

**The role of reversal learning impairment in social disinhibition following severe traumatic brain injury**

Katherine Osborne-Crowley<sup>1</sup>, Skye McDonald<sup>1</sup>, Jacqueline A. Rushby<sup>1</sup>

<sup>1</sup>School of Psychology, The University of New South Wales

Word Count: 5357 (excl. title page, abstract and references)

Abstract word count: 215

1<sup>st</sup> Author (corresponding)

Katherine Osborne-Crowley

School of Psychology,

The University of New South Wales, New South Wales, 2052, Australia

Email: k.osbornecrowley@unsw.edu.au

Phone: +61 2 9385 3590

2<sup>nd</sup> Author

Skye McDonald

School of Psychology,

The University of New South Wales, New South Wales, 2052, Australia

Email: s.mcdonald@unsw.edu.au

3<sup>rd</sup> Author

Jacqueline A. Rushby

School of Psychology,

The University of New South Wales, New South Wales, 2052, Australia

Email: j.rushby@unsw.edu.au

Osborne-Crowley, K., McDonald, S., & Rushby, J. A. (2016). Role of reversal learning impairment in social disinhibition following severe traumatic brain injury. *Journal of the International Neuropsychological Society*, 22(03), 303-313, DOI:10.1017/S1355617715001277.

## Abstract

Objective: The current study aimed to determine whether reversal learning impairments and feedback-related negativity (FRN), reflecting reward prediction error signals generated by negative feedback during the reversal learning tasks, were associated with social disinhibition in a group of participants with traumatic brain injury (TBI).

Method: Number of reversal errors on a social and a non-social reversal learning task and FRN were examined for 21 participants with TBI and 21 control participants matched for age. Participants with TBI were also divided into low and high disinhibition groups based on rated videotaped interviews.

Results: Participants with TBI made more reversal errors and produced smaller amplitude FRN's than controls. Further, participants with TBI high on social disinhibition made more reversal errors on the social reversal learning task than did those low on social disinhibition. FRN amplitude was not related to disinhibition.

Conclusions: These results suggest that impairment in the ability to update behaviour when social reinforcement contingencies change plays a role in social disinhibition after TBI. Further, the social reversal learning task used in this study may be a useful neuropsychological tool for detecting susceptibility to acquired social disinhibition following TBI. Finally, that the FRN amplitude was not associated with social disinhibition suggests that reward prediction error signals are not critical for behavioural adaptation in the social domain.

**Keywords:** brain injuries, social disinhibition, orbitofrontal cortex (OFC), reversal learning, social reinforcement, feedback-related negativity (FRN), reward prediction error

Severe traumatic brain injury (TBI) results in significant neuropsychological and psychosocial sequelae with devastating consequences both for the individual and for their family (Tate, Broe, & Lulham, 1989). However, it is the disruption to social after TBI that is often reported as being the most disabling and distressing for family and for the community (Brooks & McKinlay, 1983; McKinlay, Brooks, Bond, Martinage, & Marshall, 1981). A particularly debilitating behaviour change commonly reported after TBI is social disinhibition, which refers to “socially inappropriate verbal, physical or sexual acts which reflect a loss of inhibition or an inability to conform to social or cultural behavioural norms” (Arciniegas & Wortzel, 2014, p. 39). This inappropriate social behaviour may contribute to the well-documented trouble people with TBI have in maintaining social relationships post-injury, leading to social isolation and psychiatric illness such as depression and anxiety (Gould, Ponsford, Johnston, & Schonberger, 2011).

Socially disinhibited behaviour after TBI has been linked with damage to the orbitofrontal cortex (OFC) and its connections with other brain regions (Lipszyc et al., 2014; Namiki et al., 2008). Further, evidence from lesions studies in both humans (Barrash, Tranel, & Anderson, 2000; Blair & Cipolotti, 2000; Namiki et al., 2008) and monkeys (Butter, Mishkin, & Mirsky, 1968; Franzen & Myers, 1973; Machado & Bachevalier, 2006), as well as studies of neurodegenerative disease (Hornberger, Geng, & Hodges, 2011; Krueger et al., 2011), also consistently demonstrate an association between OFC damage and social disinhibition. The orbitofrontal region is particularly susceptible following TBI (Mattson & Levin, 1990) due to abrasion of the ventral surfaces of the frontal lobes as they scrape across the bony floor of the anterior fossa in response to the acceleration-deceleration forces associated with the trauma (Bigler, 2007). Damage to frontal white matter tracts, which connect the orbitofrontal region

with other brain regions has also been shown to be a common outcome of TBI (Kinnunen et al., 2011). Despite a general consensus in the literature that damage to the OFC mediates acquired social disinhibition, it is unknown what specific mechanism is involved.

Reversal learning impairment, or an impaired ability to update responding when reward contingencies change, is a neuropsychological hallmark of OFC damage (Schoenbaum, Takahashi, Liu, & McDannald, 2011). This well-documented deficit has generally been demonstrated using a visual discrimination test of reversal learning which involves the subject learning, based on reward and punishment, to respond to one of two visual stimuli presented, until, when a criterion level performance is reached, the reinforcement contingency is swapped without warning. Human subjects with damage to the OFC, but not those with damage outside the OFC, have been found to exhibit deficient performance on such tasks (Fellows & Farah, 2003; Hornak et al., 2004). Further, patients with frontal variant fronto-temporal dementia (fv-FTD), characterised by neurodegeneration which preferentially affects the OFC (Gregory, Serra-Mestres, & Hodges, 1999), similarly demonstrate an impairment in reversal learning (Rahman, Sahakian, Hodges, Rogers, & Robbins, 1999). Finally, people with TBI have also been found to perform poorly on reversal learning tasks (Rolls, Hornak, Wade, & McGrath, 1994). This impairment in the ability to flexibly adapt responding in an environment of changing social reinforcement contingencies may underlie acquired social disinhibition (Bachevalier & Loveland, 2006). While reversal learning impairment has been documented in people with TBI and other clinical groups with OFC damage, no studies have yet demonstrated an impairment of reversal of social reinforcement contingencies after TBI. Thus, the first aim of the current study was to determine whether participants with TBI are impaired on a social reversal learning task and whether this impairment is related to social disinhibition.

Although it is clear that the OFC is crucial for reversal learning, the precise role it plays has been the subject of debate. Schoenbaum, Roesch, Stalnaker, and Takahashi (2009) argued that the role of the orbitofrontal cortex in reversal learning behaviour is its contribution to the generation of reward prediction error signals which indicate the need for behavioural change when an outcome is worse than expected (Walsh & Anderson, 2011a). Specifically, Schoenbaum et al. (2009) suggests that the OFC provides important information about the value of the expected outcome which is used in the generation of these reward prediction error signals in the dopaminergic midbrain. Evidence from neural recording studies (Gottfried, O'Doherty, & Dolan, 2003; Hikosaka & Watanabe, 2004; Padoa-Schioppa & Assad, 2006) and behavioural studies (Izquierdo, Suda, & Murray, 2004) in animals support the role of the OFC in signalling expected outcomes. Crucially, in a reversal learning task reward prediction errors are necessary to signal the need to update behaviour when negative feedback is delivered. Thus, the current study focused also on the role of reward prediction error signals in reversal learning and socially disinhibited behaviour.

