neurotypical participants for social decisions, but the same dynamics for non-social stimuli (Gonzalez-Gadea et al., 2016). Therefore, even though there is some indication for the possible commonalities between social and non-social decision making, it is still unknown whether these common trends persist across all types of social choices and whether the similarities are preserved for the decision-making stages beyond the 'common value' conversion.

The social decision-making literature has also proposed regions, which may be linked to the processing of uniquely social phenomena, especially during early stages of value encoding. It has even been suggested that for decisions involved in an Ultimatum game, value is computed in a 3-stage process, whereby the co-player's bonus is encoded in the right TPJ and left dlPFC, followed by the computation of an effective value (i.e. the effective influence of the offer) in the right anterior insula and the final decision value is ultimately computed in the mPFC (Fukuda et al., 2019). Moreover, the amygdala has also been heavily implicated in trust decisions based on facial trustworthiness, whereas the orbitofrontal cortex has been shown to be involved in cooperation, and the TPJ has been associated with the assessment of the co-player's trustworthiness (Rilling and Sanfey, 2011). It has further been suggested that under negative affect, the activity of the TPJ during trust decisions is decreased and its connectivity to areas responsible for emotional processing such as the amygdala are disrupted (Engelmann et al., 2019). It has also been shown that the suppression of the activity in the right TPJ leads to a decrease

in the ability to infer the actions of others (Hill et al., 2017). Furthermore, there is evidence to suggest that consensus decisions are based on considerations of personal preference, group majority preference and an estimate of the persistence of each choice, which are respectively encoded in the vmPFC, TPJ and intraparietal sulcus, with the signal from these areas being integrated in the dorsal cingulate cortex (Suzuki et al., 2015). The amygdala has also been implicated in tracking the differences in money allocations between self and others and this effect has been shown to be modulated by the administration of oxytocin (Liu et al., 2019). Finally, the orbitofrontal grey matter volume has been associated with conformity tendencies (Campbell-Meiklejohn et al., 2012).

In summary, the social decision-making literature has historically focused on identifying the processes associated with uniquely social phenomena, but has recently attempted to address the potentially shared neurocomputational mechanisms across decision domains. This field has benefited from the use of economic games to extensively examine social behaviours, which have highlighted the central role of facial features in driving trustworthiness assessments. The recent use of computation models has demonstrated some of the potential common processes involved in social and non-social choices and these commonalities have even been supported by studies exploring the implementational stages of social decision making. These studies have also demonstrated potentially domain-specific activations and processing, which may correspond to the early information processing stages involved in non-social decision making. However, in order to determine whether the two types of decision making are truly based on a universal decision-making process, there is a need to also consider other types of social choices (such as ones based on social cues) and directly compare them to non-social decisions on both their algorithmic and implementational levels and across their value processing and EA stages, while attempting to ensure that their uncertainty varies across comparable scales.

## 1.6   Goals of the thesis

The literature presented above demonstrates that while recent efforts have begun to investigate if the spatiotemporal neural characteristics employed in social

decision making correspond to the ones involved in non-social choices (Behrens et al., 2008; Harris et al., 2007; Janowski et al., 2013; Konovalov et al., 2021; Krajbich et al., 2015; Tarantola et al., 2017), many questions remain unanswered. The goal of this thesis therefore is to try to examine the neurocomputational mechanisms of choices based on social uncertainty beyond the value-computation stage and to assess whether those might be a part of a universal decision-making framework. Here we specifically examine social decisions, which are based on social information - facial trustworthiness. Conversely, the non-social choices in this investigation will not involve an uncertainty manipulation based on a social component, but will instead be based on probabilistic assessments. Specifically, the social and the non-social decisions will be evaluated in the context of an economic game (i.e. the Trust game), where the outcome of the social choices will reflect subject-specific estimates of the likelihood of splitting an endowment for a set of face displays. For the non-social choices it will instead be determined by a range of reward probabilities chosen from the full probability range, thus confining the uncertainty for social and non-social decisions to the same 0 to 1 probability scale.

Recent studies have highlighted the benefits of employing a multimodal approach by means of simultaneous EEG-fMRI to study the neurocomputational characteristics of non-social decision making (e.g. Pisauro et al., 2017). This approach allows us to benefit from the high temporal resolution provided by the EEG measure and the spatial information provided by the fMRI, while gaining additional information to what each modality can offer on its own (Philiastides et al., 2021). Specifically, the combination of these imaging methods allows to exploit the internally generated trial-by-trial temporal dynamics captured by the EEG signal and use them to identify associated activations within the fMRI data, which in turn offers a more precise understanding of the spatial distribution of the activity (Sajda et al., 2011). Such EEG measures reflect latent states, which allow the identification of activations, typically unidentifiable with standalone fMRI measures informed only by the task-specific manipulations or behavioural indices (e.g. Fouragnan et al., 2015). Similarly, while modeling parameters have also been used to reflect such internally-generated states, they may be subject to misspecification, and very few studied have attempted to use sampling models that provide more than a single estimate of a mean and a standard deviation (Gluth et al., 2017; Turner et al., 2015). Therefore, the variability occurring on a by-trial basis is better captured through

the EEG data and offers stronger explanatory power than alternative measures. We used this multimodal approach while attempting to equalize the uncertainty manipulations across the social and non-domain presented in the context of a Trust game. This allowed us to identify whether the two decision domains were associated with 'common currency' computations in addition to universal evidence-to-bound processing. This study will be presented over two chapters (**Chapter 2 and Chapter 3**), with the first aiming to investigate the algorithmic characteristics of the two decision domains by examining the behavioural data and computational modelling results. The second chapter will focus on two types of EEG analyses, an EEG-informed fMRI examination and an interaction analysis, aiming to elucidate the implementational characteristics of social and non-social choices.

Specifically, the first experimental chapter (**Chapter 2**) examines whether the universal decision-making mechanism extends to social choices. It highlights the behavioural similarities across the two domains and showcases that in economic contexts both social and non-social decisions are modulated by reward probabilities and that reaction times reflect the difficulty of the choices. We were able to observe these trends across both a Pilot experiment (Fig. 1, left) as well as from the Main EEG-fMRI experiment (Fig. 1, middle). Considering that the evidence-to-bound process has been strongly implicated across various non-social decision-making domains (e.g. Polanía et al., 2014; Summerfield and Tsetsos, 2012), as explained in the previous sections of this chapter, we examined whether the algorithmic characteristic of social choices follow a similar trend. We observed that a DDM was equally well-suited to explaining the behaviour in a social and a non-social context and that the parameter estimates produced for the two decision domains were comparable, thus providing evidence in support of the notion that the two might share the same algorithmic characteristics.

Having observed similarities between the social and non-social algorithmic processes, we turned to examining the implementational characteristics of these choices in the second experimental chapter (**Chapter 3**). Initially, we used our modelling results to create model-based predictions of the EA dynamics in the two domains, which we found matched the activity captured by electrodes from centroparietal cluster previously implicated in EA processing (Gherman and Philiastides, 2015; Kelly and O'Connell, 2013; Pisauro et al., 2017; Polanía et al.,

2014; van Vugt et al., 2019) and we saw that the model-EA closely corresponded to the actual EEG signal, which reflected previously described difficulty and speed-related trends (e.g. Connolly et al., 2009; Law and Gold, 2008; Mazurek et al., 2003; Roitman and Shadlen, 2002). However, since we would use the EEG signal to identify the relevant neural site of EA, we needed to ensure the robustness of the signal. To this end, we analysed the EEG data with a multivariate technique aimed at maximizing the contribution of the relevant electrodes that capture the process of EA. Having identified this robust measure of EA, we used it to establish which brain areas corresponded to the trial-by-trial fluctuations in EA captured in the EEG signal. The variability captured in the EEG for both social and non-social sources of uncertainty corresponded to activity in the pMFC, which has previously been implicated in non-social EA (Pisauro et al., 2017) as well as in modulating the amount of evidence necessary for a decision (Forstmann et al., 2008). We also found that across the two domains, the activity in this region co-varied in a task-specific manner with regions of the human valuation system. Finally, we also observed domain-specific activations at the time of stimulus presentation. To-gether, these observations are consistent with a domain-general decision-making process similar to the one described in Fig. 3, where domain-specific information is initially encoded in early domain-relevant brain areas, then compared and con-verted into a 'common currency' in the human valuation system and accumulated for decision in regions adjacent to the motor areas responsible for carrying out the relevant choice action.

Having identified a candidate brain region responsible for EA across decision-making domains, in the final experimental chapter (**Chapter 4**) we aimed to pro-vide causal evidence for its involvement in the EA dynamics through the use of stimulation methods in a pre-registered experiment (Fig. 1, right). Specifically, we used tDCS to alter the functioning of the region and to establish the mechanistic function it might have in the **social** EA process. By focusing on the social choices, we could potentially demonstrate the capability of the social decision-making do-main to provide further insight into the universal decision-making process. We combined our stimulation method with the use of a hierarchical DDM (hDDM) to ex-amine changes to the latent variables in the model space, as any changes produced by the stimulation itself may be too subtle to appear on the behavioural level. We made use of the hierarchical version of the DDM as it makes fewer assumptions

about the distribution of the parameters and therefore might be more sensitive to the small changes produced by the stimulation (Kruschke, 2010; Wabersich and Vandekerckhove, 2014; Wiecki et al., 2013). This allowed us to create models corresponding to the current hypotheses about the potential function of this region in the context of the evidence-to-bound process - that it reflects the dynamics of the EA process (Pisauro et al., 2017) or that it modulates the amount of evidence needed for a decision (Forstmann et al., 2008). We also considered the possibility that the pMFC is responsible for both processes. Due to especially low yield rates caused by the pandemic, we were not able to collect the full sample size specified in the pre-registration, which is why we are only able to provide preliminary results for this thesis. These results so far suggest that the role of this region may be to govern multiple facets of EA, where it reflects both the EA dynamics and modulates the amount of evidence needed to make a decision. However, as the formal analyses on the parameter estimates were not significant, we conclude that more evidence is needed to determine the role of the pMFC in social decision making.

# 2 Chapter 2, Experiment 1

## 2.1 Introduction

Imagine you are playing a game of poker. You can decide whether to place a bet based on your hand and the associated likelihood that someone else on the table may be holding a better combination of cards. Alternatively, you could make your decision based on the social information produced by the other people as they announce their choices. Both options come with some degree of uncertainty - depending on the cards in your hand, your chances of winning may vary from very favorable to highly unlikely. Similarly, your ability to correctly identify the social signals unwillingly produced by the other players may vary between the different players or from round to round. Currently, however, there are still many unanswered questions about whether the uncertainty associated with the social information is processed and used in the brain in the same way as the non-social uncertainty.

The first step in examining if such a universal decision-making process exists, is to study the algorithmic characteristics of the two domains. In other words, it is necessary to investigate whether the neural computations involved in social choices are the same as the ones employed by the non-social decisions (Lockwood et al., 2020). To this end, we used a DDM in order to examine whether social choices rely on accumulation-to-bound computations like the ones implicated in non-social decisions (e.g. Kelly and O'Connell, 2013; O'Connell et al., 2012; Philiastides and Sajda, 2006; Polanía et al., 2014; Ratcliff et al., 2009). However, before we investigate the similarities between the algorithmic processes involved in social and non-social choices, it is important to use a paradigm, which tries to ensure that the only difference between the two decision domains is the source of uncertainty i.e. social and non-social (Lockwood et al., 2020). One of the difficulties in examining how similar social choices are to non-social decisions comes from finding ways of putting the social and non-social uncertainty on equal footing i.e. making sure that they vary along a comparable scale. Additionally, it is important that this uncertainty is defined through objective measurements since measures like Likert scales make use of arbitrary ranges and values. A convenient option for

investigating the non-social domain is to make use of probabilistic decision making, as the probability of receiving a reward is bound to an objective and easily interpretable 0-1 range. Therefore, to enable the fair comparison across the social and non-social domain, a similarly objective social equivalent is needed.

One potential option is to use a social feature, which is known to lead to a wide range of perceptions and actions, and embed it in the context of an economic game to define rules for obtaining a probabilistic reward distribution modulated by this feature. One such candidate feature is facial trustworthiness, since it has been shown that facial features and expressions reliably predict behaviour in economic games (Mussel et al., 2013; Van't Wout and Sanfey, 2008). Furthermore, these trustworthiness behaviours are based on unmalleable facial features (Rezlescu et al., 2012) and are preserved even when people are only shown a co-player's picture and verbally assured that the co-player's behaviour is genuine (Tzieropoulos, 2013). Making such assurances is critical since the associated neural activity has been shown to depend on the belief that one is playing against an actual human. The importance of knowing whether these trustworthiness cues may be used to inform future behaviour is further demonstrated by the observation that when participants are allowed to choose their own avatars in a computer game, they choose the more trustworthy ones and in return receive higher offers, however if these avatars are allocated by someone else, the effects disappear (Tingley, 2014). These trustworthiness-based behaviours, however, are not incidental since facial expressions signaling enjoyment have been linked to cooperation, whereas expressions of contempt have been associated with defection in a Prisoner's Dilemma (Fig. 2, b, Reed et al., 2012). Similarly, it has been shown that men with wider faces are less likely to act cooperatively or favour mutual gain (Stirrat and Perrett, 2010). Additionally, people can reliably identify individuals who are less cooperative (Tognetti et al., 2013). There may be an evolutionary incentive behind these phenomena as it would be advantageous to recognize people's intentions and predispositions based on their expressions and facial features. This is also why people tend to remember untrustworthy people with trustworthy faces better than those with untrustworthy faces (Suzuki and Suga, 2010) and why such judgments can be made relatively quickly and reliably after only 100ms (De Neys et al., 2017) as well as why these trends can be observed even in young children (Ewing et al., 2015). Nevertheless, such judgments are not universal across cultures once one

controls for facial attractiveness (Xu et al., 2012). Despite these effects not being cross-cultural, there is nevertheless real-life evidence showing that trustworthiness judgments correlate with criminal sentence severity (Wilson and Rule, 2015, 2016).

These examples show that trustworthiness assessments are vital for social interactions and that they are governed by an automatic process, in which people engage regularly to inform their decisions. Consequently, we employed facial trustworthiness as our social feature used to modulate the social uncertainty in our comparison between social and non-social choices. Specifically, we used a Trust game, in which an Investor has to decide whether to pass their endowment on to a Trustee for the chance to increase their rewards, with the outcome of that decision depending on the likelihood that the Trustee would decide to cooperate so that both parties can benefit from the interaction (Fig. 4, a, Berg et al., 1995). This allowed us to express the facial trustworthiness as the likelihood associated with the Trustee splitting the investment, thus ensuring that the social and non-social uncertainty can be measured and modulated along a comparable 0-1 range.

It is noteworthy that previous examinations comparing social and non-social decision making have employed paradigms where the social aspect involves the consideration of another person and the non-social involves the consideration of one's self (e.g. Hutcherson, Bushong, et al., 2015; Janowski et al., 2013; Lockwood et al., 2017; Tusche and Bas, 2021), which may be considered as being a more unambiguously social manipulation (as the benefactor of the choice is not always the participant) and these experiments have also argued for social-specific processing (e.g. Lockwood et al., 2021). However, we highlight that the current paradigm aimed to investigate whether the neurocomputational principles observed in these previous investigations of social choices will still be present when using a different type of social choice - where the social information is carried by a social cue. We argue that it is important to also consider such alternative types of social decisions in order to investigate the degree to which social choices are domain-specific. We further highlight that this type of social cue was chosen in an attempt to directly and fairly match the processes involved across the two decision types, as by closely matching the two domains, we believe that any differences would be more confidently attributed to the domain-specific nature of the decision and not

to differences in the types of considerations that form the basis of the decisions. Specifically, it is conceivable that in self (non-social) vs other (social) considerations, the social choice may involve the assessment of the non-social aspects of the decision (i.e. how the decision maker would appreciate the outcome of the choice) in conjunction with second-order social considerations (how well these personal considerations map on to the other individual e.g. Edelson et al., 2018) and therefore may not necessarily represent a one-to-one comparison between the two domains. We further note that later in this thesis we report that the social condition was associated with activations typically attributed to social cognition, which would not have been the case if our paradigm did not successfully employ social judgments. Moreover, the behavioural patterns which we will discuss in this chapter are unlikely to have been evident if the social manipulations included in the images was not successfully engaging the relevant social judgments. To this end we believe that while the task outlined here may not be as overtly social as in some previous investigations, it still employs social processing and would be able to offer further insight into the degree of potential overlap between social and non-social decision making.

We initially ran a behavioural pilot experiment (henceforth Pilot, Fig. 1, left panel) in order to validate our task and stimuli. Having done so, we were then able to investigate whether social and non-social decisions use the same neurocomputational principles in a simultaneous EEG-fMRI study (henceforth Main experiment, Fig. 1, middle panel), which consisted of 1) a behavioural session, 2) an online rating session and 3) an EEG-fMRI session. The behavioural sessions of the Main experiment were used to identify participants, who would engage in the task and were thus not formally analysed, apart from using the ratings to inform stimulus selection. Similarly, the online rating task was only used for stimulus selection. In the Pilot and all three sessions of the Main experiment we would show our participants a series of face displays, which they were told belonged to people, who had previously taken part in a Trust game (Fig. 4, a) as Trustees and we asked our participants to indicate the likelihood of splitting the investment associated with each face display. We told them that the people whose faces they were viewing actually participated in this game to ensure that our participants would later make use of their trustworthiness judgments. In the Pilot and all sessions of the Main experiment apart from the online rating session, we would then ask our partici-

pants to actually take part in a Trust game, where they were assigned the roles of Investors and were presented with both social and non-social trials, in which they had to make the choice between a small, but safe reward and a larger but uncertain payoff. In the non-social trials the outcome of the uncertain option would be indicated by a range of reward probabilities. The reward probabilities in the social trials would depend on the face displays and participants were told that the outcome would reflect the decisions made by that individual, whereas in reality the outcome was determined by the indirect trustworthiness rating.

These tasks allowed us to demonstrate the behavioural similarities between the social and the non-social domains, which highlighted that our paradigm enabled the social and non-social choices to be placed on equal footing, which facilitated the examination of their algorithmic characteristics. We were also able to fit the data to a DDM and showcase that the model fit the data well and reflected the main trends seen in the behavioural data and that both decision types follow similar evidence-to-bound computations, providing evidence in favour of them sharing the same algorithmic principles.

## 2.2   Materials and Methods

### 2.2.1   Participants

We recruited 41 participants for the Pilot through the University of Glasgow subject pool. Since facial perception may depend on one's race and racial history (e.g. Scott and Monesson, 2009), participants were chosen to be Caucasians, aged 18-35 to match the available face stimuli (see below). All participants had normal or corrected-to-normal vision and reported no history of psychiatric, neurological or major medical problems, and were free of psychoactive medications at the time of the study. 18 participants were removed (total n = 23) due to poor performance consisting of either chance performance or displaying the same response pattern regardless of the reward probability manipulation.

40 participants took part in the three sessions of the Main experiment involving simulataneous EEG-fMRI (Fig. 1, middle). They were selected based on the same

criteria outlined for the Pilot. Two participants were removed from all analyses due to poor behaviour during the EEG-fMRI session (one had near chance performance across all levels of reward probability in the social context, whereas the other had chosen to nearly always 'Play' across all levels of reward probability in the non-social context). For the multivariate EEG analysis, which aimed to increase the contribution of the relevant sensor information (outlined in the next chapter) seven participants were removed (total n = 31) due to poor (chance) discrimination performance. All statistics presented in relation to the Main experiment reflect the sample size from the multivariate EEG analysis. We note that originally for the univariate EEG analysis on evoked responses on individual sensors we only removed 3 additional participants (total n = 35) due to discrepancies between the DDM predictions and the behaviour, however, for consistency between the analyses, here we are presenting the dataset based on the sample identified for the multivariate analysis and we note that the statistics for the univariate EEG analysis across the two samples were consistent with each other. The study was approved by the College of Science and Engineering Ethics Committee at the University of Glasgow (300180147) and informed consent was obtained from all participants.

### 2.2.2  Stimulus creation and modulation

For the Pilot as well as the initial behavioural session (i.e. the first session in the Main EEG-fMRI experiment) we used 100 photorealistic face identities (56 female and 44 male images), which were selected from a database collected at the Institute of Neuroscience and Psychology at the University of Glasgow. The identities were chosen to be Caucasian, aged 18-35. A reverse correlation (Ahumada Jr and Lovell, 1971) based on previous ratings of the trustworthiness of the faces was used to identify features associated with higher trustworthiness scores. Afterwards, these facial features were manipulated in all faces to create different trustworthiness versions of each face using a Generative Face Grammar (Yu et al., 2012). The collection of these images, as well as the reverse correlation procedure were not done by the author of this thesis. Instead, they were collected and produced by researchers from the Schyns lab at the University of Glasgow. This facial manipulation was done in order to produce stimuli that would be likely to fall into a wide range of different trustworthiness categories i.e. highly trustworthy,

Figure 5: **Sample stimuli from a representative participant. Top row:** Social stimuli at five different subject-specific indirect trustworthiness levels, corresponding to the reward probability levels employed for the non-social stimuli. **Bottom row:** Non-social stimuli with five explicit reward probability levels superimposed on a neutral face. The photo-realistic face images were created using the procedure described in (Gill et al., 2014) and presented in Stimulus creation and modulation.



Figure 6: **Example of stimulus trustworthiness manipulation at 4 levels (trustworthiness levels from left to right: 1, 7, 14, 20).** Level 1 trustworthiness represents the lowest trustworthiness level, level 20 corresponds to the highest trustworthiness level.

highly untrustworthy and ones considered neither trustworthy nor untrustworthy. 20 trustworthiness versions were created for each face, ranging from the least (1) to the most (20) trustworthy version of the face (Fig. 6). Only one version per face identity was chosen. The face stimuli were allocated to 5 trustworthiness bands based on their trustworthiness level, which corresponded to the likelihood of split-

ting the investment given a 'Play' choice (0-0.2, 0.2-0.4, 0.4-0.6, 0.6-0.8, 0.8-1). All images were visually inspected to ensure they convey the trustworthiness level associated with the respective bands. For the non-social condition 1 neutral face, which was not included in the social condition, was used and 5 distorted images containing reward probability ranges (as text) were displayed on top of the face (Fig. 5, bottom row). These reward probability ranges were created using the R package magick (version 2.0; R Core Team, 2018). The probability ranges were blurred by convolving the images with a Laplacian with a width of 19 pixels. After-wards, Poisson noise was added and light-gray pixels were removed. The non-social face stimuli were displayed in the same manner as the social ones.

For the rating task as part of the online rating session and the EEG-fMRI session (i.e. the second two sessions of the Main experiment), we instead used 150 pho-torealistic face images based on 131 face identities (made up of 61 male images and 70 female images). 19 of the original face identity images were additionally manipulated to create extra identities, so as to increase the image sample size, which would eventually be used during the economic game as part of the EEG-fMRI session in the Main experiment. The same trustworthiness manipulation was applied to the fabricated identities. Only one version per face identity (original and fabricated) was chosen. Half of the additional 50 faces were chosen from the highly trustworthy versions of the face identities and the other half - from the highly untrustworthy versions.

We used Presentation software (Neurobehavioral Systems Inc., Albany, CA) to present the face displays as 300x400 pixel images centrally on a screen (frame rate = 60 Hz) for the rating and economic game task in both the Pilot and the behavioural session of the Main experiment. They were also presented in the same way during the economic game as part of the EEG-fMRI session in the Main experiment. However, during the Pilot and the behavioural session in the Main experiment, the stimuli were shown on a monitor 80 cm away from the participant, whereas during the EEG-fMRI session of the Main experiment, the stimuli were presented centrally via an LCD projector on a screen placed at the rear opening of the bore of the MRI scanner and viewed through a mirror mounted on the head coil (distance to screen = 95 cm). For the rating task during the online rating session and the EEG-fMRI session of the Main experiment, the face displays were presented centrally on the

screen as 3x4 images that took up 30% of the height of the screen. They were presented using the online functionality of PsychoPy3 (v3.0.7, Peirce, 2007).

### 2.2.3  Procedure for Pilot experiment



Figure 7:  **A variant of a Trust Game and procedure for EEG-fMRI session. a)** The participant (Investor) receives 1 point, which they can either 'Keep' or 'Play' for the chance of winning 2 points. If they decide to 'Play', the point is quadrupled and passed on to a Trustee, which can be social (red) or non-social (blue). The Trustee can either keep all 4 points ('Keep') or split them evenly with the participant (each receiving 2 points, 'Split'). In the social trials the probability of winning is based on the subject-specific trustworthiness score associated with the current face display, while in the non-social trials it is determined by the reward probability range displayed on a face, neutral for trustworthiness. **b)** Social (S; red outline) and non-social (NS; blue-outline) trials. Trials begin with a variable fixation cross screen (1-4 s). Following the fixation screen a stimulus is presented for up to 1.3 s, during which participants indicate their choice ('Play' or 'Keep'). Following the response, a fixation cross screen appears for the remainder of the 1.3 s.

The pilot experiment only involved a behavioural session, which was made up of two tasks - a rating task (based on 100 face identities) and an economic game (Fig. 1, left). Both of these tasks were based on a version of a Trust game (Fig. 4, a). Specifically, this variant involved one-shot Trust games between two players - Investor (Player 1) and Trustee (Player 2). The Investor is allocated 1 point per

trial and has two options. The first offers a small, but a certain reward by selecting to keep the point for that trial, thus ending the trial ('Keep' option). The second option presents a bigger, but more uncertain payoff ('Play' option). Specifically, if the Investor chooses this option, they pass the point to the Trustee. At this stage, the point is quadrupled, however it is up to the Trustee to determine how the four points should be distributed. The Trustee is also presented with two options: they either keep all four points for themselves, or they can split them evenly between the two players.

For the rating task, participants were told that the face displays belong to individuals who have previously taken part in an economic game (i.e. a Trust game like the one described above) and who were assigned the roles of Trustees (i.e. Player 2) in that economic game. To make our cover story more realistic for our participants, we took pictures of their own faces and told them that their face displays and responses would be used for similar experiments in the future. In reality, we would delete the pictures after each session. The goal of the participants during the rating task was to assess the face identities' trustworthiness by estimating the likelihood of splitting the investment for each face display (on a 0-1 scale). They were also told that they would eventually play with these face displays and were thus encouraged to use the rating task as an opportunity to familiarize themselves with the faces. When the task was presented to the participants all mentions of trustworthiness were avoided in order to sidestep the possibility of participants developing unusual strategies in the game due to social desirability biases. Furthermore, by framing the rating task in such a manner, we ensured that their subsequent decisions in the economic game would be based on the same economic considerations as the ones that they would employ for the non-social choices. This was a key feature of the design since embedding the rating in the context of a Trust game avoids the use of arbitrary explicit trustworthiness ratings (e.g. using Likert scales) and ensures a direct mapping between social and non-social choices. Furthermore, it has been suggested that such indirect trustworthiness measures may yield more ecologically valid estimates than explicit ratings (Uleman and Kressel, 2013). During each trial, a red cursor would move (from left to right) along a scale representing the likelihood of splitting the investment (Fig. 8, a). The cursor moved continuously and looped around the scale. The speed of the cursor was determined by the refresh rate of the monitor (60Hz) and moved 2 pixels to

the left every 16.6 milliseconds. The initial position of the cursor changed on a random basis. This was done in order to avoid any anchoring effects when the participants gave their responses. Responses were given by pressing a button on a Cedrus RB-740 button box (Cedrus Corporation, 2019) when the position of the cursor reflected participants' estimate of the likelihood that the person whose face was displayed would split the investment.



Figure 8: **The rating task in the Pilot and the Main experiment**. In both versions of the rating task, participants rated the face identities based on the perceived likelihood of splitting an investment in a Trust game. **a)** In the Pilot and the behavioural session of the Main experiment participants were asked to press a button once a continuously moving cursor was at the position reflecting their rating. **b)** The rating task in the online rating session and the EEG-fMRI session of the Main experiment required participants to indicate their ratings by clicking on the corresponding section on the scale.

Following the rating task, the participants would take part in the previously described game as the Investor (i.e. Player 1) in two conditions (social and non-social), each consisting of 250 trials (500 in total) presented in an interleaved fashion. On each trial, participants would see a variable fixation cross (1.25s - 1.75s). Then they would be presented with a stimulus from either the social or the non-social condition and had to choose between the 'Keep' and the 'Play' option, however the outcome of the trial depended on the condition (Fig. 7, a). In the social condition, participants were told that the probability of doubling the

reward was determined by a randomly selected response from when the person whose face display they were currently viewing played the game. In reality, the reward probability was governed by the trustworthiness estimate each participant gave during the rating task (divided into five different levels of reward probabilities (given a 'Play' choice); 0–0.2, 0.2–0.4, 0.4-0.6, 0.6–0.8 and 0.8–1). For the non-social condition, the distorted text reflected the 5 reward probability levels given a 'Play' choice. The 5 bins in the two domains (social and non-social) would also reflect the difficulty of the choices as easy trials would correspond to the two extreme bins (i.e. 0–0.2 and 0.8–1) as they contain strong information favouring either a 'Keep' or a 'Play' choice, respectively; medium difficulty trials would still favour either a 'Keep' or 'Play' response, but the uncertainty of the outcome would be higher (0.2–0.4 and 0.6–0.8, respectively) and the difficult trials would be associated with the highest uncertainty and as such would not favour one response over the other (i.e. 0.4-0.6).

We placed a background face in the non-social trials to equalize the perceptual load across domains, to enable direct comparisons between the two domains and as an attempt to guarantee that any potential differences seen between the conditions could not be attributed to bottom-up processing of the stimuli. We also distorted the numbers presenting the reward probability ranges to parallel the early encoding of the perceptual stimuli across domains in an effort to produce similar RTs and non-decision time estimates for the two domains. During the non-social trials, our participants were instructed to focus and make their choices based on the numbers displayed on the stimuli. Taken together, these design characteristics aimed to ensure that participants' decisions in the economic task would be based on the same economic considerations across the two domains. In other words, participants would be basing their responses on the reward probability associated with a 'Play' choice in both domains.

After participants made their choice (they were given up to 1.3 seconds to respond), they would be informed of the outcome of their decision via a feedback screen (0.75s). Specifically, they would see their updated score, which would be coloured based on their choice. If the participant chose to 'Keep' during that trial, the score would be coloured in blue. If they chose to 'Play' and successfully doubled their points for the trial, the updated score would be presented in green,

however, if they were unsuccessful, their score was shown in red. Finally, if they did not manage to respond within the allocated 1.3 seconds, they would see a screen informing them that they were too slow to respond. Furthermore, to motivate participants to engage with the task, we told them that in addition to their base rate payment (£6) they would receive a variable bonus (up to £4) based on the points they accumulate during the experiment. Almost all participants received £8.

### 2.2.4 Procedure for Main experiment

The Main experiment involved three sessions: 1) a behavioural session, 2) an online rating session and 3) an EEG-fMRI session (Fig. 1, middle panel). Initially participants were invited to a behaviour session in order to identify individuals who would engage with the task. The procedure for the behavioural session was the exact same as the one described for the Pilot, however based on participant feedback from the Pilot, we decreased the noise level in one of the text images showcasing the reward probability range in the non-social domain (the 0-0.2 bin). This new images was used in all subsequent versions of the economic game. To encourage participants to engage with the task we once again offered them a fixed rate (£6) and an additional payment (up to £4) based on their performance. On average, the participants who were later included in the EEG-fMRI multivariate analysis (see Participants) received £8.03 ± 0.31.

