



**University of
Zurich** UZH

**Zurich Open Repository and
Archive**

University of Zurich
University Library
Strickhofstrasse 39
CH-8057 Zurich
www.zora.uzh.ch

Year: 2014

**The feedback-related negativity (FRN) revisited: New insights into the localization,
meaning and network organization**

Hauser, Tobias U ; Iannaccone, Reto ; Stämpfli, Philipp ; Drechsler, Renate ; Brandeis, Daniel ; Walitza, Susanne ;
Brem, Silvia

DOI: <https://doi.org/10.1016/j.neuroimage.2013.08.028>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-88962>

Journal Article

Accepted Version

Originally published at:

Hauser, Tobias U; Iannaccone, Reto; Stämpfli, Philipp; Drechsler, Renate; Brandeis, Daniel; Walitza, Susanne; Brem, Silvia (2014). The feedback-related negativity (FRN) revisited: New insights into the localization, meaning and network organization. *NeuroImage*, 84:159-168.

DOI: <https://doi.org/10.1016/j.neuroimage.2013.08.028>

The Feedback-Related Negativity (FRN) revisited: New insights into the localization, meaning and network organization

Tobias U. Hauser^{a,b}, Reto Iannaccone^{a,c}, Philipp Stämpfli^d, Renate Drechsler^a, Daniel Brandeis^{a,b,e,f}, Susanne Walitza^{a,e}, Silvia Brem^a

^a University Clinics for Child and Adolescent Psychiatry (UCCAP), University of Zurich, Neumünsterallee 9, 8032 Zürich, Switzerland.

^b Neuroscience Center Zurich, University of Zurich and ETH Zurich, Switzerland

^c PhD Program in Integrative Molecular Medicine, University of Zurich, Zurich, Switzerland.

^d MR-Center of the Psychiatric University Hospital and the Department of Child and Adolescent Psychiatry, University of Zurich, 8032 Zürich, Switzerland.

^e Zurich Center for Integrative Human Physiology, University of Zurich, Zurich, Switzerland.

^f Department of Child and Adolescent Psychiatry and Psychotherapy, Central Institute of Mental Health, Medical Faculty Mannheim/ Heidelberg University, 68159 Mannheim, Germany.

Corresponding author:

Silvia Brem

University Clinics for Child and Adolescent Psychiatry (UCCAP)

University of Zurich

Neumünsterallee 9

8032 Zürich

Switzerland

sbrem@kjpd.uzh.ch

phone: +41 43 449 2760

fax: +41 43 449 2604

Abstract

Changes in response contingencies require adjusting ones assumptions about outcomes of behaviors. Such adaptation processes are driven by reward prediction error (RPE) signals which reflect the inadequacy of expectations. Signals resembling RPEs are known to be encoded by mesencephalic dopamine neurons projecting to the striatum and frontal regions. Although regions that consistently process RPEs, such as the dorsal anterior cingulate cortex (dACC), have been identified, only indirect evidence links timing and network organization of RPE processing in humans. In electroencephalography (EEG), which is well known for its high temporal resolution, the feedback-related negativity (FRN) has been suggested to reflect RPE processing. Recent studies, however, suggested that the FRN might reflect surprise, which would correspond to the absolute, rather than the signed RPE signals. Furthermore, the localization of the FRN remains a matter of debate.

In this simultaneous EEG-functional magnetic resonance imaging (fMRI) study, we localized the FRN directly using the superior spatial resolution of fMRI without relying on any spatial constraint or other assumption. Using two different single-trial approaches, we consistently found a cluster within the dACC. One analysis revealed additional activations of the salience network. Furthermore, we evaluated the effect of signed RPEs and surprise signals on the FRN amplitude. We considered that both signals are usually correlated and found that only surprise signals modulate the FRN amplitude. Last, we explored the pathway of RPE signals using dynamic causal modeling (DCM). We found that the surprise signals are directly projected to the source region of the FRN. This finding contradicts earlier theories about the network organization of the FRN, but is in line with a recent theory stating that dopamine neurons also encode surprise-like saliency signals.

Our findings crucially advance the understanding of the FRN. We found compelling evidence that the FRN originates from the dACC. Furthermore, we clarified the functional role of the FRN, and determined the role of the dACC within the RPE network. These findings should enable us to study the processing of surprise and adjustment signals in the dACC in healthy and also in psychiatric patients.

Keywords: Reward Prediction Error (RPE); Feedback-Related Negativity (FRN); dorsal Anterior Cingulate Cortex (dACC); simultaneous Electroencephalography (EEG) and functional Magnetic Resonance Imaging (fMRI); surprise; Dynamic Causal Modeling (DCM); probabilistic reversal learning.

1. Introduction

The ability to adapt to changes in the environment is essential for survival. A growing corpus of neuroscientific literature strongly suggests that adaptation involves similar learning processes in decision making, perception and movement control and is driven by (reward) prediction error (RPE) signals (Friston, 2010). RPEs signal the difference between expected and received outcomes of a behavior and are assumed to update the expectations about stimulus-outcome relations (Rescorla and Wagner, 1972). Classical reinforcement learning theories (Sutton and Barto, 1998) assume that learning is driven by signed RPEs (sRPEs). This means that the value of an object increases after a positive RPE (outcome for this object is better than expected), whereas the value decreases after a negative RPE (outcome is worse than expected). Other theories (Courville et al., 2006; Hayden et al., 2011; Pearce and Hall, 1980) focus on the aspect that learning is mainly driven by the surprisingness of an event, regardless of its valence, which can be characterized by absolute RPEs ($|RPEs|$). $|RPEs|$ are a measure for the unsigned deviance of an expected outcome, which is often called surprise (cf Hayden et al., 2011). Dopamine (DA) neurons of the mesencephalon have been found to display sRPE-like signals (Schultz et al., 1997). More recent findings show that these dopamine neurons also encode $|RPE|$ -like salience signals (Bromberg-Martin et al., 2010a, 2010b; Matsumoto and Hikosaka, 2009). These neurons are known to project to a variety of sub- and neocortical areas (Bromberg-Martin et al., 2010b; Oades and Halliday, 1987). A wealth of studies in human and non-human primates showed that these projected RPEs are also processed in the striatum, the orbitofrontal/ventromedial prefrontal cortex (OFC/vmPFC), the dorsal anterior cingulate cortex (dACC) and the dorsolateral prefrontal cortex (dlPFC) and the amygdalae (Gläscher et al., 2010; Glimcher, 2011; Haber and Knutson, 2010; Hare et al., 2008; Rutledge et al., 2010). In humans, however, we have only very limited knowledge about the temporal evolution of RPE signals in the brain, because of the poor temporal resolution of conventional neuroimaging methods, such as functional magnetic resonance imaging (fMRI; Meyer-Lindenberg, 2010).

Electroencephalography (EEG) is a method in human neuroimaging, which has an excellent temporal, but lacks in spatial resolution (Meyer-Lindenberg, 2010). In a seminal work, Holroyd and Coles (2002)

suggested an event-related potential (ERP) of the EEG, the feedback-related negativity (FRN), to reflect sRPE processing. The FRN is commonly computed as the difference wave between rewards and punishments at mid-central sites, such as the vertex electrode Cz and peaks between 200 and 300 ms after feedback onset (Miltner et al., 1997; Nieuwenhuis et al., 2004). Holroyd & Coles (2002) assumed that the neuronal populations which elicit the FRN are usually inhibited by the tonic firing of the dopamine neurons. A decrease in tonic DA, as reflected by negative sRPEs, would therefore disinhibit these neuronal populations causing the widespread neuronal activity which is detectable on the scalp as FRN. Over the last decade, various studies examined the relation between the FRN and sRPEs (for review cf Walsh and Anderson, 2012). Most approaches which evaluated the association between the FRN amplitude and sRPEs (e.g., Bellebaum and Daum, 2008; Holroyd et al., 2004; Pfabigan et al., 2011) relied on static paradigms and averaged ERPs. The averaging of multiple trials using ERPs has the advantage of increasing the physiologically poor signal-to-noise-ratio (SNR) of single-trial EEG data. Unfortunately, such approaches sometimes underestimate the dynamic, context-specific nature of RPE signals in the framework of reinforcement learning. Additional evidence which directly links RPEs with the FRN amplitudes by also considering their covariation over time, such as in single-trial analyses, would therefore be important to complement the current knowledge.