In humans, feedback-related negativity (FRN), an event related potential (ERP) component of the electroencephalogram (EEG) occurring approximately 200 to 400 ms after feedback onset, is thought to reflect reward prediction error signals (Nieuwenhuis, Holroyd, Mol, & Coles, 2004). The FRN originates at the ACC, where it is hypothesised that the reward prediction error signals are used to update behaviour such as is required in reversal learning tasks. The FRN is theorised to reflect the influence of midbrain dopaminergic reward prediction error signals on the ACC, such that a more negative FRN reflects a negative reward prediction error and a more positive FRN reflects a positive reward prediction error (Holroyd & Coles, 2002). This is evidenced by the finding that FRN amplitudes are most negative following

unpredicted non-reward and least negative following unpredicted reward, and only occur when error feedback is not expected or probable (Hajcak, Moser, Holroyd, & Simons, 2007; Holroyd & Coles, 2002; Holroyd, Krigolson, Baker, Lee, & Gibson, 2009; Holroyd, Nieuwenhuis, Yeung, & Cohen, 2003; Walsh & Anderson, 2011a, 2011b). Studies demonstrating that FRN can predict behavioural change (Cohen & Ranganath, 2007; Holroyd & Krigolson, 2007; van der Helden, Boksem, & Blom, 2010) supports the assumption that the FRN reflects the dopaminergic signalling of reward prediction errors which guide behavioural adaptation when an outcome is worse than expected. If the role of the OFC in reversal learning is its contribution to the generation of reward prediction error signals as Schoenbaum et al. (2009) suggests, it would be expected that an impaired ability to generate FRN signals to social feedback would be related to social disinhibition after TBI.

The current research compared the performance of a group of participants with TBI to a control group on both a social and a non-social reversal learning task. Feedback-related negativities elicited by negative feedback on the reversal learning tasks were also measured. In order to determine whether reversal impairments were related to social disinhibition, participants with TBI were also rated by two independent, blind-raters on their level of social disinhibition based on a video-taped interview. It was predicted that participants with TBI would make more reversal errors and have attenuated feedback-related negativities compared to controls on both the non-social and the social task. Further, if reversal learning deficits play a role in acquired social disinhibition, those TBI participants high on social disinhibition should demonstrate an impairment compared to those low on social disinhibition in the ability to update responding when social reinforcement contingencies change in the social reversal learning task. Finally, it was hypothesised that attenuated feedback-related negativity amplitudes elicited by negative

social feedback would be observed for the participants with TBI high on social disinhibition compared with those low on social disinhibition.

## Method

### Participants

Twenty-one adults (19 males) who had sustained a severe traumatic brain injury (TBI) of mean age 46.90 years ( $SD=14.54$ , range: 22 to 68) with an average of 13.10 years of formal education ( $SD=1.87$ , range: 10 to 17) participated. Participants were recruited from the outpatient records of three metropolitan brain injury units in Sydney. Included participants met the following criteria: they had sustained a severe TBI resulting in at least one day of altered consciousness (Russell & Smith, 1961), were discharged from hospital and living in the community, were proficient in English and had no substance abuse or dependence. The participants with TBI had experienced post-traumatic amnesia (PTA) ranging from 2 to 137 days ( $Mean=56.8$ ,  $SD=33.52$ ), and time post injury ranging from 3 to 46 years ( $Mean=13.90$ ,  $Median=12.0$ ,  $SD=11.09$ ). PTA scores were obtained from patient medical records, with an exception of one participant whose records were unavailable. In this case, the injury was recorded as severe because coma duration exceeded 24 hours (Corrigan, Selassie, & Orman, 2010). The participants' injuries were sustained as a consequence of motor vehicle accidents ( $n=11$ ), falls ( $n=8$ ) and assaults ( $n=2$ ). CT scans from the clinical records showed that injuries were left hemisphere focused ( $n=4$ ), right hemisphere focused ( $n=5$ ) and bilateral ( $n=11$ ). A CT scan was not available for one participant. Specific frontal lobe injuries were reported in 12 participants. However, traditional imaging technology is not a reliable indicator of orbitofrontal damage. Orbitofrontal damage has been found using high resolution MRI in patients with behavioural change despite no obvious frontal lesions detected by traditional imaging technology

(Namiki et al., 2008). Further, frontal white matter damage has been identified using diffusion tensor imaging in patients with little cortical damage evident using standard imaging (Kinnunen et al., 2011).

Control participants were 21 adults (18 males) without brain injury with a mean age of 45.29 ( $SD=13.70$ , range: 22 to 68) and an average of 14.52 years of education ( $SD= 1.69$ , range: 11 to 18). Controls were recruited from the community via online and local newspaper advertisements. The control group did not differ significantly from the TBI group with respect to age,  $t(40)=.37$ ,  $p=.712$ ,  $d=.11$ , or with respect to emotion recognition scores,  $t(40)=-1.70$ ,  $p=.097$ ,  $d=-.52$ . However, the control group did differ from the TBI group in terms of number of years of education,  $t(40)=-2.60$ ,  $p=.013$ ,  $d=-.80$  and Depression, Anxiety and Stress Scale (DASS; Lovibond & Lovibond, 1995) total score,  $t(40)=3.07$ ,  $p=.004$ ,  $d=.94$ . To address these differences between groups in analyses, years of education was entered into the behavioural analyses as a covariate since it correlated with the outcome measure. Further, emotion recognition scores were entered as a covariate as they were theoretically relevant. Table 1 provides demographic information for the TBI and control group.

Table 1 about here.

## Materials

### Reversal Learning Task.

Participants were told that they could gain points in the task by selecting symbols displayed on the screen. As in Chase, Swainson, Durham, Benham, and Cools (2011), on each trial, two different hiragana symbols appeared on the screen and participants made a selection using a left or right mouse click. Participants learned by trial and error which of these symbols was correct and which was incorrect. Selection of the correct symbol was rewarded by the



delivery of the text “You WIN 1 point!”, while selection of the incorrect symbol was punished by the delivery of the text “You LOSE 1 point” in red font. The position of the symbols on the screen was randomised. Once the participant reached a criterion level of performance, the reinforcement contingency swapped, without warning, such that the previously correct symbol became incorrect and the previously incorrect symbol became correct. The contingencies continued to switch at the beginning of each block for a total of 16 blocks. The criterion level of performance to be reached before the reinforcement contingencies were reversed differed for each block, but was between 7 and 11 consecutive correct responses. This was to prevent participants from anticipating the reversal. If an error was made, the count toward the criterion level of performance for that block began again from zero. Thus, the number of trials per block depended on the performance of the individual. Each block had a maximum of 30 trials, after which the reward contingencies reversed whether or not the participant had reached criterion. Feedback presentation was displayed for 1000ms and the inter-trial interval was 500ms. Stimuli remained on the screen until a selection was made.