If participants' behaviour in the behavioural session showed that they understood the task by displaying a parametric modulation of their 'Play' choices as a function of the Reward Probability bin, they would be invited to the EEG-fMRI session. One day prior to the EEG-fMRI session, participants would be sent a link to the online rating session. The online rating session involved only the rating task, which functioned in the same manner as the one used in the Pilot and the behavioural session, but was implemented using the online functionality of PsychoPy3 (v3.0.7, Peirce, 2007) and instead of a continuously moving cursor, participants were simply asked to use their cursor and click on the corresponding location on the scale (Fig. 8, b). Additionally, during the Online rating session, participants were exposed to 150 face displays (see Stimulus creation and manipulation) in order to increase the

stimulus sample size since following the Pilot we observed that participants were more likely to cluster their ratings around the middle of the scale. To this end we introduced the additional 50 faces, which were meant to encourage the selection of more faces in the more extreme bins. One participant did not complete the Online rating session due to software issues.

Finally, the EEG-fMRI session also involved a rating task and an economic game. The rating task was the same as the one in the online rating session. Following the rating task, participants were prepared for the EEG-fMRI recordings and entered the scanner. Inside, they played a final round of the economic game (Fig. 7). The trials were presented over 5 runs (100 trials each), lasting approximately 7 minutes (205 volumes). In each run we included a 30-second break at the middle and the end (i.e. after every 50 trials). During each trial, a jittered (1-4s, mean = 2s) fixation cross would be followed by the presentation of the face stimulus (either social or non-social). To further optimize the design (Henson, 2007) and thus increase power and ensure maximal dissociability between the two sources of uncertainty (Mumford et al., 2015), a simulation was run to produce optimal stimulus presentation and jitter lists. Based on the efficiency score of the hemodynamic response from 5000 simulated stimulus presentations and jitter lists, the top 40 lists were chosen. In the social context, we assigned each of the face identities into the five bins based on the subject-specific ratings given on the day of the EEG-fMRI session. However, their ratings across all three rating tasks involved in the Main experiment were used to identify face identities that received inconsistent scores (more than two bins apart), which we removed from the experiment (on average, 10.807 face identities were removed). With the one problematic subject, the final rating was considered twice as it was assumed that it would represent the trustworthiness judgments better than the earlier rating. On average, there were 23, 34, 32, 36 and 14 face identities across the five reward probability bins, respectively. In the non-social context, like in the behavioural sessions and the Pilots, the reward probability was presented explicitly through a probability range displayed on a face neutral for trustworthiness. For both conditions, participants had 1.3 seconds to respond and the stimulus remained on screen until the participants made their choice. If the participant's response was faster than 1.3 seconds, the fixation cross was presented again to pad the remaining time up to 1.3 seconds, so as to keep the run times consistent between participants. Participants were

asked to use a 4-button inline MR-compatible response box (Cambridge Research Systems, 2019) to make their responses. They were informed of their performance only during the breaks following every 50 trials via a 30s feedback screen showing how many points they gained in the last 50 trials. Participants were once again told that they would receive a fixed and an additional variable payment based on their performance (£16 base rate + up to £4). In reality all participants received £20 for their participation.

### 2.2.5   Pilot data analysis

The analysis of the Pilot data involved the proportion of 'Play' choices and the reaction times (RTs) gathered during the economic game. We aimed to analyse the data using two generalized linear mixed models (GLMMs), as they allow more flexibility for specifying the random effect and can thus represent the inter-participant variability better than alternative methods (like repeated measure ANOVAs, Baayen et al., 2008). However due to convergence issues with the RT data, we used a linear model for the RT analysis.

For the choice behaviour we were interested in establishing whether the probability of a 'Play' choice was predicted by the Reward probability given a 'Play' choice. We aimed to specify a maximal model including all main effects and interactions of our two predictors: i.e. domain (social vs non-social) and reward probability given a 'Play' choices (0-0.2, 0.2-0.4, 0.4-0.6, 0.6-0.8, 0.8-1). We also introduced subject-specific random slopes for the main effects and interactions to maximally reflect our design and to decrease the likelihood of a type I error (Barr et al., 2013). However, due to convergence issues we excluded the random correlations. The Reward probability given a 'Play' response predictor was divided into 4 predictor variables in order to account for the ordinal nature of the comparison and were backward difference coded (UCLA IDRE, 2011). For both analyses (choice behaviour and RT) we used deviation coding for the domain predictor by mean-centering the dummy coding assigned to the two domains.

$$logit(Play) = 0.247$$
$$+ 0.173 * domain$$
$$+ 0.704 * RewardProbability1$$
$$+ 1.943 * RewardProbability2$$
$$+ 1.656 * RewardProbability3$$
$$+ 1.339 * RewardProbability4 \tag{1}$$
$$- 0.767 * domain : RewardProbability1$$
$$+ 1.959 * domain : RewardProbability2$$
$$+ 0.795 * domain : RewardProbability3$$
$$+ 0.559 * domain : RewardProbability4$$
$$+ (1 + domain * RewardProbabilityAll || Subject)$$

For the analysis of RTs we asked if RTs can be predicted by the difficulty of the trials (1: easy trials, i.e. 0-0.2 and 0.8-1, 2: mid, i.e. 0.2-0.4 and 0.6-0.8 and 3: difficult, i.e. 0.4-0.6). The difficulty predictors were once again backward difference coded and divided into 2 predictors.

$$RT = 692.733$$
$$+ 35.304 * domain$$
$$+ 33.079 * Difficulty1$$
$$+ 27.369 * Difficulty2 \tag{2}$$
$$+ 8.532 * domain : Difficulty1$$
$$+ 26.143 * domain : Difficulty2$$

To establish the significance of our main effects and interactions, we then performed post-hoc likelihood-ratio chi square ($LR\chi^2$) model comparisons for the choice behaviour GLMM, and *F*-statistics test for the reaction times linear regression, by comparing the full model (as defined above) to models, which exclude a predictor (or a group of predictors) of interest. Here we are only reporting the

results obtained from these model comparisons. Significant interactions were further decomposed into simple effects. This was achieved by creating two sets of the data with dummy coded domain predictors. When assigning '0' to the social domain and rerunning the analysis, we obtain the simple effects for the social condition. The same was done for the non-social condition.

Both analyses were performed in R version 3.5.1 (2018-07-02) and the code was developed using RStudio version 1.1.456. The results for the choice behaviour were analysed using a binary mixed models logistic regression using the packages tidyverse (version 1.2.1) and lme4 (version 1.1.18.1). The model family was specified as binom(logit) in order to account for the binomial nature of the response and we selected the bobyca optimizer. The same packages and a simple linear regression model were used to analyse the RTs.

### 2.2.6 Choice probability calculations for the EEG-fMRI session from the Main experiment

To assess the similarity between the probabilities of 'Play' choices across the social and non-social contexts during the EEG-fMRI session, we used a conventional likelihood-ratio test implemented using MATLAB version 2018b. This test aimed to assess whether a single common distribution can account for the data better than two separate distributions for the social and non-social decisions. Specifically, we examined whether a single sigmoid curve (Weibull function) would fit the combined social and non-social choice data across the five reward probability levels as well as two separate curves (Philiastides and Sajda, 2006). We performed this separately for each participant by fitting the best single Weibull function jointly to the two data sets in addition to the individual fits. The likelihoods (L) obtained from this procedure were transformed using the following equation:

$$\lambda = -2ln\frac{\frac{1}{N}\sum_{i=1}^{N}L_i(\text{data} \mid \text{joined curve})}{\frac{1}{N}\sum_{i=1}^{N}L_i(\text{data} \mid \text{individual curves})} \tag{3}$$

where $N$ represents the number of participants and $\lambda$ is distributed as $\chi2$ with two degrees of freedom (Hoel et al., 1971). If $\lambda$ exceeds the criterion value (for $p = 0.05$), we concluded that a single function fits the data as well as two separate

domain-specific functions.

### 2.2.7 Single-trial regressions for EEG-fMRI session from the Main experiment

We examined the relationship between the reward probability (i.e. indirect trustworthiness bins and pure probability levels in the social and non-social contexts, respectively) and the probability of playing (1: 'Play', 0: 'Keep') on individual trials (Fig. 10, a) with a single-trial logistic regression analysis (separately for each participant and for each of the social and non-social trials):

$$P_{play} = [1 + e^{-(\beta_0 + \beta_1 \times y(\text{reward probability}))}]^{-1} \tag{4}$$

We also investigated the link between task difficulty (i.e. 1: easy (reward probabilities 0-0.2 and 0.8-1), 2: medium (reward probabilities 0.2-0.4 and 0.6-0.8), 3: difficult (reward probabilities 0.4-0.6)) and RTs on individual trials (Fig. 10, b) with a single-trial regression analysis (separately for each participant and for each of the social and non-social trials):

$$RT = \beta_0 + \beta_1 \times (\text{difficulty level}) \tag{5}$$

In both cases we assessed whether the regression coefficients across participants ($\beta_1$ values in Eqs. 4 and 5) came from a distribution with a mean different from zero (using separate $t$-tests). All single trial regression were performed using using MATLAB version 2018b.

### 2.2.8 Drift diffusion model

We used a version of a DDM, which involved a special case of the leaky competing accumulator model (Ornstein–Uhlenbeck process) to model EA based on the behavioral data from the EEG-fMRI experiment, following procedures described in Pisauro et al. (2017) and Polanía et al. (2014):

$$EA(t+1) = EA(t) + (\lambda EA(t) + k(evidence)dt + N(0, \sigma)) + bias(evidence = 0) \quad \textbf{(6)}$$

The decision thresholds for 'Play' and 'Keep' choices were assigned to +1 and – 1, respectively, with positive drift rates linked to reward probability levels favoring 'Play' choices, whereas negative drift rates were associated with reward probabilities favoring 'Keep' choices. To this end, in Eq. 6 the evidence corresponds to a transformed version of the original five reward probability levels, centered around zero (i.e. -0.5, -0.25, 0, 0.25, 0.5).

The free parameter $k$ modulates the evidence input, $\lambda$ signifies the acceleration to threshold and $N(0, \sigma)$ is a Gaussian noise term with standard deviation $\sigma$. We re-calculated the EA on a time increment dt = 0.001s and we examined the model choice once |EA| > boundary. We used a non-decision time free paramater (nDT) to account for early visual encoding and motor preparation. This nDT estimate was included into the total RT. For trials corresponding to the indecision point (i.e. 0 evidence, where the evidence does not favour 'Keeping' or 'Playing') we included an additional free parameter, *bias*, to capture inter-individual biases towards either 'Play' or 'Keep' choices. The RTs were split based on the choice ('Keep' or 'Play'), and we flipped the sign for the 'Keep' trials, so that all 'Keep' trials received a negative sign and then recombined the RTs into a single distribution (Voss et al., 2004). We then compared this RT distribution and participants' choice probabilities to those generated by the model. For a given set of parameter estimates, we estimated the log likelihood (LL) of the data using the following formula:

$$LL \sim \sum_{evidence=1}^{5} log(KS(RT_{data}^{evidence}, RT_{model}^{evidence}))$$
$$+ \sum_{evidence=1}^{5} log(exp(-(\frac{Pplay_{data}^{evidence} - Pplay_{model}^{evidence}}{0.01})^2)) \quad \textbf{(7)}$$

We used $KS(p, q)$ to estimate the probability that our distributions are equal, based on the Kolmogorov–Smirnov test (via the ktest2 function in MATLAB). Pplay corresponds to the fraction of 'Play' choices for each of the five evidence levels.

We used a two-step process to fit the model. We first ran the fmincon MATLAB function 20 times to obtain an initial estimate of subject-specific parameters and selected the parameters associated with the smallest LL for the next step. We then ran a grid search fitting procedure for each participant using a fine-grained parameter space around the estimates generated in the previous step. Choices and RT distributions were created for each possible combination of the four free parameters from 5000 simulated decision traces per decision domain.

## 2.3 Behavioural and Modelling results

### 2.3.1 Comparable behaviour for social and non-social trials from the Pilot study

We ran a Pilot to examine the suitability of our paradigm for comparing social and non-social decision making. We also used it as an opportunity to validate our choice of stimuli. We ran a GLMM to assess the effects of domain and reward probability given a 'Play' response on the choice behaviour. This analysis demonstrated that there was no significant difference between the social and the non-social 'Play' choices (LR$\chi^2$(1) = 0.7, $p$ = 0.403), which would be consistent with the notion of a common underlying decision-making mechanism. There was an overall significant effect of reward probability (LR$\chi^2$(4) = 138.981, $p$ < 0.001) demonstrating that the probability of a 'Play' response increased with the probability of a reward given a 'Play' choice. There was also a significant interaction between the reward probability and the condition LR$\chi^2$(4) = 39.019, $p$ < 0.001. The significant interaction was further decomposed to examine whether the reward probability had a significant effect for both decision domains. This was indeed the case for both the social (LR$\chi^2$(4) = 118.97, $p$ < 0.001) and the non-social trials (LR$\chi^2$(4) = 117.905, $p$ < 0.001). Therefore, for both social and non-social choices, the probability of 'Playing' increased as the reward probability given a 'Play' choice increased, showing that both the social and non-social domain follow the same trends, even though there was a significant difference in the strength of these effects as suggested by the significant interaction. The fact that we see this parametric modulation in the social condition and the fact that there were not many inconsistently rated

**Figure 9:** **Pilot Choice and reaction time (RT) results. a)** Participants' average Play responses for the social (red) and non-social (blue) condition. The probability of 'Play' choices increased as the probability of a reward given a 'Play' response increased (P(payoff|play)). **b)** Participants' average RTs for the social (red) and non-social (blue) condition. RTs were the highest when there was no strong evidence for or against 'Play' decisions for both conditions. Average individual behaviour is presented in grey circles

faces removed (on average, 10.807 face identities were removed) suggests that our social manipulations captured through the initial ratings were successful in identifying the social processing associated with these decisions and that participants reliably used the trustworthiness information embedded in these images.

We also investigated the effects of domain and difficulty on the RTs using a linear regression. This RT analysis demonstrated that there was a significant (albeit small = 34.516 ms) main effect of domain ($F(1, 11377) = 135.28$, $p < 0.001$), with the social responses ($676.344 \pm 95.752$ ms) being on average faster than the non-social ones ($710.861 \pm 85.227$ ms). There was also a significant main effect of difficulty ($F(2, 11377) = 114.4$, $p < 0.001$). Their interaction, was also significant ($F(2, 11377) = 8.807$, $p < 0.001$). Therefore, we further decomposed the interaction to examine

whether difficulty had a significant effect on the RTs in both decision domains. We found a significant effect of difficulty on the social ($F(1, 11377) = 32.304$, $p < 0.001$) and non-social domain ($F(1, 11377) = 90.854$, $p < 0.001$). This suggests that although this effect might have differed in strength between the two domains, both RTs were influenced by the difficulty of the trials.

### 2.3.2 Comparable behaviour for social and non-social trials from EEG-fMRI experiment



Figure 10: **EEG-fMRI behavioural Choice and reaction time (RT) results.** Social and non-social behavioral responses (red and blue circles) versus modelling performance of a drift diffusion model (black crosses) for proportion of 'Play' choices (**a**) and RTs (**b**). 'Play' responses increased with probability of reward given a 'Play' choice (P(payoff|play)) and RTs were the highest when there was no strong evidence for or against 'Play' decisions. Participant-specific behavior is presented in grey circles.

There was a positive correlation between the participants' fraction of 'Play' choices and the reward probability for both the social and non-social trials (social: $t(30) = 17.769$, $p < 0.001$; non-social: $t(30) = 4.086$, $p < 0.001$). This suggests that participants selected the riskier option more frequently when the likelihood

of receiving the higher payoff was also higher (Fig. 10, a). Crucially, we found that the choice behavior was comparable between the social and non-social trials, since a likelihood-ratio test (see Material and Methods) showed that a single sigmoid function fit the fraction of 'Play' choices (jointly across both conditions) as well as two separate functions ($\lambda(30) = 0.551$, $p = 0.759$).

We found an inverted 'V' relationship between the mean RTs and the reward probability bins for both the social and non-social trials (Fig. 10, b), reflecting a positive relationship with task difficulty (social: $t(30) = 10.024$, $p < 0.001$; non-social: $t(30) = 10.692$, $p < 0.001$). Specifically, the longest RTs were linked to the most difficult trials (reward probabilities 0.4–0.6), whereas the shortest RTs were associated with the easiest trials (reward probabilities 0–0.2 and 0.8–1) and intermediate RTs corresponded to medium difficulty trials (reward probabilities 0.4–0.6 and 0.6–0.8). Overall there was a small (41.637 ms), albeit significant difference between the social and non-social RTs ($t(30) = -3.274$, $p = 0.003$), with social trials ($M_S = 677.864$ms, $SD_S = 86.479$ ms) being on average faster than non-social ones ($M_{NS} = 719.502$ms, $SD_{NS} = 91.287$ ms).

### 2.3.3 Comparable modelling output

We found that the model-derived choice behaviour and RTs were highly correlated with the actual observed behaviour: social – fraction 'Play' Choice: r = 0.945, $t(154) = 36.464$, $p < 0.001$, RT: r = 0.754; $t(154) = 15.154$, $p < 0.001$; non-social - fraction 'Play' Choice: r = 0.968, $t(154) = 94.196$, $p < 0.001$, RT: r = 0.765; $t(154) = 14.461$, $p < 0.001$. Therefore, our modelling results were able to reflect the trends observed in our behavioural dataset in that the models would parametrically increase their proportion of 'Play' choices with the increase in reward probability and similarly, the times for decision would increase with the difficulty of the trials (Fig. 10). Consequently, the model was able to reflect the main manipulations of the design across the two decision-making domains, thus justifying features of the modeling procedure such as the parametric modulation of the drift rate parameter.

The mean parameter estimates for the social context were: ($\lambda$: $5.774 \pm 2.357$, $k$: $3.206 \pm 1.555$, $\sigma$: $0.02 \pm 0.01$, $bias$: $-0.00004 \pm 0.0004$, $nDT$: $0.336 \pm 0.09$ and the mean parameter estimates for non-social context were: ($\lambda$: $5.277 \pm 2.37$,

$k$: 2.611 ± 1.355, $\sigma$: 0.011 ± 0.006, $bias$: -0.00002 ± 0.0006, $nDT$; 0.304 ± 0.089. Most parameters did not differ significantly between the two decision domains ($\lambda$: $t$(30) = -1.3, $p$ = 0.203, $bias$: $t$(30) = 0.26, $p$ = 0.8, $k$: $t$(30) = -1.349, $p$ = 0.188, $nDT$: $t$(30) = -1.363, $p$ = 0.183), thus suggesting high degrees of similarity between the social and the non-social choice behaviour. There was nonetheless a small but significant difference in the noise term ($\sigma$: $t$(30) = -4.244, $p < 0.001$). This minor difference may be attributed to additional internal variability when processing the facial trustworthiness in the social context due to the larger number of stimuli compared to identifying the numbers in the non-social trials.

## 2.4  Discussion

Determining whether or not social choices use the same underlying neurocomputational mechanism as non-social decisions could offer a deeper understanding of human decision making and behaviour as a whole. So far, however, there have been few direct comparisons between the two where the social uncertainty is determined by a social cue, which examine decision-making processes beyond the value comparison stages and where the social and non-social uncertainty explicitly vary along comparable scales. Establishing whether the two domains share a common underlying mechanism would require the consideration of several levels of functioning, the most consequential of which are the algorithmic and the implementational levels (Lockwood et al., 2020). Therefore the aim of this chapter was to provide an initial examination of the algorithmic level used in the two decision domains. To this end, participants were asked to make economic decisions based either on an explicitly disclosed probability range or an indirect trustworthiness rating in a Pilot experiment and in the Main EEG-fMRI study. In this chapter we examined the behavioural similarities between social and non-social choices from the Pilot and the Main experiment and used a DDM to showcase the potentially shared computational principles between the two domains.

Behaviourally, the choice selection presents two pieces of evidence showcasing the similarities across the social and non-social decisions. The first one is the similarity across the two conditions in terms of 'Play' choices, demonstrated by the lack of a significant main effect in the Pilot. Even though there was a significant

interaction in the Pilot, both simple effects were significant, showcasing that the social and non-social 'Play' choices showed a common trend of scaling positively with reward probability. Additionally, for the data from the Main experiment we aimed to determine whether a single choice behaviour distribution can account for the data and we saw that a single sigmoid function fit the data as well as two separate functions. The second way in which our choice data highlights the similar trends associated with social and non-social choices is that the proportion of 'Play' choices increased with the probability of a reward given a 'Play' choice. This was demonstrated by the significant main and simple effects of P(payoff|Play) in the Pilot, whereas in the Main experiment it was shown by the positive correlation between P(payoff|Play) and proportion of 'Play' choices as well as by the significant effect in the logistic regression. This demonstrates that both conditions follow a trend, consistently seen in non-social experiments where the probability of selecting one of the alternatives scales with the available evidence for that choice (e.g. Bogacz et al., 2006; Gold and Ding, 2013; Philiastides and Sajda, 2006; Ratcliff et al., 2009). We highlight that for the social condition, we determined the reward probabilities based on the participant-specific ratings. On potential issue with this approach would be if these ratings were highly inconsistent, which would signify that the behavioural responses would not be associated with our indirect trustworthiness ratings. However we removed any displays with unreliable ratings and on average there was only a small number of faces that were removed (10.807 out of 150 displays). Furthermore, as we observed a parametric modulation of Play behaviour reflective of the ratings, this signifies that the remaining ratings were reliable and were reflective or the social processing that participants would engage in in this condition. Therefore, these results indicate that our paradigm was successful in modulating behavioural difference at the different levels of reward probability. Moreover, these comparable choices suggest that we might be able to make a direct comparison for the algorithmic process investigation with our modeling procedure.

In terms of the RTs, in both the Pilot and the Main experiment there was a small, but significant effect of domain. Nevertheless, the overall trends seen in the social and non-social reaction times were similar and consistent with previous literature. Specifically, both the social and non-social trials were associated with an inverted 'V' pattern, where trials in the middle probability range (0.4-0.6), where there

is no strong evidence indicating which choice option ('Play' or 'Keep') is the most beneficial one (thus making them more difficult), were associated with the longest reaction times. Conversely, the trials with the strongest evidence in favour of 'Keeping' (0-0.2) and 'Playing' (0.8-1, thus making them easier), were associated with the shortest reaction times. These observations were further supported by a significant main effect of difficulty in the Pilot as well as in the Main experiment. Such trends have been observed extensively in the behavioural results in previous PDM and VBDM literature (e.g. Drugowitsch et al., 2012; Philiastides and Ratcliff, 2013; Philiastides and Sajda, 2006; Pisauro et al., 2017; Ratcliff et al., 2009; Uchida et al., 2006, which have argued that the amount of available decision information (i.e. the choice uncertainty) ultimately affects the decision (Urai et al., 2017).

However, in addition to examining the purely behavioural responses, we were also able to fit the data from the Main experiment to a DDM in order to offer a mechanistic account of social and non-social choice selection. Therefore, this allowed us to go beyond the simple observation that the behavioural tendencies are shared among social and non-social decisions and instead describes the processes and computations necessary to drive choices within the two domains. We observed that we could fit such sequential sampling DDMs to both the social and non-social behaviour and we could show that the models gave rise to the same trends seen in the behaviour. This suggests that the underlying mechanistic steps that make up the DDM might correspond to the ones employed by the decision-makers. Furthermore, we saw that most of the estimated parameters did not differ significantly between the two domains. These results indicate that our participants were utilizing the domain-specific uncertainty in a comparable way to drive both their social and non-social decisions. This observation corresponds to the decision-making process model highlighted in the Introduction chapter (Fig. 3, Rangel et al., 2008; Rangel and Clithero, 2014), which suggests that the only differences between decision-making domains should be the initial appraisal of the domain-specific information, but that the subsequent processing steps should be shared across domains. These modeling results therefore give credence to the notion that this decision-making process, which might stand at the basis of various non-social decisions making domains (e.g. Krajbich and Rangel, 2011; Nunez et al., 2017; Philiastides and Ratcliff, 2013; Pisauro et al., 2017), might also reflect the algorithmic characteristics of social choices. We do note that there was a differ-

ence in the noise parameter ($\sigma$), however considering how for the social condition we included a larger number of stimuli per bin it may not be surprising that there might have been some minimal, but significant additional variability in the social EA process. Nevertheless, these results provide mechanistic evidence for the investigation of whether social decisions are a part of a universal decision-making process by suggesting that that social and non-social choices potentially share the same algorithmic computations.

While our paradigm was largely successful in that it led to prototypical results seen throughout the non-social domain, it is also worth highlighting that there were some participants who were removed due to inconsistent behaviour, which suggests that not everyone engaged with the task in a way that we expected. It is possible that these participants were not motivated to engage in the task as they would receive the base payment rate regardless of their performance. It is also possible that some of the participants might have been exposed to similar stimuli in previous experiments and therefore did not believe the cover story about the faces belonging to people who have taken part in a Trust game. Some of the removed participants also shared that they would respond based on a pattern (e.g. 1 'Play' response, followed by 2 'Keep' responses), which suggests that they either ignored the cover story and the explanation of the rules or that they were not convinced that these rules actually reflected the likelihood of the outcome. While it might have been possible to show more explicit evidence in favour of the fact that unlike those individuals, the participants included in the final sample believed the cover story through the use of a questionnaire, we also highlight that such a post-experiment questionnaire might have suffered from social desirability biases. A further issue with such a post-task examination would be the fact that more of our participants could have then guessed the true nature of the experiment and thus revealed it to other potential participants. While we fully debriefed participants, who wanted to know more about the experiment, many opted out of this option and the ones we did debrief, we asked not pass this information on to others. However it is conceivable that they did not follow these instructions, which may be why some participants might not have believed the manipulation and thus responded in unexpected ways. Regardless of the reason for the inconsistent behaviour displayed by some participants, we highlight that some of those included in the final sample would unpromptedly express that they 'hated' some of the

faces and that they thought that others were 'nice' to them. Furthermore, some participants would refuse having their faces be used in future experiments, which suggests that they believed that the faces they were seeing could have come from a previous study. We also take the behavioural patterns seen in the social condition to suggest that the indirect trustworthiness judgments were indeed used to inform the social choices. Moreover, the fMRI results presented in the next chapter, showing that the social trials were linked to activations strongly associated with social processing can be seen as further evidence for the notion that participants were extracting relevant social information from the face displays. Based on these pieces of information, we believe that the participants included in our final sample accepted our cover story and interpreted the social stimuli in a social way.

To summarise, this chapter aimed to investigate whether the algorithmic processes involved in social choices might be the same as the ones regularly linked to non-social decisions. To this end, we developed a paradigm, which aimed to vary the social and non-social uncertainty across comparable scales and to thus enable a fair comparison between the two decision domains. The observed behavioural similarities between the social and non-social choices indicated that our paradigm was suitable for the direct comparison of the two sources of uncertainty. Consequently, this allowed us to elucidate the algorithmic properties associated with social and non-social choices through our modelling procedure. Specifically, our results suggested that the two decision domains might rely on accumulation-to-bound computations consistent with the decision-making process model described in the previous chapter and therefore provide mechanistic evidence supporting the notion that social choices might use the same neurocomputational principles as non-social decisions.

# 3 Chapter 3, Experiment 2

## 3.1 Introduction

The outcomes of our decisions are rarely certain and the source of this uncertainty can vary - for some choices we may rely on the known probabilities associated with the choice alternatives, or we may use social information to guide our decisions. Recent research has begun to investigate and compare the neurocomputational principles of these two types of decision making (e.g. Janowski et al., 2013; Levy and Glimcher, 2012). However, there are still unanswered questions about the extent to which we process and use social and non-social uncertainty in the same way to guide our decisions and if these mechanisms are implemented within the same neural network.

In recent years, the examination of a potentially universal decision-making process has become more prevalent with studies demonstrating at least a partial overlap in value computations for social and non-social choices (e.g. Behrens et al., 2008; Janowski et al., 2013; Levy and Glimcher, 2012). Such observations have given rise to the notion of a 'common currency' implemented in the brain, where domain-specific information is transformed into a universal signal used for decision making (Ruff and Fehr, 2014). However, whether the exact same structures are involved in this process remains unknown since in many of these examples, the social and non-social choices are examined in separate experiments or their uncertainty is not explicitly equalized. Moreover, there is little literature focusing on whether the social and non-social information is subsequently utilized in the same manner and whether the same neural framework is used to implement the decision processes. In the previous chapter we investigated the algorithmic similarities between social and non-social choices and highlighted the common behavioural trends, prototypical of the ones seen throughout the non-social decision-making literature. We also suggested that the computations involved in these choices can be explained through a DDM and that the parameters used in these computations are mostly consistent between the two domains. These results provide evidence in favour of the notion that social choices may rely on the same algorithmic considerations as non-social decisions since our paradigm explicitly attempted to allow

only the source of the uncertainty to vary between the two domains.

In this chapter we aim to continue the investigation of the potential similarities between the social and non-social choices by focusing on their implementational characteristics to investigate whether they share a common underlying mechanism for integrating and accumulating relevant decision evidence. We aim to achieve this by examining the spatiotemporal characteristics of social and non-social decisions using simultaneous EEG-fMRI recordings (e.g. Pisauro et al., 2017). Previous studies have used EEG to investigate the temporal characteristics that underlie non-social decision making (e.g. Philiastides et al., 2014) and fMRI has been used to implicate certain brain regions in the process (e.g. Heekeren et al., 2004). However, the combination of the two modalities in a simultaneous EEG-fMRI experiment not only takes advantage of the benefits associated with each modality on its own, but their synergy allows us to obtain further insight into the underlying processes of the two decision-making domains (Sui et al., 2012). There are two main ways in which this synergy can be achieved. The first is to use them equally in a symmetric approach, which often uses machine learning techniques to combine the multimodal information into a common feature space in order to investigate the shared state-dependant processes that give rise to the EEG and fMRI observations (e.g. Valdes-Sosa et al., 2009). The second way of combining the multimodal information is through the asymmetric approach, where one of the modalities is used to inform the investigation of the other modality (Philiastides et al., 2021). Therefore one may use the trial-by-trial variability in the EEG amplitude as a predictor in a general linear model (GLM) analyses of the fMRI signal, or take advantage of the localization information provided by the fMRI data to inform the source localization of the EEG data. The EEG-informed fMRI analysis is the more common of the asymmetric approaches and it relies on identifying latent endogenous features of interest within the EEG information to explain the variance in the fMRI data with the assumption that the neural substrate that gives rise to the variation in the EEG signal also modulates the activity captured by the fMRI. This often leads to the implication of brain regions, the contribution of which would have been overlooked by standard unimodal methods (e.g. Fouragnan et al., 2015). Even though some fMRI studies use model-derived predictors to examine similar latent processes, very few go beyond the utilization of a simple mean and standard deviation estimate (Gluth et al., 2017; Turner et al., 2015) and these predictors can be influenced by

model misspecifications. Therefore, the trial-by-trial variability recorded from an endogenous source offers stronger explanatory power than alternative measures. To this end in this chapter we benefited from our simultaneously acquired EEG-fMRI recordings and employed an asymmetric approach where we used the trial-by-trial variability captured by the EEG to inform our fMRI analysis.

Initially, we utilized the insight gained from our DDM to construct model-derived estimates of the EA dynamics associated with the social and non-social choices to identify neural signatures of EA through a univariate analysis of the evoked EEG responses (e.g. Pisauro et al., 2017). However, since the reliability of the EEG signal is crucial for the EEG-informed fMRI analysis, we then took advantage of the multi-dimentional nature of the EEG signal, which enables the use of multivariate decoding methodologies to extract more reliable neural signatures by exploiting the full spatiotemporal data structure to investigate how informative the observed brain patterns are in describing the processes guiding decision making (Haynes, 2011). Specifically, the EEG data contains multi-channel noisy information, however, the information from some of these channels is more relevant for the investigated phenomenon than others and neighbouring channels may carry complementary information, which is often not taken into consideration in univariate EEG analyses. Therefore, by spatially integrating this information, we were able to increase the relevant signal-to-noise ratio and thus produce more robust single-trial estimates (Parra et al., 2005). Such advanced quantitative techniques enabled a deeper understanding of the underlying mechanisms of social and non-social choices and allowed us to reliably explore the degree of overlap of neural networks guided by the two sources of uncertainty.

