

# **Orbitofrontal lesion patients show an implicit approach bias to angry faces**

Macià Buades-Rotger<sup>\*1,2,3</sup>, Anne-Kristin Solbakk<sup>4,5,6,7</sup>, Matthias Liebrand<sup>1</sup>, Tor Endestad<sup>4,5,6,7</sup>, Ingrid Funderud<sup>4,5,7</sup>, Paul Siegwardt<sup>1</sup>, Dorien Enter<sup>3,8</sup>, Karin Roelofs<sup>3,8</sup>, Ulrike M. Krämer<sup>1,2</sup>

<sup>1</sup> Department of Neurology, University of Lübeck, Lübeck, Germany

<sup>2</sup> Department of Psychology, University of Lübeck, Lübeck, Germany

<sup>3</sup> Donders Institute for Brain, Cognition and Behaviour, Radboud University, Nijmegen, The Netherlands

<sup>4</sup> Department of Psychology, University of Oslo, Oslo, Norway

<sup>5</sup> RITMO Centre for Interdisciplinary Studies in Rhythm, Time and Motion, University of Oslo, Norway

<sup>6</sup> Oslo University Hospital, Rikshospitalet, Oslo, Norway

<sup>7</sup> Department of Neuropsychology, Helgeland Hospital, Mosjøen, Norway

<sup>8</sup> Behavioural Science Institute, Radboud University, Nijmegen, The Netherlands

\* Corresponding author:

Macià Buades-Rotger, PhD

Ratzeburger Allee 160

23562 Lübeck (Germany)

[macia.rotger@neuro.uni-luebeck.de](mailto:macia.rotger@neuro.uni-luebeck.de) ; +49 451 31017422

Acknowledgements and funding: this study was funded by the German Science Foundation (grant number KR3691/5-1). AKS, TE, and IF were supported by the Research Council of Norway through a grant (project number 240389) and through the Centres of Excellence scheme (project number 262762 RITMO). KR was supported by a VICI grant (#453-12-001) from the Netherlands Organization for Scientific Research (NWO) and a consolidator grant from the European Research Council (ERC\_CoG-2017\_772337). We thank Dr. Torstein R. Meling for help with patient recruitment and Dr. Per Kristian Hol for clinical evaluation of the patients' MRI scans. Dr. Buades-Rotger, Dr. Solbakk, Dr. Liebrand, Dr. Endestad, Dr. Funderud, Mr. Siegwardt, Ms. Enter, Dr. Roelofs, and Dr. Krämer report no biomedical financial interests or potential conflicts of interest.

## Abstract

Damage to the orbitofrontal cortex (OFC) can cause maladaptive social behavior, but the cognitive processes underlying these behavioral changes are still uncertain. Here, we tested whether patients with acquired OFC lesions show altered approach-avoidance tendencies to emotional facial expressions. Thirteen patients with focal OFC lesions and 31 age- and gender-matched healthy controls performed an implicit approach-avoidance task in which they either pushed or pulled a joystick depending on stimulus color. While controls avoided angry faces, OFC patients displayed an incongruent response pattern characterized by both increased approach and reduced avoidance of angry facial expressions. The approach bias was stronger in patients with higher self-reported impulsivity and disinhibition, and in those with larger lesions. Moreover, patients committed more errors in the task, which in turn was correlated with self-rated clinical impairment. We further used linear ballistic accumulator modelling to investigate latent parameters underlying approach-avoidance decisions. Controls displayed negative drift rates when approaching angry faces, whereas OFC lesions abolished this bias. In addition, OFC patients had weaker response drifts than controls during angry face avoidance. Finally, patients showed generally reduced variability in drift rates and shorter non-decision times, indicating impulsive and rigid decision-making. In sum, our findings suggest that OFC damage alters the pace of evidence accumulation in response to threat signals, eliminating a default, protective avoidant bias and facilitating dysfunctional approach behavior.

## Significance statement

Lesions in the orbitofrontal cortex (OFC) may alter social behavior, rendering individuals irritable or reckless. However, the precise cognitive mechanisms underlying these changes are unknown. We here examined whether OFC damage impacts how persons respond to social signals using a joystick-based task. Contrary to control participants, patients showed both increased approach to, and reduced avoidance of angry facial expressions, i.e. they were quicker to pull angry faces close and slower to push them away. Further analyses of reaction times revealed that OFC patients lack a default tendency against angry face approach, and that they show a slower decision build-up when avoiding angry faces. Thus, our findings suggest that OFC lesions reduce fearful responses to social threat signals.

## Introduction

Patients with damage to the orbitofrontal cortex (OFC) often show disruptive social behavior (Barrash et al. 2000; Blair 2004; Beer et al. 2006). OFC lesions typically impact adjacent white matter, thereby hindering OFC-amygdala cross-talk (Folloni et al. 2019) and rendering individuals more emotionally reactive (Motzkin et al. 2015). Consequently, antisocial behavior related to OFC dysfunction has been classically attributed to deficits in emotion regulation (Davidson et al. 2000). However, this view has proven difficult to reconcile with the many other functions ascribed to the OFC, such as subjective value computation (Clithero and Rangel 2014). Recent investigations hence suggest a more general *evaluative* and *generative* role for the OFC (Hiser and Koenigs 2018). According to this view, the OFC codes for the potential hedonic or threatening value of a given stimulus in order to steer the organism towards or away from it (Rudebeck and Rich 2018). In this framework the OFC is assumed to generate cognitive maps of current internal states and external sensory information, enabling the selection of the most appropriate course of action (Wilson et al. 2014; Stalnaker et al. 2015). Such a process has been termed model-based or goal-directed behavior because it operates on the basis of internal representations of oneself and the environment rather than by force of habit (Lucantonio et al. 2012).

From this rationale, it follows that antisocial behavior after OFC damage could arise from inaccurate assessment and selection processes. More specifically, OFC lesions might impair the ability to correctly predict the consequences of one's own actions in response to social signals (Rudebeck and Murray 2014), e.g., wrongly expecting rewards from approaching potential punishment cues. Nevertheless, evidence to support this tenet is scarce in humans with OFC lesions. One report suggests that OFC-damaged patients display an altered sense of personal distance, e.g., they get closer to strangers (Perry et al. 2016). Comparably, a study showed that persons with OFC lesions judge negative facial expressions (i.e., angry, disgusted, fearful and sad) as *more* approachable (Willis et al. 2010). It remains to be tested, however, whether these tendencies can be attributed to implicit biases during action selection, and whether these putative alterations are linked with actual impairments in daily functioning. Moreover, it is unclear which precise cognitive mechanisms underlie such abnormal behavioral dispositions. These are important steps in understanding how OFC-dependent disturbances in social behavior play out in everyday life.

In order to clarify these issues, we investigated whether OFC lesions lead to implicit response biases towards or away from negative, positive, or neutral facial expressions. We used a version of the

approach-avoidance task (AAT) wherein subjects have to either push or pull a joystick depending on the color (e.g. red or green) of a human face (Roelofs et al. 2010). Faces are programmed to grow or shrink in size accordingly, giving the impression that they loom closer or recede upon pulling and pushing, respectively. Hence, the AAT allows measuring implicit response tendencies to task-irrelevant features of the faces such as their emotional expression. A study with this task suggested that psychopaths lack automatic avoidance of angry faces, and that this effect was correlated with aggressiveness (von Borries et al. 2012). Following a similar rationale, we tested whether task scores correlated with patients' daily emotional behavior as measured with validated clinical scales in order to assess the clinical relevance of possible approach-avoidance biases.