Furthermore, also the localization of the FRN remains a matter of debate. Although many source estimation studies localized the FRN in the large area of the ACC (for review cf Walsh and Anderson, 2012), several studies also located the origin of the FRN in the posterior cingulate cortex (Badgaiyan and Posner, 1998; Cohen and Ranganath, 2007; Hewig et al., 2007; Müller et al., 2005; Nieuwenhuis et al., 2005b). Others (Carlson et al., 2011; Foti et al., 2011a, 2011b; Martin et al., 2009; Nieuwenhuis et al., 2005b) found the source within the basal ganglia, which is also known to crucially process RPEs. All source localization studies so far face the unsolvable inverse problem, which states that the same topographical maps measured at the scalp surface can be caused by several different sources and/or source configurations (Helmholtz, 1853). Therefore, current dipole localization studies heavily rely on prior assumptions of the experimenters, such as the number of dipoles and their (alternative) locations (Luck,

2005). Other methods, such as standardized low resolution brain electromagnetic tomography (sLORETA, Pascual-Marqui, 2002) overcome some of these problems, but show a very limited spatial resolution (see e.g., Cohen and Ranganath, 2007). Recently, multimodal imaging techniques have been introduced, such as simultaneous EEG-fMRI (Rosa et al., 2010). These methods can overcome the above mentioned limitations of EEG source localization by associating the variability of the EEG signal with the fluctuations of the fMRI (Debener et al., 2006; Huster et al., 2012; Lüchinger et al., 2012, 2011). To our knowledge, the localization of the FRN has not yet been investigated using concurrent EEG-fMRI.

An increasing number of studies contradict the assumption by Holroyd and Coles (2002) that the FRN amplitude reflects sRPEs (Alexander and Brown, 2011; Chase et al., 2010; Oliveira et al., 2007; Talmi et al., 2012). Rather, these studies suggest that the FRN reflects surprise or unexpectedness, which is known as the unsigned, absolute RPE signal (Hayden et al., 2011; Pearce and Hall, 1980). In a seminal study by Oliveira et al. (2007), the authors showed that the FRN is elicited by surprising, rather than by negative feedback. In many studies which evaluate the FRN, the negative feedback appears with a frequency well below 50% and therefore, a clear differentiation between valence (positive vs. negative feedback) and surprise is not possible. In a recent study by Chase and colleagues (Chase et al., 2010), the authors tried to link RPEs to single-trial amplitudes and found indications that larger negative sRPEs reflect larger FRN amplitudes. For positive sRPEs, the authors found a marginally significant negative relationship with FRN amplitudes (greater positive sRPEs predicted more negative amplitudes; Chase et al., 2010). This finding indicates that the FRN amplitude might be modulated by the absolute value of the RPEs rather than by sRPEs. Further evidence comes also from a functional model of the ACC by Alexander and Brown (2011), which successfully simulated the FRN deflection based on surprise-signals.

In their model, Holroyd and Coles (2002) suggested that the FRN-source lies within a widespread network of cortical and subcortical areas, such as the striatum, the amygdalae, the dlPFC and the OFC/vmPFC. They assumed that the FRN-source receives different motor programs from these areas and selects the best among them using the sRPE signal. Whether such a network is functionally plausible in the framework of FRN processing has not yet been evaluated. Furthermore, it remains unclear whether the FRN-source

directly receives dopaminergic RPE signals or whether they are transmitted via other areas such as the striatum or the vmPFC.

In order to localize the origin of the FRN, in this study, we simultaneously recorded EEG and fMRI from 15 healthy adults which performed a probabilistic reversal learning task. For the localization of the FRN, we used two complementary approaches. In the first analysis, we aimed for the maximal consistency across subjects and therefore defined the FRN latency based on the grand average. In the second analysis, we aimed for maximal individual temporal sensitivity and defined the FRN latency based on the individual difference waves. In both analyses, we found a cluster in the dorsal anterior cingulate cortex. In the second analysis, we additionally found activations in the salience network. Furthermore, we evaluated the relation between the FRN amplitude and model-derived RPEs using a single-trial approach. We found that the single-trial variability was better explained by surprise (encoded as $|RPEs|$) than by sRPEs. Additionally, we investigated the effective functional connectivity of the RPE-network surrounding the FRN source. Using dynamic causal modeling (DCM; Friston et al., 2003), we found that it is most likely that the FRN source directly receives $|RPE|$ -inputs rather than signals which are first processed in another area.

2. Material and methods

2.1 Subjects

Fifteen (9 females) healthy, right-handed adults with mean age of 25.7 years (± 2.5 SD) participated in the study. Each subject gave written informed consent, approved by the local ethics committee. The subjects received a voucher for their participation and half of the money won in the task was paid.

2.2 Task

The participants played a probabilistic reversal learning task (Fig. 1; Chase et al., 2010; Gläscher et al., 2009; Hampton and O'Doherty, 2007; Hampton et al., 2006; O'Doherty et al., 2001; Remijne et al.,

2005) of 120 trials. In each trial, the subjects had to select one of two objects (geometrical forms). They were instructed to choose the object of which they expected to get rewarded, in order to maximize their wins. One of the two presented stimuli was randomly assigned to be the correct stimulus, leading to a reward in 80%, causing a loss in 20% of the trials. As reward, the subjects received 50 Swiss Centimes (approx. USD 0.50), indicated by a framed coin. The punishment (-50 Swiss Centimes) was depicted by a crossed coin. The incorrect stimulus had a win-probability of 20% and thus resulted in punishments in 80%. After the subject chose the correct stimulus between 6 and 10 times (randomly determined; at least 3 consecutive correct responses), the reward-probabilities were reversed. The former incorrect stimulus therefore became the new correct stimulus, and vice versa. Prior to scanning, the participants were informed that the reward contingencies can occasionally switch. No information was given about the win probabilities and the frequency or other characteristics of the reversals. The participants also performed a training to familiarize with the task. During this training, no reversal occurred. The presented stimuli did not change over the course of the run. Besides the 120 trials, 40 null trials were presented (9000ms trial length). On each trial, the subjects had to select their favored stimulus within 1500ms. After selection, the chosen stimulus was highlighted until a total stimulus presentation duration of 2500ms. Subsequently, a jittered fixation cross appeared for 2750ms (2000ms - 4000ms) before the outcome (1000ms) was presented. Between two trials, a fixation cross of jittered length (mean 2750ms, 2000 - 4000ms) was shown. In order to minimize misses, late answers were punished by subtracting 100 Swiss Centimes. All subjects successfully mastered the task and selected the correct stimulus with an accuracy of 73.2% (± 4.2 SD) leading to an average total gain of 16.4 Swiss Francs (± 4.8 SD).

2.3 Reinforcement learning model

To infer the RPEs, we used a modified Rescorla-Wagner reinforcement learning model with an anticorrelated valuation system (Gläscher et al., 2009; Matsumoto et al., 2007; Sutton and Barto, 1998).

RPEs were computed as the difference between the expected value of the chosen object and the received reward on each trial:

$$RPE_t = R_t - V_t^{Chosen} \quad (1)$$

where V_t^{Chosen} denotes the value (expected outcome) of the chosen object at time t. R_t is the outcome at trial t (± 50 Swiss Centimes).

The value of each object was updated based on the RPE using anticorrelated valuation (Gläscher et al., 2009):

$$V_{t+1}^{Chosen} = V_t^{Chosen} + \alpha RPE_t \quad (2)$$

$$V_{t+1}^{Unchosen} = V_t^{Unchosen} + \alpha(-R_t - V_t^{Unchosen}) \quad (3)$$

where α describes the learning rate.

The probabilities of the choice options A and B to be selected were computed using a softmax action selection function:

$$p(A_t) = \frac{1}{1 + e^{-(V_t^A - V_t^B)\tau}} \quad (4)$$

$$p(B_t) = 1 - p(A_t) \quad (5)$$

where V_t^A denotes the value of object A at time t and τ denotes the inverse temperature of this sigmoid function.

For best fit between the selection probabilities of the model and the participants behavior, the free parameters (α , τ) were estimated using maximum log-likelihood estimates (Hampton et al., 2006):

$$\log L = \frac{\sum B_{switch} \log P_{switch}}{N_{switch}} + \frac{\sum B_{stay} \log P_{stay}}{N_{stay}} \quad (6)$$

where the behavioral component B indicates whether the subject switched on the next trial and N denotes the number of trials. P_{switch} , for example, describes the probability that the model chooses option B ($p(B_t)$), after choosing object A on the previous trial.