To this end, here we present a simultaneous EEG-fMRI experiment, during which participants made choices between small but certain rewards and larger but uncertain payoffs in interleaved social and non-social trials. In the non-social trials the outcome of the uncertain option was determined by explicit probability ranges, whereas in the social trials the outcome was dependant on subject-specific trustworthiness judgments (see Procedure for rating tasks and economic game tasks from Chapter 2). Both our EEG analyses uncovered build-up activity reflecting EA dynamics, which corresponded to the decision difficulty in our task. We then used the trial-by-trial variability captured in the slopes we identified following our

multivariate EEG analysis to create parametric fMRI predictors. This allowed us to implicate a region in the posterior-medial frontal cortex (pMFC) as the site for EA for social and non-social choices. We further showed that there was a task-dependant co-activation between the activity in the pMFC and areas in the human valuation network across both contexts. We also found domain-specific activations associated with early decision evidence encoding for each context. These results, therefore suggest that the relevant domain-specific information is converted into a 'common currency' in the human valuation system and accumulated for decision in the pMFC in a domain-general decision-making network.

## 3.2 Methods

### 3.2.1 Participants, Stimuli and Procedure

The participants, stimuli and procedures referred to within this chapter are the same as the ones described in Chapter 2.

### 3.2.2 EEG data acquisition

We collected the data using an MR-compatible EEG amplifier system (Brain Products, Germany), which we continuously recorded at 5000 Hz using Brain Vision Recorder software (Brain Products, Germany), based on the setup described in Gherman and Philiastides (2018). Specifically, we used a hardware 0.016-250 Hz band-pass to filter the data online. We used 64 Ag/AgCl scalp electrodes placed according to the 10 – 20 system. The reference electrode was built in between electrodes Fpz and Fz and the ground was located between electrodes Pz and Oz. To ensure subject safety, each electrode had in-line 10 kOhm surface-mount resistors and all leads were additionally bundled and twisted along their entire length. We ensured that the input impedance for each electrode was lowered to $< 50$ kOhm (25 KOhm average across participants). We synchronized the acquisition of EEG and MRI data (Syncbox, Brain Products, Germany) and we recorded the MR-scanner triggers separately for the subsequent offline removal of MR gradient artifacts. We also lengthened the scanner pulses to 50 $\mu$s via an in-house pulse stretcher

to facilitate the recording of the scanner triggers. We recorded the experimental event codes and participants' responses with the EEG data through the Brain Vision Recorder software. We also ensured that we positioned subjects inside the scanner so that electrodes Fp1 and Fp2 were aligned with the isocentre of the MR scanner. Finally, we secured the cabling connecting to the EEG amplifiers at the back of the bore by taping them to a cantilever beam to minimize scanner vibration artifacts.

### 3.2.3  EEG data preprocessing

The EEG data preprocessing and analysis were done using MATLAB 2018b (Mathworks, Natick, MA), following the procedures outlined in Pisauro et al. (2017) and Gherman and Philiastides (2018). Due to the magnetic induction on the EEG leads, the EEG signal recorded inside an MR scanner is contaminated with MR gradient and ballistocardiogram (BCG) artifacts. To remove the gradient artifacts, we constructed average artifact templates from sets of 70 consecutive functional volumes (centered on each volume of interest). We then subtracted the templates from the EEG signal for each functional volume of the dataset, which we were able to do because of the consistency of the artifacts over time (Allen et al., 2000). Any residual spike artifacts were removed by applying a 12 ms median filter. We also applied a 0.5 – 20 Hz band-pass filter to remove slow DC drifts and higher frequency noise. The data was also downsampled to 1000 Hz.

To remove eye blinks we asked our participants to blink repeatedly while there was a central fixation cross on screen. The timings of these events were recorded, which allowed us to use a principal component analysis to identify linear components associated with eye-blinks, which we subsequently removed from the broadband EEG data collected during the main task (Parra et al., 2005).

The removal of BCG artifacts is more challenging as they share frequency content with the EEG. To avoid loss of signal power in the EEG, for the **univariate** EEG analysis we only removed the two highest participant-specific BCG components using principal component analysis. For the **multivariate** approach, we instead created datasets with different number of principal components removed (up to 5). We determined the number of optimal principal components for each participant by maximizing classification performance along the task-relevant dimension

(see below) using cross validation (average number of BCG components across participants: $2.447 \pm 1.969$). These BCG principal components were identified after low-pass filtering at 4Hz (i.e. to extract the signal within the frequency range where BCG artifacts are typically observed) and then removed by using the sensor weightings corresponding to the relevant components and projecting them onto the broadband data and subtracting them out. Finally, we baseline corrected the data based on the 100 ms prior to the onset of the stimulus.

### 3.2.4   Univariate analysis of EEG data

To identify potential signatures of EA dynamics in our EEG signal, we performed a univariate analysis of the EEG evoked responses, similar to the one described in Pisauro et al. (2017). The social and non-social trials from 400 ms time windows preceding the response were averaged to create subject- and electrode-specific response-locked EEG traces. Trials with amplitudes higher than 2 standard deviations were removed. These 400 ms were further adjusted based on correlations with model-derived EA traces created based on subject-specific parameter estimates from our DDM analysis (see Chapter 2). This adjustment involved temporally shifting the 400 ms time window up to 200 ms further away from the response to account for motor preparation. This resulted in the identification of subject-specific best electrodes, which reflect accumulation-like activity. The robustness of the selected electrodes was tested by performing a permutation procedure. Specifically, for each trial we randomly assigned electrode labels across our data set. We performed this procedure 100 times. This allowed us to create a distribution of maximum correlation values, which are then compared to the original correlation value. If the original value was not within the top 5% of the distribution, then the participant would not be considered for the remainder of the analysis. All participants' best context-specific electrodes passed the permutation test ($p < 0.05$).

### 3.2.5  Regressions for Univariate EEG analysis

We examined the relationship between the model-derived EA traces and the subject-specific averaged best-electrode traces with a linear regression analysis (separately for each domain) in order to test if the $EA_{model}$ in the form of a single vector containing the time series of the last 400 ms of the model-derived EA traces from each participant would predict $EA_{EEG}$ i.e. a single vector containing the time series from all the time-shifted 400 ms windows of subject-specific best electrode EEG activity:

$$EA_{EEG} = \beta_0 + \beta_1 \times (EA_{model}) \tag{8}$$

We also explored the effects of trial difficulty and trial duration by selecting the easy trials (reward probability 0-0.2 and 0.8-1) and the difficult trials (reward probability 0.4-0.6) as well as splitting the data with a median split into short and long trials. We then used simple linear regressions (separately for each domain), where $EEG_{data}$ was the subject-specific time series based on the final 200 ms of the 400 ms window (since this window reflects the biggest differences across conditions) and $EA_{interval}$ was a monotonically increasing vector that matched the size of this 200 ms time window. We performed these regressions to estimate the subject-specific and domain-specific slopes associated with each of these categories (easy, hard, short, long):

$$EA_{interval} = \beta_0 + \beta_1 \times (EEG_{data}) \tag{9}$$

We then assessed whether the regression coefficients across participants ($\beta_1$ values in Eq. 9) for the easy/difficult and the short/long comparisons differed significantly (using separate paired $t$-tests). These comparisons were done separately for the two domains and using MATLAB, version 2018b.

### 3.2.6 Multivatiate analysis of EEG data

Here, we were interested in identifying robust representations of activity related to EA in the EEG data. To this end we followed the procedures described in Diaz et al. (2017), Franzen et al. (2020) and Gherman and Philiastides (2018) to perform a single-trial multivariate discriminant analysis (Parra et al., 2005; Sajda et al., 2009) to discriminate between easy (i.e. reward probabilities 0–0.2 and 0.8–1) and difficult trials (reward probabilities 0.4–0.6) in stimulus-locked EEG data. To examine the potential universality of this process we initially performed this analysis by collapsing across both social and non-social trials. We predicted that the existence of such an accumulating activity would be associated with a build-up rate proportional to the decision difficulty. In turn, this would be reflected in a gradual increase in the discriminator's performance as the traces for easy and difficult trials diverge over time in stimulus-locked data (Fig. 12, a). We did not use the medium difficulty trials (i.e. reward probabilities 0.2–0.4 and 0.6–0.8) in this discrimination process (i.e. they acted as "unseen" data), to more convincingly test for a full parametric effect on the build-up rate associated with the different decision difficulty levels (see below).

Specifically, this procedure allowed us to identify an optimal combination of EEG sensor linear weights (i.e., a spatial filter $\mathbf{w}$). We then applied these sensor weightings to the EEG data ($\mathbf{x}(t)$), thus producing a one-dimensional projection (i.e., a discriminant component $y(t)$) that discriminates between the two difficulty levels:

$$y(t) = \mathbf{w}^T\mathbf{x}(t) = \sum_{i=1}^{D} w_i x_i(t) \tag{10}$$

$D$ corresponds to the number of channels with an index $i$, and $T$ denotes the transpose of the matrix. This method allowed us to identify $\mathbf{w}$ for short (60 ms) overlapping time windows, all centred on 20 ms-interval time points, between -100 and 800 ms relative to the stimulus presentation for each subject and time window separately. By doing so, we integrated information spatially across the multidimensional sensor space, which allows us to increase signal-to-noise ratio, while still preserving the trial-by-trial variability in the relevant discriminating

component. Moreover, when we apply spatial filters (**w**'s) obtained in this way to individual trials, we create a measurement of the discriminant component amplitude for that trial, which can be used as a neural stand-in for the relevant decision activity.

We estimated the optimal discriminating spatial weighting vector **w** via a regularized Fisher discriminant analysis as described in Philiastides et al. (2014): $\mathbf{w} = \mathbf{S}_c(\mathbf{m}_2 - \mathbf{m}_1)$, where the estimated mean of the difficulty level i is denoted as $m_i$ and $S_c = 1/2(S_1 + S_2)$ represents the estimated common co-variance matrix (i.e., the average of the empirical co-variance matrices for the two difficulty levels, $S_i = 1/(n-1)\sum_{j=1}^{n}(\mathbf{x}_j - \mathbf{m}_i)(\mathbf{x}_j - \mathbf{m}_i)^T$, *n* = number of trials). As suggested by Philiastides et al. (2014), we replaced the co-variance matrices for the two difficulty levels with regularized versions of these matrices to counteract potential estimations errors: $\tilde{S}_i = (1 - \lambda)S_i + \lambda\nu\mathbf{I}$, with $\lambda \in$ [0, 1] being the regularization term and $\nu$ the average eigenvalue of the original $S_i$ (i.e., $trace(S_i)/D$, with *D* corresponding to the dimensionality of our EEG space). It should be noted that $\lambda$ = 0 gives an unregularized estimation and $\lambda$ = 1 corresponds to spherical co-variance matrices. Here, we optimized $\lambda$ for each participant using leave-one-trial-out cross validation with the following $\lambda$ values $\in$ [0, 0.01, 0.02, 0.04, 0.08, 0.16], ($\lambda$ mean $\pm$ SD: 0.067 $\pm$ 0.072).

We then computed area under a receiver operating characteristic (ROC) curve (i.e., the $A_z$ value), using a leave-one-trial-out cross-validation procedure like the ones described in Diaz et al. (2017), Franzen et al. (2020) and Gherman and Philiastides (2018), in order to quantify the performance of the discriminator (Duda, Hart, et al., 2006). Specifically, for every iteration we obtained out-of-sample discriminant component amplitudes ($y(t)$) by using N-1 trials to estimate a spatial filter (**w**), which we then applied to the remaining trials. These out-of-sample amplitudes were then used to compute the $A_z$. It should be noted that as this procedure aimed to allow us to identify endogenous representations of EA, for a small number of participants (n= 4), where there were $A_z$ peaks with similar amplitudes, we identified the peak with the spatial distribution most consistent with previous accounts of EA. As per Diaz et al. (2017), Franzen et al. (2020) and Gherman and Philiastides (2018), instead of assuming an $A_z$ = 0.5 as chance performance, we determined participant-specific $A_z$ significance thresholds through a bootstrap

analysis whereby we randomized trial labels and used another leave-one-trial-out test. We repeated this randomisation procedure 500 times (at 100 ms after the presentation of the stimulus as this procedure has been shown to produce comparable $A_z$ value outputs along the full time window (e.g. Franzen et al., 2020)) and obtained a probability distribution for $A_z$, which we used as a reference to estimate the $A_z$ value leading to a significance level of $p < 0.05$.

Like in Gherman and Philiastides (2018), we then aimed to produce the full temporal profile of the relevant discriminating components ($y(t)$) by applying the spatial filter **w** of the window associated with the highest discrimination performance (i.e. we projected the data through the 'spatial filters' leading to the most reliable discrimination) to the entire stimulus-locked window (-100 to 800 ms post-stimulus). We performed this separately for each of the social and non-social domains as well as the three difficulty conditions (easy, medium and difficult; Fig. 12, c). We then z-scored the time courses of these discriminating components separately for each participant and for each of the social and non-social domains.

We were then able to investigate the gradual build-up of EA activity leading up to the point of peak discrimination and to extract the corresponding single-trial build-up rates, which we then used in subsequent analyses. These build-up rates (or slopes) were computed through a linear regression based on the data between the onset and peak time of the accumulating activity, which we derived on a participant-specific basis (see Gherman and Philiastides, 2015). Specifically, we identified when the discriminating activity began to rise monotonically after an initial dip in the stimulus-locked data following any early evoked responses (onset time mean $\pm$ SD: $363.097 \pm 97.046$ for social trials and $376.161\text{ms} \pm 107.155$ for non-social trials) and we similarly pinpointed when the accumulation activity reached its peak.

Similar to Franzen et al. (2020), we were able to compute scalp topographies reflecting the relevant discriminating components due to the linearity of our model from Eq. 10 through a forward model:

$$\mathbf{a} = \frac{\mathbf{X}\mathbf{y}}{\mathbf{y}^T\mathbf{y}} \tag{11}$$

For convenience, here we present the EEG data **X** and discriminating components **y** in matrix and vector notation, respectively. Equation (11) showcases the electrical coupling between the discriminating component **y** and the EEG activity in **X**. Specifically, strong coupling is linked to low attenuation of the component and can be displayed as the intensity of vector **a**. We estimated forward models of the discriminating activity separately for social and non-social trials (Fig. 12, b).

### 3.2.7 Single-trial regression and correlation for multivariate EEG analysis

We aimed to establish whether the EEG signal we would use to identify the site of EA truly reflected the process of EA leading up to the decision. To this end, we predicted that the slopes we identified through our multivariate EEG analysis as reflective of the EA buildup rate (i.e. $y(t)$) should correlate with our DDM drift rate estimates (Pisauro et al., 2017; Polanía et al., 2014), which were obtained by fitting the model on participants' fraction of 'Play' choices and RTs in the previous chapter (Fig. 10). However, since the 'Keep' and 'Play' responses were mapped to +1 and –1, respectively (see section on DDM from Chapter 2 for details) and thus positive drift rates corresponded to 'Play' and negative drift rates reflected 'Keep' choices, we flipped the sign of the EEG slopes in the two reward probability levels, which support 'Keep' choices to match the DDM sign convention (i.e. $P(payoff|play) = \{0 - 0.2, 0.2 - 0.4\}$).

Additionally, we explored the association between the rate of EA from the neural data (i.e. the slopes) and behavioral performance with a single-trial logistic regression analysis. To this end, we used the trial-specific estimates of the EEG slopes of the EA signal (i.e. $y(t)$) to predict the probability of playing (1: 'Play', 0: 'Keep') on individual trials (Fig. 13, c and d). Once again, this analysis was performed separately for each participant and for each of the social and non-social trials:

$$P_{play} = [1 + e^{-(\beta_0 + \beta_1 \times y(\text{buildup rate}))}]^{-1} \tag{12}$$

In this regression analysis we assessed whether the coefficients across participants ($\beta_1$ values in Eq. 12) came from a distribution with a mean different from

zero (using separate $t$-tests).

### 3.2.8 fMRI data acquisition

For the (f)MRI acquisition, we used a Siemens 3-Tesla TIM Trio MRI scanner (Siemens, Erlangen, Germany) with an 12-channel head coil, following the setup described in Gherman and Philiastides (2018). The functional volumes were acquired with a T2*-weighted gradient echo with an echo-planar imaging sequence (32 interleaved slices, gap: 0.3 mm, voxel size: 3 x 3 x 3 mm, matrix size: 70 x 70, FOV: 210 mm, TE: 30 ms, TR: 2000 ms, flip angle: 80°). We recorded 5 experimental runs, each consisting of 205 whole-brain volumes. We then acquired phase and magnitude field maps (3 x 3 x 3 mm voxels, 32 axial slices, TR=488 ms, short TE=4.92 ms, long TE=7.38 ms) to correct for distortion of the acquired EPI images. Finally, we took a high-resolution anatomical volume using a T1-weighted sequence (192 slices, gap: 0.5 mm, voxel size: 1 x 1 x 1 mm, matrix size: 256 x 256, FOV: 256 mm, TE: 2300 ms, TR: 2.96 ms, flip angle: 9°), which we used as an anatomical reference for the functional scans.

### 3.2.9 fMRI data preprocessing

We removed the first 5 volumes per run to guarantee a steady-state fMRI and we used only the remaining 200 volumes for the analysis, which were preprocessed in a similar way as in Pisauro et al. (2017). Specifically, we employed FMRIB's Software Library (Functional MRI of the Brain, Oxford, UK) to perform head-related motion correction, slice-timing correction, high-pass filtering ($>$100s), and spatial smoothing (with a Gaussian kernel of 5 mm full-width at half maximum). The motion correction preprocessing step produced motion parameters which we then included as regressors of no interest in the general linear model (GLM) analysis (see fMRI analysis below). We used the Brain Extraction tool (BET) to perform brain extraction of the structural and functional images. The echo-planar imaging data for each participant was transformed into the subject-specific high-resolution space using a BBR (boundary-based registration) algorithm. The images were then registered to standard space (Montreal Neurological Institute, MNI) using FMRIB's

Non-linear Image Registration Tool with a resolution warp of 10 mm and 12 degrees of freedom. Finally, to correct for signal loss and geometric distortions due to B0 field inhomogeneities, B0 unwarping was used for 29 out of 31 participants. Field map images were not acquired for the remaining 2 participants.

### 3.2.10 fMRI data analysis

We used the FEAT module in FSL to perform a whole-brain statistical analysis of the functional data using a multilevel approach within the framework of a GLM (Smith et al., 2004), like in Gherman and Philiastides (2018) and Pisauro et al. (2017):

$$Y = x\beta + \epsilon = \beta_1 X_1 + \beta_2 X_2 + ... + \beta_N X_N + \epsilon \tag{13}$$

Y corresponds to a voxel's time series (with T time samples). X is a T $\times$ $N$ design matrix, the columns of which correspond to the GLM regressors (described below) convolved with a double-$\gamma$ canonical hemodynamic response function. The regression coefficients are represented by $\beta$ - a $N \times 1$ column vector and the residual error terms - by $\epsilon$ - a T $\times$ 1 column vector. The first-level analysis examined each participant's individual runs, which were then combined using a second-level analysis (fixed effects). At the third level of the analysis, we examined group-level effects using a mixed-effects model (FLAME 1), treating participants as a random effect. We also used FMRIB's improved linear model with local autocorrelation to carry out the time-series statistical analysis (Woolrich et al., 2005).

In our GLM we included 4 regressors of interest per each of the social and non-social domain (i.e. a total of 8 regressors). Specifically, for each domain we included 1) an EEG-informed regressor with a parametric amplitude modulation based on the trial-by-trial fluctuations in the EA rate (i.e. trial-wise slopes in $y(t)$), 2) a parametric regressor with individual trial RTs amplitude modulation, 3) a parametric regressor with individual trial task difficulty amplitude modulation (-1: difficult, 0: medium, 1: easy) and 4) an unmodulated regressor (i.e. all amplitudes set to 1) to account for any additional variance in the data (Fig. 14, a). We modelled all regressor events as boxcar functions and set their duration to 100

ms. We aligned the first two regressors to the time of response and the last two - to the onset of stimulus presentation. We also used the unmodulated regressors to compute standard contrast and conjunction maps between social and non-social trials. Finally, we added the motion correction parameters obtained from fMRI preprocessing (three rotations and three translations) as additional co-variates of no interest.

It is noteworthy that the EEG-derived EA slopes reflect key information of the underlying processes, which cannot be substituted by simply relying on individual RTs. This is due to the decision and motor-planning stages being subject to inter-trial variability as shown consistently in previous modelling and experimental studies (Philiastides et al., 2014; Ratcliff et al., 2009; Verdonck et al., 2021) and further demonstrated here by the fact that the RTs and the slopes were not highly correlated (social: r = -0.297, non-social: r = -0.333). Therefore, the RTs did not present redundant information or pose a major confounding factor and we thus included separate RT predictors in our fMRI analysis to absorb additional non-EA activity (for each of the social and non-social domains separately).

Furthermore, the addition of EEG-informed predictors in our design allowed us to capitalize on using actual neural signals reflecting latent variability in information processing that might be overlooked when relying on simple behavioral or model-derived indices (see Gherman and Philiastides, 2018; Pisauro et al., 2017; Sajda et al., 2009). For instance, very few studies have included single-trial parameter estimates (Gluth et al., 2017; Turner et al., 2015) with most relying on mean estimates of the relevant decision variables derived from computational models. Here, instead, we employ the rate of EA on individual trials from the slope of the accumulating activity we identified in the EEG data, thus allowing us to account for true endogenous variability in EA and avoid issues related to model estimation or misspecification when deriving build-up rates solely from behavioural fits.

Importantly, we would expect that trials with lower EA rates (which would take longer to reach the relevant decision boundary) to have larger areas (energy) under the accumulation curve (Basten et al., 2010; Hare et al., 2011; Liu and Pleskac, 2011; Pisauro et al., 2017). Therefore, the site of EA should appear to be more hemodynamically active when trials have longer integration times (Fig. 14, b).

This should result in a negative relationship between our EEG-informed EA slope predictor and the BOLD response in the EA site (Hare et al., 2011; Liu and Pleskac, 2011; Mulder et al., 2014).

### 3.2.11  fMRI analysis resampling procedure

We determined a reliable significance threshold for the fMRI data by correcting for multiple comparisons via a resampling procedure, which examines a priori statistics of the trial-wise variability in the parametrically adjusted regressors (i.e. regressors 1–3 above) in a way that trades off cluster size and maximum voxel Z-score. We ran this procedure following the procedures outlined in Fouragnan et al. (2015) and Gherman and Philiastides (2018). Specifically, for each iteration of the resampling procedure we maintained the onset and duration of the regressors, however, we shuffled their amplitude values across trials, runs and participants. The resulting regressors for each participant were therefore different since they were made up of a random sequence of regressor amplitude events. We repeated this procedure 100 times across all 3 analysis levels (run, participant, and group). This allowed us to estimate a joint threshold for the cluster size and Z-score based on the cluster outputs per shuffled regressor. We did so by constructing a null distribution for this joint threshold based on the size of all clusters larger than 10 voxels and with Z-scores larger than $|2.57|$ (i.e. considering both positive and negative correlations) across all shuffled regressors. We discovered that the largest 5% of cluster sizes exceeded 88 voxels. We then used this value as a corrected threshold for our statistical maps, which we then applied to the clusters observed in the original data (that is, Z=$\pm$2.57, minimum cluster size of 88 voxels, corrected at p=0.05).

### 3.2.12  Psychophysiological interaction analysis

We then aimed to investigate the functional connectivity between the pMFC, which we found to correlate with the trial-by-trial variability in our EEG-informed regressor and the rest of the brain. To this end we conducted a psychophysiological interaction (PPI) analysis like in Pisauro et al. (2017). We first extracted

time-series data from the group-level activation clusters in the pMFC (seed), separately for each of the social and non-social contexts. This involved identifying the relevant pMFC clusters that were situated within the supplementary motor area (SMA) and were most consistent with previous reports of EA-related activity in this region (Pisauro et al., 2017). We then back-projected these clusters from the group (standard) space into the individual participant's EPI (functional) space (by applying the inverse transformations estimated during the main registration procedure). We then used the average time-series data from the back-projected voxels, which displayed activations in the direction of the predicted EA profile as the physiological regressor in our PPI analysis.

If any brain regions showcase a task-dependant co-activation with the pMFC, then the coupling between these regions should be stronger while the process of EA unfolds and it should also reflect the difficulty of the decision. To this end, we constructed our psychological regressor as a parametric boxcar regressor, the amplitude of which corresponded to the difficulty (1 = difficult, 2 = medium, 3 = easy) and the duration of which reflected the RT of each trial. Since easier trials decrease integration times and correspondingly the overall integrated activity (that is, area under the accumulation curve; Fig. 14, b), We predicted that the relevant coupling would be negative. We finally used the threshold derived from the resampling procedure described above to correct the resulting fMRI statistical maps.

## 3.3   EEG and fMRI results

### 3.3.1   Univariate single-trial EEG analysis highlighting comparable EA dynamics in social and non-social trials

A correlation analysis was run for each participant, which allowed us to identify the subject- and context-specific best electrodes with the highest correlations with the subject-specific EA traces (mean non-social r = 0.962±0.027; mean social r = 0.962±0.025, Fig. 11, a). No significant difference in the slopes of the best electrodes was observed between the social and non-social conditions $t(30) = 0.901$, $p = 0.375$ and all participants' best electrodes passed the permutation test

**Figure 11:** **Response-locked univariate EEG results.** **a)** Mapping between the average activity captured by the best-electrodes and the predicted EA activity, demonstrating the similarity between the social and non-social traces as well as the fit between the actual and the predicted data. The model EA traces were normalized by dividing by the maximum value of the subject-specific EA trace. The baseline was then corrected by removing the first sample from the 400 ms window. The baseline for the EEG data was offset based on the first 50 ms of the 400 ms subject-specific window. We then scaled the EEG activity to correspond to the EA activity. **b)** Correlation maps displaying the average correlation between the electrode-specific EEG traces and the EA traces. **c)** Activity from participants' best electrodes divided into easy and hard trials. Trials were considered 'easy' if they contained strong evidence for or against a certain choice (i.e. reward probabilities 0–0.2 and 0.8–1). Trials were considered 'hard' if they came from the middle reward probability band (i.e. 0.4–0.6). **d)** Activity from participants' best electrodes for fast and slow trials (defined through a median RT split).

($p<0.05$). Furthermore, a regression analysis between the subject-specific model-derived EA traces and the EEG EA traces from the subject-specific best electrodes demonstrated that model EA traces were able to predict the EEG traces (social r = 0.919±0.003; $t(12398)$ = 251.429, p <0.001; non-social r = 0.933±0.002; $t(12398)$ = 307.687, p <0.001).

To further investigate the EEG signal, the traces were divided into long and short trials (based on a median split, Fig. 11, d). Simple linear regressions based on the last 200 ms of the subject-specific time windows, which is where the differences in the slopes appear to be the most prominent, indicated that the shorter RT trials were associated with significantly steeper slopes (social $t(30)$ = -5.59, $p$ < 0.001; non-social $t(30)$ = -3.514, $p$ = 0.001). This is consistent with the accumulation-to-bound framework, where for choices with higher drift rates, the time to reach the decision boundary would be shorter. These trends were found in both the social and non-social domains. Additionally, easier trials were associated with significantly steeper slopes than the slopes seen in the difficult trials in both domains (Fig. 11, c) (again based on the final 200ms of the subject-specific time windows; social $t(30)$ = 3.058, $p$ = 0.001; non-social $t(30)$ = 2.856, $p$ = 0.008). These results are therefore consistent with the EA framework as with easier trials, where the evidence in favour of a certain choice is the strongest, the accumulation process would be easier, the drift rates should be higher and the slopes would be steeper.

Finally, the group topographies reflecting the correlation between the EEG activity and the EA traces was not only comparable between the social and the non-social condition (r = 0.962, $p$ < 0.001), but it also mirrored the previously seen centroparietal cluster consistent with the P300 phenomenon and evidence-to-bound processing (Fig. 11, b, e.g. Herding et al., 2019; O'Connell et al., 2012; Pisauro et al., 2017; Polanía et al., 2014).

### 3.3.2 Multivariate single-trial EEG analysis highlighting comparable EA dynamics in social and non-social trials

Here, we took advantage of the high temporal resolution of the EEG data to identify robust signals exhibiting a gradual build-up of activity consistent with EA dynamics (e.g. Gherman and Philiastides, 2018; Pisauro et al., 2017; Polanía et al., 2014). If such signals exist, we should be able to find reliable ramp-like activity with a build-up rate proportional to the decision difficulty. To test this hypothesis, we used a single-trial multivariate linear classifier (e.g. Gherman and Philiastides, 2018; Parra et al., 2005; Sajda et al., 2009) to discriminate between easy vs. difficult trials (see Materials and Methods). By focusing on the two extreme difficulty

Figure 12: **Linear discriminant analysis results. a)** Theoretical EA build-up rates for easy (yellow) and difficult (green) trials (top) and predicted classifier performance (Az) based on the differences in the rate of EA. **b)** Average discrimination performance (using a leave-one-out cross validation procedure) for easy and difficult trials across participants and histogram of subject-specific times of maximum discrimination (top). The dashed line represents the the average Az value leading to a significance level of $p = 0.5$, which was derived through a bootstrap test. The standard errors of the mean are shown with thinner black lines. Insets: Average scalp topographies (forward models) of the discriminating activity at subject-specific peak discrimination times for the social (red outline) and non-social (blue outline) trials. **c)** Stimulus-locked average temporal profile of the discriminating activity after applying the subject-specific classification weights estimated at peak discrimination times for the three levels of decision difficulty for social (red) and non-social trials (blue). Insets: subject-specific EA onset times for social (red) and non-social trials (blue). **d)** EA onset-locked average temporal profile of the discriminating activity, for the three levels of decision difficulty for social (red) and non-social trials (blue).

trials (easy and hard), we aimed to facilitate the discrimination process, tasked with identifying the sensor weightings that reflect the process of EA. We found that the classifier's performance (Az) increased systematically over time, reflecting the potential divergence in the build-up rates associated with easy and difficult trials (Fig. 12, b). On average, the classifier's performance started increasing 400 ms after stimulus presentation (i.e. after early encoding of the relevant evidence) and reached its peak several hundred milliseconds later.

The spatial distribution of this discriminating activity (i.e. forward model; see Material and Methods) from participant-specific windows at the time of maximum discrimination between easy and difficult trials (Fig. 12, b, top) revealed comparable centroparietal topographies across social and non-social contexts (r = 0.896,

*p* < 0.01; Fig. 12, b, inset). These similarities point towards a common neural generator across the two contexts, consistent with clusters reported previously across the non-social domain (e.g. Herding et al., 2019; Kelly and O'Connell, 2013; Philiastides et al., 2014).

Having identified the participant-specific spatial weights from the time window of maximum discrimination, we applied them to an extended stimulus-locked time window, separately for social and non-social trials, to formally characterize the temporal profile of the discriminating activity (i.e. $y(t)$) for each condition separately. This was also done for the medium difficulty trials (i.e., the "unseen" data) by subjecting the relevant data through the same neural generators responsible for the original discrimination.