In addition, we scrutinized the putative cognitive mechanisms underlying altered task performance in OFC patients using Linear Ballistic Accumulator (LBA) modelling on response times (Brown & Heathcote, 2008). LBA modelling assumes that decisions arise from a sequential evidence accumulation process, the speed of which is determined by multiple latent variables (e.g., pre-existing response tendencies or shorter decision latencies) that can be quantified and compared between experimental conditions and/or groups. Previous modelling studies on an explicit version of the AAT reported relatively faster evidence accumulation in healthy subjects when threatening stimuli are to be avoided (Krypotos et al. 2015; Tipples 2019). LBA modelling might hence offer insights not captured by standard methods.

## Methods

### *Participants*

The clinical sample consisted of 13 patients with chronic (> 6 months post-injury or surgery), focal damage to the ventral prefrontal cortex (mean age=50.8 [27-62], 7 women, 12 right-handed). Lesions were predominantly located in ventromedial prefrontal brain regions, with a few lesions extending more dorsally and laterally (Fig. 1A). Etiology of the lesions was either meningioma (n=9), traumatic brain injury (n=2), oligodendroglioma (n=1), or astrocytoma (n=1). The control sample was composed of 31 age- and gender-matched neurologically healthy individuals (mean age=50.1 [43-54], 19 women, all right-handed). As previously reported, patients had normal or corrected to normal vision, showed no deficits in standard neuropsychological testing, and had no motor dysfunction of the hands. However, they reported greater difficulties in executive function, metacognition, and behavioral regulation as compared to a separate control sample (see Løvstad et al., 2012 for a complete report). All patients were recruited and measured at Oslo University Hospital and the University of Oslo, whereas

the behavioral control sample was recruited and measured at the University of Lübeck. All participants provided informed consent and the study procedures adhered to the Declaration of Helsinki. The study was approved by the ethics committee of the University of Lübeck and the Regional Committee for Medical Research Ethics - South East Norway.

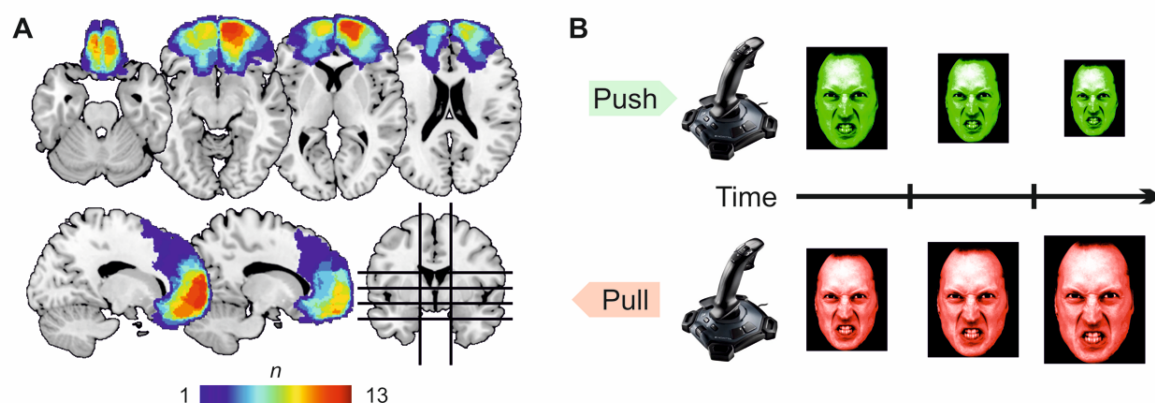
#### *Clinical scales*

Patients filled out the self-report form of the Behavior Rating Inventory of Executive Function – Adult version (BRIEF-A; Roth et al., 2005) and the Urgency, Premeditation, Perseverance, Sensation Seeking (UPPS) Impulsive Behavior Scale (Whiteside and Lynam 2001), both ad-hoc translated into Norwegian. The BRIEF-A is a standardized rating scale consisting of 75 items that tap into everyday executive functioning within the past 6 months. Internal consistency and test-retest reliability of the BRIEF-A are reportedly high and construct validity has been established in healthy and clinical populations (Waid-Ebbs et al. 2012). Since we aimed to investigate the neural control of approach-avoidance responses, for the purposes of this study we only considered the scales “Inhibit”, “Emotional Control”, and “Self-Monitor” from the BRIEF-A. The Inhibit scale measures deficits in inhibitory control and impulsivity; the Emotional Control scale assesses a person’s inability to regulate emotional responses; and the Self-Monitor scale evaluates difficulties in social or interpersonal awareness. The UPPS Impulsive Behavior Scale (Whiteside and Lynam 2001) is a 45-item self-report, assessing different facets of impulsivity on four subscales. The UPPS has been shown to display good internal consistency and construct validity (Whiteside et al. 2005). We used the total UPPS score for correlational analyses.

#### *Implicit Approach-Avoidance Task (AAT)*

Subjects performed the implicit approach-avoidance task (AAT; Fig. 1B) as previously described (Roelofs et al. 2010; von Borries et al. 2012). Stimuli were photographs (Ekman and Friesen 1976; Lundqvist et al. 1998) showing the face of one out of eight actors (four male and four female) displaying angry, happy or neutral expressions with either direct (straight) or averted (sideways) gaze. Photographs were cut out ovally and tinted red or green, amounting to a total number of 384 trials. Participants performed 18 practice trials, followed by the experimental trials. After half of the trials, subjects had a break, performed two additional practice trials to recall task demands and completed the second half. Stimuli were presented randomly, with no more than three of the same emotion-response combinations in succession.

Pictures were presented at a 1024 x 768 pixels resolution on a computer screen. We placed the joystick (Logitech Attack 3) between subject and screen to allow for comfortable pull and push movements. Participants started each trial by pressing the fire button with the index finger of the dominant hand. A face stimulus appeared in the center of the screen. Participants were instructed to ignore the facial expression and only respond to the color of the face. Half the participants had to push the joystick in response to red and pull in response to green stimuli, the other half had the opposite instruction. To visually emphasize that pull movements meant approach, and push movements meant avoidance, pictures grew or shrank in size following pull or push movements, respectively. Stimuli had a starting size of 9.5° by 13° and could shrink to a minimum of 3.5° by 4.5° when pushing or grow to a maximum of 15.5° by 20° when pulling. In practice trials, pictures remained visible after erroneous responses to allow for response correction, whereas in the task proper stimuli disappeared after they had reached minimal or maximal size. Participants were instructed to respond as quickly and accurately as possible. Importantly, trials could only be initiated once the joystick was placed back in its original centered position.



**Figure 1. A)** Lesion overlap. Warmer colors depict more overlap between patients. Peak overlap was located in  $x=4, y=58, z=-14$  (Montreal Neurological Institute space). **B)** Schematic depiction of the implicit Approach Avoidance Task (AAT). Subjects had to either push or pull a joystick in response to the color of the presented face while ignoring its facial expression (angry, happy, or neutral), gaze (direct or averted), gender (male or female), and identity (eight actors). Pushing made faces shrink in size, whereas pulling made them grow larger. The 384 trials were self-paced.