### **Social Reversal Learning Task.**

The social reversal learning task was based on that described by Kringelbach and Rolls (2003). This task ran identically to the non-social reversal learning task described above, except that the stimuli were black and white photographs of two faces with neutral expressions and the feedback consisted of a happy or angry expression of the photographed actor appearing in the place of the neutral expression. The first 8 blocks used two female faces and the second eight blocks used two male faces. The design of this task is represented in Figure 1. In this task, participants were not told that they were to gain points throughout the task but were just told to figure out which face to select at any given time. These instructions were designed to avoid the

possibility of participants applying a rule such as “a happy expression means I have gained a point” and thus to make reinforcement as close to natural social feedback as possible. The design of this task is represented in Figure 1. The order in which the participants received the social and the non-social reversal learning tasks was counterbalanced in order to minimise the impact of practice effects, since it has been suggested that reversal learning deficits disappear quickly with practice (Dias, Robbins, & Roberts, 1997; Schoenbaum, Nugent, Saddoris, & Setlow, 2002). Counterbalancing was achieved for the comparison between the TBI and control group as well as for comparison between the low disinhibition and high disinhibition group.

### **Social Disinhibition Interview Task.**

The current study used an adaptation of the self-disclosure task developed by Beer, John, Scabini, and Knight (2006). Participants were initially told that they would be asked a number of questions about themselves and their experiences, and that it was their choice how much information they wished to disclose and that they could skip any question at any time. These instructions were designed to minimise an expectation of excessive self-disclosure. Participants were then asked a series of nine questions, which included: “Tell me about an embarrassing moment you’ve had” and “Tell me about something someone has done to make you angry”. The interviews were videotaped and rated by two independent judges, blind to participant condition. Judges rated the frequency of the participants socially inappropriate behaviour on a scale of 1 to 5 (where 1 represented ‘never’ and 5 represented ‘always’) on the following items: ‘While talking with the interviewer, the participant spoke too candidly’, ‘Considering that they didn’t know the interviewer very well, the participant disclosed an inappropriate amount of information about themselves’, ‘The participant revealed more intimate details than most people would’, ‘The participant was rude’, ‘The participant made inappropriate jokes or remarks’, ‘The

participant was impatient', 'The participant did not know when to stop talking', 'The participant was critical or argumentative'. These items were based on a thorough review of literature reporting socially inappropriate behaviours displayed by individuals with damage to the OFC. The inter-rater reliability for ratings across both TBI and control groups was analysed with an intraclass coefficient (ICC) using a two factor mixed effect model. The inter-rater absolute agreement was good,  $ICC=.70$ , 95% CI [.43, .84]. The ICC was similar when looking at ratings for the TBI group alone,  $ICC=.70$ , 95% CI [.28, .87].

### **Emotion Recognition Task.**

Stimuli were 18 static images of one of four actors (two male and two female) portraying one of six emotions (happiness, surprise, sadness, anger, fear and disgust). Stimuli were still images taken from the emotion recognition task (ERT; Montagne, Kessels, De Haan, & Perrett, 2007), a computer-generated program which shows a series of 216 video clips of facial expressions across different intensities. The stimuli were developed using algorithms (Benson & Perrett, 1991) which created intermediate morphed images between a neutral face (0% emotion) and a full-intensity expression (100% emotion). Data from a study by Rosenberg, McDonald, Dethier, Kessels, and Westbrook (2014) which used the ERT video stimuli suggest that some emotions are much easier to recognise than others. Thus, in order to avoid floor and ceiling effects in recognition, 100% intensity of expression was used for fear, sadness and surprise stimuli, 80% intensity was used for anger and disgust stimuli, while 30% intensity was used for happy stimuli. Following the protocol of Heberlein, Padon, Gillihan, Farah, and Fellows (2008), participants were asked to rate the intensity of each of six emotions they detected in each stimulus. For each participant an accuracy score was derived by determining the number of trials on which participants correctly rated the expressed emotion as the most intense emotion in that

stimulus. This task was included in order to determine whether poor performance on the social reversal learning task could be explained by poor emotion recognition.

### **Procedure**

This study and its procedures were approved by the University of NSW Human Research Ethics Committee.

### **EEG Acquisition.**

EEG data was acquired using a PC-based digital signal-processing hardware and software package from Neuroscan (Compumedics, Acquire Version 4.5). Continuous EEG was recorded from 64 scalp sites using the Neuroscan Quick-cap. Signals were then filtered with a bandpass of 0.1-30 Hz, referenced to the nose and grounded by the cap electrode. Tin cup electrodes were placed 2 cm above and below the left eye, and on the outer canthus of each eye, measuring vertical (vEOG) and horizontal (hEOG) eye movements respectively. The maximum impedance was always below 5 k $\Omega$  for both EOG and cap electrodes.

### **EEG Data Analyses.**

Neuroscan Edit software (Compumedics 4.5) was used to calculate ERPs. The continuous data was bandpass filtered (0.01-30 Hz, zero-phase shift, down 24 db) and subjected to an EOG correction procedure (Semlitsch, Anderer, Schuster, & Presslich, 1986). Waveforms were segmented into epochs 200 ms pre- and 600 ms post-feedback onset. The feedback-locked data was then baseline corrected by subtracting the average activity during the 200 ms preceding the feedback onset. For each participant, difference waves were computed by subtracting the average wave for correct feedback from the average wave for error feedback. The reversal learning tasks used ensured at least 15 errors were made by each participant across a minimum of 150 trials. As is conventional in the literature, the FRN was measured base-to-peak (Hajcak, Moser, Holroyd,

& Simons, 2006; Holroyd et al., 2003; Yasuda, Sato, Miyawaki, Kumano, & Kuboki, 2004). The amplitude at the most negative peak between 200 and 500ms were derived from the individual difference waves. This large window accommodated the large variance in latency found for participants with a TBI. The FRN component was defined as the difference in an individual's difference wave between the negative peak identified and the preceding positive peak at medio-frontal channel FCZ. This electrode location was chosen because the FRN was largest at that site on examination of grand-averaged waveforms for the control group and based on previous studies showing the FRN is maximal at this medio-frontal site (Hajcak et al., 2006; Holroyd, Larsen, & Cohen, 2004; Holroyd et al., 2003). For each participant, two FRN's were derived, one for the social task and one for the non-social task. One control participant's EEG data for the social task was excluded due to faulty equipment. A task (social vs. non-social task) by group (TBI vs. control) repeated measures ANOVA was performed with FRN amplitude as the dependant variable. The FRN was not correlated with years of education nor with DASS total score for either task. Thus, no covariates were entered in this analysis. In addition, because there is evidence of laterality of processing for social information in the literature, FRN amplitude at both FC3 (over the right hemisphere) and FC4 (over the left hemisphere) was reported.