We compared this model with anticorrelated valuation with a standard Rescorla-Wagner model (Rescorla and Wagner, 1972; Sutton and Barto, 1998), which was set up similarly as the described model, simply replacing equation 3 with

$$V_{t+1}^{Unchosen} = V_t^{Unchosen} \quad (7)$$

In contrast to the anticorrelated model, the standard Rescorla-Wagner model updates only the value of the chosen object, but does not change the value of the unchosen object in any way. While a standard Rescorla-Wagner model performs well in environments with unrelated choice-options, we hypothesized that in reversal learning, where the win probabilities of both options are strongly linked, a steady update of both options would be beneficial.

Maximum $\log L$ was estimated for each subject separately using a genetic algorithm, resulting in an average learning rate α of 0.58 (± 0.16 SD) and a τ of 2.12 (± 1.1 SD).

2.4 EEG acquisition

EEG was recorded during fMRI-acquisition using MR-compatible DC-amplifiers (BrainProducts GmbH, Gilching, Germany). The signal was recorded with a sampling rate of 5kHz (filters highpass: DC, lowpass for all scalp electrodes: 250Hz, for ECG channels: lowpass 1000Hz, recording reference: Fz). The EEG-cap contained 63 scalp electrodes, one ground electrode placed at AFz, and two electrocardiogram (ECG) electrodes for cardioballistogram correction (CBC). The scalp electrodes covered the 10-20-system plus the following additional sides: FPz, AF1/2, FCz, CPz, POz, Oz, Iz, F5/6, FC1/2/3/4/5/6, FT7/8/9/10,

C1/2/5/6, CP1/2/3/4/5/6, TP7/8/9/10, P5/6, PO1/2/9/10, OI1/2, LE/RE (EOG electrodes placed lateral and below the left (LE) and right (RE) eyes). O1'2' and FP1'2' were placed 15% more laterally to Oz/FPz for more even coverage.

2.5 EEG analyses

Data preprocessing was conducted using Analyzer 2.01 (BrainProducts GmbH, Gilching, Germany). MR-artifact correction was computed using implemented sliding average subtraction (Allen et al., 2000). To eliminate pulse artifacts, we used CBC-algorithms provided in Analyzer and then inspected the correction manually. The data were then resampled (256Hz) and filtered (time constant 0.1s, low-pass cutoff 30Hz, 50Hz notch). Ocular and remaining cardioballistogram artifacts were excluded using independent component analysis (ICA; components excluded: 11.67 ± 2.64 , thereof due to cardioballistogram artifacts: $1.87 \pm .90$, the remaining excluded components reflected artifacts induced by blinks, lateral and vertical eye movements) (Jung et al., 2000). The continuous data was re-referenced to average reference (Lehmann and Skrandies, 1980) and epoched (from -100 to 700ms relative to the onset of the feedback screen). Further analyses were conducted using Matlab (7.11; MathWorks Inc., Natick, USA) and eeglab-toolbox (10.2.2.4; Delorme and Makeig, 2004).

For our subsequent single-trial FRN localization analyses, we used two different analytical approaches. The first approach (hereafter: FRN_{GA}) focused on the most consistent part of the FRN across the whole group. We therefore computed the difference wave between the averages of punishments (number of trials: 42.0 ± 3.74) and rewards (number of trials: 75.8 ± 5.2). Then, the timepoint of maximal negativity at the a-priori site of interest (Cz; Holroyd and Coles, 2002; Nieuwenhuis et al., 2004) in the grand average between 200 and 300ms after feedback-onset was determined (Fig. 2A-B). From this timepoint (223ms), the amplitudes of every single-trial were derived for each subject. The second approach (hereafter: FRN_{indDW}) takes into account that the subjects may differ in the individual ERP latencies. We derived the single-trial amplitudes at the most negative peak between 200 and 300ms after feedback-onset (Fig. 2A,C) from the individual difference waves, rather than from the grand average.

As previous literature suggested, the FRN amplitudes could either reflect signed RPEs (sRPE) as suggested by Holroyd and Coles (2002) or they could resemble surprise, which can be characterized by absolute RPEs (|RPE|; Hayden et al., 2011; Pearce and Hall, 1980). To evaluate the mere impact of both theories, we entered sRPEs as well as |RPEs| as predictors and the amplitudes of the FRN as the dependent variable in two separate linear regression analyses. The resulting beta-weights for each subject were entered into a t-test, similar to the summary statistics approach used in fMRI analyses (Holmes and Friston, 1998). Since the regression analysis of the sRPEs as well as of the |RPEs| turned out to be significant (see results) and because sRPEs and |RPEs| are correlated, we also entered both measures as predictor variables into a multiple regression analysis. The betas for each subject were then again entered into a t-test. Because such analyses only account for uniquely explained variance (Andrade et al., 1999), we can be sure that betas from this analysis were not just spurious effects of the correlated measure.

2.6 fMRI acquisition

MR acquisitions were performed using a 3 T Achieva whole-body system (Philips Medical Systems, Best, the Netherlands), equipped with a 32-element receive head coil array, which was designed for simultaneous EEG-recordings. FMRI data were recorded with an EPI-sequence (40 slices, 2.5*2.5*2.5mm voxels, 0.7mm gap, FA: 85° FOV: 240*240*127mm) optimized for minimal signal dropouts in orbitofrontal regions due to susceptibility-induced artifacts (TR: 1850ms, TE: 20ms, 15° tilted downwards of AC-PC). Additionally, T1-weighted structural images (FOV: 240*240*160mm, sagittal orientation, 1*1*1mm voxels, TR: 8.14ms, TE: 3.7ms, flip angle: 8°) were recorded in order to analyze the fMRI data with an optimized normalization procedure.

2.7 fMRI analyses

The fMRI-data were analyzed using SPM8 (<http://www.fil.ion.ucl.ac.uk/spm/>). First, the structural T1-images were segmented using new segmentation, which we then entered into the DARTEL-template

generation procedure (Ashburner, 2007). The realigned and coregistered EPI files were subsequently normalized to MNI space using the DARTEL-generated flow fields (new voxel size 1.5mm) and smoothed (9mm FWHM).

For the localization of the FRN-peak modulation, we entered the single-trial amplitudes (FRN_{GA} and FRN_{indDW} in two separate analyses) as parametric modulators at the time of feedback presentation into the first-level GLM (cf. Debener et al., 2006). Additionally, the onsets of the cue stimuli, missed trials and the movement parameters were modeled as regressors of no interest. For the second-level random-effects analysis, a significance threshold of $p < 0.05$, cluster-extent corrected for multiple comparisons (Forman et al., 1995; Slotnick et al., 2003) was used. Because single-trial EEG suffers from a poor SNR, we set the voxel height threshold at $p < 0.01$ to increase sensitivity. The Monte-Carlo simulations led to a cluster threshold of $k > 118$ (398.25mm^3) to correct for multiple comparisons. For EEG-localization studies, significance thresholds using a cluster-extent correction for multiple comparisons has an increased physiological plausibility than height-thresholds for single voxels, because EEG signals are known to be caused by large and spatially extended populations of neurons rather than small, highly active subpopulations.

To investigate the effective functional connectivity, we also analyzed the main effects of $sRPE / |RPE|$ processing. In two separate analyses, we entered two parametrically modulated regressors in the GLM: $sRPEs / |RPEs|$ as parametric modulators at the time of feedback presentation and the value of the chosen stimulus as parametric modulation at time of stimulus presentation. Additionally, we entered the six movement parameters as regressors of no interest to control for remaining movement artifacts. Misses were entered as regressors of no interest. For group level analyses, the same significance level was used as in the EEG-fMRI-analysis. From these analyses, we derived the localizations and time series for the dynamic causal modeling (DCM) analysis (see below). Figures and activation tables can be found in the supporting material.