### Behavioral data analysis

Reaction times (RT) were recorded as time from stimulus onset until the first joystick movement. We excluded incorrect trials as well as those with RT shorter than 150ms or longer than 1000ms, and extracted mean log-transformed RT per cell (as in Bertsch et al., 2018). We then ran an analysis of variance (ANOVA) on the resulting values with within-subject factors emotion (happy, neutral, angry), actor gender (male, female), gaze (left, right, and direct), movement (pull or push), and the between-



subject factor group (OFC vs healthy controls) using the *ez* package (version 4.4-0). We modelled all relevant task factors as in previous studies with the implicit AAT (Roelofs et al. 2010; von Borries et al. 2012). In order to control for multiple testing, we applied a False Discovery Rate (FDR) correction as recommended for exploratory ANOVAs (Cramer et al. 2016). Color and condition were counterbalanced across participants (green=pull for one half, green=push for the other half) and are thus controlled for by design. We inspected significant effects with post-hoc t-tests.

Due to the relatively low and unevenly distributed number of errors, we simply compared the mean error rate between groups using a Welch's t-test, which is robust to unequal variances and uneven sample sizes (Ruxton, 2006). Subsequently, we computed Pearson correlation coefficients between AAT scores (between-condition differences in RT and overall error rates) and each of the four clinical scales. We assessed the robustness of significant correlations with bootstrap resampling to obtain 95% bias-corrected accelerated confidence intervals (BCa CI) with 10000 iterations using the *bootstrap* package (version 2019.5). We performed all analyses described in this section in R (version 3.6.1) running on R Studio (version 1.1.423).

#### *Linear ballistic accumulator (LBA) modelling of reaction times*

We subsequently implemented Linear Ballistic Accumulator (LBA) modelling on reaction time data (Brown and Heathcote 2008). LBA models assume that decisions stem from a sequential evidence accumulation process (Fig. 3A). Evidence for each response option is gathered linearly by a separate accumulator, which races against the other/s until one of them reaches a decision threshold. Evidence accumulation starts after a variable period of non-decision time and its speed is given by the drift rate, which is sampled from a normal distribution. The standard deviation of this distribution constitutes what we here label drift noise, i.e., variability in the pace of evidence accumulation. In addition, the accumulators might begin each trial from a different starting point, which is drawn from a uniform distribution. Therefore, a response option will be taken more quickly if starting point and decision threshold are nearer, if the drift rate is higher and less variable, and if the non-decision time is shorter. LBA models are akin to the now-popular drift diffusion models (DDM), but are simpler and more tractable computationally and thus well-suited for the relatively low amount of trials available in the present dataset (see Heathcote and Hayes, 2012, for a detailed empirical comparison between LBA and DDM).

Here, we fitted a series of LBA models with two accumulators (approach and avoidance) and four parameters: decision threshold, starting point, drift rate, and drift noise. We tested a total of 16 models in which a given combination of these parameters was allowed to vary between the six

experimental conditions of interest: pull angry, pull happy, pull neutral, push angry, push happy, and push neutral. We could not test for a modulation of experimental condition on non-decision time because models including this effect failed to converge in most subjects. See Table 1 for a summary of all models. We fitted each model on each participant's reaction time data using full information maximum likelihood estimation as implemented in the *glba* package version 0.2. We used raw RT excluding errors and responses quicker than 150ms or slower than 1s. For model comparison and inspected which model yielded the lowest Bayesian Information Criterion (BIC) values across participants. BIC is a standard goodness of fit measure that penalizes for model complexity (Raftery 1995; Burnham and Anderson 2004). Our model fitting and comparison approach is highly comparable to that of a recent DDM study on social approach-avoidance decisions (Mennella et al. 2020). Afterwards, we simulated data per group using the *rlba()* function and the average parameter estimates from the winning model. Finally, we compared the parameters of the winning model between groups with independent-samples Welch t-tests. We used R (version 3.6.1) running on R Studio (version 1.1.423) for all analyses in this section.

#### *Neuroimaging data acquisition and analysis*

Structural brain volumes were recorded at the Intervention center at Oslo University hospital (Norway) on a Philips Ingenia 3-T scanner. We acquired structural images with a T1-weighted 3D turbo gradient-echo sequence with the following settings: repetition time (TR)=1.900ms, echo time (TE)=2.23ms, flip angle=8°, voxel size=1mm<sup>3</sup>, field-of-view (FOV)=256x256mm. Members of the team at the University of Oslo, trained in lesion reconstruction, manually delineated lesion masks on each patient's anatomical images. We normalized these masks as recommended for lesioned brains (Ripollés et al. 2012) and created lesion overlap maps using MRIcron (Rorden and Brett 2000). We also inspected whether lesion size was linked with reaction times and error rates in the task. We correlated lesion size with behavioral parameters showing a group difference in the AAT and obtained the 95% bootstrapped CIs with 10000 iterations using the *bootstrap* R package to assess these effects' robustness.

## **Results**

### *Approach-Avoidance Task (AAT) results*

In our primary analysis of reaction times we observed main effects of group ( $F_{1,42}=11.92$ ,  $p=.001$ ,  $pFDR=.013$ ) and emotion ( $F_{2,84}=6.89$ ,  $p=.001$ ,  $pFDR=.010$ ) which were qualified by an emotion x movement interaction that did not survive multiple comparison correction ( $F_{2,84}=4.36$ ,  $p=.015$ ,  $pFDR=.083$ ), and, crucially, by a group x emotion x movement interaction ( $F_{2,84}=12.64$ ,  $p<.001$ ,  $pFDR<.001$ ). In order to dissect the latter three-way interaction, we computed the difference between