## Results

### Behavioural Results

Emotion recognition, DASS, disinhibition and reversal learning scores for both groups are outlined in Table 2. Correlations between these variables are provided in Table 3.

Table 2 about here.

Table 3 about here.

A 2 x 2 (task x group) repeated measures ANCOVA was conducted with number of

reversal errors as the dependant variable. The analysis revealed a significant main effect of group,  $F(1,40)=9.54$ ,  $p=.004$ ,  $\eta^2=.19$ , such that controls ( $M=17.64$ ,  $SE=1.54$ ) made fewer errors than did participants with TBI ( $M=24.36$ ,  $SE=1.54$ ). Group differences remained with the addition of years of education and emotion recognition as a covariate,  $F(1,38)=4.081$ ,  $p=.05$ , indicating that these variables were not important factors in this effect. Mean reversal errors for both groups and both tasks are shown in Figure 2. There was no significant main effect of task,  $F(1,40)=.02$ ,  $p=.892$ , and no significant interaction,  $F(1,40)=.14$ ,  $p=.709$ .

Social disinhibition ratings were not normally distributed in the TBI group, with a significant positive skewness of 3.08 ( $SE=.37$ ,  $p<.05$ ; Cramer & Howitt, 2004). To provide a meaningful metric based on these ratings individuals were categorised as low ( $n=10$ ) on social disinhibition if they received the lowest possible social disinhibition rating of 8. They were categorised as high ( $n=11$ ) on social disinhibition if they received a score of 9 or above. These two groups did not differ with regards to age ( $p=.396$ ), years of education ( $p=.369$ ), post-traumatic amnesia ( $p=.758$ ), time since injury ( $p=.731$ ) or DASS total score ( $p=.921$ ). Figure 3 shows reversal errors on both tasks for TBI participants high on social disinhibition and TBI participants low on social disinhibition. A repeated measures 2 x 2 (task x disinhibition) ANCOVA with number of reversal errors as the dependant variable revealed a trend toward a task by disinhibition interaction,  $F(1,19)=4.02$ ,  $p=.059$ ,  $\eta^2=.18$ . This result was significant when years of education and emotion recognition were added as covariates,  $F(1,17)=7.48$ ,  $p=.014$ ,  $\eta^2=.31$ . Because an a priori hypothesis was made about a specific relationship between the social reversal learning task and social disinhibition, univariate ANOVA's were carried out to determine whether differences between groups existed for each task separately. These analyses revealed that participants high on social disinhibition ( $M=29.18$ ,  $SD=11.04$ ) made significantly

more errors than those low on social disinhibition ( $M=19.80$ ,  $SD=4.66$ ) on the social reversal learning task,  $F(1,21)=9.23$ ,  $p=.007$ ,  $\eta^2=.34$ , but not on the non-social task,  $F(1,21)=.001$ ,  $p=.971$ .

### EEG Results

Figure 4 displays mean correct and incorrect waveforms, as well the difference waves (FRN), at electrode FCZ for each group and each task. Figure 5 displays the variance (SEM) contributing to the correct and incorrect wave forms for both groups and for both tasks. The repeated measures 2 x 2 (task x group) ANOVA with FRN amplitude as the dependant variables revealed a significant main effect of group,  $F(1,39)=8.97$ ,  $p=.005$ ,  $\eta^2=.19$ , such that controls ( $M=8.85$ ,  $SE=.85$ ) had higher FRN amplitudes than did the TBI group ( $M=5.29$ ,  $SE=.83$ ). There was also a main effect of task,  $F(1,39)=10.80$ ,  $p=.002$ ,  $\eta^2=.22$ , such that FRN amplitudes were higher in the social task ( $M=8.63$ ,  $SE=.92$ ) than in the non-social task ( $M=5.51$ ,  $SE=.57$ ). There was no significant interaction,  $F(1,39)=1.13$ ,  $p=.295$ .

In order to determine whether these results were affected by the inclusion of more correct trials than incorrect in the analysis, a separate analysis was run with equal number of trials. The above analysis was re-run on randomly selected 15 correct and 15 incorrect trials for each participant and each task and results remained the same. There was a significant group effect,  $F(1,39)=12.14$ ,  $p=.001$ ,  $\eta^2=.24$ , and a significant task effect,  $F(1,39)=4.98$ ,  $p=.031$ ,  $\eta^2=.11$ , but no interaction,  $F(1,39)=.79$ ,  $p=.378$ .

Figure 6 depicts the FRN difference wave at FC3 (left hemisphere), FCZ (central) and FC4 (right hemisphere) and shows that the FRN was larger over the right hemisphere compared to central and left hemisphere sites for the social task. A repeated measures 3 (electrode: FC3, FCZ, FC4) x 2 (task) ANOVA revealed a significant electrode by task interaction,

$F(2,80)=10.09, p<.001$ . Follow-up tests of simple effects revealed that there was a main effect of electrode for the social task,  $F(2,80)=16.42, p<.001$ , but not for the non-social task,  $F(2,82)=1.25, p=.291$ . For the social task, pairwise comparisons with Bonferroni correction revealed that the FRN difference wave at FC4 was greater than at FC3 ( $M_{diff}=1.92, p<.001$ ) but not different than at FCZ ( $M_{diff}=.63, p=.168$ ).

Finally, using only the TBI group, a repeated measures 2 x 2 (task x disinhibition) ANOVA with FRN amplitude as the dependant variable revealed no significant effect of task,  $F(1,19)=3.51, p=.076$ , no significant main effect of disinhibition,  $F(1,19)=.588, p=.453$ , and no significant interaction,  $F(1,19)=.07, p=.789$ .

## Discussion

The current study aimed to determine whether reversal learning deficits play a role in acquired social disinhibition after TBI by comparing performance of a group of people with TBI and a control group on a social and a non-social reversal learning task. As predicted, the TBI group made significantly more reversal errors across both versions of the reversal learning task than did controls, demonstrating an impaired ability to update behaviour when reinforcement contingencies change. Although reversal learning impairment has been previously demonstrated in a brain-injured sample (Rolls et al., 1994), the current study was the first to show that TBI participants are also impaired at reversing responding when social reinforcement contingencies change. Further, the current study found that TBI participants high on social disinhibition performed more poorly on the social reversal learning task than did those low on social disinhibition. This is consistent with Rolls et al. (1994) report of a reversal learning deficit in TBI patients who displayed socially inappropriate behaviours as reported by caregivers. The current research, however, is the first to demonstrate that reversal learning impairment is



associated with social disinhibition observed in an experimental setting. Further, this result could not be explained by poor emotion recognition in the high social disinhibition group. Together, these findings suggest that an inability to reverse social reinforcement contingencies may contribute to inappropriate social responding after TBI. Further, the current results suggest that the social reversal learning task may be a useful neuropsychological tool for detecting susceptibility to social disinhibition after TBI. This is significant because past research has been unable to identify neuropsychological predictors of social disinhibition, often reporting that disinhibited individuals perform normally on neuropsychological tests (Cicerone & Tanenbaum, 1997; Damasio, Grabowski, Frank, Galaburda, & Damasio, 1994).