2.8 Effective connectivity analysis using dynamic causal modeling (DCM)

We analyzed the effective functional connectivity using dynamic causal modeling (DCM10) in order to evaluate the network surrounding the origin of the FRN. We selected four additional striatal and prefrontal regions, which were suggested by Holroyd and Coles (2002) and are known to process RPEs. These regions were also found to be active in our fMRI analyses of the sRPEs and the |RPEs| (Fig. S1; Table S1). Our regions-of-interest (ROIs) consisted of the vmPFC/OFC, bilateral striatum, bilateral amygdalae and the right dlPFC. To extract the time series of the ROIs, we created spheres at the peak of the group level activation. A radius of 4mm was selected for the FRN_{indDW} -cluster (MNI: $x = 15, y = 18, z = 42$), left putamen (MNI: $x = -24, y = -12, z = -3$), right putamen (MNI: $x = 29, y = 0, z = 3$), left amygdala (MNI: $x = -21, y = -9, z = -18$) and right amygdala (MNI: $x = 23, y = -10, z = -15$). A radius of 8mm was chosen for the vmPFC (MNI: $x = -3, y = 56, z = -3$) and dlPFC (MNI: $x = 47, y = 9, z = 41$), due to a spatially larger activation in the fMRI data (cf Fig. S1). The peaks of the group activations for the striatum, amygdalae and the vmPFC were derived from the sRPE analysis. The peak of the dlPFC was determined from the analysis of the |RPEs|. The same analysis was also conducted using the peak of the FRN_{GA} (MNI: $x = 11, y = 24, z = 24$) analysis.

As inputs, we entered the sRPEs as well as the |RPEs|, based on the assumption that both may mimic mesencephalic dopamine activity. Phasic activity of dopamine neurons in the ventral tegmental area is known to reflect sRPE activity (Bromberg-Martin et al., 2010b; Glimcher, 2011; Schultz et al., 1997). More recently, it has been suggested that mesencephalic dopamine neurons not only represent sRPEs, but also |RPE|-like salience signals (Bromberg-Martin et al., 2010a, 2010b; Matsumoto and Hikosaka, 2009). Furthermore, it has been hypothesized that the dACC may receive sRPE- as well as |RPE|-like dopaminergic inputs directly from the mesencephalon (Bromberg-Martin et al., 2010b). Such a direct input contradicts the assumptions made by Holroyd and Coles (2002), who suggested that the sRPE-inputs are projected via the basal ganglia.

The problem for the analysis of such complex models is the enormously large model space. We therefore decided to perform a stepwise model-exploration approach. First (models 1-21), we constructed a set of models which had only one region (bilateral ROIs were combined) which received input. This region

projected to all other regions. As inputs, we chose either sRPEs, |RPEs|, or both. We additionally explored whether no connectivity between the regions or full connectivity of all regions performed better. The comparison of these models revealed a clear winner (see results). As a second step (models 22-48), we took the winning model and explored whether additional inputs through other regions improved the model fit. We also evaluated whether reciprocal or unidirectional connections between the regions which received inputs performed better.

The model comparison was conducted using Bayesian random-effects group analysis (Stephan et al., 2009). To estimate the coupling parameters (Table S3), Bayesian parameter averaging (Kasess et al., 2010) of the winning model was conducted.

3. Results

3.1 Reinforcement learning model

For probabilistic reversal learning tasks, it was suggested to use Rescorla-Wagner models with an anticorrelated valuation system (Gläscher et al., 2009; Hampton et al., 2006; Matsumoto et al., 2007). Using a random-effects Bayesian model selection approach (Stephan et al., 2009), the model with the anticorrelated valuation showed a clearly better fit than the standard Rescorla-Wagner model (exceedance probability $p > .999$).

3.2 Localization of the FRN

Due to major limitations of source estimation models for EEG data, we aimed to localize the source of the FRN using simultaneously acquired EEG and fMRI. EEG-informed fMRI approaches do not rely on prior spatial constraints and profit from the superior spatial resolution of fMRI. In the FRN_{GA} analysis, which focuses on the most consistent FRN signal across subjects, we found one single cluster, located at the dACC (Table 1; Fig. 2D). In the FRN_{indDW} analysis, which gives more freedom to detect individual FRN peaks, we also found a cluster in the dACC (Fig. 2E). Furthermore, additional regions of the salience

network (Menon and Uddin, 2010), such as the anterior insula, the dlPFC and the posterior parietal cortex (PPC) were also activated (Table 1, Fig. 2F).

3.3 The influence of RPEs on the FRN amplitude

We evaluated whether the FRN amplitudes reflect sRPE processing or whether they signal surprise, encoded as |RPEs|. Using a summary statistics approach, we found a significant effect of the sRPEs on the FRN amplitude (FRN_{GA}: $t(14) = 2.1$, $p < .05$, one-tailed; FRN_{indDW}: $t(14) = 4.2$, $p < .001$). However, the analysis using |RPEs| resulted in a significant effect as well (FRN_{GA}: $t(14) = -4.2$, $p < .001$; FRN_{indDW}: $t(14) = -6.1$, $p < .001$). In tasks where subjects learn about values of objects, sRPEs and |RPEs| are usually correlated. Linear regression analyses with only one of the correlated variables as predictor are not able to take this correlation into account. This means that a significant effect for the first variable (e.g. sRPE) could also be driven by the second, correlated variable (e.g. |RPE|). We therefore entered both variables into a multiple regression analysis, where the correlated aspects of both variables end up in the error term (Andrade et al., 1999). In this analysis, the |RPEs| remained significant (FRN_{GA}: $t(14) = -4.3$, $p < .001$; FRN_{indDW}: $t(14) = -4.9$, $p < .001$), whereas the sRPEs did not remain significant (FRN_{GA}: $t(14) = -.8$, $p > .05$; FRN_{indDW}: $t(14) = -.2$, $p > .05$). This finding shows that, in contrast to the |RPEs|, the sRPEs do not uniquely explain variance. A similar result was obtained when we used a model selection approach (see supplementary materials). This relation is also visible in Figure 3, which shows an inverted u-shaped relation between sRPEs and amplitudes. An example for the relation between the single-trial amplitudes and the |RPEs| is shown in Figure S2 for two subjects.

3.4 Effective connectivity of the FRN source

We wanted to examine, from which area the RPE signals are projected to the source of the FRN in the dACC. We selected four regions which were suggested by Holroyd and Coles (2002) to be involved in the FRN network: striatum, vmPFC, amygdala and dlPFC. These areas are known to be crucially involved during RPE processing and they show strong functional and structural connectivities with the dACC. Furthermore, we evaluated whether the origin of the FRN might receives sRPE or |RPE| signals directly,

both of them potentially reflecting inputs from dopaminergic midbrain regions (Bromberg-Martin et al., 2010b). Among the 21 basic models (Table S2, Fig. S5) which we compared, the model which assumes that the |RPEs| are projected directly to the FRN source (model 17) clearly outperformed the other models (exceedance probability $p = .62$; Fig. 4). We then compared the winning model with modified versions of the winning model. None of these outperformed the winning model 17 (Fig. S3). The same result was obtained when we used the FRN-cluster-peak of the FRN_{GA} instead of the FRN_{indDW} analysis (Fig. S4).

4. Discussion

The feedback-related negativity (FRN) is a component of the EEG which occurs 200 to 300 ms after feedback onset over mid-central scalp sites (Miltner et al., 1997; Nieuwenhuis et al., 2004). A decade ago, Holroyd and Coles (2002) suggested that the FRN might reflect a temporal aspect of reward prediction error (RPE) processing and originates from the anterior cingulate cortex. Over the last ten years, a variety of studies tried to localize the FRN (for review cf Walsh and Anderson, 2012). Several studies localized its origin along the anterior cingulate cortex (Bellebaum and Daum, 2008; Gruendler et al., 2011; Hewig et al., 2007; Mathewson et al., 2008; Miltner et al., 1997; Potts et al., 2006; Zhou et al., 2010). Others found the source rather in the posterior cingulate cortex (Badgaiyan and Posner, 1998; Cohen and Ranganath, 2007; Hewig et al., 2007; Müller et al., 2005; Nieuwenhuis et al., 2005b) or even in the basal ganglia (Carlson et al., 2011; Foti et al., 2011a, 2011b; Martin et al., 2009; Nieuwenhuis et al., 2005b). Besides the inconsistent localizations, the methods used in these studies face several limitations when trying to overcome the unsolvable inverse problem (Helmholtz, 1853). Most source localization studies depend on strong prior constraints, such as the number of possible dipoles (Luck, 2005). Furthermore, the spatial resolution of these methods is often quite limited. The simultaneous acquisition of EEG and fMRI has been demonstrated to be an adequate method to localize single-trial EEG-fluctuations using the superior spatial resolution of fMRI without relying on prior constraints (Debener et al., 2006; Huster et al., 2012; Lüchinger et al., 2012, 2011; Rosa et al., 2010). Both approaches to associate the single-trial FRN

amplitudes with the fMRI signal consistently yielded one single cluster located in the dorsal anterior cingulate cortex (dACC). In the FRN_{indDW} analysis, we additionally found areas of the salience network (Menon and Uddin, 2010) activated. Taken together, our findings provide compelling evidence that the FRN truly originates from the dACC, rather than from posterior cingulate or even striatal areas.

