

Learning affective values for faces is expressed in amygdala and fusiform gyrus

Predrag Petrovic, Raffael Kalisch, Mathias Pessiglione, Tania Singer, and Raymond J. Dolan

Wellcome Trust Centre for Neuroimaging, University College of London, 12 Queen Square, London, WC1N 3BG, UK

To monitor the environment for social threat humans must build affective evaluations of others. These evaluations are malleable and to a high degree shaped by responses engendered by specific social encounters. The precise neuronal mechanism by which these evaluations are constructed is poorly understood. We tested a hypothesis that conjoint activity in amygdala and fusiform gyrus would correlate with acquisition of social stimulus value. We tested this using a reinforcement learning algorithm, Q-learning, that assigned values to faces as a function of a history of pairing, or not pairing, with aversive shocks. Behaviourally, we observed a correlation between conditioning induced changes in skin conductance response (SCR) and subjective ratings for likeability of faces. Activity in both amygdala and fusiform gyrus (FG) correlated with the output of the reinforcement learning algorithm parameterized by these ratings. In amygdala, this effect was greater for averted than direct gaze faces. Furthermore, learning-related activity change in these regions correlated with SCR and subjective ratings. We conclude that amygdala and fusiform encode affective value in a manner that closely approximates a standard computational solution to learning.

Keywords: amygdala; fusiform gyrus; reinforcement learning; face stimuli; affective value

INTRODUCTION

There is an evolutionary benefit in tracking potential sources of threat in our social environment. For example, a specific identity may predict an aversive intention based upon previous associations with that individual. This may be expressed in low-level learning (Buchel *et al.*, 1998; Olsson *et al.*, 2005) as well as in how we experience others (Gottfried and Dolan, 2004; Singer *et al.*, 2004, 2006). Thus, the brain needs to continually update how we evaluate specific individuals depending on past emotional associations. The mechanisms supporting the encoding of these explicit evaluations are poorly understood.

Neurobiological research on the human amygdala has focused on fear-learning, as expressed in classical conditioning, and processing of fear-related stimuli. Responses to such stimuli are largely dependent on amygdala and include change in autonomic output and fight-or-flight behavioural repertoires (Davis and Whalen, 2001; Phelps, 2006). However, the amygdala also influences cortical activity either through direct anatomical connections (Amaral *et al.*, 2003) or indirect connectivity (Kapp *et al.*, 1994).

The fusiform gyrus (FG) plays a critical role in processing faces (Kanwisher *et al.*, 1997; Haxby *et al.*, 2000) and its activity is modulated by emotional expressions (Vuilleumier *et al.*, 2001; Vuilleumier and Pourtois, 2006). Activity in FG is augmented, in conjunction with amygdala activation,

when a face expresses fear even when attentional processes are directed elsewhere (Vuilleumier *et al.*, 2001). A causal link to amygdala is supported by observations that patients with damage to this structure do not show emotion induced modulation (Vuilleumier *et al.*, 2004). Neutral faces may also be associated to positive or negative values both through social (Singer *et al.*, 2004) or fear learning (Gottfried and Dolan, 2004).

While it is known that amygdala is critical in fear learning where autonomic responses were the dependent measure (Buchel *et al.*, 1998; Tranel, 2000; Cheng *et al.*, 2003; Knight *et al.*, 2005; Olsson *et al.*, 2005; Carter *et al.*, 2006) the computational processes involved in learning explicit evaluation of faces has not been described. Here, we tested whether the change in value associated with aversive learning is implemented by a reinforcement learning algorithm implemented in amygdala and FG.

To alter face evaluations, we presented subjects with faces that were, or were not, paired with shocks while we simultaneously recorded brain activity using fMRI. Subjects rated how sympathetic they perceived the faces, before and after conditioning, a rating that provided a measure of how the affective value of faces changed with conditioning. Value learning for faces was assessed using a formal reinforcement model, Q-learning (Sutton and Barto, 1998). Our hypothesis was that affective values would change after conditioning, with faces paired with shock being rated as less sympathetic than faces not paired with shocks, and the change would correlate with the magnitude of change in skin conductance responses (SCRs). Crucially, we also hypothesized that the output of our reinforcement learning model would correlate in time with activity in amygdala and FG. In other words,

Received 25 June 2007; Accepted 9 January 2008

Advance Access publication 8 February 2008

We thank Stefan Kloppel for helpful discussions during preparation of the article and the staff at the Wellcome Trust Centre for Neuroimaging. This work is supported by a Wellcome Trust Programme Grant to R.J.D., Hjärnfonden and Vetenskapsrådet.

Correspondence should be addressed to Predrag Petrovic, Wellcome Trust Centre for Neuroimaging, University College of London, 12 Queen Square, London, WC1N 3BG, UK. E-mail: predrag.petrovic@ki.se.