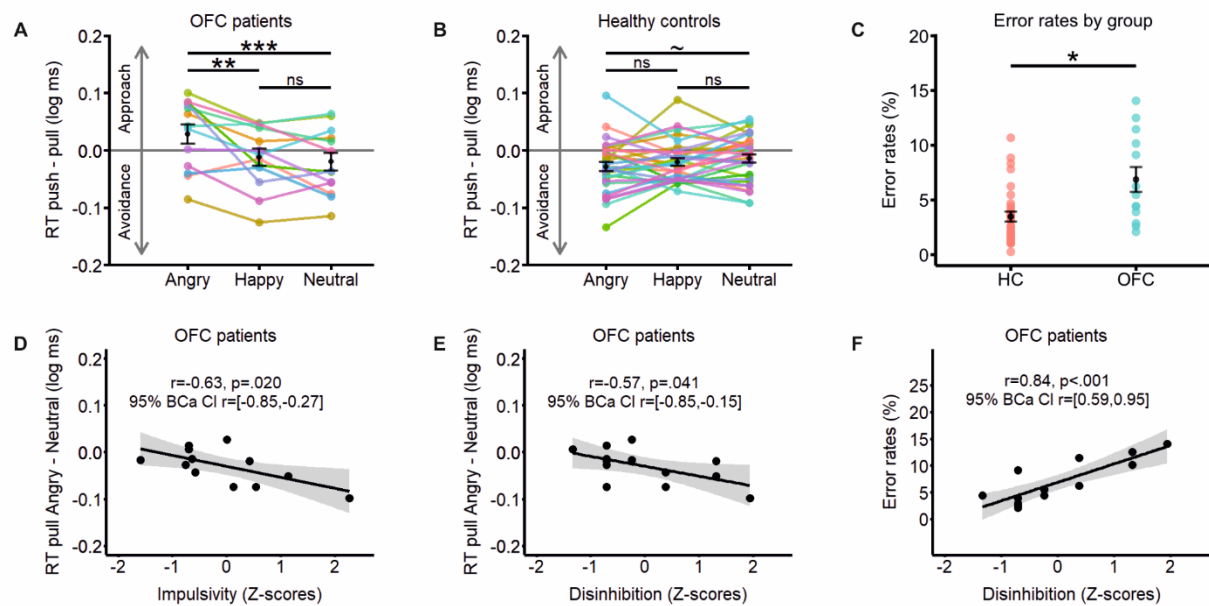


push and pull (i.e., approach minus avoidance) for each emotion and inspected for differences between emotion categories in each group, following previous work (Roelofs et al. 2010; von Borries et al. 2012). As shown in Fig. 2A, OFC patients showed a stronger approach bias toward angry relative to both happy ( $t_{12}=3.17$ ,  $p=.008$ ) and neutral faces ( $t_{12}=4.32$ ,  $p<.001$ ), with no difference between happy and neutral faces ( $p=.416$ ). In comparison (Fig. 2B), controls showed a trend-level avoidant bias for angry relative to neutral faces ( $t_{30}=1.75$ ,  $p=.089$ ), with no further differences between categories (all  $p>.272$ ). Thus, OFC patients were generally slower when pushing angry faces away relative to pulling them close.

In order to ascertain whether these effects were predominantly driven by approach or avoidance, we computed the difference in reaction times between emotions separately for push and pull movements in each group. Regarding approach movements, OFC patients were faster to pull angry relative to neutral ( $t_{12}=3.66$ ,  $p=.003$ ) but not happy faces ( $p=.107$ ). Controls showed no between-emotion differences in pull movements (all  $p>.278$ ). For avoidance movements, OFC patients were slower to push angry relative to happy faces ( $t_{12}=2.88$ ,  $p=.013$ ) but comparably fast when pushing angry and neutral ones ( $p=.284$ ). Controls were quicker to push angry as compared to neutral faces ( $t_{30}=2.27$ ,  $p=.030$ ) but not happy ones ( $p=.605$ ). Therefore, controls specifically showed avoidance of angry in comparison with neutral expressions. In contrast, OFC patients showed increased approach of angry relative to neutral faces, and reduced avoidance of angry as compared to happy ones. We used these significant between-emotion differences for later correlation analyses, as they index the increased threat approach (pull angry minus pull neutral) and reduced threat avoidance (push angry minus push happy) demonstrated by OFC patients.

Additionally, there was an emotion x gaze interaction across the whole sample ( $F_{4,168}=4.63$ ,  $p=.001$ ,  $pFDR=.011$ ). We computed the difference in reaction times between direct and averted gaze and compared between emotions over all participants to further investigate this effect. The interaction was driven by slower reactions to directly-gazing neutral faces relative to happy ( $t_{43}=3.09$ ,  $p=.003$ ) and, at trend level, angry ones ( $t_{43}=1.81$ ,  $p=.077$ ).

We subsequently compared error rates between groups. Although both groups performed the task well, OFC patients committed about twice as many errors ( $6.87\pm1.13\%$ ) than healthy controls ( $3.47\pm0.46\%$ ),  $t_{16.14}=2.77$ ,  $p=.013$  (Fig. 1C).



**Figure 2.** **A)** Patients with orbitofrontal cortex (OFC) lesions showed an approach bias (reaction times [RT] for push minus pull) towards angry relative to happy and neutral faces. **B)** Healthy controls (HC) showed no bias in either direction, with a trend towards avoidance of angry relative to neutral faces. **C)** OFC patients made more errors than HC. **D)** Shorter RT for pull angry minus pull neutral trials were linked with greater self-rated impulsivity in OFC patients. **E)** Shorter RT for pull angry minus pull neutral trials were correlated with greater disinhibition in OFC patients. **F)** Error rates were correlated with all clinical self-reports in OFC patients, including greater self-rated disinhibition.  $\sim p < .1$ ,  $*p < .05$ ,  $**p < .01$ ,  $***p < .001$ .

# *Correlations between task scores and clinical scales*

We then inspected for associations between clinical scales and task-derived scores, with the aim of testing the clinical relevance of approach-avoidance biases as measured with the AAT. The approach bias for angry minus neutral faces was linked with increased self-reported impulsivity (Fig. 1D;  $r = -.63$ ,  $p = .020$ , 95% BCa CI =  $[-.86, -.28]$ ), and greater disinhibition (Fig. 1E;  $r = -.57$ ,  $p = .041$ , 95% BCa CI =  $[-.85, -.15]$ ), but there were no correlations with either of the other two clinical scales, or between the angry push minus happy push difference and any of the scales (all  $p > .160$ ). Error rates were positively associated with all clinical scales, namely impulsivity ( $r = .57$ ,  $p = .040$ , 95% BCa CI =  $[.05, .85]$ ), disinhibition (Fig. 1F;  $r = .84$ ,  $p < .001$ , 95% BCa CI =  $[.59, .95]$ ), emotional control ( $r = .59$ ,  $p = .033$ , 95% BCa CI =  $[.23, .78]$ ), and self-monitoring ( $r = .70$ ,  $p = .007$ , 95% BCa CI =  $[.25, .90]$ ).

# *Correlations between lesion size and AAT scores*

Subsequently, we tested whether task-derived response biases were linked with lesion size. Patients with larger lesions were quicker to approach angry relative to neutral faces ( $r = -.72$ ,  $p = .004$ , 95% BCa CI =  $[-.90, -.33]$ ). Lesion size was not correlated with the push angry minus push happy difference ( $p = .374$ ) or with error rates ( $p > .663$ ). Lesion extension was thus exclusively associated with threat approach, but not with the reduced threat avoidance and increased error rates displayed by OFC patients.

## Linear Ballistic Accumulator (LBA) modelling results

Next, we turned to Linear Ballistic Accumulator (LBA) modelling in order to uncover which latent decision parameters might account for OFC patients' response patterns. We provide the complete list of models in Table 1. The winning model assumed that emotional expression and movement modulated drift rates exclusively. This model had the lowest BIC across subjects (median BIC=-646.62, k=10 free parameters) and was the best-fitting model in all 13 OFC patients as well as in 28/31 control participants. According to model-comparison guidelines (Raftery 1995; Burnham and Anderson 2004), the evidence for this model can be considered substantial relative to the two next best-fitting (and slightly more complex) models, one assuming an effect of emotional expression on drift rate and drift noise (median BIC=-629.92, k=15 free parameters), and one in which emotional expression impacted drift rate and decision threshold (median BIC=-629.57, k=15 free parameters). Further, the winning model could reproduce reaction times in pull angry trials with a precision of around ~30-50ms across successive simulations for both OFC patients (example mean simulated data=495ms; mean real data=544ms) and control participants (example mean simulated data=593ms; mean real data=624ms).