The FRN is assumed to reflect sRPE processing in order to select the optimal among several concurring motor programs (Holroyd and Coles, 2002). A variety of studies tried to associate the FRN with sRPEs within the static framework of averaged ERP-measures (e.g., Bellebaum and Daum, 2008; Cohen and Ranganath, 2007; Holroyd et al., 2004; Pfabigan et al., 2011). Averaged ERPs, however, face the difficulty that they are not sensitive to account for trial-by-trial modulation of RPEs and the corresponding dynamic changes in the ERP amplitudes such as those occurring in reinforcement learning tasks. Over the last years, accumulating evidence queries the linear relation between sRPEs and the FRN (Chase et al., 2010; Oliveira et al., 2007; Talmi et al., 2012). It has been suggested that the FRN might reflect a signal of surprise, which can be encoded as absolute RPEs ($|RPE|$; Hayden et al., 2011). In this study, we used a single-trial analysis approach to evaluate the relation between the model-derived RPEs and the FRN amplitude. We found that the FRN is modulated by $|RPEs|$ rather than by sRPEs. In feedback-dependent learning tasks, such as ours, sRPEs are usually correlated with $|RPEs|$. In our study, sRPEs also significantly modulated the FRN amplitude. However, by using multiple regressions as well as a model selection approach, we were able to show that this effect was driven by the correlated $|RPEs|$. This fact could explain why other studies found significant associations between the FRN and sRPEs. Our finding therefore advocates that the FRN reflects surprise rather than sRPEs.

Support for our finding that the FRN reflects surprise rather than sRPEs also comes from recent studies on error processing. Already Holroyd and Coles (2002) assumed that the FRN reflects similar processes as the error-related negativity (ERN; Falkenstein et al., 1990; Gehring et al., 1993). More recently, Gentsch et al. (2009) found that the ERN and the FRN can be explained by very similar topographical components and are therefore most probably caused by a similar neural generator. The same group (Wessel et al., 2012) also found that the ERN topography and neural network closely resembles those of surprise-related

signals. Strikingly, this network also strongly resembles the network which we found in our FRN_{indDW} analysis. Furthermore, the source of the ERN has been localized to the dACC using simultaneous EEG-fMRI (Debener et al., 2005). Taken together, the findings strongly support the suggestion of a common surprise-network including the dACC.

The dACC is well known to be a crucial area within the reward network (Beckmann et al., 2009; Haber and Knutson, 2010; Rushworth et al., 2011, 2004). Although it is known to process RPEs (Matsumoto et al., 2007), its exact role is not yet clear. Besides its role in tracing uncertainty (Behrens et al., 2007), the dACC is assumed to play a crucial role in response evaluation and behavioral adaptation (Williams et al., 2004). This assumption is supported by the finding that the dACC has a highly predictive value for behavioral adaptation processes (Hampton and O'Doherty, 2007). Because the dACC is densely connected to motor regions (Beckmann et al., 2009; Paus, 2001) as well as to areas which process RPEs, such as the dIPFC (Bates and Goldman-Rakic, 1993; Pandya et al., 1981), vmPFC (Morecraft and Van Hoesen, 1998), the amygdalae (Beckmann et al., 2009; Morecraft et al., 2007) or the striatum (Haber et al., 2006; Kunishio and Haber, 1994), several authors suggested that the dACC receives inputs from these areas to compare behavioral options and choose the optimal among them (Hare et al., 2011; Holroyd and Coles, 2002). However, others suggested that the dACC acts more as an independent module receiving reward and RPE-signals directly from its dense dopaminergic innervations (Lindvall et al., 1974; Oades and Halliday, 1987) rather than receiving these signals from other areas of the sRPE network (Alexander and Brown, 2010; Silvetti et al., 2012, 2011). Whether these dopaminergic inputs reflect sRPEs, |RPEs|, or even a mixture of both, remains unknown (Bromberg-Martin et al., 2010b). Using dynamic causal modeling (Friston et al., 2003), we explored whether the source of the FRN receives direct inputs of sRPEs, |RPE|, or both. We also evaluated whether these signals are projected to the dACC directly or via another region, as Holroyd and Coles (2002) suggested. As preceding regions, we considered the striatum, vmPFC, amygdalae and the dIPFC due to their dense functional and structural connections (Beckmann et al., 2009; Paus, 2001). In our DCM analysis, we found that a model which showed direct |RPE|-inputs to the FRN source outperformed all other models which assumed that input signals (sRPEs or |RPEs|) are

first transformed in another region before being forwarded to the FRN. This finding contradicts the assumption made by Holroyd and Coles (2002), which assumed that the dACC receives the sRPEs via the striatum. However, our findings are in line with the recently stated theory that dopaminergic neurons also encode |RPE|-like salience signals (Bromberg-Martin et al., 2010a, 2010b; Matsumoto and Hikosaka, 2009) which are most probably also projected to the anterior cingulate cortex (Bromberg-Martin et al., 2010b).

Although our findings clearly increase knowledge about the FRN and further the debate about the localization, the relation to RPEs and the functional network, we see some limitations of our findings. First, the FRN is defined very differently across studies and therefore its characterization is not always undisputed. In the first study on the FRN by Miltner et al. (1997), the authors found a negative displacement after about 250ms for the negative feedback. To better characterize the divergence, the authors computed the difference wave between the positive and negative feedbacks. In later studies this difference wave-approach has very often been pursued (e.g., Baker and Holroyd, 2011; Cohen and Ranganath, 2007; Holroyd and Coles, 2002; Nieuwenhuis et al., 2005b; Walsh and Anderson, 2011). However, other studies also tried to characterize the FRN by computing the amplitude difference between two subsequent deflections in the ERPs of a single feedback condition (cf Yeung and Sanfey, 2004). When looking across FRN-studies, the peak of the FRN deflection is caused by different aspects of underlying ERPs for rewards and punishments. For example, while Gehring and Willoughby (2002) found negative going deflections for both feedbacks to compute the difference, Miltner et al. (1997) did not find a negative or negative going deflection for rewards. In our study, we defined the FRN peak, which was then used for single-trial analysis, as the most negative deflection of the difference wave at the central electrode Cz between 200 and 300ms after feedback, similar to Miltner et al. (1997) and Holroyd and Coles (2002). In the FRN_{GA} analysis, the peak of our difference wave was mainly caused by a decreased positive deflection for punishments rather than an increased negative going peak as reported by Gehring and Willoughby (2002). Given that a large number of studies defined rather long intervals to analyze the FRN (e.g., Baker and Holroyd, 2011; Holroyd and Coles, 2002; Miltner et al., 1997; Nieuwenhuis et al.,

2005a, 2005b), our analyses (mainly the FRN_{GA} analysis) may not reflect the FRN as a whole, but rather an early peak of a more extended FRN interval. Such an early peak in the FRN, however, has been found in other studies as well (e.g., Baker and Holroyd, 2011; Gehring and Willoughby, 2002; Holroyd et al., 2004; Miltner et al., 1997; Oliveira et al., 2007; Walsh and Anderson, 2011). The fact that our analysis captured mainly the early part of the FRN might also explain its slightly more posterior negative maximum in the FRN_{GA} analysis. Second, in this study, we used a single-trial approach for our analyses. It is well known that single-trial EEG analyses suffer from a low signal-to-noise ratio. We therefore chose a rather liberal voxel-level threshold for our fMRI analyses. With choosing this threshold, we wanted to ascertain that we do not oversee FRN sources due to a too conservative threshold. Still, we corrected for multiple comparisons using an adjusted cluster-extent threshold. Third, we restricted our single-trial analysis only to the FRN-ERP. One could also have used a more exploratory approach and evaluated the relation at every time point after feedback (cf Philiastides et al., 2010). However, our goal was to determine the relation between the RPEs and the FRN and we therefore strictly limited our analysis to the FRN deflection. Last, in our DCM analysis, we used the model-derived sRPEs and |RPEs| as inputs assuming that these signals might reflect dopaminergic activity of the mesencephalon. We decided to do so because fMRI usually struggles with acquiring adequate signals from midbrain regions. Furthermore, it is assumed that dopaminergic neurons encode either sRPE or |RPE| signals (Bromberg-Martin et al., 2010b). However, we cannot be sure that our inputs really reflect dopaminergic activity. Furthermore, due to the intractable full model space of our DCM analysis, we decided to perform a stepwise exploratory search. With this approach we do not evaluate highly complex models. Despite these limitations, our multimodal approach allowed us to directly associate the FRN single-trial amplitude with RPE surprise signals originating in the dACC which seems to directly receive |RPEs| without being processed in other regions before.