This resulted in a gradual build-up of activity akin to a process of EA in both social and non-social trials (Fig. 12, c; top: social, bottom: non-social). Like the classifier performance, the neural activity began to rise around 400 ms after stimulus presentation in both the social and non-social trials, with the build-up rate reflecting the decision difficulty. Importantly, the build-up rate from the medium difficulty trials appeared between the two extreme conditions used to train the classifier, establishing a fully parametric effect across the three difficulty levels ($F(2, 90)=16.88$, $p<0.001$ for the social condition, $F(2, 90) = 26.76$, $p < 0.001$ for the non-social condition, post-hoc $t$-tests, all $p < 0.001$).

Finally, we identified participant-specific EA onset times as the time point at which the discriminating activity began to rise monotonically after an initial dip in activity following any early evoked responses present in the data (Fig. 12, c; insets). Since we observed some inter-individual variability in these onset times, we predicted that by re-aligning the relevant signals to the participant-specific EA onset times we would see a more pronounced depiction of the underlying process of EA at the population level, which was indeed the case (Fig. 12, d; top: social, bottom: non-social).

**Figure 13:** **Examining the link between EA dynamics and behaviour. a), b)** Positive correlation between subject-specific EA slopes (**y**(**t**)), i.e. build-up rate) for each of the five levels of $P(payoff|play)$ and DDM estimates of drift rate for both the social **a)** and non-social **b)** contexts. We normalized the EEG-derived EA signal **y**(**t**) used to derive the EA slopes per trial to factor out effects unrelated to the EA processing (e.g. attentional drifts). Black circles represent group averages. **c), d)** Positive correlation between EA slopes and the probability of playing (Eq. 7) for both the social **c)** and non-social **d)** contexts. For visualization purposes we divided the EA slopes into five bins. The black curves were created by fitting Eq. 7 to individual trials.

## Examining the link between EA dynamics and behaviour

Considering the importance of using a robust representation of the EA dynamics in our EEG-informed fMRI analysis, we aimed to further establish that our EEG signals from the multivariate EEG analysis reveal the process of EA. To this end we carried out two additional analyses to examine whether these signals reflect our participants' behavioral performance. We predicted that the build-up rate of these signals should correlate with drift rate estimates obtained from a DDM (Pisauro et al., 2017; Polanía et al., 2014) fit on participants' fraction of 'Play' choices and RTs (Fig. 10, see details in Chapter 2). We found robust correlations

between the slopes of the EEG and drift rates, across both social (Fig. 13, a; r = 0.653, *p* < 0.001) and non-social trials (Fig. 13, b; r = 0.709, *p* < 0.001). Note that we flipped the sign of the EEG slopes in the two reward probability levels which support 'Keep' choices (i.e. $P(payoff|play) = \{0 - 0.2, 0.2 - 0.4\}$) to align them to the EA slopes (see Methods).

We also used a logistic regression where the EA slopes directly predicted the probability of 'Playing' on individual trials (once again flipping the sign of the EEG slopes in the two lowest reward probability levels supporting 'Keep' choices). We expected that easy choices (both 'Keep' and 'Play') would be linked to high negative and high positive EA rates, respectively. We also predicted that medium difficulty trials would be associated with intermediate magnitude slopes and that difficult trials would be reflected by slopes near zero. Indeed, the EEG slopes significantly predicted the eventual probability of 'Playing' for both the social (Fig. 13, c; $t(30)$ = 7.582, $p < 0.001$) and non-social trials (Fig. 13, d; $t(30)$ = 8.173, $p < 0.001$).

We also tested whether the trial-by-trial changes in the EEG slopes may simply reflect fluctuations in attention (as it varies over the course of the experiment). To this end, we ran a linear serial autoregression model predicting the EEG-derived EA slope in the current trial from the slopes from the previous four trials, individually for all participants. This analysis, however only accounted for a very small portion of the overall variance in the EEG slopes (social: $R^2$ = 0.02, non-social: $R^2$ = 0.019), indicating the negligible serial autocorrelation in slopes between neighbouring trials.

### 3.3.3  pMFC reflects domain-general EA and shows functional connectivity with regions of the human valuation system

Both EEG analyses described above revealed the comparable EA dynamics in social and non-social choices, which may indicate that they share a common neural generator. To further examine the existence of a common EA site, we constructed EEG-informed fMRI predictors by taking advantage of the endogenous trial-by-trial variability in the slope of EA activity, which we derived through our multivariate EEG analysis (Fig. 14, a) and accounted for task-specific additional variance

**Figure 14:** **EEG-informed fMRI analysis. a)** Our fMRI GLM model contained two parametric boxcar regressors at the time of response reflecting the trial-by-trial variability in the EA slope (separately for the social and non-social condition). We accounted for the variance associated with other task-related processes by introducing 3 additional regressors per decision domain (6 in total): UNMOD - unmodulated boxcar regressors at the onset of the stimuli (amplitude set to 1) , DIFF - parametric boxcar regressors of task difficulty (-1: hard, 0: mid, 1: easy) at stimulus onset, and RT – parametric boxcar regressors with trial-by-trial RT modulations at the time of response. The duration of all regressors was set to 100 ms. **b)** Example EA activity with different slopes shown via coloured arrows. When these EA traces are convolved with a hemodynamic response function (HRF) we should see that longer integration times (larger gray areas under EA traces) would produce higher predicted fMRI activity. Correspondingly, we would expect higher predicted fMRI activity for shallower than for steeper EA slopes. **c)** In line with the predictions from **b)** we saw a negative association between the EEG-informed fMRI predictors of the EA slopes and the activity in the posterior medial-frontal cortex (pMFC) for both social and non-social trials. **d)** The activity in the pMFC co-varied in a task-dependant fashion with regions of the human valuation system i.e. clusters in the medial prefrontal cortex and the posterior cingualte cortex. All clusters shown correspond to mixed-effects activations that survived a $|Z| > 2.57$ threshold and a cluster-size correction ($P < 0.05$) using a resampling procedure (minimum cluster size = 88 voxels; see Materials and Methods). A complete lists of activations is available in Tables 1 and 2. vmPFC: ventromedial profrontal cortex; dmPFC: dorsomedial prefrontal cortex; vPCC: ventral posterior cingualte cortex. The theoretical panels are partially adapted from Pisauro et al. (2017).

| Region | Hemisphere | BA | Peak MNI coordinate d (mm) | | | Voxels | Z Value (peak) |
|---|---|---|---|---|---|---|---|
| | | | X | Y | Z | | |
| **_EEG social_** | | | | | | | |
| posterior medial-frontal cortex | B | 24 | 0 | -16 | 56 | 192 | -3.49 |
| pre/postcentral gyrus | L | 3\4 | -40 | -16 | 52 | 1028 | -4.45 |
| | R | 3\4 | 56 | -16 | 50 | 140 | -3.57 |
| | R | 3\4 | 30 | -18 | 62 | 113 | -3.32 |
| **_EEG non-social_** | | | | | | | |
| posterior medial-frontal cortex | R | 24 | 8 | 8 | 44 | 965 | -4.39 |
| pre/postcentral gyrus | L | 3\4 | -38 | -22 | 50 | 1686 | -4.43 |
| | R | 3\4 | 48 | 8 | 36 | 393 | -3.7 |
| | R | 3\4 | 34 | -4 | 52 | 1355 | -4.68 |
| | R | 3\4 | 60 | -16 | 46 | 95 | -3.25 |
| lateral occipital complex | R | 17 | 28 | -86 | -10 | 433 | -3.75 |
| | L | 17 | -36 | -74 | -12 | 579 | -4.1 |
| superior temporal gyrus | R | 41 | 48 | -22 | 16 | 163 | -3.84 |
| | L | 41 | -58 | 4 | 12 | 162 | -3.96 |
| inferior frontal gyrus | R | 44 | 56 | 12 | 0 | 97 | -3.78 |
| superior parietal lobule | R | 7 | 20 | -62 | 50 | 105 | -3.26 |
| insular cortex | R | 16 | 34 | 24 | 8 | 91 | -3.52 |

**Table 1:** List of brain activations associated with the social and non-social EEG predictors.

by including RT, task-difficulty and unmodulated predictors for both decision domains (Fig. 14, a). It should be noted that even though the RTs were not highly correlated with our EA slopes as the decision and motor-planning stages are associated with inter-trial variability (Philiastides et al., 2014; Ratcliff et al., 2009; Verdonck et al., 2021), we still included them as nuisance predictors. Importantly, we would expect the accumulator region(s) to show a negative association with our EEG-informed predictors (Hare et al., 2011; Liu and Pleskac, 2011; Mulder et al., 2014), given that trials with lower EA rates have larger areas under the accumulation curve as they take longer to reach the decision boundary (Basten et al., 2010; Hare et al., 2011; Liu and Pleskac, 2011, Fig. 14, b). We found a region in the posterior medial frontal cortex (pMFC), consistent with previous observations (Pisauro et al., 2017) that satisfied this negative correlation criterion (Fig. 14, c) in both social and non-social domains, thus being in accordance with the notion of a universal EA neural framework. We observed that the social choices cluster was more posterior compared to non-social choices cluster, however, there were no domain-

specific activations surviving the direct EA slope predictor contrast between the two domains.

We then aimed to test whether the pMFC was coupled in a task-dependent way with brain regions encoding the relevant decision evidence. We were also interested in establishing whether such a coupling would highlight domain-general or domain-specific neural representations of the decision evidence. This involved running separate PPI analyses for each of the social and non-social trials, using the domain-specific pMFC clusters as seeds and the trial-specific task difficulty as the psychological predictor (see Materials and Methods). Once again, we predicted that the relevant task-specific co-activating regions should show a negative relationship with the pMFC, since easier trials should lead to shorter integration times and thus to a decreased overall integrated activity under the EA trace (Fig. 14, b).

Through this PPI analysis we found that the pMFC coupled negatively (by task difficulty) with regions of the valuation system (posterior cingulate cortex (PCC) as well as in dorso- and ventro-medial prefrontal cortex (dmPFC/vmPFC)) for both decision domains (Fig. 15). These regions along the medial wall of the brain have consistently been implicated in converting value signals into a 'common neural currency' subsequently used in EA (Pearson et al., 2014; Piva et al., 2019; Rangel and Hare, 2010) and recent resting state connectivity reports have highlighted their negative relationship with the pMFC (Neubert et al., 2015). Similar to the EA clusters, the social activations were situated more posterior relative to the non-social ones, consistent with previous reports about the organization within the prefrontal cortex (Chib et al., 2009; Clithero and Rangel, 2014; Smith et al., 2010). However, once again, there were no domain-specific activations surviving the direct contrast between social and non-social PPI predictors. These findings therefore support the notion that relevant decision evidence may be converted into a 'common currency' in the human valuation system and then accumulated in the pMFC to reach a decision.

### 3.3.4 Difficulty and early domain-specific processing

In our experiment, the task difficulty was determined as a function of the reward probability following a 'Play' choice, with easy trials corresponding to reward

| Region | Hemisphere | BA | Peak MNI coordinate d (mm) | | | Voxels | Z Value (peak) |
|---|---|---|---|---|---|---|---|
| | | | X | Y | Z | | |
| **PPI social** | | | | | | | |
| dorsomedial prefrontal cortex | | 10 | -4 | 48 | 26 | 276 | -3.5 |
| ventromedial prefrontal cortex | L | 10 | -2 | 52 | 2 | 133 | -3.19 |
| posterior cingulate cortex | R | 7 | 2 | -70 | 24 | 113 | -3.54 |
| middle frontal gyrus | R | 9 | 26 | 24 | 40 | 818 | -3.73 |
| | L | 9 | -26 | 36 | 42 | 526 | -3.8 |
| pre/postcentral gyrus | R | 3\4 | 20 | -24 | 68 | 687 | -4.23 |
| | R | 3\4 | 56 | 0 | 22 | 114 | -3.24 |
| | L | 3\4 | -50 | -14 | 38 | 366 | -3.65 |
| frontal pole | R | 10 | 22 | 58 | -10 | 496 | -3.8 |
| inferior frontal gyrus | R | 45 | 56 | 22 | 14 | 108 | -3.65 |
| | L | 45 | -52 | 16 | 14 | 97 | -3.46 |
| **PPI non-social** | | | | | | | |
| dorsomedial prefrontal cortex | L | 10 | -4 | 52 | 24 | 135 | -3.39 |
| ventromedial prefrontal cortex | R | 10 | -8 | 62 | 0 | 134 | -3.74 |
| posterior cingulate cortex | R | 7 | -4 | -66 | 22 | 329 | -3.88 |
| posterior medial-frontal cortex | R | 24 | 2 | -26 | 62 | 209 | -3.81 |
| pre/postcentral gyrus | R | 3\4 | 38 | -26 | 62 | 2379 | -5.02 |
| | L | 3\4 | -44 | -12 | 28 | 686 | -4.16 |
| | L | 3\4 | -20 | -26 | 70 | 259 | -3.77 |
| middle frontal gyrus | R | 9 | 24 | 24 | 40 | 858 | -4.15 |
| | L | 9 | -30 | 34 | 38 | 475 | -4.56 |
| superior temporal gyrus | R | 19 | 62 | -50 | 18 | 1826 | -4.08 |
| supramarginal gyrus | L | 19 | -58 | -46 | 42 | 301 | -3.82 |
| middle temporal gyrus | L | 19 | -54 | -56 | 12 | 97 | -3.21 |

Table 2: List of brain activations associated with the social and non-social PPI analyses.

probability bins with strong evidence in favour of an alternative (both 'Keep' and 'Play', i.e. 0-0.2 and 0.8-1), medium difficulty trials associated with reward probability bins favouring an alternative, but not as strongly as the easy trials (i.e. 0.2-0.4 and 0.6-0.8) and difficult trials with ambiguous reward probabilities, which do not clearly favour either alternative (i.e. 0.4-0.6). Therefore, the task difficulty predictor should correlate positively (i.e. easy > difficult) with areas responsible for encoding choice value. It should also correlate negatively (i.e. difficult > easy) with regions of the attentional network linked to the processing of task demands. Indeed, the conjunction of the two task difficulty regressors revealed positive correlations with regions of the valuation system such as the vmPFC, ventral stria-
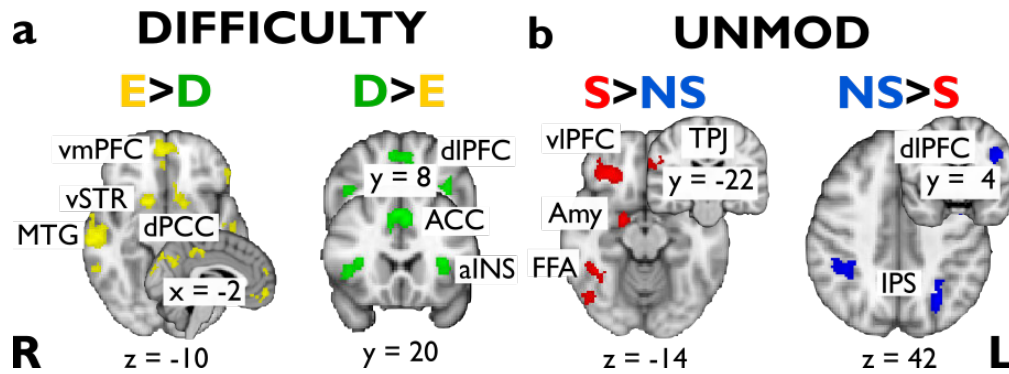
**Figure 15:** **Task difficulty and domain-specific fMRI activations. a)** Activations showing greater BOLD response for easy (E) than difficult (D) trials (yellow) and those exhibiting higher response for difficult compared to easy trials (green). These activations arise from the conjunction of the two task difficulty regressors (DIFF) for social and non-social trials in Fig. 14, **a)**. **b)** Activations showing greater BOLD response for social (S) than non-social (NS) trials (red) and those exhibiting higher response for non-social compared to social trials (blue). These activations arise from the contrast of the two unmodulated regressors (UNMOD) in Fig. 14, **a)**. All clusters represent mixed-effects activations that survived $|Z| > 2.57$ and that were cluster-corrected ($P < 0.05$) using a resampling procedure with a minimum cluster size of 88 voxels (see Materials and Methods). The complete list of activations is shown in Tables 3 and 4. vlPFC: ventrolateral prefrontal cortex; Amy: amygdala; FFA: fusiform face area; TPJ: temporoparietal junction; IPS: intraparietal sulcus; vmPFC: ventromedial prefrontal cortex; vSTR: ventral striatum; MTG: medial temporal gyrus; dPCC: dorsal posterior cingulate cortex; dlPFC: dorsolateral prefrontal cortex; ACC: anterior cingulate cortex; aINS: anterior insula.

tum and the PCC (Clithero and Rangel, 2014; Domenech et al., 2018, Fig. 15, a). Conversely, regions previously implicated in the encoding of uncertainty and attentional control such as the anterior cingulate cortex, lateral prefrontal cortex and anterior insula (Grinband et al., 2008; Monosov, 2017; Philiastides and Sajda, 2007) displayed a negative correlations with our difficulty regressors (Fig. 15, a). The fact that we could identify these prototypical results from across the decision-making literature demonstrates the validity of our fMRI analysis design and suggests that it offers a reliable account of the relevant experimental manipulations.

We further examined the contrast between the unmodulated predictors for the social vs. non-social trials to identify brain areas responsible for the processing of early domain-specific representations (i.e. prior to their conversion to a 'common currency', Fig. 14, a). The regions showing higher activations during the social trials were the right fusiform gyrus, right amygdala and right ventrolateral prefrontal cortex, which all belong to the so-called 'face network' (Skelly and Decety, 2012)

| Region | Hemisphere | BA | Peak MNI coordinated (mm) | | | Voxels | Z Value (peak) |
|---|---|---|---|---|---|---|---|
| | | | X | Y | Z | | |
| **Easy>Hard** | | | | | | | |
| ventromedial prefrontal cortex | R | 10 | 6 | 54 | -8 | 323 | 3.75 |
| dorsomedial prefrontal cortex | R | 10 | 12 | 44 | 14 | 213 | 3.81 |
| striatum | R | 25 | 16 | 10 | -10 | 142 | 3.88 |
| | L | 25 | -14 | 8 | -12 | 111 | 3.54 |
| posterior cingulate cortex | L | 7 | -2 | -28 | 46 | 198 | 3.55 |
| lateral occipital complex | L | 17 | -36 | -74 | 38 | 3369 | 4.72 |
| | L | 17 | -14 | -96 | 24 | 1244 | 4.12 |
| middle temporal gyrus | R | 21 | 64 | -20 | -4 | 2723 | 4.66 |
| | L | 21 | -58 | 2 | -18 | 113 | 3.48 |
| superior parietal lobule | R | 7 | 22 | -54 | 72 | 433 | 3.75 |
| | L | 7/1 | -22 | -52 | 72 | 329 | 3.66 |
| inferior frontal gyrus | L | 45 | -56 | 22 | 14 | 363 | 3.8 |
| middle frontal gyrus | L | 9 | -30 | 26 | 46 | 154 | 3.74 |
| frontal pole | R | 10 | 16 | 68 | 14 | 89 | 3.37 |
| **Hard>Easy** | | | | | | | |
| anterior cingulate cortex | L | 24 | -2 | 16 | 48 | 604 | -5.09 |
| anterior insula | R | 13 | 32 | 20 | 10 | 183 | -3.76 |
| | L | 13 | -32 | 22 | 8 | 148 | -3.91 |
| dorsolateral prefrontal cortex | L | 46 | -38 | 4 | 24 | 228 | -4.11 |
| | R | 46 | 44 | 8 | 24 | 90 | -3.34 |

Table 3:    List of brain activations associated with the Difficulty predictors.

as well as the right temporoparietal junction (TPJ, Fig. 15, b). These findings suggest that the social trials required the processing of facial features necessary for mentalizing and examining the opponent's intentions, respectively (Cerniglia et al., 2019). We also found that the lateral intraparietal cortex bilaterally, as well as the left dorsolateral preforntal cortex (dlPFC, Fig. 15, b), which have been previously implicated in encoding risk and reward probabilities in non-social contexts (Burke and Tobler, 2011; Daw et al., 2006; Smith et al., 2009) displayed stronger activations for non-social trials. We argue that these patterns showcase that our paradigm was successful in capturing social and non-social processing as we would not expect to see these domain-specific differences if our participants did not rely on social information in the social trials and if the non-social trials were not associated with reward probability considerations. These results, combined with the trends seen in the behavioural results and the observation that in the social condition only a small number of inconsistent faces was removed (on average 10.807 face display removed, with on average 23, 34, 32, 36 and 14 faces per bin, respectively) suggests that our paradigm was largely successful in capturing the social and non-social variability involved in these decisions.

## 3.4 Discussion

Establishing whether social decisions employ the same universal decision-making mechanism involved in non-social choices could offer deeper insight into human behaviour. However, despite recent attempts to investigate whether the two decision-making domains employ the same processes across the various decision-making levels, many questions still remain. The previous chapter highlighted the potential algorithmic similarities between social and non-social choices. In this chapter, we aimed to further explore their implementational underpinnings and used simultaneous EEG-fMRI to examine whether the spatiotemporal neural characteristics involved in social choices correspond to the ones underlying non-social decisions.

To this end, we initially aimed to uncover signatures of EA dynamics through a univariate analysis of the EEG evoked responses. As a result, we used our model estimates from the previous chapter to create model-derived EA traces, which

| Region | Hemisphere | BA | Peak MNI coordinated (mm) | | | Voxels | Z Value (peak) |
|---|---|---|---|---|---|---|---|
| | | | X | Y | Z | | |
| *Social>Non-Social* | | | | | | | |
| amygdala | R | 25 | 18 | -2 | -16 | 105 | 4.48 |
| fusiform gyrus | R | 37 | 40 | -48 | -16 | 107 | 5.18 |
| temporoparietal junction | R | 13 | 64 | -42 | 26 | 285 | 3.39 |
| ventrolateral prefrontal cortex | R | 11 | 30 | 34 | -14 | 359 | 5.42 |
| lateral occipital complex | R | 17 | 48 | -60 | -2 | 1027 | 5.26 |
| | L | 17 | -48 | -74 | 4 | 113 | 3.77 |
| superior parietal lobule | L | 7 | -6 | -58 | 68 | 647 | 4.22 |
| inferior frontal gyrus | R | 45 | 52 | 32 | 2 | 246 | 4.13 |
| occipital pole | L | 18 | -14 | -102 | 16 | 117 | 3.55 |
| *Non-Social>Social* | | | | | | | |
| intraparietal sulcus | L | 5 | -28 | -62 | 46 | 225 | -4.26 |
| | R | 5 | 38 | -42 | 42 | 136 | -3.63 |
| dorsolateral prefrontal cortex | L | 4 | -38 | 0 | 32 | 240 | -4.12 |
| lateral occipital complex | R | 17 | 24 | -86 | -6 | 111 | -3.82 |

**Table 4:** List of brain activations associated with the Unmodulated predictors.

we compared to our EEG data. This comparison revealed group topographies, which were not only comparable between the social and the non-social condition, but they also reflected the previously seen centroparietal cluster observed in VBDM (Pisauro et al., 2017), PDM (Gherman and Philiastides, 2015; Mostert et al., 2015; O'Connell et al., 2012; Philiastides et al., 2014; Twomey et al., 2016) and in memory-based choices (van Vugt et al., 2019). It has also been shown that positive amplitudes recorded from these electrode sites are associated with the subjective impression of the available evidence and that this signal seems to be linked to confidence and (pre)motor choice-specific signals (Herding et al., 2019), thus showing the relevance of this cluster to EA dynamics. Interestingly, the observation that the two topographies were highly correlated also provided an early indication as to the possibility that these signals might arise from the same cortical source.

The model-derived EA traces also allowed us to identify subject-specific best electrodes, which we then examined to further investigate these signature of EA. Firstly, we established that there was no significant difference between the social and non-social average best electrode traces. Furthermore, slopes across both do-

mains reflected the speed and the difficulty of the trial - with shorter reaction times and easier trials being linked to steeper slopes than longer reaction times and more difficult trials - which have been benchmark observations in a number of decision-making studies (e.g. Connolly et al., 2009; Law and Gold, 2008; Mazurek et al., 2003; Roitman and Shadlen, 2002). These trends are once again consistent with the decision-making process model from the Introduction (Fig. 3). Therefore, these EEG results demonstrate that our data was able to capture signs of EA and even offer evidence in favour of a universal decision-making mechanism. Furthermore, by identifying these EA dynamics through correlations with our model predictions, these EEG results provide a neural validation of our modeling by suggesting that it was able to reflect the same neural computations that gave rise to these evoked responses.

Although our model-informed EEG analysis was successful in demonstrating that EA dynamics were captured in our EEG data, we needed to ensure that we obtain a robust representation of these dynamics in order to utilize them to uncover the site of EA. We therefore took advantage of the multidimensional nature of our EEG signal and we spatially integrated this information to increase the relevant signal-to-noise ratio and thus produce more robust single-trial estimates of EA (Parra et al., 2005). This procedure gave rise to even more convincing representations of the trends seen in the univariate EEG analysis, with more localized centroparietal clusters, which were still comparable between the two domains, and even more pronounced parametric effects of our difficulty modulation across both the social and non-social choices. To further validate these signals we demonstrated that the EA slopes we derived from this multivariate EEG analysis corresponded to the drift-rate estimates from our modelling and that they were able to predict our participants' choice behaviour. Having validated these slopes we then used them as internally-generated indices of EA and utilized their trial-wise variability to identify a region in the pMFC as responsible for this process across both choice contexts, thus providing empirical evidence in favour of a common domain-general integration of evidence for the decision occurring in the pMFC.

This finding - that social and non-social choices share an EA site - is not only in favour of the 'common currency' schema, but also provides an extension to this framework. Specifically, the 'common currency' schema typically examines

whether the motivational processes for social and non-social choices are implemented in the same regions, with different connectivity patterns for the social and non-social trials, and it often focuses exclusively on considering the value processing stages of decision making (Ruff and Fehr, 2014). Here we show evidence that the commonalities between social and non-social decisions extend to the embodiment of the choice since the EA processing was contained within the same pMFC structure, which encompasses pre-motor areas, even though the social cluster was more posterior than the non-social cluster. However, since there were no significant clusters associated with the direct contrast of these two activations, we suggest that the two decision domains were largely processed the same way in this region.

These results are therefore in accordance with views postulating that when decisions are made under time pressure, they are embodied in the same sensorimotor areas responsible for carrying out the choice - a notion which has found support through a number of neuroimaging studies (e.g. Donner et al., 2009; Filimon et al., 2013). Specifically, such (pre)motor areas have been shown to reflect the updating of the relevant information (Kolling et al., 2016) and inactivating these areas has been found to lead to major behavioural impairments in simple discrimination tasks (Peixoto et al., 2021; Wu et al., 2020). Additionally, a recent resting state MRI experiment linked the centroparietal activity with the dorsal frontoparietal network where the pMFC is located (Brosnan et al., 2020). Specifically, it suggested that the centroparietal build up rate is associated with the connectivity between the structures of this network and with the speed of decisions. In other words, the higher the build up, the stronger the connection and the faster the decisions. These studies linking (pre)motor areas to EA dynamics have led to the development of computational accounts proposing that (pre)motor areas are responsible for an active 'motor accumulation', which receives integrated information from non-embodied regions and proceeds to transfer this information to the relevant motor system once its internal threshold is reached (Steinemann et al., 2018; Verdonck et al., 2021).

We also found further implementational evidence in favour of the 'common currency' schema through our PPI analysis which showed task-specific co-activation between the pMFC and areas of the human valuation system along the medial wall

- specifically, the vm/dmPFC and the PCC, all of which have previously been implicated in value processing both in the social (Fehr and Krajbich, 2014; Gusnard et al., 2001; Heatherton et al., 2006; Lee and Seo, 2016; Mar, 2011; Suzuki and O'Doherty, 2020) and the non-social domain (Coricelli and Nagel, 2009; Grueschow et al., 2015; Nagel et al., 2018). We nevertheless found a degree of anterior-posterior dissociation linked to the prefrontal activations, consistent with recent findings (Kolling et al., 2021). Specifically, we found that the social activations were more posterior (Ferrari et al., 2016; Jenkins and Mitchell, 2011; Lieberman et al., 2019) relative to the non-social cluster (Chib et al., 2009; Clithero and Rangel, 2014; Smith et al., 2010). However, since the direct contrast between the two PPI activations did not show any significant clusters, this suggests that the two decision domains were processed in a largely comparable way within these structures.

Further in accordance with the 'common currency' schema, which suggests that the domain-specific information is converted into a 'common currency' (Hutcherson, Montaser-Kouhsari, et al., 2015; Lim et al., 2013; Ruff and Fehr, 2014) as well as a recent review, which suggests that complex social choices involve value computations from various social sources ultimately converted into a common value signal along the medial wall of the prefrontal cortex (Suzuki and O'Doherty, 2020), we found domain-specific activations at the time of stimulus presentation. Specifically, the social trials were associated with activations in the 'face network' i.e. the FFA, the amygdala and vlPFC (most of which displayed a right lateralization), all of which have been implicated in face identification and processing (Garvert et al., 2014; Vuilleumier et al., 2004). These trials were also linked to activity in the right TPJ, which has been shown to be involved in social cognition and mentalizing, which require the extraction of the relevant information carried by the opponent/co-player's face (Morishima et al., 2012; Van Overwalle, 2009). On the other hand, for the non-social trials, which required the consideration of the likelihood of receiving a reward, presented explicitly during the task, we saw activations in the lateral intraparietal cortex which has been shown to reflect pure reward probabilities both in humans (Daw et al., 2006; Wu et al., 2015) and in primates (Burke and Tobler, 2011; Sugrue et al., 2004).

This experiment also offered prototypical results seen across the non-social

decision-making literature, thus demonstrating the validity of our task since it successfully produced the desired behavioural and neural responses we were interested in. These results also justified our choice of analysis, as they were able to uncover these well-known neural patterns. For instance, in relation to our difficulty manipulation we found that the easy trials were linked to areas, which have typically been associated with value encoding and risk assessments, such as the striatum (Yacubian et al., 2007), mPFC and the PCC (Levy et al., 2010). We also found that the more difficult choices were associated with brain areas typically linked to uncertainty, in keeping with a recent view suggesting that the choice uncertainty affects decisions (Urai et al., 2017). These areas included the ACC, which has been shown to respond to unknown vs known probabilities (Hayden et al., 2011), the dlPFC which is known to respond to stress (Bogdanov and Schwabe, 2016) as well as the insula, which has been shown to be associated with the assessment of risk (Levy et al., 2010; Mohr et al., 2010), feelings, empathy and uncertainty (Singer et al., 2009).

Finally, in addition to providing initial evidence in favour of a shared mechanism for social and non-social choices spanning multiple decision-making levels, this experiment fits in with the wider literature suggesting the potential universality of decision making as stipulated by standard economic utility models (Caplin and Glimcher, 2014; Von Neumann and Morgenstern, 2007). Specifically, within the non-social domain it has shown that VBDM follow the same integrative mechanisms as PDM (Krajbich et al., 2010; Mormann et al., 2010) demonstrated by the fact that they show similar centroparietal clusters (Pisauro et al., 2017) and that parietal gamma frequencies capture similar aspects of EA across VBDM and PDM (Polanía et al., 2014). Therefore, even though there has not been a direct comparison investigating the implementational levels between these two non-social sub-domains, these studies provide initial evidence suggesting that they might share the same algorithmic processes. As a result, our current results provide an additional piece of evidence supporting the potential existence of a universal decision-making process across its valuation and EA stages.