The current study also measured feedback-related negativity amplitudes evoked by negative feedback in both the non-social and social reversal learning tasks. FRN's are thought to reflect dopaminergic midbrain reward prediction error signals, which drive the updating of reinforcement contingencies and thus the updating of behaviour (Holroyd & Coles, 2002). Participants with TBI had attenuated FRN amplitudes compared with controls across both tasks, indicating an impaired ability to generate reward prediction error signals when negative social and non-social feedback is encountered. Consistent with this, previous research has shown that people with TBI did not differentiate reward from non-reward at an electrophysiological level (Larson, Kelly, Stigge-Kaufman, Schmalfluss, & Perlstein, 2007). Together these findings suggest that people with TBI are impaired at reward processing and thus at signalling when a predicted reward has not been delivered. This impairment in reward prediction error signalling was not, however, related to social disinhibition. This finding is contrary to the hypothesis that FRN amplitudes reflecting social reward prediction error signals drive changes in behaviour to enable adaptive and context appropriate social behaviour. It suggests that while these signals

may be important in indicating when social feedback is worse than was expected, they may not necessarily correlate with updated behaviour. In fact, while some studies have found a link between FRN amplitude and the updating of behaviour (Cohen & Ranganath, 2007; Holroyd & Krigolson, 2007; van der Helden et al., 2010), other studies have demonstrated that FRN's are generated when no behavioural adaptation is required (Gehring & Willoughby, 2004; Luu, Tucker, Derryberry, Reed, & Poulsen, 2003), suggesting that the FRN is not necessarily a signal used for learning. Thus, social reward prediction errors may not constitute sufficient information upon which to base a decision to change behaviour.

Since the FRN has been widely reported to be maximal centrally, the right hemisphere lateralisation of the FRN in the social task, illustrated in Figure 6, warrants discussion. Another study has similarly found a right-hemisphere lateralised 'social FRN' elicited by unfair offers from other 'players' in a computerised game (Boksem & De Cremer, 2010). Gehring and Willoughby (2004) have suggested that lateralised contributing activity could result in a lateralised FRN. The right hemisphere lateralisation of social FRNs, then, is in line with a pattern of literature documenting right hemisphere lateralisation of social reward processing (Demaree, Everhart, Youngstrom, & Harrison, 2005). For example, right hemisphere dominance has been found for processing of negative emotional expressions (Adolphs, Damasio, Tranel, & Damasio, 1996; Nakamura et al., 1999) and in responding to negative social feedback (Kaplan & Zaidel, 2001). Thus, the right hemisphere lateralisation of the FRN produced by negative social feedback in the current study likely results from right hemisphere dominance of negative social feedback processing.

A couple of limitations of the current study must be considered when interpreting the results. The TBI group had a slightly higher probability of experiencing error feedback in the

reversal learning tasks than did controls. It is well established that a larger amplitude FRN is produced by less probable events (Sambrook & Goslin, 2015). This is because the more a reward comes to be expected, the greater the reward prediction error signal will be when the reward is not delivered. In the current study, the control group experienced error feedback on 11.5% of trials on average, while the TBI group experienced error feedback on 13.7% of trials. This seems a trivial difference in terms of participant's perceptions of the probability of error feedback and is unlikely to be the source of group differences. Even so, future research should attempt to replicate this finding using a paradigm which equates number of errors as a percentage of total trials. Further, despite ample evidence to suggest that reversal learning impairment and social disinhibition stem from OFC damage, the current study cannot confirm the origins of observed impairments in the TBI group. The use of high resolution imaging technology in combination with the measures used here could clarify these findings.

In summary, the current research found increased reversal errors and decreased FRN amplitudes elicited by error feedback in participants with TBI when compared with controls across both a social and a non-social reversal learning task. Further, participants with TBI high on social disinhibition made more errors on the social reversal learning task than did those low on social disinhibition, supporting the hypothesis that reversal learning impairments underlie acquired social disinhibition after TBI. Attenuated FRN amplitudes in people with TBI indicate an impairment in feedback monitoring, possibly driven by an inability to differentiate reward from non-reward at an electrophysiological level. This impairment was not found to be a feature of socially disinhibited individuals specifically, though, suggesting that reward prediction error signals are not critical for behavioural adaptation in the social domain.

#### **Acknowledgements**

445           We express our gratitude to people with traumatic brain injuries who participated in the  
446 studies reported here as well as to our community control participants who gave willingly of  
447 their time. The authors have no competing or conflicts of interest to report.

## References

- Adolphs, R., Damasio, H., Tranel, D., & Damasio, A. R. (1996). Cortical systems for the recognition of emotion in facial expressions. *The Journal of Neuroscience*, 16(23), 7678-7687.
- Arciniegas, D. B., & Wortzel, H. S. (2014). Emotional and behavioral dyscontrol after traumatic brain injury. *Psychiatric Clinics of North America*, 37(1), 31-53. doi: 10.1016/j.psc.2013.12.001
- Bachevalier, J., & Loveland, K. A. (2006). The orbitofrontal–amygdala circuit and self-regulation of social–emotional behavior in autism. *Neuroscience & Biobehavioral Reviews*, 30(1), 97-117. doi: 10.1016/j.neubiorev.2005.07.002
- Barrash, J., Tranel, D., & Anderson, S. W. (2000). Acquired personality disturbances associated with bilateral damage to the ventromedial prefrontal region. *Developmental Neuropsychology*, 18(3), 355-381. doi: 10.1207/S1532694205Barrash
- Beer, J. S., John, O. P., Scabini, D., & Knight, R. T. (2006). Orbitofrontal cortex and social behavior: integrating self-monitoring and emotion-cognition interactions. *Journal of Cognitive Neuroscience*, 18(6), 871-879. doi: 10.1162/jocn.2006.18.6.871
- Benson, P. J., & Perrett, D. I. (1991). Perception and recognition of photographic quality facial caricatures: Implications for the recognition of natural images. *European Journal of Cognitive Psychology*, 3(1), 105-135. doi: 10.1080/09541449108406222
- Bigler, E. D. (2007). Anterior and middle cranial fossa in traumatic brain injury: Relevant neuroanatomy and neuropathology in the study of neuropsychological outcome. *Neuropsychology*, 21(5), 515-531. doi: 10.1037/0894-4105.21.5.515 17784800