5. Conclusions

The aims of our simultaneous EEG-fMRI study were (1) to localize the origin of the FRN overcoming the limitations of previous source localization studies; (2) to determine the relation between the FRN amplitude, sRPEs and surprise signals (encoded as |RPEs|); (3) to explore the reward network to determine the signaling pathway of RPEs. We consistently localized the FRN in the dACC, using the superior spatial resolution of fMRI without relying on operator dependent priors, such as spatial constraints on the source distribution. We were also able to disentangle the effect of signed and absolute RPEs on the FRN amplitude. Our findings clearly demonstrate that the FRN reflects a surprise rather than a sRPE signal. Our analysis also showed that the correlation between these signals might also explain why others found a linear relation between sRPEs and the FRN. Using effective connectivity measures, we found that the FRN source region most probably receives direct |RPE| inputs rather than signals which are first processed in other areas of the RPE-network. These findings make the FRN an ideal temporal component to study surprise-based behavior evaluation in the dACC. Such knowledge may not only advance the field of cognitive neuroscience, but also aid understanding of the disturbed dACC-related adaptation and decision making processes in various neuropsychiatric disorders such as addiction, obesity, obsessive-compulsive disorder or attention deficit/hyperactivity disorder (Cubillo et al., 2012; Kishida et al., 2010; Rangel et al., 2008; Schultz, 2011; Sharp et al., 2012).

Acknowledgments

We thank M. Piccirelli and J. Kronschnabel for their helpful inputs on MR-acquisition and analysis. We thank the reviewers for proposing new, complementary and fruitful analyses (e.g., FRN_{indDW}). This study was supported by the Swiss National Science Foundation (No. 320030_130237) and the Hartmann Müller Foundation (No. 1460).

Conflict of interest

S. Walitza received speakers honoraria from Eli Lilly, Janssen-Cilag and AstraZeneca in the last five years. The other authors declare no competing financial interests.

References

- Alexander, W.H., Brown, J.W., 2010. Competition between learned reward and error outcome predictions in anterior cingulate cortex. *NeuroImage* 49, 3210–3218.
- Alexander, W.H., Brown, J.W., 2011. Medial prefrontal cortex as an action-outcome predictor. *Nat Neurosci* 14, 1338–1344.
- Allen, P.J., Josephs, O., Turner, R., 2000. A method for removing imaging artifact from continuous EEG recorded during functional MRI. *NeuroImage* 12, 230–9.
- Andrade, A., Paradis, A.L., Rouquette, S., Poline, J.B., 1999. Ambiguous results in functional neuroimaging data analysis due to covariate correlation. *NeuroImage* 10, 483–486.
- Ashburner, J., 2007. A fast diffeomorphic image registration algorithm. *NeuroImage* 38, 95–113.
- Badgaiyan, R.D., Posner, M.I., 1998. Mapping the cingulate cortex in response selection and monitoring. *NeuroImage* 7, 255–260.
- Baker, T.E., Holroyd, C.B., 2011. Dissociated roles of the anterior cingulate cortex in reward and conflict processing as revealed by the feedback error-related negativity and N200. *Biol Psychol* 87, 25–34.
- Bates, J.F., Goldman-Rakic, P.S., 1993. Prefrontal connections of medial motor areas in the rhesus monkey. *J Comp Neurol* 336, 211–228.
- Beckmann, M., Johansen-Berg, H., Rushworth, M.F.S., 2009. Connectivity-based parcellation of human cingulate cortex and its relation to functional specialization. *J. Neurosci.* 29, 1175–1190.
- Behrens, T.E.J., Woolrich, M.W., Walton, M.E., Rushworth, M.F.S., 2007. Learning the value of information in an uncertain world. *Nat Neurosci* 10, 1214–1221.
- Bellebaum, C., Daum, I., 2008. Learning-related changes in reward expectancy are reflected in the feedback-related negativity. *Eur J Neurosci* 27, 1823–1835.
- Bromberg-Martin, E.S., Delazer, M., Hikosaka, O., 2010a. Distinct tonic and phasic anticipatory activity in lateral habenula and dopamine neurons. *Neuron* 67, 144–155.
- Bromberg-Martin, E.S., Matsumoto, M., Hikosaka, O., 2010b. Dopamine in motivational control: rewarding, aversive, and alerting. *Neuron* 68, 815–834.
- Carlson, J.M., Foti, D., Mujica-Parodi, L.R., Harmon-Jones, E., Hajcak, G., 2011. Ventral striatal and medial prefrontal BOLD activation is correlated with reward-related electrocortical activity: a combined ERP and fMRI study. *NeuroImage* 57, 1608–1616.
- Chase, H.W., Swainson, R., Durham, L., Benham, L., Cools, R., 2010. Feedback-related negativity codes prediction error but not behavioral adjustment during probabilistic reversal learning. *J Cogn Neurosci* 23, 936–946.
- Cohen, M.X., Ranganath, C., 2007. Reinforcement learning signals predict future decisions. *J Neurosci* 27, 371–378.
- Courville, A.C., Daw, N.D., Touretzky, D.S., 2006. Bayesian theories of conditioning in a changing world. *Trends in Cognitive Sciences* 10, 294–300.
- Cubillo, A., Halari, R., Smith, A., Taylor, E., Rubia, K., 2012. A review of fronto-striatal and fronto-cortical brain abnormalities in children and adults with Attention Deficit Hyperactivity Disorder (ADHD) and new evidence for dysfunction in adults with ADHD during motivation and attention. *Cortex* 48, 194–215.
- Debener, S., Ullsperger, M., Siegel, M., Engel, A.K., 2006. Single-trial EEG-fMRI reveals the dynamics of cognitive function. *Trends Cogn Sci* 10, 558–563.
- Debener, S., Ullsperger, M., Siegel, M., Fiehler, K., von Cramon, D.Y., Engel, A.K., 2005. Trial-by-trial coupling of concurrent electroencephalogram and functional magnetic resonance imaging identifies the dynamics of performance monitoring. *J Neurosci* 25, 11730–11737.
- Delorme, A., Makeig, S., 2004. EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *J Neurosci Methods* 134, 9–21.
- Falkenstein, M., Hohnsbein, J., Hoormann, J., Blanke, L., 1990. Effects of errors in choice reaction tasks on the ERP under focused and divided attention., in: Brunia, C., Gaillard, A., Kok, A. (Eds.), *Psychophysiological Brain Research*. Tilburg University Press, Tilburg, the Netherlands, pp. 192–195.