In summary, in this chapter we aimed to elucidate the fundamental implementational processes involved in decision making. Specifically, we found that the social and non-social information, encoded in domain-specific areas, is then con-

verted into a 'common currency' in the value circuit along the medial wall of the human brain.  This information is then accumulated in the pMFC - a pre-motor structure adjacent to the relevant motor area - for both decision domains.  These results therefore extend the scope of the 'common currency' schema and demonstrate that the commonalities between domains go beyond the valuation stages of decision making by specifically highlighting the potential embodied nature of the process.  Overall, these findings offer detailed implementational insight into the omnipresent strategic choices involved in every day life.

# 4 Chapter 4, Experiment 2

## 4.1 Introduction

In the last two chapters we investigated whether the neurocomputational principles of choices made based on social information match the ones identified within the non-social decision-making literature. We found evidence suggesting that the two might share common algorithmic principles (i.e. the mechanistic principles involved in these choices) and that they might also be implemented within the same general brain structures, thus providing support in favour of a universal decision-making mechanism across decision-making domains. In the previous chapter we indicated the key role the pMFC might play in the implementation of social choices. However, these results were only correlational in nature. Therefore, they were only suggesting that the activity in the pMFC corresponds to the endogenous activity reflective of an EA process. In order to thoroughly investigate the role of the pMFC, stimulation methods are needed to establish the causal role of the regions in the identified process. While continuing to compare social and non-social choices in this investigation might have allowed us to examine the degree of similarity between the social and non-social decision processing even further by studying whether the pMFC is causally involved in both process, here we focused instead on a different aspect of this investigation. Specifically, we still aimed to offer further implementational understanding of the function of the pMFC by causally implicating it in EA, but we also strived to offer further algorithmic insight by asking which aspect of the EA dynamics this region is involved in. In order to maximise our capabilities of answering these questions we opted to increase the number of trials, but we also wanted to avoid including more experimental sessions than would be conceivable that our participants would agree to take part in. To this end, we had to focus on one domain and we chose to examine social decision making in order to demonstrated the capacity of this domain to offer further insight into basic neurocomputational principles of decision making, which typically are initially investigated in non-social choices.

The results from the previous chapter suggest that the pMFC (a region of the cortex extending from the SMA to the mid-cingulate cortex) reflects the rate of

EA - the drift rate at which a stochastic accumulator integrates information until it reaches one of the two boundaries that correspond to the choice alternatives. We implicated the pMFC in this process by demonstrating that its activity matched the slopes we identified from the EEG signal reflective of these EA dynamics. This notion is consistent with recent EEG-fMRI work demonstrating that the activity in the pMFC might represent the full temporal dynamics of the process of EA itself (Pisauro et al., 2017). However, one of the main structures making up the pMFC - the SMA - has previously been thought to be responsible for boundary adjustments. In other words, in the context of the general accumulation-to-bound framework, it has been suggested that this region modulates the amount of information needed to make a choice (Bogacz, Wagenmakers, et al., 2010; Forstmann et al., 2008). This view has emerged from literature varying the goal of the task by making participants either focus on making speedy or accurate choices. This speed-accuracy trade-off literature has shown that when tasked to make accurate choices, people take longer, which in computational models is reflected in increased boundary separation. Conversely, when people make quick responses, the distance between the boundaries is reduced and therefore an erroneous choice is more likely to occur. To establish whether the pMFC simply reflects the rate of EA, modulates the boundary separation or if it is involved in both processes, here we used transcranial direct current stimulation (tDCS) in a pre-registered experiment to target the pMFC and provide causal evidence for the functioning of this brain region in EA.

tDCS is a widely used neurostimulation method, which non-invasively and temporarily changes the excitability of a brain region (anodal - increases; cathodal - decreases excitability) to examine its function and relation to behaviour (Nitsche et al., 2008). This stimulation method has garnered some scepticism regarding its efficacy with regards to the unreliability of some of the published effects (Horvath et al., 2015, but also see Antal et al., 2015) and because of the frequent violations of underlying anatomical consistency assumptions (Kim et al., 2014). Nevertheless, tDCS remains popular since the inconsistent results might partially be attributed to the variability in design (Horvath et al., 2015), choice of participants (Berryhill et al., 2014) and the inconsistent protocols (Wiethoff et al., 2014). Furthermore, tDCS has been successfully used to investigate social decision making before. For example, it has been employed to investigate if brain stimulation to the right lateral prefrontal cortex would lead to changes in social norm compliance

in both voluntary and sanction-induced scenarios (Ruff et al., 2013). The neural stimulation led to changes in both types of compliance, however it affected them differently - anodal tDCS led to greater compliance in sanction-induced scenarios, whereas the cathodal stimulation led to greater compliance in the voluntary scenario. Even though the stimulation changed the rate of compliance, it did not alter the understanding of the social norms, nor that of the sanctions. This study therefore demonstrated the usefulness of using a combination of anodal and cathodal stimulation to examine social choices and shows that the comparison between these stimulation methods may be used to uncover the exact functioning of the pMFC.

Specifically, here we used anodal and cathodal tDCS stimulation (Fig. 1, right) to test three main hypotheses: 1) the pMFC reflects the rate of EA, 2) the pMFC reflects boundary adjustments and 3) the pMFC governs multiple EA processes, whereby it reflects EA dynamics as well as boundary adjustment (Fig. 16). This was achieved by fitting the behavioural performance of our participants to a hierarchical DDM (hDDM; Wabersich and Vandekerckhove, 2014) and examining the effects of the stimulation through our hDDM parameters. Even though we also investigated the effects of the stimulation on behavioural indices (Proportion 'Play' choices, RT data), this experiment primarily made use of the hDDM's latent variables derived by considering both choice and RT data together. This joint consideration of both behavioural indices in model space offers a more sensitive measure of subtle changes in behaviour, and the model definitions also provide a mechanistic understanding of the changes induced by the stimulation.

We first examined the hypothesis most in line with the observations from the previous chapter and recent literature (Pisauro et al., 2017) - that the pMFC modulates the rate of EA (Fig. 16 left). If this hypothesis is correct, the cathodal stimulation should interfere with EA dynamics and would make the decision process noisier. This should be reflected in lower drift rates and/or higher drift rate variability (e.g. Polanía et al., 2015) as specified by a DDM. As participants might already be performing optimally, the anodal stimulation might not have any facilitatory effects and thus may not result in any parameter changes. However, if this condition reduces the noise in the process, we may still observe that it reduces drift rate variability. Nevertheless, we should observe a relative difference
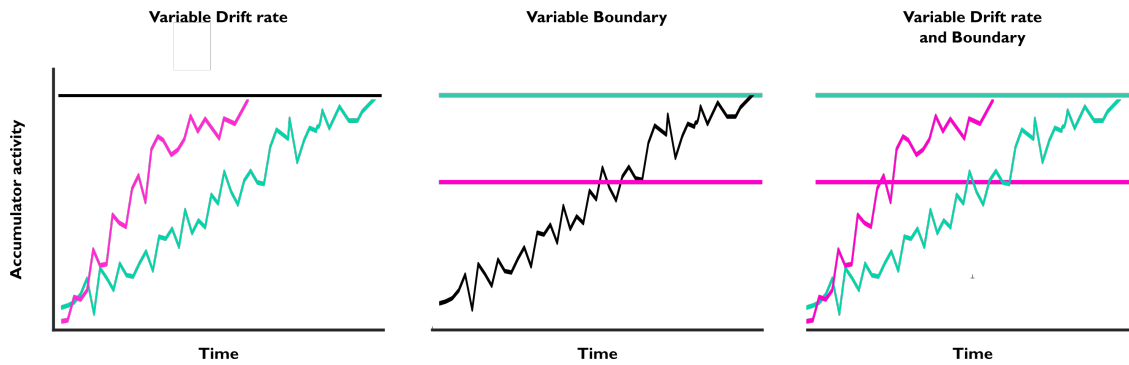
Figure 16: **Graphical representation of the Drift rate, Boundary and Combined hypotheses for the function of the pMFC.** The figure shows three ways in which the anodal (pink) and cathodal (turquoise) stimulation might affect the pMFC. The **left** panel illustrates an example in which the two stimulation conditions are associated with differences in drift rate, with the cathodal stimulation being linked to a lower drift rate than the anodal condition. The **middle** panel illustrates an example in which the two stimulation methods lead to variability in the boundary parameter, with the anodal condition being linked to a lower boundary than the cathodal condition. Finally, the **right** panel illustrates the possibility that the stimulation might lead to changes in the drift rate as well as to the boundary separation.

in the drift rate and/or drift rate variability parameter estimates of the two stimulation conditions with the anodal stimulation being associated with higher drift rates or less drift rate variability compared to the cathodal condition. We also considered whether the pMFC is responsible for threshold adjustments (Fig. 16, middle, e.g. Bogacz, Wagenmakers, et al., 2010; Forstmann et al., 2008). This hypothesis posits that the anodal stimulation should lead to lower boundary separation and the cathodal stimulation should produce higher threshold values in the hDDM. Finally, we examined the possibility that the pMFC is responsible for both EA and boundary adjustment (Fig. 16, right). Therefore, we would expect to see changes in both the drift rate/drift rate variability and the boundary parameters across the stimulation sessions as detailed above.

It is worthwhile noting that our hypotheses refer to the relative differences between the anodal and cathodal conditions as we are not including a sham condition. Nonetheless, these differences are still able to elucidate the functioning of the pMFC in social choices and allow us to provide a mechanistic account for this process by demonstrating that changes to the excitability of the region might lead to changes in the underlying computations as evidenced by the changed parameters. Our preliminary model comparison results were most consistent with the

third hypothesis implicating the pMFC as an accumulator region by suggesting that it reflects both the rate of accumulation, but also adjusts the boundary separation based on how much information is needed for a decision. However, the formal comparisons of the parameters across stimulation sessions were not significant, suggesting that more evidence may be needed to elucidate the role of this region in social decision making.

## 4.2 Materials and Methods

This experiment, its hypotheses and methods have been pre-registered:

(https://osf.io/y4jsk)

### 4.2.1 Sample size and participants

We calculated the sample size (45 typical individuals) by referring to a previous experiment, which used a DDM to asses the functioning of a brain area in decision making using stimulation methods (Philiastides et al., 2011). Since the experiment on which we are basing our sample size used transcranial magnetic stimulation, we penalized their achieved effect size by 0.1, as we would expect tDCS to produce smaller effects (Priori et al., 2009). We thus calculated the necessary sample size using the following parameters: effect size = 0.455, alpha = 0.05, power = 0.9. However, due to low yield rates during the pandemic, it was not possible to reach this sample size. Therefore, the results included in this chapter only serve as an illustration of the possible conclusions that could be drawn once the full sample size is achieved. However, these conclusions should be considered with caution as the trends outlined in this chapter may change once the full dataset is collected. As a result, for this experiment we collected data from only 24 participants. Two of these participants were excluded as they did not meet the criterion that their behavioural performance during the anodal condition should represent a sigmoid showing a parametric effect on the behaviour reflective of the ratings (See exclusion criteria bellow). We also note that one of our exclusion criteria included the removal of participants with higher nDT estimates in the anodal condition than the

last 100 trials from the training session. This was true for 2 participants, however as these differences were very small (i.e. 0.027 and 0.007s) and given the already limited number of participants, for the purposes of this thesis we are retaining their data. Therefore, the analysis included in this chapter is based on a sample of 22 (Collected between 02/06/2022 and 22/02/2022; analysis completion date: 08/03/2022).

Participants were recruited through the University of Glasgow subject pool and through word of mouth. We invited both male and female, right-handed, Caucasian participants (to avoid biases in facial perception e.g. Scott and Monesson, 2009) aged 18-35, with no neurological conditions, no tDCS contraindications and with normal/corrected to normal vision to take part in the experiment. We originally used a fixed payment rate of £30 for the participation in the experiment, but due to issues with recruitment, we increased our payment rate to £40. However, participants were originally told that they would receive £9 per tDCS session for their participation and up to an additional £9 based on their performance (i.e. up to between £18-£36 in total). For the increased rate, we told them that they would receive £18 per tDCS session and up to £5 based on their performance (i.e. up to between £36 and £46). We did not tell our participants about the fixed rate in order to encourage them to engage with the task.

The experiment included the following exclusion criteria:

Presence of tDCS contraindications (Rossi et al., 2009); Low tolerance for the tDCS stimulation, which requires the termination of a session; absence of a sigmoid curve after visually inspecting the online behavioural data or all responses falling within less than 0.3 of each other; absence of a sigmoid curve after visual inspection of the anodal condition; participants not paying attention as evidenced by increased nDT estimates during the anodal condition relative to the last 100 trials of the training block, as the longer early sensory processing might be indicating reduced top-down influences of attention; highly inconsistent ratings (>50 items rated more than 1.5 bins apart, calculated as follows: if one of the 2 ratings falls within the two extreme bins (0-0.3 or 0.7-1), the 1.5 bin critical difference would be set to 0.35 and if not, it will be set to 0.2; Insufficient ratings in a reward probability bin (i.e. 0 faces in a bin following the removal of inconsistent

ratings); Different baselines across the two conditions determined by whether a single sigmoid curve based on the combined choice behaviour across the two sessions explains the data as well as two separate curves.

### 4.2.2 Stimuli

The stimuli included in this experiment are the same 150 photorealistic face displays based on 131 face identities with 19 artificially generated faces as the ones used in the online rating session and the EEG-fMRI session as part of the Main EEG-fMRI experiment, however, in this experiment only the social displays were used.

### 4.2.3 Procedure for online pre-tests

Participants who applied for the experiment were sent a link to an online rating task and a shortened version of the economic game (100 trials, Fig. 1, right). We did this in order to ensure that the participants we invited for the tDCS sessions understood and were willing to engage with the task in that their behaviour reflected their ratings. The rating task was the same as the one employed during the online rating session from the Main EEG-fMRI experiment. Participants saw the 150 face displays, which were presented centrally on a screen as 3x4 images that took up 30% of the height of the screen using the online functionality of PsychoPy3 (v2020.1, Peirce, 2007, Fig. 8, b). Underneath each face identity was a scale ranging from 0-100%. Participants were once again told that the face identities belonged to individuals who took part in a previous experiment investigating social attitudes, whereupon they were asked to participate in a version of the Trust game as Trustees (Berg et al., 1995, Fig. 4, a). The goal of the participants was to select a likelihood (in percent) that the person whose face display they were viewing would return half of an investment (2 points) that was given to them by selecting the section of the scale that corresponds to the perceived likelihood of splitting.

The online version of the economic game involved 100 trials of a Trust game, which worked that same way as in the previous experiment, but without the non-

social trials and without a variable fixation cross. Specifically, during each trial, a fixation cross (1.25s) was followed by the presentation of a face stimulus centrally on the screen, which was displayed in the same way as in the rating task described above. On each trial participants had the option between keeping a point ('Keep') or playing ('Play') to double their profit. They would receive 1 point if they chose the 'Keep' option. Conversely, if they selected the 'Play' option, the probability of doubling the reward (P(reward|Play)) was guided by the likelihood estimate each participant gave during the rating task for each face identity. Participants were, however, once again told that the outcome depended on trials selected at random from when the individual whose face they were viewing played the game. The stimulus remained on screen until the participant made a response or for a maximum of 1.3 seconds. Since the online task was done remotely, participants were asked to use the arrow keys on their machine to make their responses. Finally, a feedback screen would show their accumulated score after every trial for 0.75 + 1.3 - response seconds.

Only participants with satisfactory behavioural performance (see exclusion criteria) were invited to the first tDCS session. One day before the first session participants were sent the online task again, which allowed us to identify stimuli with inconsistent ratings (more than 1.5 bins apart), which were subsequently removed. Following the second online task, the face stimuli were also divided into 2 sets - one for each session. This was done to avoid possible learning effects carrying over from one session to the next.

### 4.2.4  Procedure for lab sessions

At the start of each tDCS session we would place the electrode over the pMFC on the participant's head (See tDCS setup). This was done in order to give the participants a chance to get used to having an electrode attached to their scalp, so as to minimize new experiences during the tDCS stimulation. Each tDCS session consisted of two parts - a training segment and a main tDCS segment. We presented both using PsychoPy3 (v2020.1.3, Peirce, 2007), with the face displays shown as 3x4 images that took up 30% of the height of the screen (refresh rate = 60Hz). At the start of the training segment, participants were shown 5 screens summarizing how

they rated the stimuli that would be presented during that session (grouped into 5 categories: highly unlikely/somewhat unlikely/ambiguous/somewhat likely/highly likely to split an investment). These 5 categories reflected the face rating bins: 0-0.3, 0.3-0.433, 0.433-0.566, 0.566-0.7, 0.7-1. We changed the bin sizes for this experiment in an attempt to equalize the number of faces across bins. On average, per session there were 16, 14, 10, 16, 10 unique face identities across the five reward probability bins, respectively and on average we would remove 16.8 out of the 150 faces across both sessions. Following the summary pages, participants took part in the training segment, in which they played a Trust game similar to the economic game they played online (400 trials, separated into 4 blocks - 100 trials per block, with 4 short 24.6-second breaks after each block). However, they would not receive a feedback screen following every trial and instead they would only see a fixation cross padding up until the 1.3 s mark. This was done in order to minimize the overall run time of the task and to keep the total duration length maximally consistent across participants. We included the training segment to ensure performance stabilisation and to minimize additional learning effects or changes in choice strategy for the main tDCS segment. The number of trials in the training segment is based on the data from the EEG-fMRI session from the Main experiment presented in the previous chapter. Specifically, we found that after fitting a Weibull function to the proportion of 'Play' choices as a function of perceived trustworthiness in the subject-specific social choices, the slope of the sigmoid stabilized on average after 170 trials. We then increased the number of trials to 400 to enable all participants' slopes to stabilize.

Following the training session, we would place the second electrode on the participant's cheek - to give them an opportunity to take a break from the task. Participants would then receive one of the tDCS stimulation conditions (either anodal or cathodal - the allocation of which was counterbalanced across participants) and participate in the main tDCS segment. The main segment was virtually identical to the training, but it had a longer duration of 30 minutes (648 trials). This included 6 blocks - 108 trials each, with 6 short 24.6-second breaks after each block. The tDCS stimulation would be active during the first 20 minutes of the main tDCS segment. However, participants were not told that the stimulation would end at that time and that they would continue performing the economic task for the remaining 10 minutes.

We chose this mixed approach (which includes a stimulation and a no stimulation period), firstly, due to the inconsistent accounts regarding the tDCS after-effects duration, with some accounts claiming they last up to several hours after minimal stimulation (Nitsche et al., 2008), others showing peak effects occurring 10 (Carlsen et al., 2015) or 20 (Bender et al., 2017) minutes after the end of the stimulation and some proposing that the tDCS effects do not persist past the end of the stimulation (Tremblay et al., 2016). Secondly, this approach is further justified by the inconsistent literature regarding whether the stimulation should occur exclusively prior to the task (e.g. Bender et al., 2017; Carlsen et al., 2015) or whether it should be active ('on') during the task (e.g. Greinacher et al., 2019). Consequently, our main analysis focused on trials starting after the first 10 minutes of the stimulation in order to minimize the effects of unexpected tDCS sensations at the start of the stimulation, which may be distracting to our participants, and to also allow for enough time for the tDCS effects to arise. The potential sensation effects are also why this experiment does not include a sham condition, since it has been suggested that participants can reliably detect when the stimulation is 'on' (Greinacher et al., 2019). Finally, our choice of stimulation duration was governed by accounts suggesting that stimulation lengths exceeding 20 minutes could result in the opposite effect than what is anticipated (Monte-Silva et al., 2013).

During the breaks, participants were asked three questions related to the stimulation sensation. Specifically, we asked them: "Is the stimulation on?", "Are you experiencing any tingling sensation?" and "Are you experiencing any burning sensation?". The first question required a 'yes' or 'no' answer given by the left or the right arrow key, respectively. The other two questions required a rating from (1) 'Not at all' to (5) 'Very strongly', given by the respective number key. These tDCS-related questions were also presented during the training so that the participants could familiarize themselves with the questions, the response buttons and how long they had to give their answers. At the end of the breaks - after responding to the three tDCS questions, participants would see their aggregated score (both in the training and main tDCS segments).

Following the first session, participants would be invited for a second session no earlier than 48h following the first one (5.909 days on average). During the second tDCS session, participants were shown their rating summaries for the face

identities they would encounter during the second session. They would then participate in the training segment and would subsequently receive the other stimulation condition, while performing the main economic segment for a final time. Everything apart from the stimulation condition was kept identical between stimulation sessions. Following the end of the second stimulation session participants were also asked in which session they thought they performed better in order to get an indication whether they could accurately perceive the effects of the anodal vs the cathodal stimulation. They were also allowed to indicate that they could not perceive any differences between the two sessions.

### 4.2.5 tDCS setup

This experiment involved anodal and cathodal tDCS stimulation +/- 2 mA. We selected the electrode placement based on a SimNIBS simulation (Fig. 17), in which we aimed to optimize the current intensity going through the region, while keeping the pMFC electrode close to the pMFC. In the anodal condition we would place the active electrode (circular, 3.4 cm diameter) between Cz and CP1 from the 10-20 international electrode placement to target the region identified by (Pisauro et al., 2017) and the EEG-fMRI experiment presented in the previous chapter. The return electrode (circular, 4.5 cm diameter) would be placed on the right cheek, on the intersection of a line going between the lower tip on the right ear and the nose, and FT8 and FT10 on the 10-20 international system. We would switch the active-passive allocation of the electrodes for the cathodal condition. We kept the stimulation 'on' during the first 20 minutes of the task. This included a 30-second ramp up at the start and end of the stimulation (19 minutes of full activity). We would space the two sessions out by a minimum of 48h and we counterbalanced the stimulation types over the two sessions across participants.

### 4.2.6 Hierarchical DDMs specifications

Our behavioural data consisted of RTs and choice behaviour (probability of a 'Play' choice, i.e. (P(Play|choice))) measured during the two stimulation sessions. The probability of a 'Play' choice is determined for each reward probability given
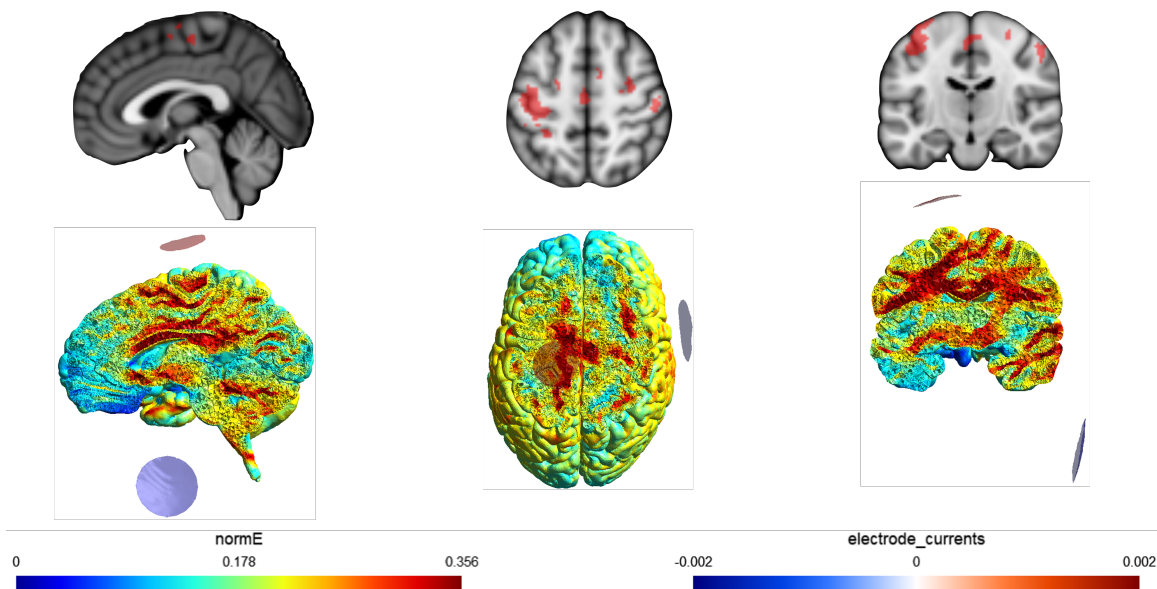
117

**Figure 17:** **fMRI activations from the pMFC for social choices and SimNIBS output.** This graphic presents the set up used for the anodal condition, with the active electrode presented in red and the return electrode presented in blue.

'Play' bin (P(reward|Play)): very unlikely (0-0.3), somewhat unlikely (0.3-0.433), ambiguous (0.433-0.566), somewhat likely (0.566-0.7) and very unlikely (0.7-1) to split an investment, and per subject. Therefore, we calculated a P(Play) from the 20 minutes of interest based on the data from the 10th minute of the stimulation until the end of the economic task.

We used conventional $\chi^2$ tests based on the last 200 trials of the two training segments to determine which participants would be eligible for the main analysis. Specifically, we examined whether a single Weibull function would fit the combined data from both sessions as well as using separate curves for each session in order to assess if there were any baseline differences between the two sessions. We performed these $\chi^2$ test on individual subject data in a similar way as in the chapter 2, but instead of basing the $\lambda$ calculation on the average values, we calculated them per participant, to establish whether there were behavioural differences in the anodal and cathodal 'Play' choices. All collected participants, satisfied this criterion (mean $\lambda$ = 0.978, $p$ = 0.623).

Since tDCS is associated with small effects and has often led to unreliable behavioural results (Horvath et al., 2015; Priori et al., 2009), our main analysis focused on the parameter effects seen in the model space as these differences might
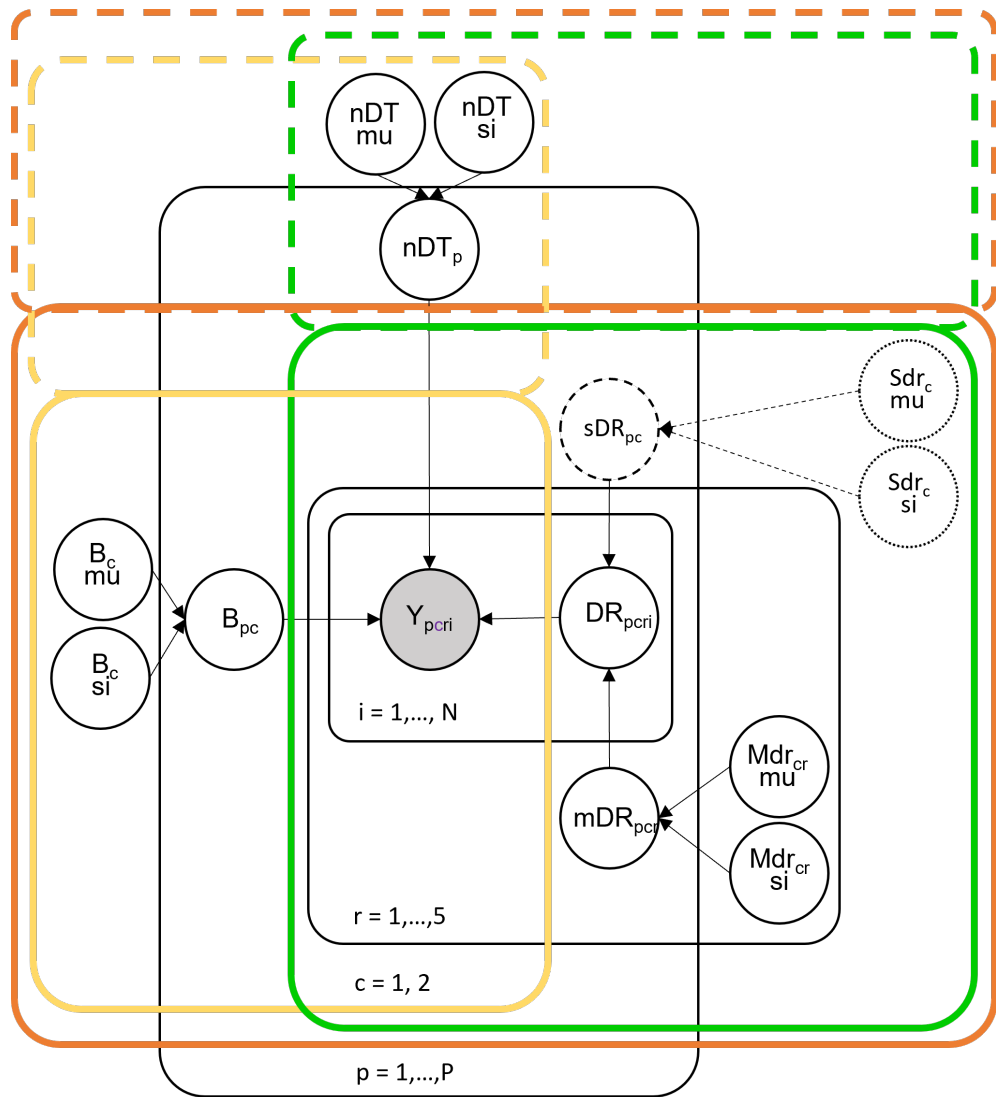
**Figure 18:** **Graphical representation of hDDM parameters.** Round white nodes represent continuous random variables ($nDT$ = non-decision time; $B$ = boundary; $sDR$ = drift rate variability; $mDR$ = mean drift rate; $DR$ = drift rate). The gray node ($Y$) represents the observed behaviour (RTs and choice behaviour). Parameters are modeled as random variables produced through means ($mu$) and variances ($si$). The rectangular plates signify that the random variables included in the plate are allowed to vary according to the stimulation condition. The colours on the plates denote the different models (yellow - variable boundary; green - variable drift rate/drift rate variability; orange - combined model with variable boundary and drift rate/drift rate variability). The dashed lines on the plates represent the three control analyses, where the $nDT$ is allowed to vary between the anodal condition of the main tDCS segment and the associated training segment according to the original three main models. The dashed nodes and lines signify that originally we allowed the $sDR$ to vary as a continuous random variable with a mean = $Sdr.mu$ and a variance = $Sdr.si$, However due to convergence issues, we fixed $sDR$ = 0.25 across both conditions and for all participants. p = participant, c = condition, r = reward probability bin, i = trial

be more sensitive to the subtle changes produced by the stimulation, since we are taking advantage of the joint choice-RT space. We fit the behavioural data from the last 20 minutes of the main economic game to three DDM families (Fig. 18), based on our pre-specified hypotheses (Fig. 16) using a hierarchical DDM (hDDM) implemented through the JAGS wiener module, via the R-package 'rjags' (version 4.10) interface in a similar fashion to Franzen et al. (2020). hDDMs assume the same stochastic accumulation-to-bound process as traditional DDMs and similarly produce estimates of internal components of processing. However, unlike most traditional DDMs, which either predict that subject-specific parameters are entirely independent from each other (by fitting the data to individual subjects) or that the parameters are the same across subjects (by fitting only to the group), the hierarchical framework simultaneously estimates parameters on an individual and on a group level i.e. at different hierarchical levels (Wiecki et al., 2013). This hDDM allows the specification of a drift rate (rate/slope of EA), boundary (threshold for decision), bias towards an alternative (starting point) and non-decision time (nDT). The specified parameters are adjusted through Markov chain Monte Carlo (MCMC) sampling so as to maximize the summed log-likelihood of the choice behaviour and RTs, provided by the Wiener first-passage time distribution. Such a hierarchical Bayesian framework for DDM parameter estimation involves updating the prior distributions for the model parameters based on likelihood of the data given the model, which ultimately produces a posterior distribution (Kruschke, 2010; Wabersich and Vandekerckhove, 2014; Wiecki et al., 2013). One major benefit of the Bayesian hDDM is its flexibility in allowing the specification of additional variables to feed into the main model parameters, thus facilitating a more complete mechanistic view of the decision-making process. Furthermore, the posterior distributions can be directly used to assess the quality of the parameter estimation and more stable parameter estimates for individual participants can be estimated since all obserevations are assumed to be drawn from a group (Wiecki et al., 2013).