- 470 Blair, R. J. R., & Cipolotti, L. (2000). Impaired social response reversal 'A case of acquired  
471 sociopathy'. *Brain*, 123(6), 1122-1141. doi: 10.1093/brain/123.6.1122
- 472 Boksem, M. A., & De Cremer, D. (2010). Fairness concerns predict medial frontal negativity  
473 amplitude in ultimatum bargaining. *Social Neuroscience*, 5(1), 118-128. doi:  
474 10.1080/17470910903202666
- 475 Brooks, N., & McKinlay, W. (1983). Personality and behavioural change after severe blunt head  
476 injury - a relative's view. *Journal of Neurology, Neurosurgery & Psychiatry*, 46(4), 336-  
477 344. doi: 10.1136/jnnp.46.4.336
- 478 Butter, C. M., Mishkin, M., & Mirsky, A. F. (1968). Emotional responses toward humans in  
479 monkeys with selective frontal lesions. *Physiology & Behavior*, 3(2), 213-215. doi:  
480 10.1016/0031-9384(68)90087-5
- 481 Chase, H. W., Swainson, R., Durham, L., Benham, L., & Cools, R. (2011). Feedback-related  
482 negativity codes prediction error but not behavioral adjustment during probabilistic  
483 reversal learning. *Journal of Cognitive Neuroscience*, 23(4), 936-946. doi:  
484 10.1162/jocn.2010.21456
- 485 Cicerone, K. D., & Tanenbaum, L. N. (1997). Disturbance of social cognition after traumatic  
486 orbitofrontal brain injury. *Archives of Clinical Neuropsychology*, 12(2), 173-188. doi:  
487 10.1093/arclin/12.2.173
- 488 Cohen, M. X., & Ranganath, C. (2007). Reinforcement learning signals predict future decisions.  
489 *The Journal of Neuroscience*, 27(2), 371-378. doi: 10.1523/JNEUROSCI.4421-06.2007
- 490 Corrigan, J. D., Selassie, A. W., & Orman, J. A. L. (2010). The epidemiology of traumatic brain  
491 injury. *The Journal of Head Trauma Rehabilitation*, 25(2), 72-80. doi:  
492 10.1097/HTR.0b013e3181ccc8b4

- 493 Cramer, D., & Howitt, D. (2004). *The Sage dictionary of statistics: A practical resource for*  
494 *students in the social sciences.*: Thousand Oaks: Sage.
- 495 Damasio, H., Grabowski, T., Frank, R., Galaburda, A. M., & Damasio, A. R. (1994). The return  
496 of Phineas Gage - Clues about the brain from the skull of a famous patient. *Science*,  
497 *264*(5162), 1102-1105. doi: 10.1126/science.8178168
- 498 Dias, R., Robbins, T., & Roberts, A. C. (1997). Dissociable forms of inhibitory control within  
499 prefrontal cortex with an analog of the Wisconsin Card Sort Test: restriction to novel  
500 situations and independence from “on-line” processing. *The Journal of Neuroscience*,  
501 *17*(23), 9285-9297.
- 502 Fellows, L. K., & Farah, M. J. (2003). Ventromedial frontal cortex mediates affective shifting in  
503 humans: evidence from a reversal learning paradigm. *Brain*, *126*(8), 1830-1837. doi:  
504 10.1093/brain/awg180
- 505 Franzen, E., & Myers, R. (1973). Neural control of social behavior: prefrontal and anterior  
506 temporal cortex. *Neuropsychologia*, *11*(2), 141-157. doi: 10.1016/0028-3932(73)90002-  
507 X
- 508 Gehring, W. J., & Willoughby, A. R. (2004). Are all medial frontal negativities created equal?  
509 Toward a richer empirical basis for theories of action monitoring. *Errors, Conflicts, and*  
510 *the Brain. Current Opinions on Performance Monitoring*, 14-20.
- 511 Gottfried, J. A., O'Doherty, J., & Dolan, R. J. (2003). Encoding predictive reward value in  
512 human amygdala and orbitofrontal cortex. *Science*, *301*(5636), 1104-1107. doi:  
513 10.1126/science.1087919
- 514 Gould, K. R., Ponsford, J. L., Johnston, L., & Schonberger, M. (2011). Relationship between  
515 psychiatric disorders and 1-year psychosocial outcome following traumatic brain injury.

- 516 *Journal of Head Trauma Rehabilitation*, 26(1), 79-89. doi:  
517 10.1097/Htr.0b013e3182036799
- 518 Gregory, C. A., Serra-Mestres, J., & Hodges, J. R. (1999). Early Diagnosis of the Frontal Variant  
519 of Frontotemporal Dementia: How Sensitive Are Standard Neuroimaging and  
520 Neuropsychologic Tests? *Cognitive and Behavioral Neurology*, 12(2), 128-135.
- 521 Hajcak, G., Moser, J. S., Holroyd, C. B., & Simons, R. F. (2006). The feedback-related  
522 negativity reflects the binary evaluation of good versus bad outcomes. *Biological*  
523 *Psychology*, 71(2), 148-154. doi: 10.1016/j.biopsycho.2005.04.001
- 524 Hajcak, G., Moser, J. S., Holroyd, C. B., & Simons, R. F. (2007). It's worse than you thought:  
525 The feedback negativity and violations of reward prediction in gambling tasks.  
526 *Psychophysiology*, 44(6), 905-912. doi: 10.1111/j.1469-8986.2007.00567.x
- 527 Heberlein, A. S., Padon, A. A., Gillihan, S. J., Farah, M. J., & Fellows, L. K. (2008).  
528 Ventromedial Frontal Lobe Plays a Critical Role in Facial Emotion Recognition. *Journal*  
529 *of Cognitive Neuroscience*, 20(4), 721-733. doi: 10.1162/jocn.2008.20049
- 530 Hikosaka, K., & Watanabe, M. (2004). Long-and short-range reward expectancy in the primate  
531 orbitofrontal cortex. *European Journal of Neuroscience*, 19(4), 1046-1054. doi:  
532 10.1111/j.0953-816X.2004.03120.x
- 533 Holroyd, C. B., & Coles, M. G. H. (2002). The neural basis of human error processing:  
534 Reinforcement learning, dopamine, and the error-related negativity. *Psychological*  
535 *Review*, 109(4), 679-709. doi: 10.1037/0033-295X.109.4.679 12374324
- 536 Holroyd, C. B., & Krigolson, O. E. (2007). Reward prediction error signals associated with a  
537 modified time estimation task. *Psychophysiology*, 44(6), 913-917. doi: 10.1111/j.1469-  
538 8986.2007.00561.x