**Table 1:** Summary of Linear Ballistic Accumulator models tested

Modulated parameters in model	K	Median BIC
Threshold, starting point, drift rate, drift noise	25	-546.16
Threshold, starting point, drift rate	20	-575.52
Threshold, starting point, drift noise	20	-75.84
Threshold, drift rate, drift noise	20	-606.32
Starting point, drift rate, drift noise	20	-606.47
Threshold, starting point	15	-117.53
Threshold, drift rate	15	-629.57
Threshold, drift noise	15	-119.90
Starting point, drift rate	15	-624.78
Starting point, drift noise	15	-157.82
Drift rate, drift noise	15	-629.92
Threshold	10	-161.38
Starting point	10	-107.24
<b>Drift rate</b>	<b>10</b>	<b>-646.62</b>
Drift noise	10	-172.87
Null model	5	-169.83

K: number of free parameters; BIC: Bayesian Information Criterion. The model marked in **bold** had the best fit to the data across participants.

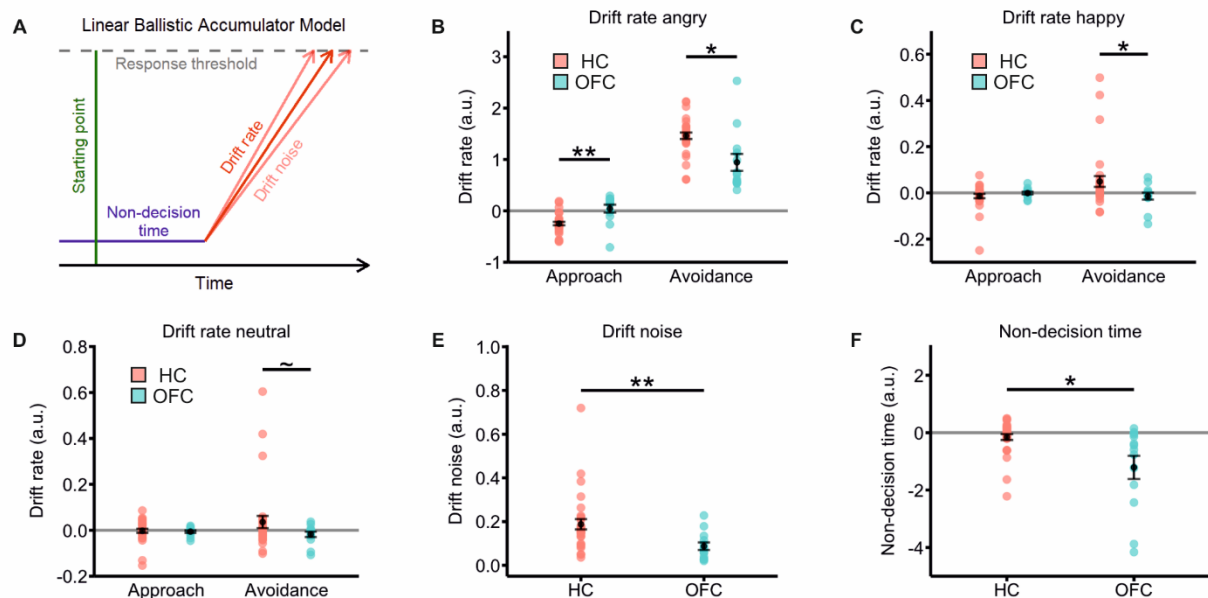
We subsequently tested whether any of the LBA model parameters differed between groups. Controls had negative response drifts when pulling angry faces close, whereas the mean value for this parameter was centered around zero in OFC patients (Fig. 3B, left;  $t_{17.46}=3.51$ ,  $p=.002$ ). OFC patients also showed lower drift rates than control participants when pushing angry faces away (Fig. 3B, right;  $t_{15.56}=2.92$ ,  $p=.010$ ). Therefore, response drifts in OFC patients were weaker when avoiding angry faces and relatively less negative (i.e. centered around null) when approaching them. OFC patients also had

reduced response drifts when pushing happy faces away (Fig. 2C, right;  $t_{41.99}=2.29$ ,  $p=.026$ ), but not when pulling them close (Fig. 2C, left;  $p=.258$ ). This pattern was also present at trend level for neutral expressions (Fig. 2D; avoid:  $t_{39.10}=1.85$ ,  $p=.070$ ; approach:  $p=.723$ ). Thus, OFC patients had generally lower drift rates than controls during avoidance movements, especially for angry faces. Regarding the remaining parameters, the patient group displayed reduced drift noise (Fig. 2E;  $t_{41.11}=3.42$ ,  $p=.001$ ; HC:  $0.18 \pm 0.02$ , OFC:  $0.08 \pm 0.01$ ), and non-decision times (Fig. 2F;  $t_{13.58}=2.54$ ,  $p=.023$ ; HC:  $-.15 \pm .10$ , OFC:  $-1.21 \pm .40$ ). There were no group differences in decision threshold ( $p=.126$ ) or starting point ( $p=.364$ ). Hence, evidence accumulation began earlier and was less variable across conditions in OFC patients.

**Table 2:** Group-wise means and standard errors of free parameters from the winning model

Parameter	HC	OFC
Drift rate pull angry**	$-.24 \pm .03$	$.04 \pm .07$
Drift rate pull happy	$-.01 \pm .009$	$-.001 \pm .005$
Drift rate pull neutral	$-.002 \pm .008$	$-.005 \pm .005$
Drift rate push angry*	$1.46 \pm .06$	$.94 \pm .16$
Drift rate push happy*	$.04 \pm .02$	$-.01 \pm .01$
Drift rate push neutral	$.03 \pm .02$	$-.01 \pm .01$
Drift noise**	$.18 \pm .02$	$.08 \pm .01$
Starting point	$.13 \pm .02$	$.08 \pm .03$
Threshold	$.87 \pm .14$	$1.56 \pm .39$
Non-decision time*	$-.15 \pm .10$	$-1.21 \pm .40$

HC: healthy controls; OFC: orbitofrontal cortex patients. Asterisks denote significant between-group differences in parameter estimates at \* $p < .05$  or \*\* $p < .01$ .



**Figure 3.** **A)** Schematic depiction of a Linear Ballistic Accumulator (LBA) model, which operationalizes decisions as the result of a sequential evidence accumulation process. The model assumes separate, competing accumulators for each response option, with faster decisions when the response threshold is lower, starting point is higher, non-decision time is shorter, and the drift towards a given option is stronger and less variable (i.e. higher drift rate and lower drift noise). We estimated the parameters from each participant's reaction time distribution with a maximum likelihood algorithm. **B)** Orbitofrontal cortex (OFC) patients showed less negative (i.e. around zero) drift rates than healthy controls (HC) when pulling angry faces close (left), and lower drift rates when pushing angry faces away (right). **C)** OFC patients had lower drift rates than HC when pushing happy faces away. **D)** OFC patients displayed trend-level lower drift rates than HC when avoiding neutral faces. **E)** OFC patients had lower drift noise. **F)** OFC patients showed shorter non-decision times. A.u.: arbitrary units. ~ $p < .1$ , \* $p < .05$ , \*\* $p < .01$ .

In a final exploratory analysis, we tested for associations between LBA parameters and clinical scales in OFC patients as done with reaction times. We limited these analyses to drift rates for threat approach (pull angry minus pull neutral) and threat avoidance (push angry minus push happy), as these were the same contrasts that we computed for correlations with reaction times. There were no associations between either score and any of the clinical scales (all  $p > .234$ ).