The first of the three model families (Fig. 16 and Fig. 18) aimed to address the variable drift rate hypothesis and thus allowed the drift rate ($mDR$) and/or drift rate variability ($sDR$) to vary with the stimulation conditions. For this model we kept the boundary parameter ($B$) fixed across stimulation conditions. The second model family reflected the boundary adjustment hypothesis, which involved

us varying $B$ according to the stimulation condition. Here, $mDR$ and $sDR$ did not vary per stimulation condition. The final model family combined the previous two and we allowed $B$, $mDR$ and $sDR$ to vary across stimulation conditions. Across all three model families (Fig. 18), $mDR$ and $sDR$ varied with respect to the probability of reward given a 'Play' choice. We drew the parameters' means and variance from non-informative uniform distributions ($Mdr.mu \sim$ U(-5.00, 5.00), $Mdr.si \sim$ U(0.0001, 3.00); $Sdr.mu \sim$ U(0.0100, 7.00); $Sdr.si \sim$ U(0.0001, 1); $B.mu \sim$ U(0.0100, 3.00); $B.si \sim$ U(0.0001, 2.00); $nDT.mu \sim$ U(0.0100, 0.70); $nDT.si \sim$ U(0.0001, 0.25)) and we fixed the starting point to 0.5. Due to convergence issues (i.e. $\hat{R} > 1.05$) we had to use a fixed parameter for $sDR$ = 0.25. We chose these values based on models run on the simultaneous EEG-fMRI dataset. We then created the parameter estimates by drawing from normal distributions with these estimates serving as means and by converting the variances into standard deviations. For all three models, we followed the procedure outlined in Franzen et al. (2020) and used 5 Markov chains, each consisting of 5500 samples (created using the 'coda' R package (version 0.19-1)). Following standard, MCMC protocols (Wabersich and Vandekerckhove, 2014), we treated the first 500 samples as 'burn-in' since these early samples may be influenced by the starting value and we thus discarded them. Furthermore, since neighbouring samples are likely to be highly correlated, we additionally subsampled ('thinned') our data by a factor of 50. We based our subsequent parameter estimations and analyses on the remaining samples.

To assess the convergence of our chains, we used the Gelman Rubin Rhat ($\hat{R}$) statistic, which assesses within- and between-chain variance. We concluded that the traces converged if the $\hat{R} < 1.05$. This was true for all parameters after we fixed the sDR parameter. We also performed formal model comparisons in order to arbitrate between the different hypotheses using the deviance information criterion (DIC), which is widely used to assess model fits and for model comparisons, because it takes into consideration both the goodness-of-fit of the model and its complexity (Spiegelhalter et al., 2002). The better performing models with the highest likelihood and least degrees of freedom are associated with the lowest DIC values.

For the winning model, we directly tested for differences between the relevant

parameter(s) using a paired *t*-test between the anodal and cathodal conditions. Since our design does not include a sham condition, the stimulation effects we observed were only relative and without a clear indication of whether one condition drives the effect more than the other. Nevertheless, these relative differences were sufficient to provide a mechanistic account of the pMFC. However, since our formal statistical analysis was not significant, we conclude that more evidence is needed in order to determine the role of the pMFC in EA.

### 4.2.7 Control analyses: Learning effects, attention

We additionally wanted to examine potential learning effects by comparing participants' choice behaviour (P(Play|choice)) from the first session versus the second session. Given that we counterbalanced the session order across participants and that we included extensive training, we did not expect to observe any significant effects following this control analysis. To this end, we used linear regressions like the ones described for the Choice probability calculations section in Chapter 2. We ran these regressions separately for the first and the second stimulation session as well as for each participant. We further fitted the main three models, but allowed the parameters to vary depending on the session number to further investigate potential differences arising in the parameter space.

To avoid the possibility that the observed effects are due to changes in attention, which would potentially result in longer early processing during low attention sessions, we also considered models with variable nDTs and investigated potential changes between the anodal stimulation segment and the last 100 trials from the associated training segment. Specifically, we might expect the cathodal condition to cause the participants to perform worse, however if the choice behaviour is worse during the anodal condition, it might signal that participants were not paying attention, which should also be reflected in bigger nDT estimates.

### 4.2.8 Exploratory analyses: temporal, RT/choice data and sensation effects

In addition to our main analysis, for which we specified a time window during the last 20 minutes of the main tDCS segment, we also ran the same three sets

of hDDMs (as described above) after adjusting the analysis window, both in terms of its onset relative to the start of the main tDCS segment as well as in regards to its duration. The aim of these additional analyses was to investigate potential differences between the sessions during three additional time windows, for which we included trials from: 1) the full duration of the economic game (30 minutes), 2) the stimulation section of the tDCS segment (first 20 minutes) and 3) after the end of the stimulation (last 10 minutes).

While our main analysis focused on the latent DDM variables, we also examined potential relative differences between the two stimulation sessions, by relying on the formally tested behavioural data (choice behaviour and RTs) using linear regressions akin to the ones described in the Choice probability calculations section in Chapter 2. We ran these regressions separately for the anodal and the cathodal stimulation conditions as well as for each participant. We also performed conventional $\chi^2$ tests between the anodal and cathodal trials from the main analysis window to assess whether a single Weibull function would fit the combined data from both sessions as well as using separate curves for each session. These tests based on the purely behavioural measurements were done in order to examine whether the trends emerging from the modelling were also observable on the behavioural level.

We also compared the responses to the sensation questionnaires between the anodal and cathodal conditions to assess whether there were differences in the duration or the severity of the tDCS-related sensations. We specifically examined whether the strength and the duration of these sensations, or the beliefs regarding whether the session involved anodal or cathodal stimulation had any effect on the model parameters. In relation to the duration of the tDCS-related sensations, we examined whether the total number of blocks during which the participants thought the stimulation was 'on' differed between the two sessions. In regards to the sensation severity, we averaged the severity scores for the tingling and burning sensations per participant separately for the two sensation types and then averaged those values in order to get an aggregate measure of the sensation severity for each participant. We then compared these averaged values between the two sessions (anodal and cathodal). We also considered differences in the duration and severity of the tDCS-related sensations arising during the first half of the

main economic game, while the stimulation was still 'on', and the second part of the stimulation, during most of which the stimulation was switched off (the stimulation was still 'on' during only the first block of the second half). Finally, we examined the modeling parameters from the winning main model in relation to how the participants assessed their own performance. In other words, we rearranged the estimated parameters in accordance with the sessions in which the participants thought they performed better and the ones in which they thought they performed worse. We investigated both the boundary and drift rate parameters in relation to these 'better' and 'worse' categories. Since 8 participants could not perceive any differences in their own performance, we randomly labeled one of their stimulation sessions as the 'better' one. We performed this assignment procedure 1000 times, while running paired *t*-tests based on the 'better' vs 'worse' sessions and we took note of the associated *t*-values. We then calculated the average t-values for the boundary and the drift rate parameters and the associated p-values.

We also planned to assess cases in which the behaviour during the training sessions was inconsistent. We planned to do so in two ways - the first one was to directly correct the behavioural output, by re-adjusting the RTs and choice behaviour in the second session. This readjustment would be based on the differences found during the final 200 trials of the training between the two stimulation sessions. This would have been achieved by subtracting the value of the mean RT and choice difference per bin between the two sessions. The second version of this readjustment was meant to take place directly on the parameter space, by adjusting the parameters themselves by the mean difference per bin in the estimates. All other aspects of the analyses would have been kept the same. However, we did not perform these additional analyses as there were no participants with inconsistent training segment performance.

## 4.3 Modelling and additional results

### 4.3.1 Modelling results

Our main analysis focused on examining potential differences in the DDM space, since these latent variables may be more sensitive to subtle changes produced by
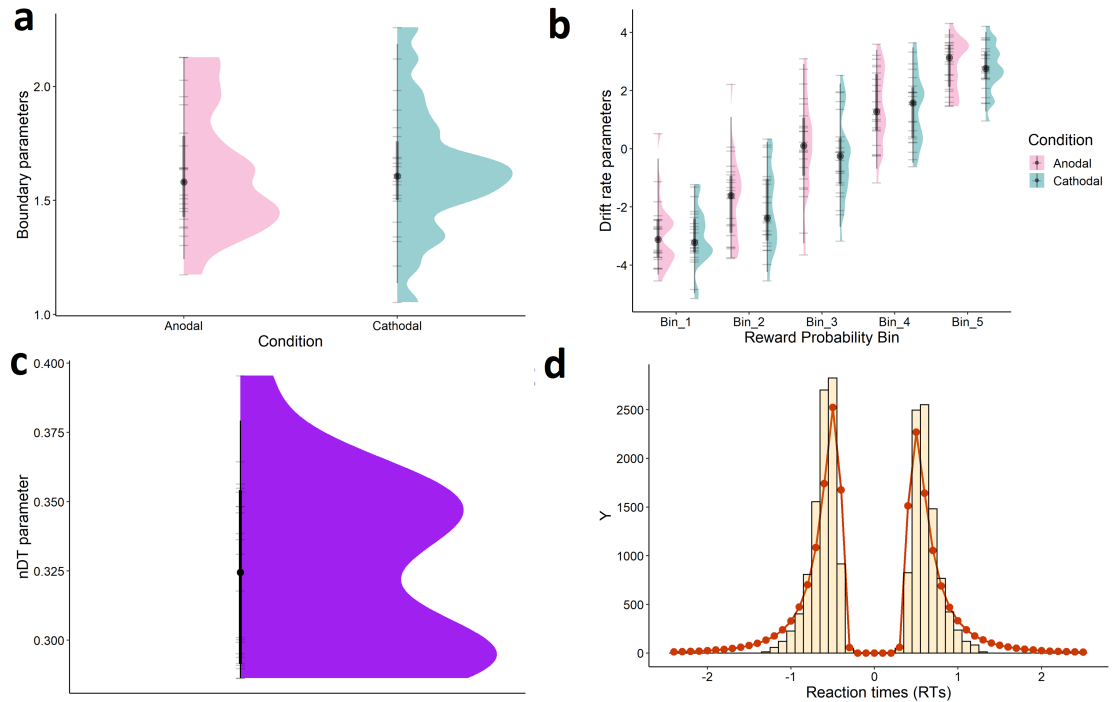
**Figure 19:** **Modelling results. a)** Subject-specific boundary parameter estimates raincloud plots for the anodal (pink) and cathodal (turquoise) conditions showing slightly lowered boundary estimates for the anodal condition. **b)** Subject-specific drift rate parameter estimates raincloud plots per reward probability bin for the anodal and cathodal conditions, showing that on average the anodal condition was associated with slightly higher drift rates. **c)** Subject-specific nDT parameter estimates raincloud plot, showing that the nDTs were very closely distributed in the 0.28-0.4 range. **d)** Histogram of the observed RT distributions ('Keep' choices were assigned negative RTs) and hDDM model fits (orange line) demonstrating that the model fit our data well. In all raincloud plots each participant estimate is presented as a horizontal line. The black circle shows the median value and the thicker line shows the interquartile range and the thinner line shows the 95% data range. The shaded area shows the density distribution of the estimates.

the stimulation by taking advantage of the combined choice-RT space. We therefore employed three different hypothesis-driven hierarchical variants of the traditional DDM, which has been consistently used to examine rapid decision making and offers mechanistic insight into the underlying processes (Ratcliff and McKoon, 2008; Wiecki et al., 2013).

Traditional hDDMs employ a number of parameters (drift rate, boundary separation, nDT and choice bias) to describe task performance. By obtaining such hDDM parameter estimates for both the anondal and cathodal stimulation condition we were able to identify the internal changes produced by the stimulation conditions and consequently establish the role of the stimulated region (pMFC) in

social decision making.

We specifically created three model families based on predictions for the functioning of the pMFC (Fig. 16 and 18). The first one exclusively varied the drift rate parameter while keeping the boundary constant across conditions to examine the variable drift rate account postulating that the activity in the pMFC reflects the rate of EA. The second model kept the drift rate constant across stimulation conditions, but varied the boundary separation so as to allow us to examine the boundary adjustment hypothesis, which suggests that the role of the pMFC is to determine the amount of information needed for a choice. Finally, we combined both models by varying both the drift rate and the boundary parameters across stimulation conditions in order to examine the possibility that the pMFC may be an accumulator region reflecting both the rate of EA and the necessary amount of information needed for the choice. For all of these models we allowed the drift rate to vary with the reward probability given a 'Play' choice. We initially allowed the drift rate variability to vary as well, but due to convergence issues, we fixed it to 0.25 (see Methods).

We used the DIC to arbitrate between these models and we established that the best-fitting model was the one, which allowed for both the boundary and the drift rate to vary across stimulation conditions (DIC Boundary: -10763, Drift rate: -11384, Combined: -11760). The fact that the model that best explains the data, shows the necessity to draw the boundary and drift rate parameters from separate distributions, suggests that varying only one of these parameters is not enough to capture the variability in the trial-wise choice-RT data. We found that on average, the anodal condition was associated with higher drift rates (by 0.191) and lower boundary parameters (by 0.003). This therefore suggests that the anodal condition, which is meant to increase the excitability of the underlying region made the information integration clearer, leading to higher drift rates. Additionally, since the read-out was less ambiguous, the boundary was also lowered since less evidence was needed to make a choice. These results, therefore suggest that the pMFC might be causally implicated in both reflecting the rate of EA and adjusting the boundary separation in correspondence with the amount of information necessary for a decision. We also formally examined the differences between the parameters using paired $t$-tests, however we did not find significant effects between

the boundary parameters across the two stimulation sessions ($t(21) = -0.053$, $p =$ 0.958), and similarly there was no significant difference between the drift rate parameters ($t(109) = 1.786$, $p = 0.077$). These results therefore suggest that although the best performing model allowed for both parameters to vary between stimulation conditions, there was not a significant difference between these parameter estimates, potentially due to the limited sample size presented in this thesis as opposed to the calculated necessary samples size (achieved effect size: 0.011 and power = 0.05 for boundary test; achieved effect size: 0.17 and power = 0.425 for drift rate test).

### 4.3.2   Exploratory and Control results



Figure 20:   **Behavioural results**. Across the two stimulation sessions (anodal - pink, cathodal - turquoise) we observed typical and comparable behaviour. **a)** Specifically, we saw that the P(Play|Reward Probability) increased as the reward probability increased. **b)** Although not fully symmetrical, we saw a general difficulty effect in the RTs with easy trials, where the reward probability was strongly in favour of a choice option ('Play' or 'Keep'), were associated with shorter reaction times than more difficult trials, where the reward probability was not clearly favouring a choice.

In addition to our main analysis aiming to investigate differences between the model parameters produced by the two stimulation conditions, we also performed a number of control and exploratory analyses in order to assess whether such differences exist even on the purely behavioural level, to examine whether our data was consistent with previously seen trends within the decision-making literature and to account for the effects of attention, learning, analysis window and tDCS-related sensations. On the purely behavioural level, we saw that participants' 'Play' behaviour for both the anodal and cathodal sessions reflected the reward probability given a 'Play' choice (Anodal: $t(21)$ = -14.5, $p < 0.001$; Cathodal: $t(21)$ = -12.303, $p < 0.001$). Specifically, we found that they selected the riskier option more frequently as the likelihood of receiving the higher payoff increased. As we see this parametric modulation and considering that across the 2 sessions we only removed a small portion of the 150 faces (16.8) per participant and on average, per session there were 16, 14, 10, 16, 10 unique face identities across the five reward probability bins, respectively, we consider this to be evidence that the ratings we were using as indirect trustworthiness measures were largely consistent and captured relevant social information displayed in our social stimuli. We also observed that these probabilities did not differ significantly between the two stimulation conditions ($t(21)$ = 1.488, $p$ = 0.14). We additionally used a likelihood-ratio test (see Chapter 2, Choice probability calculation) fit to each participant separately to show that a single sigmoid function fit the fraction of 'Play' choices (jointly across both stimulation conditions) as well as two separate functions for the anodal and the cathodal condition (average $\lambda(21)$ = 1.023, $p$ = 0.6). We also found similar RT trends to what we observed in Chapter 2. They once again were mostly consistent with an inverted 'V' relationship with reward probability reflective of the difficulty of the trial (Anodal: $t(21)$ = 3.948, p <0.001; Cathodal: $t(21)$ = 6.056, $p < 0.001$) and those did not differ significantly between the two stimulation sessions (RT : $t(21)$ = 0.121, $p$ = 0.904). We therefore observed that the RTs were longer for more difficult trials where there is no strong evidence to indicate the appropriate choice ('Play' or 'Keep'). These results suggest that as we predicted, although the behavioural data managed to reflect typical trends seen across the decision-making literature, they did not reflect the effects of our stimulation manipulation, potentially due to the typically small effects produced by tDCS, which are often not identifiable through simple behavioural tests (Horvath et al., 2015;

Priori et al., 2009).

We also investigated potential learning effects in relation to the session number and found similar behavioural trends across the two sessions. Once again, participants' 'Play' choices increased with the reward probability for both the first and the second session (First: $t(21)$ = -13.227, $p$ < 0.001, Second: $t(21)$ = -13.246, $p$ < 0.001). There was also no significant difference between the two ($t(21)$ = 1.552, $p$ = 0.124). The RT trends were preserved between the first and the second sessions (First: $t(21)$ = 5.479, $p$ < 0.001; Second: $t(21)$ = 4.019, $p$ = 0.001) with no significant difference dependant on the session number ($t(21)$ = -0.696, $p$ = 0.488). We also examined the effects of session number on the parameter space. We ran the same three main models as the ones in our main analysis, but this time we allowed the boundary and drift rate parameters to vary depending on the session number. Once again, the best-performing model was the one where both the boundary and the drift rate were allowed to vary (DIC Boundary: -10723; Drift rate: -11284; Combined: -11695). We then proceeded to perform formal $t$-tests between the two sets of parameters and found no significant differences: Boundary $t(21)$ = 0.525; $p$ = 0.605), Drift rate $t(109)$= 1.513, $p$ = 0.133). These results therefore suggest that any difference observed in this experiment were likely not driven by the session number.

Furthermore, to investigate the effects of the stimulation sensations, we compared the number of times a participant said that they believed the stimulation was 'on' across the two sessions. We found no significant differences between the anodal and cathodal conditions ($t(21)$ = -0.591, $p$ = 0.561). We similarly looked at the severity of the sensations by first averaging their ratings for the tingling and burning sensations separately and then averaging those two values per participant. Here, however, we did find a significant difference ($t(21)$ = -2.17, $p$ = 0.042), with the cathodal condition being associated with slightly higher severity scores (by 0.326) even though on average the impedance was slightly lower, but not significantly so, during the cathodal condition (Anodal: 3.568 k$\Omega$, Cathodal: 3.355 k$\Omega$, $t(21)$ = 1.305, $p$ = 0.206). We also found a significant difference in the duration of the tDCS duration estimations during the first and second half of the main segment after aggregating across the two stimulation sessions ($t(43)$ = 2.466, $p$ = 0.018) with the first half being associated with a higher number of 'Yes' re-

sponses. We also found that the difference in the duration between the first and second half of the main segment did not differ significantly between the two stimulation conditions ($t(21) = -0.271$, $p = 0.789$). We also found a significant difference between the first and second half of the trials in regards to the sensation severity after collapsing between stimulation sessions ($t(43) = 8.997$, $p < 0.001$), with the average sensations in the first half being stronger than the ones in the second half by 0.765. Once again the differences between the first and second half did not differ significantly between the stimulation sessions ($t(21) = -1.7986$, $p = 0.087$). These results suggest that participants were able to identify when the stimulation was 'on', in accordance with previous literature (Greinacher et al., 2019).

To further investigate potential effects related to participants' beliefs about the stimulation condition, we took the boundary and drift rate estimates from the winning model in our main analysis and regrouped them based on participants' beliefs into a session in which they thought they performed 'better' and a session in which they believed they performed 'worse'. As 8 participants could not perceive any difference between the two sessions, we assigned their 'better' and 'worse' sessions randomly 1000 times and redid a $t$-test for each permutation. This allowed us to calculate an average $t$-value for the boundary and drift rate differences between the 'better' and 'worse' session. We found no difference in the boundary parameters $t(21) = 0.669$, $p = 0.511$ or the drift rate parameters $t(109) = 1.598$, $p = 0.113$. We also note that only 6 participants replied that they performed better in the anodal condition, which was meant to increase the excitability of the pMFC and thus the signal read-out and make our participants perform better. Therefore, the differences between the stimulation sessions were likely not caused by participants' beliefs about their performance.

To investigate the effect of attention we ran all three main hDDMs on the data from the anodal tDCS segment (once again based on the last 20 minutes of the task) as well as on the last 100 trials of the training prior to the anodal segment. In these models we also allowed the nDT parameter to vary between the two segments (training, anodal). We did so under the assumption that if participants were not paying attention to the task, their initial stimulus encoding might take longer. As outlined in the Sample size and participants subsection, we also intended to use this analysis to identify and remove participants who were likely not paying

attention (as evidenced by their nDT estimates). We found two participants with higher nDTs during the anodal segment, however as these differences were very small (0.027 and 0.007s) and considering the already limited samples size, we retained the data from these participants for the purposes of this thesis. We once again compared the three models and found that the combined model explained the data the best (DIC Boundary: -7911, Drift rate: -8070, Combined: -8114). We then formally compared the nDT parameters between the training and the anodal condition with paired $t$-test and found a significant difference ($t$(21) = -5.255, $p <$ 0.001). Although there was a significant effect, we saw that the nDTs were longer in the training than in the anodal condition (by 0.043), suggesting that our participants might have paid more attention during the anodal segment than in the associated training segment. However, we also note that it is also possible that these results might be indicating that participants were not as engaged during the training or that there might have been residual subtle learning effects occurring towards the end of the training segment.

Finally, we varied the windows of our analysis in order to establish the robustness of our results. We selected three additional windows: across the full main economic game (30 min), during the whole stimulation (first 20 min) and following the stimulation (last 10 min). We reran the three main analysis models during these three windows and found that the results were consistent with our main analysis: full duration (DIC Boundary: -15680; Drift rate: -16583, Combined: -17069), during stimulation (DIC Boundary: -11157; Drift rate: -11618, Combined: -11963), post-stimulation (DIC Boundary: -5731; Drift rate: -6068 , Combined: -6361). These results showed that across these windows, the winning model was the combined model allowing for both the boundary and the drift rate to vary across stimulation conditions, thus demonstrating that our main findings were not biased by our choice of analysis window. Once again, across all windows, the boundary and the drift rate parameter did not show any significant differences.

## 4.4  Discussion

The previous two chapters demonstrated that social choices might use the same algorithmic and implementational mechanisms employed by non-social decisions.

Specifically, we showed that similar to non-social choices, social information is encoded in domain-specific regions, converted into a 'common currency' in the human valuation system and then integrated according to an accumulation-to-bound process in the pMFC. In this chapter we aimed to establish a causal and mechanistic role of the pMFC in social decision making by using anodal and cathodal stimulation to enhance and diminish (respectively) the activity in the region. Specifically, we aimed to understand whether the pMFC is involved in reflecting the rate of EA as suggested by a recent account (Pisauro et al., 2017) or if it is involved in boundary adjustments as suggested by the literature investigating the speed-accuracy trade-off (Bogacz, Wagenmakers, et al., 2010; Forstmann et al., 2008). We also considered the possibility that under time constraints, the pMFC regulates both the rate of EA as well as the amount of information necessary for a decision.

Our model comparison results were consistent with the third hypothesis - that the pMFC may act as an accumulator region reflecting both the changes to the rate of EA as well as the amount of information needed for a decision. This observation could have implications for the embodied nature of social choices since the anodal condition, which is meant to increase the excitability of the pMFC (relative to the cathodal condition) was linked to higher drift rates and a lower boundary. This could signify that the increased excitability enhanced the decision signal, which also led to a decreased need to deliberate on the decision outcome. Such observations would, therefore, implicate the pMFC in reflecting the read-out quality of the decision signal in addition to modulating how much of this information is needed prior to committing to a decision.

It is noteworthy however, that the formal comparisons between the boundary and drift rate parameters across stimulation conditions were not significant. This may be attributed to the fact that the tDCS literature has historically been highly inconstant with some of the main assumptions at the basis of this experiment - that the anodal stimulation is meant to increase excitability in the stimulated region and the cathodal is meant to decrease excitability - having been contested (e.g. Moliadze et al., 2018). Moreover, while the simulated effect of the tDCS (Fig. 17) was spatially dispersed and could have accounted for potential heterogeneity in participant-specific anatomical features, it is conceivable that the anatomy of some participants may not have been consistent with the tDCS stimulation. Ad-

ditionally, as the placebo effects associated with tDCS have recently been placed at the forefront of tDCS research (e.g. Brunoni et al., 2014; Greinacher et al., 2019), it is also important to note that the cathodal stimulation was rated as more strongly associated with negative tDCS sensations. Therefore, it is conceivable that the effects we observe here may have been caused by these sensations and it may not be possible to resolve what the source of these differences is. Our preliminary results nevertheless showed that the stimulation session order did not affect the boundary and drift rate parameters and we also found that these parameters were likely not affected by participants' beliefs about their own performance. Taken together these results suggest that more research is needed in order to determine the true function of the pMFC in social decision making. However, considering how this experiment aimed to optimise the setup by including a stimulation simulation, a more informed analysis window and by avoiding the often easily discernible sham condition (Brunoni et al., 2014; Greinacher et al., 2019), it is possible that some of the uncertainty about the potentially meaningfulness of the current results may be resolved with the collection of the full sample size. Consequently, the rest of this discussion section will explore the potential contributions that this experiment could offer for the examination of the role of the pMFC in social decision making once the full dataset is collected.

Specifically, the role of this pre-motor region and one of its main sub-regions - the SMA - has primarily been investigated under the framework of the speed-accuracy trade-off. Consequently, these boundary effects have only been investigated when the task requires major changes to the participants' goal - from focusing on accurate to speedy choices (Bogacz, Wagenmakers, et al., 2010; Forstmann et al., 2008). Therefore, here, we offered to extend our current understanding of these boundary modulations by examining if they are still as prominent when participants are not required to consider the nature of their responses. This experiment may also offer additional insight in relation to recent computational accounts, which propose the existence of a 'motor accumulator' (Verdonck et al., 2021). Specifically, this view proposes that the motor accumulator receives information from a non-embodied accumulator region and proceeds to transfer this information to the motor system once it reaches its own threshold. This is consistent with observation showing that when the responses are known, the relevant (pre)motor areas are active, but when they are not, the accumulation activity is

primarily reserved to the frontal cortex (Filimon et al., 2013). Therefore, given that the location of the pMFC is adjacent to the relevant motor area responsible for the execution of the choice, as well as the observation that motor preparation partially overlaps with EA, which seems to be reflected in this region, the pMFC may act as a motor accumulator by adjusting the rate of EA and the boundary separation, at least in the scenarios involving the use of the motor cortex to perform the outcome of the choice.

Moreover, examining decision embodiment within the context of the social domain may have repercussions for the decision-making literature as a whole, especially after having shown evidence in favour of the notion of a universal decision-making process in the previous two chapters. Specifically, this experiment showcases the potential for paradigms examining social phenomena to offer further mechanistic insight into the universal decision-making processes by providing a range of new potential paradigms defined though the game-theory framework.

Another key aspect of the experiment presented in this chapter is the combined use of tDCS and modelling. Specifically, relying on simple behavioural measurements may not be sufficient for the investigation of the processes underlying decision making (Horvath et al., 2015; Priori et al., 2009). Instead by relying on parameters derived through a hDDM we offered to elucidate the underlying mechanisms that are responsible for social decisions and to provide potential explanations for the role of the pMFC in this process. These latent variables are more appropriate for the investigation of the functioning of the pMFC since they combine the available behavioural output in a joint choice-RT space in such a way so as to detect more subtle changes in behaviour. Furthermore, by using a hierarchical model instead of a traditional DDM, we offer a more robust representation of the underlying process since we do not assume that our participant's parameters are fully independent, nor do we assume that they do not differ at all. Thus this approach allows the parameters to vary between participants while maintaining that they arise from a common across-participant distribution (Wiecki et al., 2013). Moreover, by using a stimulation method, we attempt to provide causal evidence for the involvement of this region in social decision making by showing that changes to its excitability directly result in changes in the EA processing. Additionally, the stimulation setup we chose - in which we both increased and decreased the ex-

citability in the pMFC, without relying on a sham condition (which has been shown to be correctly identifiable by participants e.g. Greinacher et al., 2019) - allows us to alter the functioning of this region in such as way so as to be maximally informative about the underlying processes by examining the relative differences between the two conditions.

Consequently, even though previous attempts to use tDCS to study the role of the SMA in decision making have not managed to find meaningful findings (e.g. de Hollander et al., 2016) and the fact that our preliminary results were not able to find any significant differences associated with the two stimulation conditions, our experiment may still offer mechanistic insight for the role of the pMFC in social decision making once the full dataset is collected. This is due to the fact that this experiment holds an advantage over previous examinations of the SMA using tDCS due to the choice of modelling procedure, which may be more sensitive to subtle changes in behaviour, in addition to our more informed electrode placement (Fig. 17), targeting the region more reliably.

In summary, here we offered preliminary results, which were not able to provide a strong account for the role of the pMFC in social decision making potentially due to the insufficient sample size. Nevertheless, this experiment has the potential to offer further insight into the embodied nature of decision making by involving (pre)motor structures used to express the relevant choice. To this end, we used models, which make fewer assumptions about the distribution from which the parameters were drawn and which are more sensitive to minor changes in behaviour. Furthermore, our use of stimulation methods could allow us to causally implicate the pMFC in the process of EA and potentially demonstrate the key role of this region in social decision making. Finally, by examining the neurocomputational function of this region in the social domain, we highlight the potential of social decision-making paradigms to be used to examine the underlying basis principles of decision making.

# 5   Chapter 5, General Discussion

## 5.1   Key findings

Social choices are an ever-present part of our lives. From the inconsequential choices about whether to trust the salesperson when they say that a piece of clothing looks great on us, to the potentially life-altering decision about whether to marry someone, social decisions require the careful consideration of the options we have available. In the introduction of this thesis, we highlighted that until recently the investigation of such social decisions has primarily focused on their uniquely social qualities and to a lesser extent on their basic neurocomputational mechanisms (Ruff and Fehr, 2014). Conversely, the non-social decision-making domain has historically put a emphasis on examining these basic principles and has highlighted several key stages in the decision-making process - value encoding, value comparison, EA, action selection and value updating (Fig. 3, Rangel et al., 2008) - and has identified the accumulation-to-bound mechanism as a key driver for these decisions (Kelly and O'Connell, 2013; O'Connell et al., 2012; Polanía et al., 2014). Recent work has proposed that social and non-social value may be processed within the same neural structures (e.g. Janowski et al., 2013; Lockwood et al., 2016) or in adjacent sub-areas (Harris et al., 2007) and that common processes may be involved in social and non-social choices (Krajbich et al., 2015). However, often times the social and non-social uncertainty are not explicitly put on equal footing (e.g. Harris et al., 2007; Janowski et al., 2013) and the data is sometimes collected across different samples/studies (e.g. Krajbich et al., 2015; Tarantola et al., 2017). Similarly, most of this research has focused on scenarios in which the social/non-social aspect arises not from the source of uncertainty, but from varying the outcome receiver (e.g. Lockwood et al., 2016) as well as on the value computation stages of decision making. Therefore, more research is needed in order to determine whether social and non-social value are processed in a comparable way in a 'common currency' framework (Ruff and Fehr, 2014) and whether any potential similarities extend across other stages of the decision-making process such as whether social choices are embodied in the same way as non-social decisions. Recently, a new framework (Lockwood et al., 2020) has been developed for the systematic comparison of social and non-social choices based on Marr's framework

(Marr and Poggio, 1979), which presents the necessary steps to establish whether a process is social-specific or universal. This framework calls for the comparison across both the algorithmic and implementational levels and advocates for the use of computational modelling and stimulation methods. To this end, this thesis aimed to provide some insight into the specificity of social decision-making processes based on social cues, whether they are governed by the same mechanisms as non-social choices and to investigate the potentially embodied nature of these choices.