- 539 Holroyd, C. B., Krigolson, O. E., Baker, R., Lee, S., & Gibson, J. (2009). When is an error not a  
540 prediction error? An electrophysiological investigation. *Cognitive, Affective, &*  
541 *Behavioral Neuroscience*, 9(1), 59-70. doi: 10.3758/CABN.9.1.59
- 542 Holroyd, C. B., Larsen, J. T., & Cohen, J. D. (2004). Context dependence of the event-related  
543 brain potential associated with reward and punishment. *Psychophysiology*, 41(2), 245-  
544 253. doi: 10.1111/j.1469-8986.2004.00152.x
- 545 Holroyd, C. B., Nieuwenhuis, S., Yeung, N., & Cohen, J. D. (2003). Errors in reward prediction  
546 are reflected in the event-related brain potential. *Neuroreport*, 14(18), 2481-2484. doi:  
547 10.1097/01.wnr.0000099601.41403.a5
- 548 Hornak, J., O'doherty, J., Bramham, J., Rolls, E. T., Morris, R., Bullock, P., & Polkey, C. (2004).  
549 Reward-related reversal learning after surgical excisions in orbito-frontal or dorsolateral  
550 prefrontal cortex in humans. *Journal of Cognitive Neuroscience*, 16(3), 463-478. doi:  
551 10.1162/089892904322926791
- 552 Hornberger, M., Geng, J., & Hodges, J. R. (2011). Convergent grey and white matter evidence of  
553 orbitofrontal cortex changes related to disinhibition in behavioural variant frontotemporal  
554 dementia. *Brain*, 134(9), 2502-2512. doi: 10.1093/brain/awr173
- 555 Izquierdo, A., Suda, R. K., & Murray, E. A. (2004). Bilateral orbital prefrontal cortex lesions in  
556 rhesus monkeys disrupt choices guided by both reward value and reward contingency.  
557 *The Journal of Neuroscience*, 24(34), 7540-7548. doi: 10.1523/JNEUROSCI.1921-  
558 04.2004
- 559 Kaplan, J. T., & Zaidel, E. (2001). Error monitoring in the hemispheres: the effect of lateralized  
560 feedback on lexical decision. *Cognition*, 82(2), 157-178. doi: 10.1016/S0010-  
561 0277(01)00150-0

- 562 Kinnunen, K. M., Greenwood, R., Powell, J. H., Leech, R., Hawkins, P. C., Bonnelle, V., . . .  
563 Sharp, D. J. (2011). White matter damage and cognitive impairment after traumatic brain  
564 injury. *Brain*, 134(2), 449-463. doi: 10.1093/brain/awq347
- 565 Kringelbach, M. L., & Rolls, E. T. (2003). Neural correlates of rapid reversal learning in a  
566 simple model of human social interaction. *Neuroimage*, 20(2), 1371-1383. doi:  
567 10.1016/S1053-8119(03)00393-8
- 568 Krueger, C. E., Laluz, V., Rosen, H. J., Neuhaus, J. M., Miller, B. L., & Kramer, J. H. (2011).  
569 Double dissociation in the anatomy of socioemotional disinhibition and executive  
570 functioning in dementia. *Neuropsychology*, 25(2), 249-259. doi: 10.1037/a0021681
- 571 Larson, M. J., Kelly, K. G., Stigge-Kaufman, D. A., Schmalfluss, I. M., & Perlstein, W. M.  
572 (2007). Reward context sensitivity impairment following severe TBI: an event-related  
573 potential investigation. *Journal of the International Neuropsychological Society*, 13(04),  
574 615-625.
- 575 Lipszyc, J., Levin, H., Hanten, G., Hunter, J., Dennis, M., & Schachar, R. (2014). Frontal white  
576 matter damage impairs response inhibition in children following traumatic brain injury.  
577 *Archives of Clinical Neuropsychology*, 29(3), 289-299. doi: 10.1093/arclin/acu004
- 578 Lovibond, P. F., & Lovibond, S. H. (1995). The structure of negative emotional states:  
579 Comparison of the Depression Anxiety Stress Scales (DASS) with the Beck Depression  
580 and Anxiety Inventories. *Behaviour Research and Therapy*, 33(3), 335-343. doi:  
581 10.1016/0005-7967(94)00075-U
- 582 Luu, P., Tucker, D. M., Derryberry, D., Reed, M., & Poulsen, C. (2003). Electrophysiological  
583 responses to errors and feedback in the process of action regulation. *Psychological*  
584 *Science*, 14(1), 47-53. doi: 10.1111/1467-9280.01417 12564753

- 585 Machado, C. J., & Bachevalier, J. (2006). The impact of selective amygdala, orbital frontal  
586 cortex, or hippocampal formation lesions on established social relationships in rhesus  
587 monkeys (*Macaca mulatta*). *Behavioral Neuroscience*, 120(4), 761-786. doi:  
588 10.1037/0735-7044.120.4.761
- 589 Mattson, A. J., & Levin, H. S. (1990). Frontal lobe dysfunction following closed head injury. A  
590 review of the literature. *The Journal of Nervous and Mental Disease*, 178(5), 282-291.
- 591 McKinlay, W., Brooks, N., Bond, M., Martinage, D., & Marshall, M. (1981). The short-term  
592 outcome of severe blunt head injury as reported by relatives of the injured persons.  
593 *Journal of Neurology, Neurosurgery & Psychiatry*, 44(6), 527-533. doi:  
594 10.1136/jnnp.44.6.527
- 595 Montagne, B., Kessels, R. P. C., De Haan, E. H. F., & Perrett, D. I. (2007). The emotion  
596 recognition task: A paradigm to measure the perception of facial emotional expressions at  
597 different intensities. *Perceptual and Motor Skills*, 104(2), 589-598. doi:  
598 10.2466/Pms.104.2.589-598
- 599 Nakamura, K., Kawashima, R., Ito, K., Sugiura, M., Kato, T., Nakamura, A., . . . Fukuda, H.  
600 (1999). Activation of the right inferior frontal cortex during assessment of facial emotion.  
601 *Journal of Neurophysiology*, 82(3), 1610-1614.
- 602 Namiki, C., Yamada, M., Yoshida, H., Hanakawa, T., Fukuyama, H., & Murai, T. (2008). Small  
603 orbitofrontal traumatic lesions detected by high resolution MRI in a patient with major  
604 behavioural changes. *Neurocase*, 14(6), 474-479. doi: 10.1080/13554790802459494
- 605 Nieuwenhuis, S., Holroyd, C. B., Mol, N., & Coles, M. G. (2004). Reinforcement-related brain  
606 potentials from medial frontal cortex: origins and functional significance. *Neuroscience*  
607 & *Biobehavioral Reviews*, 28(4), 441-448. doi: 10.1016/j.neubiorev.2004.05.003