- Forman, S.D., Cohen, J.D., Fitzgerald, M., Eddy, W.F., Mintun, M.A., Noll, D.C., 1995. Improved assessment of significant activation in functional magnetic resonance imaging (fMRI): Use of a cluster-size threshold. *Magn Reson Med* 33, 636–647.
- Foti, D., Weinberg, A., Dien, J., Hajcak, G., 2011a. Event-related potential activity in the basal ganglia differentiates rewards from nonrewards: Temporospacial principal components analysis and source localization of the feedback negativity. *Hum Brain Mapp* 32, 2207–2216.
- Foti, D., Weinberg, A., Dien, J., Hajcak, G., 2011b. Event-related potential activity in the basal ganglia differentiates rewards from nonrewards: Response to commentary. *Hum Brain Mapp* 32, 2267–2269.
- Friston, K.J., 2010. The free-energy principle: a unified brain theory? *Nat Rev Neurosci* 11, 127–138.
- Friston, K.J., Harrison, L., Penny, W., 2003. Dynamic causal modelling. *NeuroImage* 19, 1273–1302.
- Gehring, W.J., Goss, B., Coles, M.G.H., Meyer, D.E., Donchin, E., 1993. A neural system for error detection and compensation. *Psychological Science* 4, 385–390.
- Gehring, W.J., Willoughby, A.R., 2002. The medial frontal cortex and the rapid processing of monetary gains and losses. *Science* 295, 2279–2282.
- Gentsch, A., Ullsperger, P., Ullsperger, M., 2009. Dissociable medial frontal negativities from a common monitoring system for self- and externally caused failure of goal achievement. *Neuroimage* 47, 2023–2030.
- Gläscher, J., Daw, N., Dayan, P., O’Doherty, J.P., 2010. States versus rewards: Dissociable neural prediction error signals underlying model-based and model-free reinforcement learning. *Neuron* 66, 585–595.
- Gläscher, J., Hampton, A.N., O’Doherty, J.P., 2009. Determining a role for ventromedial prefrontal cortex in encoding action-based value signals during reward-related decision making. *Cereb Cortex* 19, 483–495.
- Glimcher, P.W., 2011. Understanding dopamine and reinforcement learning: the dopamine reward prediction error hypothesis. *Proc Natl Acad Sci U.S.A.* 108, 15647–15654.
- Gruendler, T.O.J., Ullsperger, M., Huster, R.J., 2011. Event-related potential correlates of performance-monitoring in a lateralized time-estimation task. *PLoS ONE* 6, e25591.
- Haber, S.N., Kim, K.-S., Maily, P., Calzavara, R., 2006. Reward-related cortical inputs define a large striatal region in primates that interface with associative cortical connections, providing a substrate for incentive-based learning. *J Neurosci* 26, 8368–8376.
- Haber, S.N., Knutson, B., 2010. The reward circuit: linking primate anatomy and human imaging. *Neuropsychopharmacology* 35, 4–26.
- Hampton, A.N., Bossaerts, P., O’Doherty, J.P., 2006. The role of the ventromedial prefrontal cortex in abstract state-based inference during decision making in humans. *J Neurosci* 26, 8360–8367.
- Hampton, A.N., O’Doherty, J.P., 2007. Decoding the neural substrates of reward-related decision making with functional MRI. *Proc Natl Acad Sci U.S.A.* 104, 1377–1382.
- Hare, T.A., O’Doherty, J.P., Camerer, C.F., Schultz, W., Rangel, A., 2008. Dissociating the role of the orbitofrontal cortex and the striatum in the computation of goal values and prediction errors. *J Neurosci* 28, 5623–5630.
- 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. *Proc. Natl. Acad. Sci. U.S.A.* 108, 18120–18125.
- 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. *J. Neurosci.* 31, 4178–4187.
- Helmholtz, H., 1853. Ueber einige Gesetze der Vertheilung elektrischer Ströme in körperlichen Leitern, mit Anwendung auf die thierisch-elektrischen Versuche. *Annalen der Physikalischen Chemie* 89, 211 – 233.
- Hewig, J., Trippe, R., Hecht, H., Coles, M.G.H., Holroyd, C.B., Miltner, W.H.R., 2007. Decision-making in Blackjack: an electrophysiological analysis. *Cereb. Cortex* 17, 865–877.

- Holmes, A.P., Friston, K.J., 1998. Generalisability, random effects and population inference. *NeuroImage* 7, S754.
- Holroyd, C.B., Coles, M.G.H., 2002. The neural basis of human error processing: reinforcement learning, dopamine, and the error-related negativity. *Psychol Rev* 109, 679–709.
- Holroyd, C.B., Larsen, J.T., Cohen, J.D., 2004. Context dependence of the event-related brain potential associated with reward and punishment. *Psychophysiology* 41, 245–253.
- Huster, R.J., Debener, S., Eichele, T., Herrmann, C.S., 2012. Methods for simultaneous EEG-fMRI: an introductory review. *J. Neurosci.* 32, 6053–6060.
- Jung, T.-P., Makeig, S., Westerfield, M., Townsend, J., Courchesne, E., Sejnowski, T.J., 2000. Removal of eye activity artifacts from visual event-related potentials in normal and clinical subjects. *Clin Neurophysiol* 111, 1745–1758.
- Kasess, C.H., Stephan, K.E., Weissenbacher, A., Pezawas, L., Moser, E., Windischberger, C., 2010. Multi-subject analyses with dynamic causal modeling. *NeuroImage* 49, 3065–3074.
- Kishida, K.T., King-Casas, B., Montague, P.R., 2010. Neuroeconomic approaches to mental disorders. *Neuron* 67, 543–554.
- Kunishio, K., Haber, S.N., 1994. Primate cingulo-striatal projection: limbic striatal versus sensorimotor striatal input. *J Comp Neurol* 350, 337–356.
- Lehmann, D., Skrandies, W., 1980. Reference-free identification of components of checkerboard-evoked multichannel potential fields. *Electroencephalogr Clin Neurophysiol* 48, 609–621.
- Lindvall, O., Björklund, A., Moore, R.Y., Stenevi, U., 1974. Mesencephalic dopamine neurons projecting to neocortex. *Brain Res* 81, 325–331.
- Lüchinger, R., Michels, L., Martin, E., Brandeis, D., 2011. EEG-BOLD correlations during (post-)adolescent brain maturation. *NeuroImage* 56, 1493–1505.
- Lüchinger, R., Michels, L., Martin, E., Brandeis, D., 2012. Brain state regulation during normal development: Intrinsic activity fluctuations in simultaneous EEG-fMRI. *NeuroImage* 60, 1426–1439.
- Luck, S.J., 2005. *An Introduction to the event-related potential technique*. MIT Press, Cambridge, MA.
- Martin, L.E., Potts, G.F., Burton, P.C., Montague, P.R., 2009. Electrophysiological and hemodynamic responses to reward prediction violation. *Neuroreport* 20, 1140–1143.
- Mathewson, K.J., Dywan, J., Snyder, P.J., Tays, W.J., Segalowitz, S.J., 2008. Aging and electrocortical response to error feedback during a spatial learning task. *Psychophysiology* 45, 936–948.
- Matsumoto, M., Hikosaka, O., 2009. Two types of dopamine neuron distinctly convey positive and negative motivational signals. *Nature* 459, 837–841.
- Matsumoto, M., Matsumoto, K., Abe, H., Tanaka, K., 2007. Medial prefrontal cell activity signaling prediction errors of action values. *Nat. Neurosci.* 10, 647–656.
- Menon, V., Uddin, L.Q., 2010. Saliency, switching, attention and control: a network model of insula function. *Brain Struct Funct* 214, 655–667.
- Meyer-Lindenberg, A., 2010. From maps to mechanisms through neuroimaging of schizophrenia. *Nature* 468, 194–202.
- Miltner, W.H.R., Braun, C.H., Coles, M.G.H., 1997. Event-related brain potentials following incorrect feedback in a time-estimation task: Evidence for a “generic” neural system for error detection. *J Cogn Neurosci* 9, 788–798.
- Morecraft, R.J., McNeal, D.W., Stilwell-Morecraft, K.S., Gedney, M., Ge, J., Schroeder, C.M., Van Hoesen, G.W., 2007. Amygdala interconnections with the cingulate motor cortex in the rhesus monkey. *J Comp Neurol* 500, 134–165.
- Morecraft, R.J., Van Hoesen, G.W., 1998. Convergence of limbic input to the cingulate motor cortex in the rhesus monkey. *Brain Res Bull* 45, 209–232.
- Müller, S.V., Möller, J., Rodriguez-Fornells, A., Münte, T.F., 2005. Brain potentials related to self-generated and external information used for performance monitoring. *Clin Neurophysiol* 116, 63–74.