In Chapters 2 and 3 we aimed to compare social and non-social choices across the algorithmic and implementational levels (respectively). Specifically, we developed an experimental design, which aimed to enable the social and non-social uncertainty to vary across comparable scales by embedding the two decisions in the context of a Trust game. The social uncertainty was determined by extracting indirect subject-specific trustworthiness judgments that reflect the likelihood of splitting a reward. These trustworthiness judgments were matched to the explicit probabilities of a reward that reflected the non-social uncertainty. In Chapter 2, we saw behavioural and modelling results, which showed that both social and non-social conditions were associated with trends typically observed within the non-social literature (e.g. Bogacz et al., 2006; Gold and Ding, 2013; Philiastides and Sajda, 2006; Ratcliff et al., 2009) in that the probability of 'Playing' positively correlated with reward probability given a 'Play' choice and that the RTs reflected the difficulty of the trial with more difficult trials leading to longer RTs. We also found that both the social and the non-social choices could be modelled with a DDM and that they follow an 'accumulation-to-bound' process. Additionally, we observed that the parameters associated with these two conditions were largely comparable. Therefore, the results from Chapter 2 suggest that our paradigm might have been successful in putting the social and non-social uncertainty on the same scale and that the two decision domains might share a common algorithmic mechanism.

We then aimed to investigate whether social and non-social choices rely on the same brain structures to carry out their decision-making computations. To this end we first aimed to identify neural signatures of EA in our EEG data by using our model predictions to examine whether such EA dynamics occur in our EEG

signal. We found that the model-predicted EA traces corresponded to EEG activity arising from a centroparietal cluster consistently associated with EA dynamics (e.g. Gherman and Philiastides, 2018; Kelly and O'Connell, 2013; Pisauro et al., 2017) and that the dynamics of the electrodes with the highest correlations with the model traces showed typical difficulty and speed modulations (e.g. Connolly et al., 2009; Law and Gold, 2008; Mazurek et al., 2003; Roitman and Shadlen, 2002), thus suggesting that our EEG signal managed to capture EA dynamics. However, since our EEG-informed fMRI analysis relied on the accurate representation of the EA dynamics, we used a multivariate analysis to increase the signal-to-noise ratio of the relevant signal. Specifically, we used a single-trial LDA based on the easy vs difficult trials to produce a discriminator reflective of EA. We then projected our EEG data through this discriminator to examine the full temporal dynamics of this process and identified single-trial slopes reflective of the trial-by-trial variability in EA. We were able to validate these slopes by showing that they correlated with our drift rate estimates from our modelling procedure and that they were able to predict the choice behaviour of our participants. We then used these slopes as endogenous predictors of EA in our EEG-informed fMRI analysis and implicated the pMFC as the potential site for EA for both social and non-social choices. An additional PPI analysis showed that the activity in the pMFC co-varied in a task-dependant manner with regions of the human valuation system consistent with the notion that domain-specific information is converted into a 'common currency' in these value-processing regions and then accumulated for a decision in the pMFC - a region adjacent to the motor region responsible for carrying out the decision - thus showing the potentially embodied nature of both social and non-social choices.

In Chapter 4 we specifically focused on social choices in order to further investigate the role of the pMFC by varying the excitability of the region through anodal and cathodal tDCS stimulation. We only managed to produce preliminary results due to slow recruitment during the pandemic, which point to the pMFC having a function beyond merely being responsible for modulating the rate of EA or for boundary adjustments, but that it might reflect both EA processes. This is consistent with recent views proposing the existence of a second motor accumulator, which triggers responses once its internal boundary is reached (Verdonck et al., 2021). However, as the formal statistical analysis on the drift rate and boundary parameters was not significant, we note that additional evidence is needed in or-

der to causally implicate the pMFC in social decision making and to establish its exact algorithmic mechanism.

## 5.2 Limitations and future directions

Despite this thesis offering a detailed examination of the spatiotemporal neural characteristics of social decision making, this work has limitations, which need to be acknowledged. The first limitation relates to the fact that throughout this thesis we have only considered one type of social decision i.e. choices based on a social cue. As we highlighted in the Introduction section of this thesis, there are other types of social decisions (Ruff and Fehr, 2014), which need to be examined in order to determine if our observations about the universality of the decision-making process extend to other social decision-making domains. Furthermore, as some of the other subtypes of social decision-making (such as making decisions on someone else's behalf), may have a more overt social component, it may be beneficial to examine whether the degree of social involvement may also be associated with the strength of potential domain-specific processing. Similarly, throughout this thesis we have exclusively examined choices occurring under time constraints and as highlighted by the literature examining the trade-off between making speedy and accurate choices (Bogacz, Wagenmakers, et al., 2010; Forstmann et al., 2008), changing the conditions under which the decision is made, can lead to changes in how these choices are made. Therefore, once again, it is important for future research to examine whether the consistencies between social and non-social choices would be preserved once the goal of the decision is changed.

Furthermore, the main results presented in this thesis come from the simultaneous EEG-fMRI experiment (i.e. the Main experiment), which later informed our tDCS investigation. This approach has multiple benefits such as the opportunity to exploit the high temporal resolution of the EEG signal while providing concrete implementational information about where in the brain information is integrated to guide a choice. However, we relied on the EEG measurement to inform our fMRI analysis in order to identify the site of EA and due to the additive nature of recording electrical signals from the surface of the scalp, these signals may have contained information arising from other brain areas, impeding the identification

of the true relevant signal. Consequently, our EEG predictors may have been influenced by factors unrelated to the process of EA on a trial-by-trial basis. While we used methods to ensure the robustness of the relevant signals, which identified a structure that has been previously implicated in decision making (Pisauro et al., 2017), future research will need to validate our current findings and test their soundness and could potentially also explore the oscillatory phenomena associated with these choices in addition to the broadband signal examined in this thesis (e.g. Polanía et al., 2014).

Even though our tDCS experiment was carefully designed to include a more informed analysis window and stimulation setup, a further limitation of this thesis stems from the fact that the tDCS literature has been associated with certain inconsistencies about what kinds of effects the stimulation might have on cognition, which may explain why our formal statistical analysis did not lead to significant results. Additionally, some of the stimulation by-products, such as differences in the perceived stimulation sensations, might make it impossible to attribute behavioural or modelling changes to stimulation-driven neural changes. Furthermore, to avoid adding further experimental sessions to the experiment, we did not include an initial anatomical scan and while our tDCS setup allowed for some underlying anatomical inconsistencies between our participants, it is possible that the anatomical characteristics of some of them might not have allowed for the stimulation to lead to any cognitive changes, thus providing another potential explanation for the non-significant results. However, it is also conceivable the the collection of the full sample size may reduce some of the uncertainty associated with these potential pitfalls. Another limitation associated with our tDCS experiment relates to the fact that our conclusion were further limited by the models we considered. Specifically, it is conceivable that another model exists, which may provide a better explanation for the functioning of the pMFC. Nevertheless, the models we included were governed by the existing literature and offer to reconcile the existing discrepant views regarding the role of the pMFC in EA.

Despite these limitations, this thesis offers detailed insight about whether social decision making is domain-specific or if it falls within a universal decision-making framework, by providing support for the latter option. However, there are still aspect of the decision-making process, which have not been addressed and it

remains to be uncovered whether the similarities between social and non-social choices are preserved across these other facets of decision making, such as the potential shared motivational aspects between social and non-social behaviours. Specifically, it is conceivable that since social goal-oriented actions are effortful, deciding whether to engage in them might be subject to the same effort and motivation-based considerations as non-social actions and might be governed by the same neural substrates. However, recent accounts argue otherwise and have suggested that such effort-based computations may be associated with domain-specific processing (Contreras-Huerta et al., 2020; Lockwood et al., 2017). For instance, recently it has been suggested that the anterior cingulate gyrus uniquely encodes effort costs for prosocial actions, whereas other areas such as the ventral tegmental area and the ventral insua have been linked to self-related effort processing Lockwood et al., 2021. Such observations offer further insight into the degree of overlap between social and non-social decisions by suggesting that the motivational processes may be domain-specific and this would have implications for the notion of a domain-general decision-making process. However, future research may aim to elucidate whether some circumstances may drive more comparable processing across domains, such as when other types of social decisions are involved and to thus establish whether the differences between social and non-social motivation are categorical or if there is a potential continuum. Nevertheless, determining the extent of the shared motivational processing across the social and non-social domains may enable a deeper understanding of states like apathy marked by reduced goal-driven behaviours.

Another vital aspect of decision making, which has been heavily investigated in recent years (e.g. Gherman and Philiastides, 2015, 2018), is the role of confidence. It has been shown that when there is no explicit feedback in relation to performance, confidence (i.e. the internal belief that the selected choices will be correct) modulates how much information is needed for a choice to be made, thus affecting the boundary separation for the decision (Desender, Boldt, Verguts, et al., 2018). Moreover, it has been suggested that these confidence judgments arise from the same EA processes (Gherman and Philiastides, 2015) observed in the experiments outlined in this thesis. It has also been shown that they involve the human valuation system, with the vmPFC in particular (Gherman and Philiastides, 2018) as well as the dmPFC when considering the reliability of social information

(De Martino et al., 2017), both of which were also implicated in the Main EEG-fMRI experiment presented here. Additionally, it has been suggested that confidence predicts information-seeking tendencies (Desender, Boldt, and Yeung, 2018) and that it affects when we seek out social information (Toelch et al., 2014). Since confidence has been shown to be involved in social choices and it relies on brain regions, which are also implicated in these social decisions, it is important to examine how confidence emerges and affects the decision-making processes across the decision-making domains in addition to comparing its effects on information-seeking strategies across different decision types.

The examination of the potentially universal decision-making mechanism would also be incomplete without investigating how social and non-social values are updated over multiple decisions and whether these learning processes are implemented in the same way in the brain. To this end, recent proposals have been made as to how to study learning in social situations (Lockwood and Klein-Flügge, 2021). The importance of this decision-making stage is highlighted by the fact that learning has been shown to affect the decision-selection stage itself and not early value encoding (Diaz et al., 2017). Specifically, in the non-social domain, a meta analysis of fMRI experiments investigating the neural mechanisms underlying the prediction error signal - the signal that reflects the difference between what was expected and what occurred, which drives how we update our values and beliefs - has elucidated how two learning valence systems are involved in how we make decisions (Fouragnan et al., 2018). The first one processes negative outcomes and is associated with areas responsible for alertness and switching behaviour such as the midcingulate cortex, the thalamus and the dorsolateral prefrontal cortex. The second one preferentially responds to positive outcomes and is linked to higher activation in reward-processing areas like the ventral striatum and vmPFC. Additionally, a third network has also been suggested to be responsible for encoding surprise, the processing of which overlaps temporally with the valence signals and which is associated with activations in the anterior cingulate cortex, anterior insula and dorsal striatum (Fouragnan et al., 2017). These findings therefore demonstrate the importance of understanding the mechanisms at play when we make choices, however considering the evidence from this thesis it is also important to consider their generalizability to other decision-making domains.

Recently, a number of studies have investigated the application of reinforcement learning models, developed in the non-social domain (e.g. Rescorla, 2008; Sutton and Barto, 1998), to social scenarios and have demonstrated the key role of the mPFC in these processes (e.g. Diaconescu et al., 2017; Lockwood et al., 2016; Lockwood and Wittmann, 2018; Wittmann et al., 2016). There have also been studies, which have demonstrated that in most aspects, social and non-social learning are the same, however unlike non-social choices, social decision are affected by additional social information such as knowing the popularity of a certain option (Tarantola et al., 2017). Others have also shown that the ventral striatum reflects reward-prediction errors for both social and non-social choices (e.g. Behrens et al., 2008; Lockwood et al., 2016). Such observations support the notion of a universal decision-making mechanism, however in order to robustly determine if learning is the same across the two domains further investigations may be needed as these effects are often exclusively studied in self-other choices and therefore it remains unknown whether these observation are preserved for other types of social choices.

Finally, the design used in our experiments provides a blueprint for further investigations of the 'common currency' schema by considering other types of decision making, for instance in decisions based on linguistic (such as in lexical decision tasks, e.g. Wagenmakers et al., 2004) or multimodal information (e.g. Franzen et al., 2020). It also offers a way to investigate the mechanisms involved in neurodivergent processing. For instance, it has been shown that adults who score high on the autistic quotient do not show any impairments when it comes to forming trustworthiness judgments (Hooper et al., 2019). However, they do display a lack of association between their responses in economic games. This suggests that they do not make use of these judgments as part of their strategy while participating in economic games. Therefore, by employing a design similar to the one used in this thesis it may be possible to examine the exact algorithmic and implementational differences that occur between people on the autistic spectrum and neurotypical individuals.

## 5.3 Conclusion

In this thesis we offered an examination of the spatiotemporal characteristics of decision based on social information in three studies. For this examination, we compared social choices to non-social decisions to establish whether they share the same neurocomputational underpinnings and we attempted to provide causal evidence for the involvement of the pMFC in social choices. We found evidence suggesting that social decisions employ comparable algorithmic and implementational processes to the ones used in non-social choices and that both might rely on accumulation-to-bound processes embedded in the pMFC. We further showcase a tDCS experiment, which has the potential to elucidate whether the pMFC reflects the continuous changes in EA dynamics, whether it determines the optimal amount of information necessary for a decision or if it reflects both EA processes. These experiments offer detailed insight into a potentially universal decision-making architecture and provide an example for the systematic comparison across decision-making domains.

# References

Adolphs, R. (2010). Conceptual challenges and directions for social neuroscience. *Neuron*, *65*(6), 752-767.

Ahumada Jr, A., & Lovell, J. (1971). Stimulus features in signal detection. *The Journal of the Acoustical Society of America*, *49*(6B), 1751–1756.

Allen, P. J., Josephs, O., & Turner, R. (2000). A method for removing imaging artifact from continuous eeg recorded during functional mri. *Neuroimage*, *12*(2), 230-239.

Antal, A., Keeser, D., Priori, A., Padberg, F., & Nitsche, M. (2015). Conceptual and procedural shortcomings of the systematic review "evidence that transcranial direct current stimulation (tdcs) generates little-to-no reliable neurophysiologic effect beyond mep amplitude modulation in healthy human subjects: A systematic review" by horvath and co-workers. *Brain Stimulation: Basic, Translational, and Clinical Research in Neuromodulation*, *8*(4), 846–849.

Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed random effects for subjects and items. *Journal of memory and language*, *59*(4), 390–412.

Bach, D. R., Hulme, O., Penny, W. D., & Dolan, R. J. (2011). The known unknowns: Neural representation of second-order uncertainty, and ambiguity. *Journal of Neuroscience*, *31*(13), 4811–4820.

Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of memory and language*, *68*(3), 255–278.

Basten, U., Biele, G., Heekeren, H. R., & Fiebach, C. J. (2010). How the brain integrates costs and benefits during decision making. *Proceedings of the National Academy of Sciences*, *107*(50), 21767-21772.

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

Bender, A. D., Filmer, H. L., & Dux, P. E. (2017). Transcranial direct current stimulation of superior medial frontal cortex disrupts response selection during proactive response inhibition. *Neuroimage*, *158*, 455–465.

Berg, J., Dickhaut, J., & McCabe, K. (1995). Trust, reciprocity, and social history. *Games and economic behavior*, *10*(1), 122–142.

Berryhill, M. E., Peterson, D. J., Jones, K. T., & Stephens, J. A. (2014). Hits and misses: Leveraging tdcs to advance cognitive research. *Frontiers in psychology*, *5*, 800.

Bogacz, R., Brown, E., Moehlis, J., Holmes, P., & Cohen, J. D. (2006). The physics of optimal decision making: A formal analysis of models of performance in two-alternative forced-choice tasks. *Psychological review*, *113*(4), 700.

Bogacz, R., Hu, P. T., Holmes, P. J., & Cohen, J. D. (2010). Do humans produce the speed–accuracy trade-off that maximizes reward rate? *The Quarterly Journal of Experimental Psychology*, *63*(5), 863–891.

Bogacz, R., Wagenmakers, E.-J., Forstmann, B. U., & Nieuwenhuis, S. (2010). The neural basis of the speed–accuracy tradeoff. *Trends in neurosciences*, *33*(1), 10–16.

Bogdanov, M., & Schwabe, L. (2016). Transcranial stimulation of the dorsolateral prefrontal cortex prevents stress-induced working memory deficits. *Journal of Neuroscience*, *36*(4), 1429–1437.

Brosnan, M. B., Sabaroedin, K., Silk, T., Genc, S., Newman, D. P., Loughnane, G. M., Fornito, A., O'Connell, R. G., & Bellgrove, M. A. (2020). Evidence accumulation during perceptual decisions in humans varies as a function of dorsal frontoparietal organization. *Nature human behaviour*, *4*(8), 844–855.

Brunoni, A. R., Schestatsky, P., Lotufo, P. A., Benseñor, I. M., & Fregni, F. (2014). Comparison of blinding effectiveness between sham tdcs and placebo sertraline in a 6-week major depression randomized clinical trial. *Clinical Neurophysiology*, *125*(2), 298–305.

Burke, C. J., & Tobler, P. N. (2011). Coding of reward probability and risk by single neurons in animals. *Frontiers in neuroscience*, *5*, 121.

Busemeyer, J. R., & Diederich, A. (2014). Estimation and testing of computational psychological models. *Neuroeconomics* (pp. 49–61). Elsevier.

Camerer, C. F. (2003). Behavioral game theory: Plausible formal models that predict accurately. *Behavioral and Brain Sciences*, *26*(2), 157–158.

Campbell-Meiklejohn, D. K., Kanai, R., Bahrami, B., Bach, D. R., Dolan, R. J., Roepstorff, A., & Frith, C. D. (2012). Structure of orbitofrontal cortex predicts social influence. *Current Biology*, *22*(4), R123–R124.

Caplin, A., & Glimcher, P. W. (2014). Basic methods from neoclassical economics. *Neuroeconomics* (pp. 3–17). Elsevier.

Carlsen, A. N., Eagles, J. S., & MacKinnon, C. D. (2015). Transcranial direct current stimulation over the supplementary motor area modulates the preparatory activation level in the human motor system. *Behavioural brain research*, *279*, 68–75.

Cerniglia, L., Bartolomeo, L., Capobianco, M., Lo Russo, S. L. M., Festucci, F., Tambelli, R., Adriani, W., & Cimino, S. (2019). Intersections and divergences between empathizing and mentalizing: Development, recent advancements by neuroimaging and the future of animal modeling. *Frontiers in behavioral neuroscience*, *13*, 212.

Chen, F., & Krajbich, I. (2018). Biased sequential sampling underlies the effects of time pressure and delay in social decision making. *Nature communications*, *9*(1), 1–10.

Chib, V. S., Rangel, A., Shimojo, S., & O'Doherty, J. P. (2009). Evidence for a common representation of decision values for dissimilar goods in human ventromedial prefrontal cortex. *Journal of Neuroscience*, *29*(39), 12315–12320.

Cisek, P., & Kalaska, J. F. (2005). Neural correlates of reaching decisions in dorsal premotor cortex: Specification of multiple direction choices and final selection of action. *Neuron*, *45*(5), 801–814.

Clithero, J. A., & Rangel, A. (2014). Informatic parcellation of the network involved in the computation of subjective value. *Social cognitive and affective neuroscience*, *9*(9), 1289–1302.

Connolly, P. M., Bennur, S., & Gold, J. I. (2009). Correlates of perceptual learning in an oculomotor decision variable. *Journal of Neuroscience*, *29*(7), 2136–2150.

Contreras-Huerta, L. S., Pisauro, A., & Apps, M. A. (2020). Effort shapes social cognition and behaviour: A neuro-cognitive framework. *Neuroscience & Biobehavioral Reviews*.

Coricelli, G., & Nagel, R. (2009). Neural correlates of depth of strategic reasoning in medial prefrontal cortex. *Proceedings of the National Academy of Sciences*, *106*(23), 9163–9168.

Daw, N. D., O'doherty, J. P., Dayan, P., Seymour, B., & Dolan, R. J. (2006). Cortical substrates for exploratory decisions in humans. *Nature*, *441*(7095), 876–879.

De Martino, B., Bobadilla-Suarez, S., Nouguchi, T., Sharot, T., & Love, B. C. (2017). Social information is integrated into value and confidence judgments according to its reliability. *Journal of Neuroscience*, *37*(25), 6066–6074.

De Neys, W., Hopfensitz, A., & Bonnefon, J.-F. (2017). Split-second trustworthiness detection from faces in an economic game. *Experimental Psychology*.

de Hollander, G., Labruna, L., Sellaro, R., Trutti, A., Colzato, L. S., Ratcliff, R., Ivry, R. B., & Forstmann, B. U. (2016). Transcranial direct current stimulation does not influence the speed–accuracy tradeoff in perceptual decision-making: Evidence from three independent studies. *Journal of cognitive neuroscience*, *28*(9), 1283–1294.

Desender, K., Boldt, A., Verguts, T., & Donner, T. H. (2018). Post-decisional sense of confidence shapes speed-accuracy tradeoff for subsequent choices. *bioRxiv*, 466730.

Desender, K., Boldt, A., & Yeung, N. (2018). Subjective confidence predicts information seeking in decision making. *Psychological science*, *29*(5), 761–778.

Diaconescu, A. O., Mathys, C., Weber, L. A., Kasper, L., Mauer, J., & Stephan, K. E. (2017). Hierarchical prediction errors in midbrain and septum during social learning. *Social cognitive and affective neuroscience*, *12*(4), 618–634.

Diaz, J. A., Queirazza, F., & Philiastides, M. G. (2017). Perceptual learning alters post-sensory processing in human decision-making. *Nature Human Behaviour*, *1*(2), 1–9.

Ding, L., & Gold, J. I. (2010). Caudate encodes multiple computations for perceptual decisions. *Journal of Neuroscience*, *30*(47), 15747–15759.

Ding, L., & Gold, J. I. (2012). Neural correlates of perceptual decision making before, during, and after decision commitment in monkey frontal eye field. *Cerebral Cortex*, *22*(5), 1052–1067.

Domenech, P., Redouté, J., Koechlin, E., & Dreher, J.-C. (2018). The neurocomputational architecture of value-based selection in the human brain. *Cerebral Cortex*, *28*(2), 585–601.

Donner, T. H., Siegel, M., Fries, P., & Engel, A. K. (2009). Buildup of choice-predictive activity in human motor cortex during perceptual decision making. *Current Biology*, *19*(18), 1581–1585.

Drugowitsch, J., Moreno-Bote, R., Churchland, A. K., Shadlen, M. N., & Pouget, A. (2012). The cost of accumulating evidence in perceptual decision making. *Journal of Neuroscience*, *32*(11), 3612–3628.

Duda, R. O., Hart, P. E. et al. (2006). *Pattern classification*. John Wiley & Sons.

Edelson, M. G., Polania, R., Ruff, C. C., Fehr, E., & Hare, T. A. (2018). Computational and neurobiological foundations of leadership decisions. *Science*, *361*(6401).

Engelmann, J. B., Meyer, F., Ruff, C. C., & Fehr, E. (2019). The neural circuitry of affect-induced distortions of trust. *Science advances*, *5*(3), eaau3413.

Ennew, C., & Sekhon, H. (2007). Measuring trust in financial services: The trust index. *Consumer Policy Review*, *17*(2), 62.

Ewing, L., Caulfield, F., Read, A., & Rhodes, G. (2015). Perceived trustworthiness of faces drives trust behaviour in children. *Developmental Science*, *18*(2), 327–334.

Fehr, E., & Krajbich, I. (2014). Social preferences and the brain. *Neuroeconomics* (pp. 193–218). Elsevier.

FeldmanHall, O., & Shenhav, A. (2019). Resolving uncertainty in a social world. *Nature human behaviour*, *3*(5), 426–435.

Feng, C., Luo, Y.-J., & Krueger, F. (2015). Neural signatures of fairness-related normative decision making in the ultimatum game: A coordinate-based meta-analysis. *Human brain mapping*, *36*(2), 591–602.

Ferrari, C., Lega, C., Vernice, M., Tamietto, M., Mende-Siedlecki, P., Vecchi, T., Todorov, A., & Cattaneo, Z. (2016). The dorsomedial prefrontal cortex plays a causal role in integrating social impressions from faces and verbal descriptions. *Cerebral cortex*, *26*(1), 156–165.

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.

Forstmann, B. U., Dutilh, G., Brown, S., Neumann, J., Von Cramon, D. Y., Ridderinkhof, K. R., & Wagenmakers, E.-J. (2008). Striatum and pre-sma facilitate decision-making under time pressure. *Proceedings of the National Academy of Sciences*, *105*(45), 17538–17542.

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.

Fouragnan, E., Queirazza, F., Retzler, C., Mullinger, K. J., & Philiastides, M. G. (2017). Spatiotemporal neural characterization of prediction error valence and surprise during reward learning in humans. *Scientific reports*, *7*(1), 1–18.

Fouragnan, E., Retzler, C., Mullinger, K., & Philiastides, M. G. (2015). Two spatiotemporally distinct value systems shape reward-based learning in the human brain. *Nature communications*, *6*(1), 1–11.

Fouragnan, E., Retzler, C., & Philiastides, M. G. (2018). Separate neural representations of prediction error valence and surprise: Evidence from an fmri meta-analysis. *Human brain mapping*, *39*(7), 2887–2906.

Franzen, L., Delis, I., De Sousa, G., Kayser, C., & Philiastides, M. G. (2020). Auditory information enhances post-sensory visual evidence during rapid multisensory decision-making. *Nature communications*, *11*(1), 1–14.

Fukuda, H., Ma, N., Suzuki, S., Harasawa, N., Ueno, K., Gardner, J. L., Ichinohe, N., Haruno, M., Cheng, K., & Nakahara, H. (2019). Computing social value conversion in the human brain. *Journal of Neuroscience*, *39*(26), 5153–5172.

Garvert, M. M., Friston, K. J., Dolan, R. J., & Garrido, M. I. (2014). Subcortical amygdala pathways enable rapid face processing. *Neuroimage*, *102*, 309–316.

Gherman, S., & Philiastides, M. G. (2015). Neural representations of confidence emerge from the process of decision formation during perceptual choices. *Neuroimage*, *106*, 134–143.

Gherman, S., & Philiastides, M. G. (2018). Human vmpfc encodes early signatures of confidence in perceptual decisions. *Elife*, *7*, e38293.

Gill, D., Garrod, O. G., Jack, R. E., & Schyns, P. G. (2014). Facial movements strategically camouflage involuntary social signals of face morphology. *Psychological science*, *25*(5), 1079–1086.

Gluth, S., Hotaling, J. M., & Rieskamp, J. (2017). The attraction effect modulates reward prediction errors and intertemporal choices. *Journal of Neuroscience*, *37*(2), 371–382.

Gold, J. I., & Ding, L. (2013). How mechanisms of perceptual decision-making affect the psychometric function. *Progress in neurobiology, 103*, 98–114.

Gold, J. I., & Shadlen, M. N. (2007). The neural basis of decision making. *Annu. Rev. Neurosci., 30*, 535–574.

Gonzalez-Gadea, M. L., Sigman, M., Rattazzi, A., Lavin, C., Rivera-Rei, A., Marino, J., Manes, F., & Ibanez, A. (2016). Neural markers of social and monetary rewards in children with attention-deficit/hyperactivity disorder and autism spectrum disorder. *Scientific reports, 6*(1), 1–11.

Grecucci, A., Giorgetta, C., Van't Wout, M., Bonini, N., & Sanfey, A. G. (2013). Reappraising the ultimatum: An fmri study of emotion regulation and decision making. *Cerebral cortex, 23*(2), 399–410.

Greinacher, R., Buhôt, L., Möller, L., & Learmonth, G. (2019). The time course of ineffective sham-blinding during low-intensity (1 ma) transcranial direct current stimulation. *European Journal of Neuroscience, 50*(8), 3380–3388.

Grinband, J., Wager, T. D., Lindquist, M., Ferrera, V. P., & Hirsch, J. (2008). Detection of time-varying signals in event-related fmri designs. *Neuroimage, 43*(3), 509–520.

Grueschow, M., Polania, R., Hare, T. A., & Ruff, C. C. (2015). Automatic versus choice-dependent value representations in the human brain. *Neuron, 85*(4), 874–885.

Gusnard, D. A., Akbudak, E., Shulman, G. L., & Raichle, M. E. (2001). Medial prefrontal cortex and self-referential mental activity: Relation to a default mode of brain function. *Proceedings of the National Academy of Sciences, 98*(7), 4259–4264.

Güth, W., Schmittberger, R., & Schwarze, B. (1982). An experimental analysis of ultimatum bargaining. *Journal of economic behavior & organization, 3*(4), 367–388.

Hackel, L. M., & Amodio, D. M. (2018). Computational neuroscience approaches to social cognition. *Current Opinion in Psychology, 24*, 92–97.

Hare, T. A., Schultz, W., Camerer, C. F., O'Doherty, J. P., & Rangel, A. (2011). Transformation of stimulus value signals into motor commands during simple choice. *Proceedings of the National Academy of Sciences, 108*(44), 18120–18125.

Harlé, K. M., Chang, L. J., van't Wout, M., & Sanfey, A. G. (2012). The neural mechanisms of affect infusion in social economic decision-making: A mediating role of the anterior insula. *Neuroimage*, *61*(1), 32–40.

Harris, L. T., McClure, S. M., Van Den Bos, W., Cohen, J. D., & Fiske, S. T. (2007). Regions of the mpfc differentially tuned to social and nonsocial affective evaluation. *Cognitive, Affective, & Behavioral Neuroscience*, *7*(4), 309–316.

Harsanyi, J. C. (1967). Games with incomplete information played by "bayesian" players, i-iii part i. the basic model. *Management science*, *14*(3), 159–182.

Harsanyi, J. C. (1968). Games with incomplete information played by "bayesian" players part ii. bayesian equilibrium points. *Management Science*, *14*(5), 320–334.

Hayden, B. Y., Heilbronner, S. R., Pearson, J. M., & Platt, M. L. (2011). Surprise signals in anterior cingulate cortex: Neuronal encoding of unsigned reward prediction errors driving adjustment in behavior. *Journal of Neuroscience*, *31*(11), 4178–4187.

Haynes, J.-D. (2011). Multivariate decoding and brain reading: Introduction to the special issue.

Heatherton, T. F., Wyland, C. L., Macrae, C. N., Demos, K. E., Denny, B. T., & Kelley, W. M. (2006). Medial prefrontal activity differentiates self from close others. *Social cognitive and affective neuroscience*, *1*(1), 18–25.

Heekeren, H. R., Marrett, S., Bandettini, P. A., & Ungerleider, L. G. (2004). A general mechanism for perceptual decision-making in the human brain. *Nature*, *431*(7010), 859–862.

Heekeren, H. R., Marrett, S., & Ungerleider, L. G. (2008). The neural systems that mediate human perceptual decision making. *Nature reviews neuroscience*, *9*(6), 467–479.

Henson, R. (2007). Efficient experimental design for fmri. *Statistical parametric mapping: The analysis of functional brain images*, 193–210.

Herding, J., Ludwig, S., von Lautz, A., Spitzer, B., & Blankenburg, F. (2019). Centro-parietal eeg potentials index subjective evidence and confidence during perceptual decision making. *NeuroImage*, *201*, 116011.