- 608 Padoa-Schioppa, C., & Assad, J. A. (2006). Neurons in the orbitofrontal cortex encode economic  
609 value. *Nature*, 441(7090), 223-226. doi: 10.1038/nature04676
- 610 Rahman, S., Sahakian, B. J., Hodges, J. R., Rogers, R. D., & Robbins, T. W. (1999). Specific  
611 cognitive deficits in mild frontal variant frontotemporal dementia. *Brain*, 122(8), 1469-  
612 1493. doi: 10.1093/brain/122.8.1469
- 613 Rolls, E. T., Hornak, J., Wade, D., & McGrath, J. (1994). Emotion-related learning in patients  
614 with social and emotional changes associated with frontal lobe damage. *Journal of*  
615 *Neurology, Neurosurgery & Psychiatry*, 57(12), 1518-1524. doi:  
616 10.1136/jnnp.57.12.1518
- 617 Rosenberg, H., McDonald, S., Dethier, M., Kessels, R. P. C., & Westbrook, R. F. (2014). Facial  
618 Emotion Recognition Deficits following Moderate-Severe Traumatic Brain Injury (TBI):  
619 Re-examining the Valence Effect and the Role of Emotion Intensity. *Journal of the*  
620 *International Neuropsychological Society*, 20(10), 994-1003. doi:  
621 10.1017/S1355617714000940
- 622 Sambrook, T. D., & Goslin, J. (2015). A neural reward prediction error revealed by a meta-  
623 analysis of ERPs using great grand averages. *Psychological Bulletin*, 141(1), 213-235.  
624 doi: 10.1037/bul0000006
- 625 Schoenbaum, G., Nugent, S. L., Saddoris, M. P., & Setlow, B. (2002). Orbitofrontal lesions in  
626 rats impair reversal but not acquisition of go, no-go odor discriminations. *Neuroreport*,  
627 13(6), 885-890.
- 628 Schoenbaum, G., Roesch, M. R., Stalnaker, T. A., & Takahashi, Y. K. (2009). A new perspective  
629 on the role of the orbitofrontal cortex in adaptive behaviour. *Nature Reviews*  
630 *Neuroscience*, 10(12), 885-892. doi: 10.1038/nrn2753

- 631 Schoenbaum, G., Takahashi, Y., Liu, T. L., & McDannald, M. A. (2011). Does the orbitofrontal  
632 cortex signal value? *Annals of the New York Academy of Sciences*, 1239(1), 87-99. doi:  
633 10.1111/j.1749-6632.2011.06210.x
- 634 Semlitsch, H. V., Anderer, P., Schuster, P., & Presslich, O. (1986). A solution for reliable and  
635 valid reduction of ocular artifacts, applied to the P300 ERP. *Psychophysiology*, 23(6),  
636 695-703. doi: 10.1111/j.1469-8986.1986.tb00696.x
- 637 Tate, R. L., Broe, G. A., & Lulham, J. M. (1989). Impairment after severe blunt head-injury - the  
638 results from a consecutive series of 100 patients. *Acta Neurologica Scandinavica*, 79(2),  
639 97-107. doi: 10.1111/j.1600-0404.1989.tb03719.x
- 640 van der Helden, J., Boksem, M. A., & Blom, J. H. (2010). The importance of failure: feedback-  
641 related negativity predicts motor learning efficiency. *Cerebral Cortex*, 20(7), 1596-1603.  
642 doi: 10.1093/cercor/bhp224
- 643 Walsh, M. M., & Anderson, J. R. (2011a). Learning from delayed feedback: Neural responses in  
644 temporal credit assignment. *Cognitive, Affective, & Behavioral Neuroscience*, 11(2), 131-  
645 143. doi: 10.3758/s13415-011-0027-0
- 646 Walsh, M. M., & Anderson, J. R. (2011b). Modulation of the feedback-related negativity by  
647 instruction and experience. *Proceedings of the National Academy of Sciences*, 108(47),  
648 19048-19053. doi: 10.1073/pnas.1117189108
- 649 Yasuda, A., Sato, A., Miyawaki, K., Kumano, H., & Kuboki, T. (2004). Error-related negativity  
650 reflects detection of negative reward prediction error. *Neuroreport*, 15(16), 2561-2565.  
651  
652

Table 1

*Means, standard deviations, ranges and results of group comparisons for the TBI and comparison groups*

Table 2

*Correlations between demographic variables, emotion functioning, disinhibition, emotion recognition and reversal learning across the TBI and control group (N=42)*

*Figure 1.* Design of the social reversal learning task.

*Figure 2.* Mean number of errors on the social and the non-social reversal learning tasks for the TBI and control group.

*Figure 3.* Mean number of errors on the social and the non-social reversal learning tasks for TBI participants with high (n=11) and low (n=10) social disinhibition.

*Figure 4.* Average waveforms for the TBI and control group for correct and incorrect trials as well as the difference waveform. Waveforms for the non-social reversal learning task can be seen in the left panels and for the social reversal learning task in the right panels.

*Figure 5.* Variance (SEM) contributing to the correct and incorrect wave forms for both groups and for both tasks.

*Figure 6.* Feedback-related negativity at electrodes FC3, FCZ and FC4 for the non-social task for (a) the

677 control group and (b) the TBI group, as well as for the social task for (c), the control group and (d) the  
678 TBI group.

679

Table 1

*Means, standard deviations, ranges and results of group comparisons for demographic variables*

	Mean (SD), Range		Diff ( <i>p</i> )	Cohen's <i>d</i>
	TBI ( <i>N</i> =21)	Control ( <i>N</i> =21)		
Demographics				
PTA (days)	56.80 (33.52), 2-137			
Time Since Injury (years)	13.90 (11.09), 3-46			
Age	46.90 (14.54), 22-68	45.29 (13.70), 22-68	.712	.11
Years of education	13.10 (1.87), 10-17	14.52 (1.69), 11-18	.013*	-.80



Table 2.

*Means, standard deviations, ranges and results of group comparisons for experimental variables*

	Mean (SD), Range		Diff ( <i>p</i> )	Cohen's <i>d</i>
	TBI ( <i>N</i> =21)	Control ( <i>N</i> =21)		
Emotion Recognition	10.71 (2.72), 4-16	12.05 (2.36), 6-15	.097	.52
DASS Total	30.52 (6.66), 6-108	11.42 (12.56), 0-42	.004**	.94
Disinhibition	10.02 (3.20), 8-20	8.69 (.94), 8-11.5	.075	.57
Reversal Learning				
Non-Social Reversal Errors	24.00 (13.30), 15-64	17.81 (2.62), 14-25	.043*	.65
Social Reversal Errors	24.71 (9.68), 16-52	17.48 (1.69), 15-21	.002**	1.07

Table 3.

*Correlations between demographic and experimental variables across the TBI and control group (N=42)*

	Age	Years of Education	DASS Total Score	Disinhibition	Emotion Recognition	Non-Social Reversal Errors	Social Reversal Errors
Demographics							
Age		-.026	.238	-.039	-.208	.072	.140
Years of Education			-.198	.015	.153	-.272	-.325*
DASS Total Score				.447**	-.066	.197	.169
Disinhibition					-.030	.064	.242
Emotion Recognition						-.314*	-.266
Reversal Learning							
Non-Social Reversal Errors							.515**
Social Reversal Errors							

*Note.* \*Significant at  $p < .05$ . \*\* Significant at  $p < .001$ .