- Nieuwenhuis, S., Holroyd, C.B., Mol, N., Coles, M.G.H., 2004. Reinforcement-related brain potentials from medial frontal cortex: origins and functional significance. *Neurosci Biobehav Rev* 28, 441–448.
- Nieuwenhuis, S., Nielen, M.M., Mol, N., Hajcak, G., Veltman, D.J., 2005a. Performance monitoring in obsessive-compulsive disorder. *Psychiatry Res* 134, 111–122.
- Nieuwenhuis, S., Slagter, H.A., Geusau, V., Alting, N.J., Heslenfeld, D.J., Holroyd, C.B., 2005b. Knowing good from bad: differential activation of human cortical areas by positive and negative outcomes. *Eur J Neurosci* 21, 3161–3168.
- O’Doherty, J.P., Kringelbach, M.L., Rolls, E.T., Hornak, J., Andrews, C., 2001. Abstract reward and punishment representations in the human orbitofrontal cortex. *Nat Neurosci* 4, 95–102.
- Oades, R.D., Halliday, G.M., 1987. Ventral tegmental (A10) system: neurobiology. 1. Anatomy and connectivity. *Brain Res.* 434, 117–165.
- Oliveira, F.T.P., McDonald, J.J., Goodman, D., 2007. Performance monitoring in the anterior cingulate is not all error related: expectancy deviation and the representation of action-outcome associations. *J Cogn Neurosci* 19, 1994–2004.
- Pandya, D.N., Van Hoesen, G.W., Mesulam, M.-M., 1981. Efferent connections of the cingulate gyrus in the rhesus monkey. *Exp Brain Res* 42, 319–330.
- Pascual-Marqui, R.D., 2002. Standardized low-resolution brain electromagnetic tomography (sLORETA): technical details. *Methods Find Exp Clin Pharmacol* 24 Suppl D, 5–12.
- Paus, T., 2001. Primate anterior cingulate cortex: Where motor control, drive and cognition interface. *Nat Rev Neurosci* 2, 417–424.
- Pearce, J.M., Hall, G., 1980. A model for Pavlovian learning: variations in the effectiveness of conditioned but not of unconditioned stimuli. *Psychol Rev* 87, 532–552.
- Pfabigan, D.M., Alexopoulos, J., Bauer, H., Sailer, U., 2011. Manipulation of feedback expectancy and valence induces negative and positive reward prediction error signals manifest in event-related brain potentials. *Psychophysiology* 48, 656–664.
- Philiastides, M.G., Biele, G., Vavatzanidis, N., Kazzer, P., Heekeren, H.R., 2010. Temporal dynamics of prediction error processing during reward-based decision making. *NeuroImage* 53, 221–232.
- Potts, G.F., Martin, L.E., Burton, P., Montague, P.R., 2006. When things are better or worse than expected: the medial frontal cortex and the allocation of processing resources. *J Cogn Neurosci* 18, 1112–1119.
- Rangel, A., Camerer, C., Montague, P.R., 2008. A framework for studying the neurobiology of value-based decision making. *Nat Rev Neurosci* 9, 545–556.
- Remijnse, P.L., Nielen, M.M.A., Uylings, H.B.M., Veltman, D.J., 2005. Neural correlates of a reversal learning task with an affectively neutral baseline: an event-related fMRI study. *NeuroImage* 26, 609–618.
- Rescorla, R.A., Wagner, A.R., 1972. A theory of Pavlovian conditioning: variations in the effectiveness of reinforcement and nonreinforcement, in: Black, A.H., Prokasy, W.F. (Eds.), *Classical Conditioning II: Current Research and Theory*. Appleton-Century-Crofts, New York, pp. 64–99.
- Rosa, M.J., Daunizeau, J., Friston, K.J., 2010. EEG-fMRI integration: a critical review of biophysical modeling and data analysis approaches. *J Integr Neurosci* 9, 453–476.
- Rushworth, M.F.S., Noonan, M.P., Boorman, E.D., Walton, M.E., Behrens, T.E.J., 2011. Frontal Cortex and Reward-Guided Learning and Decision-Making. *Neuron* 70, 1054–1069.
- Rushworth, M.F.S., Walton, M.E., Kennerley, S.W., Bannerman, D.M., 2004. Action sets and decisions in the medial frontal cortex. *Trends Cogn Sci* 8, 410–417.
- Rutledge, R.B., Dean, M., Caplin, A., Glimcher, P.W., 2010. Testing the reward prediction error hypothesis with an axiomatic model. *J Neurosci* 30, 13525–13536.
- Schultz, W., 2011. Potential vulnerabilities of neuronal reward, risk, and decision mechanisms to addictive drugs. *Neuron* 69, 603–617.
- Schultz, W., Dayan, P., Montague, P.R., 1997. A neural substrate of prediction and reward. *Science* 275, 1593–1599.

- Sharp, C., Monterosso, J., Montague, P.R., 2012. Neuroeconomics: a bridge for translational research. *Biol Psychiatry* 72, 87–92.
- Silvetti, M., Seurinck, R., Verguts, T., 2011. Value and prediction error in medial frontal cortex: integrating the single-unit and systems levels of analysis. *Front Hum Neurosci* 5, 75.
- Silvetti, M., Seurinck, R., Verguts, T., 2012. Value and prediction error estimation account for volatility effects in ACC: A model-based fMRI study. *Cortex*.
- Slotnick, S.D., Moo, L.R., Segal, J.B., Hart Jr., J., 2003. Distinct prefrontal cortex activity associated with item memory and source memory for visual shapes. *Cogn Brain Res* 17, 75–82.
- Stephan, K.E., Penny, W.D., Daunizeau, J., Moran, R.J., Friston, K.J., 2009. Bayesian model selection for group studies. *NeuroImage* 46, 1004–1017.
- Sutton, R.S., Barto, A.G., 1998. Reinforcement learning: An introduction. MIT Press.
- Talmi, D., Fuentemilla, L., Litvak, V., Duzel, E., Dolan, R.J., 2012. An MEG signature corresponding to an axiomatic model of reward prediction error. *NeuroImage* 59, 635–645.
- Walsh, M.M., Anderson, J.R., 2011. Modulation of the feedback-related negativity by instruction and experience. *PNAS* 108, 19048–19053.
- Walsh, M.M., Anderson, J.R., 2012. Learning from experience: event-related potential correlates of reward processing, neural adaptation, and behavioral choice. *Neurosci Biobehav Rev* 36, 1870–1884.
- Wessel, J.R., Danielmeier, C., Morton, J.B., Ullsperger, M., 2012. Surprise and error: common neuronal architecture for the processing of errors and novelty. *J. Neurosci.* 32, 7528–7537.
- Williams, Z.M., Bush, G., Rauch, S.L., Cosgrove, G.R., Eskandar, E.N., 2004. Human anterior cingulate neurons and the integration of monetary reward with motor responses. *Nat Neurosci* 7, 1370–1375.
- Yeung, N., Sanfey, A.G., 2004. Independent coding of reward magnitude and valence in the human brain. *J. Neurosci.* 24, 6258–6264.
- Zhou, Z., Yu, R., Zhou, X., 2010. To do or not to do? Action enlarges the FRN and P300 effects in outcome evaluation. *Neuropsychologia* 48, 3606–3613.

Figure Legends

Fig. 1. Probabilistic reversal learning task. The participants had to learn which of the two presented stimuli is associated with a higher reward probability. One of the stimuli was assigned with a reward probability of 80% whereas the other had a reward probability of 20%. After several trials, the reward probabilities were reversed and the participants had to detect and adjust to these changes based on the feedback they received.

Fig. 2. Localization of the feedback-related negativity (FRN). (A) Grand average waveform of punishment (red) – reward (green) trials, revealing the FRN with its maximal amplitude after 223ms (bold red line). (B) The statistical map illustrates the potential field distribution of the FRN at the peak of the FRN in the FRN_{GA} analysis (B) and at the peaks of the individual difference waves between 200 and 300 ms in the FRN_{indDW} analysis (C). (D) The FRN_{GA} analysis revealed one single cluster within the dACC. The FRN_{indDW} analysis also showed a significant cluster in the dACC (E). Additionally, other areas of the salience network were also found to be active (F).

Fig. 3. FRN amplitudes are driven by $|RPEs|$ rather than by $sRPEs$. The illustration of the relation between $sRPEs$ and the single-trial amplitudes indicates the effect of $|RPEs|$ (± 1 SEM).

Fig. 4. Effective connectivity of the RPE-signaling pathway. (A) The exceedance probability of model 17 clearly outperformed all other models. Please note: exceedance probabilities are relative values which sum up to 1 over the complete model space. (B) The winning model of the DCM analysis reveals that $|RPEs|$ are directly projected to the source region of the FRN.