## Discussion

Maladaptive social behavior is common after orbitofrontal cortex (OFC) damage (Davidson et al. 2000), but the neurocognitive processes underlying these symptoms remain elusive. Here, we tested whether patients with acquired OFC lesions show altered automatic responses to emotional facial expressions. OFC patients displayed both reduced avoidance of, and increased approach to angry faces. Modelling of reaction times revealed relatively slower evidence accumulation when avoiding angry faces in OFC patients relative to controls. Moreover, patients lacked the negative response drifts that controls showed during approach of angry expressions. OFC patients further evinced less variable and earlier-starting evidence accumulation. The approach bias in OFC patients was associated with self-reported clinical measures of impulsive and disinhibited behavior. Patients also committed more errors, which was in turn correlated with greater self-reported impulsivity, disinhibition, problems in emotional control, and worse self-monitoring. Finally, larger lesions were linked with a relatively more pronounced approach bias to angry faces, but not with error rates or avoidance biases. All in all, these findings suggest that OFC damage can precipitate maladaptive behavior by altering the implicit processing of threatening social information during action selection.

### *OFC lesions increase approach and reduce avoidance of threatening stimuli*

Our findings expand on a previous report indicating that OFC-damaged individuals report negative facial expressions to be more approachable (Willis et al. 2010). Here, we showed that this translates into observable, automatic motor behavior, such that OFC patients were quicker to actively approach angry faces (i.e., pull them towards themselves), but slower to avoid them (i.e., push them away). Reduced implicit avoidance of angry faces has been reported in psychopathic offenders (von Borries et al. 2012), who also display dampened physiological reactivity to threatening distractors (Newman et al. 2010). Therefore, both lower threat aversion and enhanced threat approach seem to be at play in populations showing disruptive social behavior.

The present results broadly converge with clinical (Blair 2004), volumetric (Chester et al. 2017), and functional (Beyer et al. 2015; Gilam et al. 2015) studies asserting that the OFC is essential for the

regulation of aggressive urges. However, our data further indicate that the OFC does not merely suppress automatic impulses but rather directs the course of approach-avoidance reactions, in line with recent proposals (Hiser and Koenigs 2018; Rudebeck and Rich 2018), and with the well-known association between damage to this region and disadvantageous decision-making (Koenigs and Tranel 2007). Given that the OFC is involved in the anticipation and evaluation of actions related to certain stimuli (Wilson et al. 2014), we suggest that OFC dysfunction gives rise to an altered processing of threat signals. Specifically, it might be that OFC damage compromises the prediction of behavioral outcomes associated with potentially punishing stimuli, i.e., tagging angry faces as neutral or even potentially rewarding (Rudebeck and Murray 2014). These abnormal value forecasts can in turn enable the impulsive, rule-breaking behavior that characterizes the sequelae of OFC lesions.

In line with the latter statement, approach towards angry relative to neutral faces was linked with greater self-reported disinhibition and impulsive behavior. Paralleling our results, it has been reported that patients with borderline personality disorder, who regularly engage in antagonistic and aggressive behavior, also show an approach bias to angry faces (Bertsch et al. 2018) and comparable levels of impulsivity and self-reported anger as those of OFC patients (Berlin et al. 2005). Similarly, healthy individuals with high trait anger are quicker to approach angry relative to happy faces (Veenstra et al. 2017). The current results thus provide further evidence that threat signals might act as appetitive stimuli for individuals with externalizing symptomatology (Chester 2017), and further add that OFC lesions might precipitate such dysfunctional evaluation processes.

Of note, the response tendencies observed in OFC patients were independent of gaze direction. This pattern deviates from previous studies reporting group-specific approach-avoidance biases exclusively for directly-gazing angry faces (Roelofs et al. 2010; von Borries et al. 2012). Hence, the present findings tentatively suggest that OFC lesions might be associated with reduced sensitivity to gaze direction. We did find, however, that straight-looking neutral faces were linked with slower reaction times across the whole sample irrespective of movement type. The latter observation insinuates that neutral expressions, due to their inherent ambiguity (Blasi et al. 2009), are more thoroughly evaluated when they are directed to oneself.

Importantly, OFC patients performed generally worse in the approach-avoidance task (AAT) than controls. This is largely in line with previous findings on the role of the OFC and lateral frontal pole in controlling social approach-avoidance behavior (Roelofs et al. 2009; Volman et al. 2011). Here, subjects committed more errors than controls in an implicit version of the AAT, which is suggestive of



difficulties in ignoring task-irrelevant stimulus features. This observation concurs with other studies in showing that OFC patients are more susceptible to distraction by to-be-ignored stimulus characteristics (Mäki-Marttunen et al. 2017; Kuusinen et al. 2018), and agrees with the general idea that OFC damage hinders the implementation of goal-directed behavior (Rudebeck and Rich 2018). Moreover, error rates were associated with greater self-reported impulsivity and disinhibition in OFC patients, as well as with worse emotional control and self-monitoring. Such findings speak for the predictive validity of the AAT and support its potential usefulness for assessing emotional dysfunction in neurological patients (Fricke and Vogel 2020).

#### *OFC lesions affect latent decision parameters*

We used Linear Ballistic Accumulator (LBA) modelling to delve deeper into the decision processes underlying approach-avoidance responses in OFC patients. These analyses indicated that emotional facial expressions modulated drift rates (i.e., the speed of evidence accumulation after a stimulus appears) but no other parameters. These findings extend previous drift diffusion modelling work using an explicit version of the AAT in which emotional expressions impacted not only drift rates but also response thresholds and non-decision times (Tipples 2019). Hence, the influence of emotional expressions on latent decision variables may be less pronounced when facial expressions are to be ignored. The present data do however fully dovetail previous modelling studies in that response drifts were maximal when threatening stimuli were to be avoided (Krypotos et al. 2015; Tipples 2019). Our results complement these findings by showing that angry faces automatically bias evidence accumulation towards avoidance even in the absence of explicit response contingencies.

Between-group comparisons of model parameters revealed profound differences between OFC patients and control participants. OFC patients showed near-zero drift rates when approaching (i.e., pulling) angry faces, whereas healthy controls showed negative values in this parameter. OFC lesions might thus eliminate a default bias against threat approach. In addition, we observed weaker response drifts during avoidance responses (i.e., push movements) in patients relative to controls. The group difference in this parameter was strongest for angry facial expressions but also present to a lesser extent in happy and neutral trials. Evidence accumulation leading to avoidance decisions is hence more sluggish in OFC patients, and especially so in the presence of angry facial expressions. Therefore, the incongruent approach behavior often observed in OFC patients (Willis et al. 2010; Perry et al. 2016) might be partly attributable to an altered evidence accumulation process in response to threat signals. Specifically, evidence accumulation in OFC patients seems to lack a bias against threat approach and

is slower when threatening stimuli are to be avoided. In control participants, in contrast, the positive drift rates when pushing angry faces away might have outweighed the negative drifts when pulling them close, resulting in threat avoidance. These observations agree with the idea that the OFC encodes the currently relevant state-space (Wilson et al. 2014; Stalnaker et al. 2015). Angry facial expressions should, on the basis of previous experience, evoke a representation of possible negative outcomes and thereby facilitate avoidance, as seen in control participants. This negative outcome representation is abolished after OFC lesions, presumably producing the observed alterations in evidence accumulation and the resulting abnormal approach-avoidance tendencies.