Hill, C. A., Suzuki, S., Polania, R., Moisa, M., O'doherty, J. P., & Ruff, C. C. (2017). A causal account of the brain network computations underlying strategic social behavior. *Nature neuroscience*, *20*(8), 1142–1149.

Hoel, P., Port, S., & Stone, C. (1971). Introduction to probability theory. series in statistics.

Hooper, J. J., Sutherland, C. A., Ewing, L., Langdon, R., Caruana, N., Connaughton, E., Williams, N., Greenwell-Barnden, J., & Rhodes, G. (2019). Should i trust you? autistic traits predict reduced appearance-based trust decisions. *British Journal of Psychology, 110*(4), 617–634.

Horvath, J. C., Forte, J. D., & Carter, O. (2015). Evidence that transcranial direct current stimulation (tdcs) generates little-to-no reliable neurophysiologic effect beyond mep amplitude modulation in healthy human subjects: A systematic review. *Neuropsychologia, 66*, 213–236.

Houser, D., & McCabe, K. (2014). Experimental economics and experimental game theory. *Neuroeconomics* (pp. 19–34). Elsevier.

Huettel, S. A., Stowe, C. J., Gordon, E. M., Warner, B. T., & Platt, M. L. (2006). Neural signatures of economic preferences for risk and ambiguity. *Neuron, 49*(5), 765–775.

Hutcherson, C. A., Bushong, B., & Rangel, A. (2015). A neurocomputational model of altruistic choice and its implications. *Neuron, 87*(2), 451–462.

Hutcherson, C. A., Montaser-Kouhsari, L., Woodward, J., & Rangel, A. (2015). Emotional and utilitarian appraisals of moral dilemmas are encoded in separate areas and integrated in ventromedial prefrontal cortex. *Journal of Neuroscience, 35*(36), 12593–12605.

Izuma, K., Saito, D. N., & Sadato, N. (2008). Processing of social and monetary rewards in the human striatum. *Neuron, 58*(2), 284–294.

Janowski, V., Camerer, C., & Rangel, A. (2013). Empathic choice involves vmpfc value signals that are modulated by social processing implemented in ipl. *Social cognitive and affective neuroscience, 8*(2), 201–208.

Jenkins, A. C., & Mitchell, J. P. (2011). Medial prefrontal cortex subserves diverse forms of self-reflection. *Social neuroscience, 6*(3), 211–218.

Johnson, E. J., & Ratcliff, R. (2014). Computational and process models of decision making in psychology and behavioral economics. *Neuroeconomics*, 35–47.

Kelly, S. P., & O'Connell, R. G. (2013). Internal and external influences on the rate of sensory evidence accumulation in the human brain. *Journal of Neuroscience, 33*(50), 19434–19441.

Kim, J.-H., Kim, D.-W., Chang, W. H., Kim, Y.-H., Kim, K., & Im, C.-H. (2014). Inconsistent outcomes of transcranial direct current stimulation may originate from anatomical differences among individuals: Electric field simulation using individual mri data. *Neuroscience letters*, *564*, 6–10.

Klein-Flügge, M. C., & Bestmann, S. (2012). Time-dependent changes in human corticospinal excitability reveal value-based competition for action during decision processing. *Journal of neuroscience*, *32*(24), 8373–8382.

Kolling, N., Braunsdorf, M., Vijayakumar, S., Bekkering, H., Toni, I., & Mars, R. B. (2021). Constructing others' beliefs from one's own using medial frontal cortex. *Journal of Neuroscience*.

Kolling, N., Wittmann, M. K., Behrens, T. E., Boorman, E. D., Mars, R. B., & Rushworth, M. F. (2016). Value, search, persistence and model updating in anterior cingulate cortex. *Nature neuroscience*, *19*(10), 1280–1285.

Konovalov, A., Hill, C., Daunizeau, J., & Ruff, C. C. (2021). Dissecting functional contributions of the social brain to strategic behavior. *Neuron*, *109*(20), 3323–3337.

Krajbich, I., Armel, C., & Rangel, A. (2010). Visual fixations and the computation and comparison of value in simple choice. *Nature neuroscience*, *13*(10), 1292–1298.

Krajbich, I., Hare, T., Bartling, B., Morishima, Y., & Fehr, E. (2015). A common mechanism underlying food choice and social decisions. *PLoS computational biology*, *11*(10), e1004371.

Krajbich, I., & Rangel, A. (2011). Multialternative drift-diffusion model predicts the relationship between visual fixations and choice in value-based decisions. *Proceedings of the National Academy of Sciences*, *108*(33), 13852–13857.

Kruschke, J. K. (2010). What to believe: Bayesian methods for data analysis. *Trends in cognitive sciences*, *14*(7), 293–300.

Kuhn, H. (1950). Extensive games. *Proceedings of the National Academy of Sciences*, *36*, 570–576.

Law, C.-T., & Gold, J. I. (2008). Neural correlates of perceptual learning in a sensory-motor, but not a sensory, cortical area. *Nature neuroscience*, *11*(4), 505–513.

Lee, A. J., Wright, M. J., Martin, N. G., Keller, M. C., & Zietsch, B. P. (2017). Facial trustworthiness is associated with heritable aspects of face shape. *Adaptive Human Behavior and Physiology*, *3*(4), 351–364.

Lee, D. (2008). Game theory and neural basis of social decision making. *Nature neuroscience*, *11*(4), 404–409.

Lee, D., & Seo, H. (2016). Neural basis of strategic decision making. *Trends in neurosciences*, *39*(1), 40–48.

Lepora, N. F., & Pezzulo, G. (2015). Embodied choice: How action influences perceptual decision making. *PLoS computational biology*, *11*(4), e1004110.

Levy, D. J., & Glimcher, P. W. (2012). The root of all value: A neural common currency for choice. *Current opinion in neurobiology*, *22*(6), 1027–1038.

Levy, I., Snell, J., Nelson, A. J., Rustichini, A., & Glimcher, P. W. (2010). Neural representation of subjective value under risk and ambiguity. *Journal of neurophysiology*, *103*(2), 1036–1047.

Lieberman, M. D., Straccia, M. A., Meyer, M. L., Du, M., & Tan, K. M. (2019). Social, self,(situational), and affective processes in medial prefrontal cortex (mpfc): Causal, multivariate, and reverse inference evidence. *Neuroscience & Biobehavioral Reviews*, *99*, 311–328.

Lim, S.-L., O'Doherty, J. P., & Rangel, A. (2011). The decision value computations in the vmpfc and striatum use a relative value code that is guided by visual attention. *Journal of Neuroscience*, *31*(37), 13214–13223.

Lim, S.-L., O'Doherty, J. P., & Rangel, A. (2013). Stimulus value signals in ventromedial pfc reflect the integration of attribute value signals computed in fusiform gyrus and posterior superior temporal gyrus. *Journal of Neuroscience*, *33*(20), 8729–8741.

Liu, T., & Pleskac, T. J. (2011). Neural correlates of evidence accumulation in a perceptual decision task. *Journal of neurophysiology*, *106*(5), 2383–2398.

Liu, Y., Li, S., Lin, W., Li, W., Yan, X., Wang, X., Pan, X., Rutledge, R. B., & Ma, Y. (2019). Oxytocin modulates social value representations in the amygdala. *Nature neuroscience*, *22*(4), 633–641.

Lockwood, P., Wittmann, M., Nili, H., Matsumoto-Ryan, M., Abdurahman, A., Cutler, J., Husain, M., & Apps, M. (2021). Distinct neural representations for prosocial and self-benefitting effort. *bioRxiv*.

Lockwood, P. L., Apps, M. A., & Chang, S. W. (2020). Is there a 'social'brain? implementations and algorithms. *Trends in Cognitive Sciences*.

Lockwood, P. L., Apps, M. A., Valton, V., Viding, E., & Roiser, J. P. (2016). Neurocomputational mechanisms of prosocial learning and links to empathy. *Proceedings of the National Academy of Sciences*, *113*(35), 9763–9768.

Lockwood, P. L., Hamonet, M., Zhang, S. H., Ratnavel, A., Salmony, F. U., Husain, M., & Apps, M. A. (2017). Prosocial apathy for helping others when effort is required. *Nature human behaviour*, *1*(7), 1–10.

Lockwood, P. L., & Klein-Flügge, M. C. (2021). Computational modelling of social cognition and behaviour—a reinforcement learning primer. *Social Cognitive and Affective Neuroscience*, *16*(8), 761–771.

Lockwood, P. L., & Wittmann, M. K. (2018). Ventral anterior cingulate cortex and social decision-making. *Neuroscience & Biobehavioral Reviews*, *92*, 187–191.

Mar, R. A. (2011). The neural bases of social cognition and story comprehension. *Annual review of psychology*, *62*, 103–134.

Marr, D., & Poggio, T. (1979). A computational theory of human stereo vision. *Proceedings of the Royal Society of London. Series B. Biological Sciences*, *204*(1156), 301–328.

Mazurek, M. E., Roitman, J. D., Ditterich, J., & Shadlen, M. N. (2003). A role for neural integrators in perceptual decision making. *Cerebral cortex*, *13*(11), 1257–1269.

McBride, J., Sumner, P., & Husain, M. (2018). Masked primes evoke partial responses. *Quarterly Journal of Experimental Psychology*, *71*(6), 1431–1439.

Mohr, P. N., Biele, G., & Heekeren, H. R. (2010). Neural processing of risk. *Journal of Neuroscience*, *30*(19), 6613–6619.

Moliadze, V., Lyzhko, E., Schmanke, T., Andreas, S., Freitag, C. M., & Siniatchkin, M. (2018). 1 ma cathodal tdcs shows excitatory effects in children and adolescents: Insights from tms evoked n100 potential. *Brain Research Bulletin*, *140*, 43–51.

Monosov, I. E. (2017). Anterior cingulate is a source of valence-specific information about value and uncertainty. *Nature communications*, *8*(1), 1–12.

Monosov, I. E. (2020). How outcome uncertainty mediates attention, learning, and decision-making. *Trends in neurosciences*.

Monte-Silva, K., Kuo, M.-F., Hessenthaler, S., Fresnoza, S., Liebetanz, D., Paulus, W., & Nitsche, M. A. (2013). Induction of late ltp-like plasticity in the human motor cortex by repeated non-invasive brain stimulation. *Brain stimulation*, *6*(3), 424–432.

Morishima, Y., Schunk, D., Bruhin, A., Ruff, C. C., & Fehr, E. (2012). Linking brain structure and activation in temporoparietal junction to explain the neurobiology of human altruism. *Neuron*, *75*(1), 73–79.

Mormann, M. M., Malmaud, J., Huth, A., Koch, C., & Rangel, A. (2010). The drift diffusion model can account for the accuracy and reaction time of value-based choices under high and low time pressure. *Judgment and Decision making*, *5*(6), 437–449.

Morriss, J., Gell, M., & van Reekum, C. M. (2019). The uncertain brain: A co-ordinate based meta-analysis of the neural signatures supporting uncertainty during different contexts. *Neuroscience & Biobehavioral Reviews*, *96*, 241–249.

Mostert, P., Kok, P., & de Lange, F. P. (2015). Dissociating sensory from decision processes in human perceptual decision making. *Scientific reports*, *5*(1), 1–13.

Mulder, M., Van Maanen, L., & Forstmann, B. (2014). Perceptual decision neuro-sciences –a model-based review. *Neuroscience*, *277*, 872–884.

Mumford, J. A., Poline, J.-B., & Poldrack, R. A. (2015). Orthogonalization of regressors in fmri models. *PloS one*, *10*(4), e0126255.

Mussel, P., Göritz, A. S., & Hewig, J. (2013). The value of a smile: Facial expression affects ultimatum-game responses. *Judgment and Decision Making*, *8*(3).

Nagel, R., Brovelli, A., Heinemann, F., & Coricelli, G. (2018). Neural mechanisms mediating degrees of strategic uncertainty. *Social Cognitive and Affective Neuroscience*, *13*(1), 52–62.

Neubert, F.-X., Mars, R. B., Sallet, J., & Rushworth, M. F. (2015). Connectivity reveals relationship of brain areas for reward-guided learning and decision making in human and monkey frontal cortex. *Proceedings of the national academy of sciences*, *112*(20), E2695–E2704.

Nitsche, M. A., Cohen, L. G., Wassermann, E. M., Priori, A., Lang, N., Antal, A., Paulus, W., Hummel, F., Boggio, P. S., Fregni, F., et al. (2008). Transcranial

direct current stimulation: State of the art 2008. *Brain stimulation, 1*(3), 206–223.

Nunez, M. D., Vandekerckhove, J., & Srinivasan, R. (2017). How attention influences perceptual decision making: Single-trial eeg correlates of drift-diffusion model parameters. *Journal of mathematical psychology, 76,* 117–130.

O'Connell, R. G., Dockree, P. M., & Kelly, S. P. (2012). A supramodal accumulation-to-bound signal that determines perceptual decisions in humans. *Nature neuroscience, 15*(12), 1729–1735.

Parra, L. C., Spence, C. D., Gerson, A. D., & Sajda, P. (2005). Recipes for the linear analysis of eeg. *Neuroimage, 28*(2), 326–341.

Paulus, M. P., & Frank, L. R. (2006). Anterior cingulate activity modulates nonlinear decision weight function of uncertain prospects. *Neuroimage, 30*(2), 668–677.

Pearson, J. M., Watson, K. K., & Platt, M. L. (2014). Decision making: The neuroethological turn. *Neuron, 82*(5), 950–965.

Peirce, J. W. (2007). Psychopy—psychophysics software in python. *Journal of neuroscience methods, 162*(1-2), 8–13.

Peixoto, D., Verhein, J. R., Kiani, R., Kao, J. C., Nuyujukian, P., Chandrasekaran, C., Brown, J., Fong, S., Ryu, S. I., Shenoy, K. V., et al. (2021). Decoding and perturbing decision states in real time. *Nature, 591*(7851), 604–609.

Philiastides, M. G., Auksztulewicz, R., Heekeren, H. R., & Blankenburg, F. (2011). Causal role of dorsolateral prefrontal cortex in human perceptual decision making. *Current biology, 21*(11), 980–983.

Philiastides, M. G., Biele, G., & Heekeren, H. R. (2010). A mechanistic account of value computation in the human brain. *Proceedings of the National Academy of Sciences, 107*(20), 9430–9435.

Philiastides, M. G., Heekeren, H. R., & Sajda, P. (2014). Human scalp potentials reflect a mixture of decision-related signals during perceptual choices. *Journal of Neuroscience, 34*(50), 16877–16889.

Philiastides, M. G., & Ratcliff, R. (2013). Influence of branding on preference-based decision making. *Psychological science, 24*(7), 1208–1215.

Philiastides, M. G., Ratcliff, R., & Sajda, P. (2006). Neural representation of task difficulty and decision making during perceptual categorization: A timing diagram. *Journal of Neuroscience*, *26*(35), 8965–8975.

Philiastides, M. G., & Sajda, P. (2006). Temporal characterization of the neural correlates of perceptual decision making in the human brain. *Cerebral cortex*, *16*(4), 509–518.

Philiastides, M. G., & Sajda, P. (2007). Eeg-informed fmri reveals spatiotemporal characteristics of perceptual decision making. *Journal of Neuroscience*, *27*(48), 13082–13091.

Philiastides, M. G., Tu, T., & Sajda, P. (2021). Inferring macroscale brain dynamics via fusion of simultaneous eeg-fmri. *Annual Review of Neuroscience*, *44*, 315–334.

Pisauro, M. A., Fouragnan, E., Retzler, C., & Philiastides, M. G. (2017). Neural correlates of evidence accumulation during value-based decisions revealed via simultaneous eeg-fmri. *Nature communications*, *8*(1), 1–9.

Piva, M., Velnoskey, K., Jia, R., Nair, A., Levy, I., & Chang, S. W. (2019). The dorsomedial prefrontal cortex computes task-invariant relative subjective value for self and other. *Elife*, *8*, e44939.

Polanía, R., Krajbich, I., Grueschow, M., & Ruff, C. C. (2014). Neural oscillations and synchronization differentially support evidence accumulation in perceptual and value-based decision making. *Neuron*, *82*(3), 709–720.

Polanía, R., Moisa, M., Opitz, A., Grueschow, M., & Ruff, C. C. (2015). The precision of value-based choices depends causally on fronto-parietal phase coupling. *Nature communications*, *6*(1), 1–10.

Polanía, R., Woodford, M., & Ruff, C. C. (2019). Efficient coding of subjective value. *Nature neuroscience*, *22*(1), 134–142.

Priori, A., Hallett, M., & Rothwell, J. C. (2009). Repetitive transcranial magnetic stimulation or transcranial direct current stimulation? *Brain stimulation*, *2*(4), 241–245.

Rangel, A., Camerer, C., & Montague, P. R. (2008). A framework for studying the neurobiology of value-based decision making. *Nature reviews neuroscience*, *9*(7), 545–556.

Rangel, A., & Clithero, J. A. (2014). The computation of stimulus values in simple choice. *Neuroeconomics*, 125–148.

Rangel, A., & Hare, T. (2010). Neural computations associated with goal-directed choice. *Current opinion in neurobiology*, *20*(2), 262–270.

Rapoport, A., Chammah, A. M., & Orwant, C. J. (1965). *Prisoner's dilemma: A study in conflict and cooperation* (Vol. 165). University of Michigan press.

Ratcliff, R. (1978). A theory of memory retrieval. *Psychological review*, *85*(2), 59.

Ratcliff, R., & McKoon, G. (2008). The diffusion decision model: Theory and data for two-choice decision tasks. *Neural computation*, *20*(4), 873–922.

Ratcliff, R., Philiastides, M. G., & Sajda, P. (2009). Quality of evidence for perceptual decision making is indexed by trial-to-trial variability of the eeg. *Proceedings of the National Academy of Sciences*, *106*(16), 6539–6544.

Reed, L. I., Zeglen, K. N., & Schmidt, K. L. (2012). Facial expressions as honest signals of cooperative intent in a one-shot anonymous prisoner's dilemma game. *Evolution and Human Behavior*, *33*(3), 200–209.

Rescorla, R. (2008). Rescorla-wagner model. *Scholarpedia*, *3*(3), 2237.

Rezlescu, C., Duchaine, B., Olivola, C. Y., & Chater, N. (2012). Unfakeable facial configurations affect strategic choices in trust games with or without information about past behavior. *PloS one*, *7*(3), e34293.

Rilling, J. K., King-Casas, B., & Sanfey, A. G. (2008). The neurobiology of social decision-making. *Current opinion in neurobiology*, *18*(2), 159–165.

Rilling, J. K., & Sanfey, A. G. (2011). The neuroscience of social decision-making. *Annual review of psychology*, *62*, 23–48.

Roitman, J. D., & Shadlen, M. N. (2002). Response of neurons in the lateral intraparietal area during a combined visual discrimination reaction time task. *Journal of neuroscience*, *22*(21), 9475–9489.

Rossi, S., Hallett, M., Rossini, P. M., Pascual-Leone, A., of TMS Consensus Group, S., et al. (2009). Safety, ethical considerations, and application guidelines for the use of transcranial magnetic stimulation in clinical practice and research. *Clinical neurophysiology*, *120*(12), 2008–2039.

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

Ruff, C. C., Ugazio, G., & Fehr, E. (2013). Changing social norm compliance with noninvasive brain stimulation. *Science*, *342*(6157), 482–484.

Sajda, P., Philiastides, M. G., Heekeren, H., & Ratcliff, R. (2011). Linking neuronal variability to perceptual decision making via neuroimaging. In M. Ding & D.

Glanzman (Eds.), *The dynamic brain: An exploration of neuronal variability and its functional significance* (pp. 214-232). Oxford University Press.

Sajda, P., Philiastides, M. G., & Parra, L. C. (2009). Single-trial analysis of neuroimaging data: Inferring neural networks underlying perceptual decision-making in the human brain. *IEEE reviews in biomedical engineering, 2,* 97–109.

Sanfey, A. G. (2007). Social decision-making: Insights from game theory and neuroscience. *Science, 318*(5850), 598–602.

Scott, L. S., & Monesson, A. (2009). The origin of biases in face perception. *Psychological Science, 20*(6), 676–680.

Shafritz, K. M., Bregman, J. D., Ikuta, T., & Szeszko, P. R. (2015). Neural systems mediating decision-making and response inhibition for social and nonsocial stimuli in autism. *Progress in Neuro-Psychopharmacology and Biological Psychiatry, 60,* 112-120.

Shenhav, A., Straccia, M. A., Botvinick, M. M., & Cohen, J. D. (2016). Dorsal anterior cingulate and ventromedial prefrontal cortex have inverse roles in both foraging and economic choice. *Cognitive, Affective, & Behavioral Neuroscience, 16*(6), 1127-1139.

Shenhav, A., Straccia, M. A., Cohen, J. D., & Botvinick, M. M. (2014). Anterior cingulate engagement in a foraging context reflects choice difficulty, not foraging value. *Nature neuroscience, 17*(9), 1249-1254.

Singer, T., Critchley, H. D., & Preuschoff, K. (2009). A common role of insula in feelings, empathy and uncertainty. *Trends in cognitive sciences, 13*(8), 334–340.

Skelly, L. R., & Decety, J. (2012). Passive and motivated perception of emotional faces: Qualitative and quantitative changes in the face processing network. *PloS one, 7*(6), e40371.

Smith, B. W., Mitchell, D. G., Hardin, M. G., Jazbec, S., Fridberg, D., Blair, R. J. R., & Ernst, M. (2009). Neural substrates of reward magnitude, probability, and risk during a wheel of fortune decision-making task. *Neuroimage, 44*(2), 600–609.

Smith, D. V., Hayden, B. Y., Truong, T.-K., Song, A. W., Platt, M. L., & Huettel, S. A. (2010). Distinct value signals in anterior and posterior ventromedial prefrontal cortex. *Journal of Neuroscience, 30*(7), 2490-2495.

Smith, S. M., Jenkinson, M., Woolrich, M. W., Beckmann, C. F., Behrens, T. E., Johansen-Berg, H., Bannister, P. R., De Luca, M., Drobnjak, I., Flitney, D. E., et al. (2004). Advances in functional and structural mr image analysis and implementation as fsl. *Neuroimage, 23*, S208–S219.

Spiegelhalter, D. J., Best, N. G., Carlin, B. P., & Van Der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the royal statistical society: Series b (statistical methodology), 64*(4), 583–639.

Steinemann, N. A., O'Connell, R. G., & Kelly, S. P. (2018). Decisions are expedited through multiple neural adjustments spanning the sensorimotor hierarchy. *Nature communications, 9*(1), 1–13.

Stirrat, M., & Perrett, D. I. (2010). Valid facial cues to cooperation and trust: Male facial width and trustworthiness. *Psychological science, 21*(3), 349–354.

Sugrue, L. P., Corrado, G. S., & Newsome, W. T. (2004). Matching behavior and the representation of value in the parietal cortex. *science, 304*(5678), 1782–1787.

Sui, J., Adali, T., Yu, Q., Chen, J., & Calhoun, V. D. (2012). A review of multivariate methods for multimodal fusion of brain imaging data. *Journal of neuroscience methods, 204*(1), 68–81.

Summerfield, C., & Tsetsos, K. (2012). Building bridges between perceptual and economic decision-making: Neural and computational mechanisms. *Frontiers in neuroscience, 6*, 70.

Sutton, R. S., & Barto, A. G. (1998). Reinforcement learning: An introduction mit press. *Cambridge, MA, 22447*.

Suzuki, A., & Suga, S. (2010). Enhanced memory for the wolf in sheep's clothing:: Facial trustworthiness modulates face-trait associative memory. *Cognition, 117*(2), 224–229.

Suzuki, S., Adachi, R., Dunne, S., Bossaerts, P., & O'Doherty, J. P. (2015). Neural mechanisms underlying human consensus decision-making. *Neuron, 86*(2), 591–602.

Suzuki, S., & O'Doherty, J. P. (2020). Breaking human social decision making into multiple components and then putting them together again. *Cortex, 127*, 221–230.

Tarantola, T., Kumaran, D., Dayan, P., & De Martino, B. (2017). Prior preferences beneficially influence social and non-social learning. *Nature communications*, *8*(1), 1–14.

Tingley, D. (2014). Face-off: Facial features and strategic choice. *Political Psychology*, *35*(1), 35-55.

Toelch, U., Bruce, M. J., Newson, L., Richerson, P. J., & Reader, S. M. (2014). Individual consistency and flexibility in human social information use. *Proceedings of the Royal Society B: Biological Sciences*, *281*(1776), 20132864.

Tognetti, A., Berticat, C., Raymond, M., & Faurie, C. (2013). Is cooperativeness readable in static facial features? an inter-cultural approach. *Evolution and Human Behavior*, *34*(6), 427–432.

Tom, S. M., Fox, C. R., Trepel, C., & Poldrack, R. A. (2007). The neural basis of loss aversion in decision-making under risk. *Science*, *315*(5811), 515–518.

Tremblay, S., Larochelle-Brunet, F., Lafleur, L.-P., El Mouderrib, S., Lepage, J.-F., & Théoret, H. (2016). Systematic assessment of duration and intensity of anodal transcranial direct current stimulation on primary motor cortex excitability. *European Journal of Neuroscience*, *44*(5), 2184–2190.

Tsetsos, K., Pfeffer, T., Jentgens, P., & Donner, T. H. (2015). Action planning and the timescale of evidence accumulation. *PloS one*, *10*(6), e0129473.

Turner, B. M., Van Maanen, L., & Forstmann, B. U. (2015). Informing cognitive abstractions through neuroimaging: The neural drift diffusion model. *Psychological review*, *122*(2), 312.

Tusche, A., & Bas, L. M. (2021). Neurocomputational models of altruistic decision-making and social motives: Advances, pitfalls, and future directions. *Wiley Interdisciplinary Reviews: Cognitive Science*, *12*(6), e1571.

Twomey, D. M., Kelly, S. P., & O'Connell, R. G. (2016). Abstract and effector-selective decision signals exhibit qualitatively distinct dynamics before delayed perceptual reports. *Journal of Neuroscience*, *36*(28), 7346–7352.

Tzieropoulos, H. (2013). The trust game in neuroscience: A short review. *Social neuroscience*, *8*(5), 407–416.

Uchida, N., Kepecs, A., & Mainen, Z. F. (2006). Seeing at a glance, smelling in a whiff: Rapid forms of perceptual decision making. *Nature Reviews Neuroscience*, *7*(6), 485-491.

Uleman, J. S., & Kressel, L. M. (2013). A brief history of theory and research on impression formation.

Urai, A. E., Braun, A., & Donner, T. H. (2017). Pupil-linked arousal is driven by decision uncertainty and alters serial choice bias. *Nature communications*, *8*(1), 1–11.

Valdes-Sosa, P. A., Sanchez-Bornot, J. M., Sotero, R. C., Iturria-Medina, Y., Aleman-Gomez, Y., Bosch-Bayard, J., Carbonell, F., & Ozaki, T. (2009). Model driven eeg/fmri fusion of brain oscillations. *Human brain mapping*, *30*(9), 2701–2721.

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

Van't Wout, M., & Sanfey, A. G. (2008). Friend or foe: The effect of implicit trustworthiness judgments in social decision-making. *Cognition*, *108*(3), 796–803.

van Vugt, M. K., Beulen, M. A., & Taatgen, N. A. (2019). Relation between centroparietal positivity and diffusion model parameters in both perceptual and memory-based decision making. *Brain research*, *1715*, 1–12.

Verdonck, S., Loossens, T., & Philiastides, M. G. (2021). The leaky integrating threshold and its impact on evidence accumulation models of choice response time (rt). *Psychological Review*, *128*(2), 203.

Von Neumann, J., & Morgenstern, O. (2007). *Theory of games and economic behavior*. Princeton university press.

von Lautz, A., Herding, J., & Blankenburg, F. (2019). Neuronal signatures of a random-dot motion comparison task. *NeuroImage*, *193*, 57–66.

Voss, A., Rothermund, K., & Voss, J. (2004). Interpreting the parameters of the diffusion model: An empirical validation. *Memory & cognition*, *32*(7), 1206–1220.

Vuilleumier, P., Richardson, M. P., Armony, J. L., Driver, J., & Dolan, R. J. (2004). Distant influences of amygdala lesion on visual cortical activation during emotional face processing. *Nature neuroscience*, *7*(11), 1271–1278.

Wabersich, D., & Vandekerckhove, J. (2014). Extending jags: A tutorial on adding custom distributions to jags (with a diffusion model example). *Behavior research methods*, *46*(1), 15–28.

Wagenmakers, E.-J., Steyvers, M., Raaijmakers, J. G., Shiffrin, R. M., Van Rijn, H., & Zeelenberg, R. (2004). A model for evidence accumulation in the lexical decision task. *Cognitive Psychology*, *48*(3), 332–367.

Wan, X., Cheng, K., & Tanaka, K. (2015). Neural encoding of opposing strategy values in anterior and posterior cingulate cortex. *Nature neuroscience*, *18*(5), 752–759.

Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). Hddm: Hierarchical bayesian estimation of the drift-diffusion model in python. *Frontiers in neuroinformatics*, *7*, 14.

Wiethoff, S., Hamada, M., & Rothwell, J. C. (2014). Variability in response to transcranial direct current stimulation of the motor cortex. *Brain stimulation*, *7*(3), 468–475.

Williams, E. H., Bilbao-Broch, L., Downing, P. E., & Cross, E. S. (2020). Examining the value of body gestures in social reward contexts. *NeuroImage*, *222*, 117276.

Wilson, J. P., & Rule, N. O. (2015). Facial trustworthiness predicts extreme criminal-sentencing outcomes. *Psychological science*, *26*(8), 1325–1331.

Wilson, J. P., & Rule, N. O. (2016). Hypothetical sentencing decisions are associated with actual capital punishment outcomes: The role of facial trustworthiness. *Social Psychological and Personality Science*, *7*(4), 331–338.

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

Woolrich, M. W., Behrens, T. E. J., Beckmann, C. F., & Smith, S. M. (2005). Mixture models with adaptive spatial regularization for segmentation with an application to fmri data. *IEEE transactions on medical imaging*, *24*(1), 1–11.

Wu, S.-W., Delgado, M. R., & Maloney, L. T. (2015). Gambling on visual performance: Neural correlates of metacognitive choice between visual lotteries. *Frontiers in neuroscience*, *9*, 314.

Wu, Z., Litwin-Kumar, A., Shamash, P., Taylor, A., Axel, R., & Shadlen, M. N. (2020). Context-dependent decision making in a premotor circuit. *Neuron*, *106*(2), 316–328.

Wunderlich, K., Rangel, A., & O'Doherty, J. P. (2009). Neural computations underlying action-based decision making in the human brain. *Proceedings of the National Academy of Sciences*, *106*(40), 17199–17204.

Xu, F., Wu, D., Toriyama, R., Ma, F., Itakura, S., & Lee, K. (2012). Similarities and differences in chinese and caucasian adults' use of facial cues for trustworthiness judgments. *PLoS One*, *7*(4), e34859.

Yacubian, J., Sommer, T., Schroeder, K., Gläscher, J., Braus, D. F., & Büchel, C. (2007). Subregions of the ventral striatum show preferential coding of reward magnitude and probability. *Neuroimage*, *38*(3), 557–563.

Yu, H., Garrod, O. G., & Schyns, P. G. (2012). Perception-driven facial expression synthesis. *Computers & Graphics*, *36*(3), 152–162.

Zaghloul, K. A., Blanco, J. A., Weidemann, C. T., McGill, K., Jaggi, J. L., Baltuch, G. H., & Kahana, M. J. (2009). Human substantia nigra neurons encode unexpected financial rewards. *Science*, *323*(5920), 1496–1499.