In addition, OFC patients displayed relatively shorter non-decision times and lower drift rate variability irrespective of experimental condition. This implies that approach-avoidance decision processes start earlier and are more rigid in OFC patients as compared to control participants. The lower non-decision times are in consonance with the generally speeded responding and higher error rates incurred by OFC patients, as well as with the enhanced impulsivity often observed in OFC-damaged individuals (Berlin et al. 2004, 2005). On the other hand, the reduced drift rate variability observed in patients parallels the deficits in goal-directed behavior subsequent to OFC damage, i.e., a failure to update stimulus value resulting in perseverative responses (Rudebeck et al. 2013; Rudebeck and Murray 2014). Importantly, we observed no group differences in starting point or decision threshold, indicating that the approach bias observed in OFC patients is likely due to post-stimulus processing rather than to pre-existing response tendencies. Taken together, LBA results suggest that damage to the OFC might lead to rapid and invariant evidence accumulation, which is in turn slower when avoiding threatening stimuli but relatively faster when approaching these signals.

### *Limitations*

The cross-sectional nature of the design, along with the reduced sample size common in studies with focal lesion patients (Motzkin et al. 2015; Pujara et al. 2016), constrain the generalizability of the present results. Special caution should be exercised regarding the correlations: even though we used bootstrapping to assess their robustness, the ability of the implicit AAT to track interindividual differences is uncertain due to the lack of data on this instrument's reliability (Hedge et al. 2018). In general, effect sizes from discovery studies such as the present one should be assumed to be inflated until replication or follow up studies permit a more precise estimation of the true effect (Wilson et al. 2020). It should also be noted that some lesions affected medial and anterior portions of the prefrontal cortex, and damage in these regions has been linked with reduced punishment sensitivity (Gläscher et al. 2019). In

partial agreement with this finding, we observed that threat approach (but not threat avoidance) was more pronounced in patients with larger lesions. Nonetheless, the strongest lesion overlap was located in ventromedial aspects. Lesion-symptom mapping in larger patient samples is needed to clarify the regional specificity of the observed effects (Gläscher et al. 2019). Finally, due to time constraints, we were not able to measure patients' explicit emotion recognition abilities, which are sometimes (Heberlein et al. 2008) but not always (Willis et al. 2010) impaired in OFC patients. This limitation is minimized by the fact that the task did not require emotion recognition to be performed.

### Conclusion

The present study provides insight on how OFC dysfunction impacts the processing of threatening information during approach-avoidance decisions. This was manifested in altered evidence accumulation in response to threatening stimuli in combination with markers of premature and inflexible decision-making. Intervention programs to improve social functioning in OFC patients might therefore benefit from a focus on correctly interpreting and reacting to emotional information as well as on ameliorating impulsivity (Levine et al. 2008). In sum, our study demonstrates that OFC damage can steer individuals towards maladaptive approach behavior by biasing the automatic evaluation of threat signals.

### References

- Barrash J, Tranel D, Anderson SW. 2000. Acquired Personality Disturbances Associated With Bilateral Damage to the Ventromedial Prefrontal Region. *Dev Neuropsychol.* 18:355–381.
- Beer JS, John OP, Scabini D, Knight RT. 2006. Orbitofrontal cortex and social behavior: integrating self-monitoring and emotion-cognition interactions. *J Cogn Neurosci.* 18:871–879.
- Berlin HA, Rolls ET, Iversen SD. 2005. Borderline Personality Disorder, Impulsivity, and the Orbitofrontal Cortex. *Am J Psychiatry.* 162:2360–2373.
- Berlin HA, Rolls ET, Kischka U. 2004. Impulsivity, time perception, emotion and reinforcement sensitivity in patients with orbitofrontal cortex lesions. *Brain.* 127:1108–1126.
- Bertsch K, Roelofs K, Roch PJ, Ma B, Hensel S, Herpertz SC, Volman I. 2018. Neural correlates of emotional action control in anger-prone women with borderline personality disorder. *J Psychiatry Neurosci JPN.* 43:161–170.
- Beyer F, Münte TF, Göttlich M, Krämer UM. 2015. Orbitofrontal Cortex Reactivity to Angry Facial Expression in a Social Interaction Correlates with Aggressive Behavior. *Cereb Cortex.* 25:3057–3063.
- Blair RJR. 2004. The roles of orbital frontal cortex in the modulation of antisocial behavior. *Brain Cogn.* 55:198–208.
- Blasi G, Hariri AR, Alce G, Taurisano P, Sambataro F, Das S, Bertolino A, Weinberger DR, Mattay VS. 2009. Preferential amygdala reactivity to the negative assessment of neutral faces. *Biol Psychiatry.* 66:847–853.
- Brown SD, Heathcote A. 2008. The simplest complete model of choice response time: linear ballistic accumulation. *Cognit Psychol.* 57:153–178.
- Burnham KP, Anderson DR. 2004. Multimodel Inference: Understanding AIC and BIC in Model Selection. *Sociol Methods Res.* 33:261–304.
- Chester DS. 2017. The Role of Positive Affect in Aggression. *Curr Dir Psychol Sci.* 26:366–370.
- Chester DS, Lynam DR, Milich R, DeWall CN. 2017. Physical aggressiveness and gray matter deficits in ventromedial prefrontal cortex. *Cortex.* 97:17–22.

- Clithero JA, Rangel A. 2014. Informatic parcellation of the network involved in the computation of subjective value. *Soc Cogn Affect Neurosci*. 9:1289–1302.
- Cramer AOJ, van Ravenzwaaij D, Matzke D, Steingrover H, Wetzels R, Grasman RPPP, Waldorp LJ, Wagenmakers E-J. 2016. Hidden multiplicity in exploratory multiway ANOVA: Prevalence and remedies. *Psychon Bull Rev*. 23:640–647.
- Davidson RJ, Putnam KM, Larson CL. 2000. Dysfunction in the Neural Circuitry of Emotion Regulation—A Possible Prelude to Violence. *Science*. 289:591–594.
- Ekman P, Friesen WV. 1976. *Pictures of Facial Affect*. Consult Psychol Press.
- Folloni D, Sallet J, Khrapitchev AA, Sibson N, Verhagen L, Mars RB. 2019. Dichotomous organization of amygdala/temporal-prefrontal bundles in both humans and monkeys. *eLife*. 8:e47175.
- Fricke K, Vogel S. 2020. How interindividual differences shape approach-avoidance behavior: Relating self-report and diagnostic measures of interindividual differences to behavioral measurements of approach and avoidance. *Neurosci Biobehav Rev*. 111:30–56.
- Gilam G, Lin T, Raz G, Azrielant S, Fruchter E, Ariely D, Hendler T. 2015. Neural substrates underlying the tendency to accept anger-infused ultimatum offers during dynamic social interactions. *NeuroImage*. 120:400–411.
- Gläscher J, Adolphs R, Tranel D. 2019. Model-based lesion mapping of cognitive control using the Wisconsin Card Sorting Test. *Nat Commun*. 10:20.
- Heathcote A, Hayes B. 2012. Diffusion versus linear ballistic accumulation: different models for response time with different conclusions about psychological mechanisms? *Can J Exp Psychol Rev Can Psychol Exp*. 66:125–136.
- Heberlein AS, Padon AA, Gillihan SJ, Farah MJ, Fellows LK. 2008. Ventromedial Frontal Lobe Plays a Critical Role in Facial Emotion Recognition. *J Cogn Neurosci*. 20:721–733.
- Hedge C, Powell G, Sumner P. 2018. The reliability paradox: Why robust cognitive tasks do not produce reliable individual differences. *Behav Res Methods*. 50:1166–1186.
- Hiser J, Koenigs M. 2018. The Multifaceted Role of the Ventromedial Prefrontal Cortex in Emotion, Decision Making, Social Cognition, and Psychopathology. *Biol Psychiatry*. 83:638–647.
- Koenigs M, Tranel D. 2007. Irrational economic decision-making after ventromedial prefrontal damage: evidence from the Ultimatum Game. *J Neurosci Off J Soc Neurosci*. 27:951–956.
- Kryptos A-M, Beckers T, Kindt M, Wagenmakers E-J. 2015. A Bayesian hierarchical diffusion model decomposition of performance in Approach–Avoidance Tasks. *Cogn Emot*. 29:1424–1444.
- Kuusinen V, Cesnaite E, Peräkylä J, Ogawa KH, Hartikainen KM. 2018. Orbitofrontal Lesion Alters Brain Dynamics of Emotion-Attention and Emotion-Cognitive Control Interaction in Humans. *Front Hum Neurosci*. 12.
- Levine B, Turner GR, Stuss DT. 2008. Rehabilitation of frontal lobe functions. In: *Cognitive neurorehabilitation: Evidence and application*, 2nd ed. New York, NY, US: Cambridge University Press. p. 464–486.
- Løvstad M, Funderud I, Endestad T, Due-Tønnessen P, Meling TR, Lindgren M, Knight RT, Solbakk AK. 2012. Executive functions after orbital or lateral prefrontal lesions: neuropsychological profiles and self-reported executive functions in everyday living. *Brain Inj*. 26:1586–1598.
- Lucantonio F, Stalnaker TA, Shaham Y, Niv Y, Schoenbaum G. 2012. The impact of orbitofrontal dysfunction on cocaine addiction. *Nat Neurosci*. 15:358–366.
- Lundqvist D, Flykt A, Öhman A. 1998. The Karolinska directed emotional faces (KDEF). CD ROM Dep Clin Neurosci Psychol Sect Karolinska Institutet. 91:630.
- Mäki-Marttunen V, Kuusinen V, Peräkylä J, Ogawa KH, Brause M, Brander A, Hartikainen KM. 2017. Greater Attention to Task-Relevant Threat Due to Orbitofrontal Lesion. *J Neurotrauma*. 34:400–413.
- Mennella R, Vilarem E, Grèzes J. 2020. Rapid approach-avoidance responses to emotional displays reflect value-based decisions: Neural evidence from an EEG study. *NeuroImage*. 222:117253.
- Motzkin JC, Philippi CL, Wolf RC, Baskaya MK, Koenigs M. 2015. Ventromedial Prefrontal Cortex Is Critical for the Regulation of Amygdala Activity in Humans. *Biol Psychiatry*. 77:276–284.
- Newman JP, Curtin JJ, Bertsch JD, Baskin-Sommers AR. 2010. Attention moderates the fearlessness of psychopathic offenders. *Biol Psychiatry*. 67:66–70.
- Perry A, Lwi SJ, Verstaen A, Dewar C, Levenson RW, Knight RT. 2016. The role of the orbitofrontal cortex in regulation of interpersonal space: evidence from frontal lesion and frontotemporal dementia patients. *Soc Cogn Affect Neurosci*. 11:1894–1901.
- Pujara MS, Philippi CL, Motzkin JC, Baskaya MK, Koenigs M. 2016. Ventromedial Prefrontal Cortex Damage Is Associated with Decreased Ventral Striatum Volume and Response to Reward. *J Neurosci*. 36:5047.
- Raftery AE. 1995. Bayesian Model Selection in Social Research. *Sociol Methodol*. 25:111–163.

- Ripollés P, Marco-Pallarés J, de Diego-Balaguer R, Miró J, Falip M, Juncadella M, Rubio F, Rodriguez-Fornells A. 2012. Analysis of automated methods for spatial normalization of lesioned brains. *NeuroImage*. 60:1296–1306.
- Roelofs K, Minelli A, Mars RB, van Peer J, Toni I. 2009. On the neural control of social emotional behavior. *Soc Cogn Affect Neurosci*. 4:50–58.
- Roelofs K, Putman P, Schouten S, Lange W-G, Volman I, Rinck M. 2010. Gaze direction differentially affects avoidance tendencies to happy and angry faces in socially anxious individuals. *Behav Res Ther*. 48:290–294.
- Rorden C, Brett M. 2000. Stereotaxic Display of Brain Lesions. *Behav Neurol*. 12.
- Roth R, Isquith P, Gioia G. 2005. Behavior Rating Inventory of Executive Function - Adult Version (BRIEF-A). Lutz, FL: Psychological Assessment Resources.
- Rudebeck PH, Murray EA. 2014. The Orbitofrontal Oracle: Cortical Mechanisms for the Prediction and Evaluation of Specific Behavioral Outcomes. *Neuron*. 84:1143–1156.
- Rudebeck PH, Rich EL. 2018. Orbitofrontal cortex. *Curr Biol*. 28:R1083–R1088.
- Rudebeck PH, Saunders RC, Prescott AT, Chau LS, Murray EA. 2013. Prefrontal mechanisms of behavioral flexibility, emotion regulation and value updating. *Nat Neurosci*. 16:1140–1145.
- Stalnaker TA, Cooch NK, Schoenbaum G. 2015. What the orbitofrontal cortex does not do. *Nat Neurosci*. 18:620–627.
- Tipples J. 2019. Recognising and reacting to angry and happy facial expressions: a diffusion model analysis. *Psychol Res*. 83:37–47.
- Veenstra L, Schneider IK, Bushman BJ, Koole SL. 2017. Drawn to danger: trait anger predicts automatic approach behaviour to angry faces. *Cogn Emot*. 31:765–771.
- Volman I, Roelofs K, Koch S, Verhagen L, Toni I. 2011. Anterior Prefrontal Cortex Inhibition Impairs Control over Social Emotional Actions. *Curr Biol*. 21:1766–1770.
- von Borries AKL, Volman I, de Bruijn ERA, Bulten BH, Verkes RJ, Roelofs K. 2012. Psychopaths lack the automatic avoidance of social threat: Relation to instrumental aggression. *Psychiatry Res*. 200:761–766.
- Waid-Ebbs JK, Wen P-S, Heaton SC, Donovan NJ, Velozo C. 2012. The item level psychometrics of the behaviour rating inventory of executive function-adult (BRIEF-A) in a TBI sample. *Brain Inj*. 26:1646–1657.
- Whiteside SP, Lynam DR. 2001. The Five Factor Model and impulsivity: using a structural model of personality to understand impulsivity. *Personal Individ Differ*. 30:669–689.
- Whiteside SP, Lynam DR, Miller JD, Reynolds SK. 2005. Validation of the UPPS impulsive behaviour scale: a four-factor model of impulsivity. *Eur J Personal*. 19:559–574.
- Willis ML, Palermo R, Burke D, McGrillen K, Miller L. 2010. Orbitofrontal cortex lesions result in abnormal social judgements to emotional faces. *Neuropsychologia*. 48:2182–2187.
- Wilson BM, Harris CR, Wixted JT. 2020. Science is not a signal detection problem. *Proc Natl Acad Sci*. 117:5559–5567.
- Wilson RC, Takahashi YK, Schoenbaum G, Niv Y. 2014. Orbitofrontal cortex as a cognitive map of task space. *Neuron*. 81:267–279.