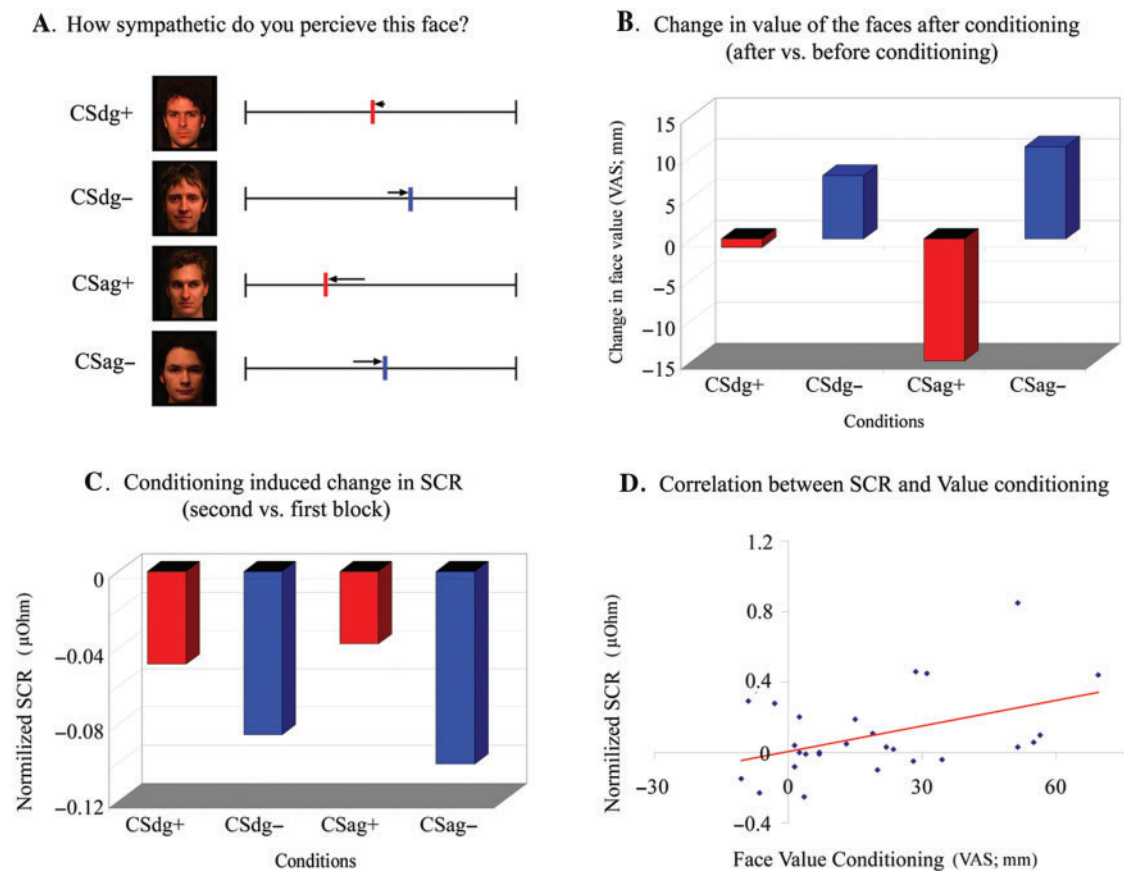


Fig. 1 Relation between learning new affective values and change in autonomic responses after conditioning of faces. (A) Two faces, one with direct gaze (CSdg+) and one with averted gaze (CSag+), were paired with shocks in a standard Pavlovian conditioning paradigm. Two matched faces were never paired with shocks. Subjects rated how sympathetic they perceived faces on a visual analogue scale before and after learning providing an index of value (the arrows show the change in value). (B) Average values of sympathy ratings (VAS 0–100 mm) decreased for faces paired with shock (CSdg+ and CSag+; change of value after shock shown in red), while the average values for the faces with direct (CSdg-) and averted gaze (CSag-), not paired with shock, increased (change of value after shock shown in blue) when pre-conditioning ratings were compared with post-conditioning ratings. The effect of conditioning on face rating was significant (Wilcoxon Signed Ranks Test, P -value <0.005). (C) The average SCR were more expressed for CSdg+ and CSag+ vs. CSdg- and CSag- in the second block compared with the first block, and a repeated measurement ANOVA showed a significant time by conditioning interaction effect ($F=4.681$, P -value <0.05). Here, we show the change in SCR in the second vs the first block in order to compare it with the change in value (B). (D) The learning effect for affective value (conditioning induced change in sympathy ratings for CS- vs CS+; VAS 0–100 mm) correlated with the learning expressed in SCR (Spearman's rho, correlation coefficient = 0.446, P -value <0.05) indicating a functional relationship between autonomic reactivity and how subjects explicitly evaluated the faces.

newly learned values encoded in these structures would support the changed affective evaluation measured post-conditioning.

METHODS

Subjects

Thirty right-handed healthy male subjects were included in the present study that was approved by the local ethical committee (Department of Neurology and Neurosurgery, UCL, London, UK). The subjects had no history of mental or psychiatric disorder. Before the study, the subjects signed a written informed consent. Three subjects were excluded from the analysis, two because they showed a high degree of drowsiness and informed us that they several times had their eyes closed and fallen asleep in the scanner in-between the shocks, and one because of large movement

artefacts in his scanning data. The data from the remaining 27 subjects (age 18–36 years) were used for the analysis.

Experimental design

During fMRI acquisition the subjects were shown four different faces, two with direct gaze and two with averted gaze (Figure 1A; previously used in George *et al.*, 2001). Direct gaze has been shown to engage amygdala and the FG to a higher degree than averted gaze in a neutral situation indicating a higher processing of social information and identity (Kawashima *et al.*, 1999; George *et al.*, 2001; Haxby *et al.*, 2002). Given that gaze direction has different social meanings and is associated differently with the amygdala–fusiform network, we predicted that evaluative learning would differ for averted vs direct gaze.

Two faces (CS+) (one with direct gaze and one with averted gaze) were paired with an electric shock (UCS) to the right hand with a 50% contingency. The other two faces were never paired with the electric shock (CS−). Thus, there were six conditions in total, CS+ with direct gaze paired with a shock (CSdg+UCS, $n = 15$), CS+ with averted gaze paired with a shock (CSag+UCS, $n = 15$), CS+ with direct gaze unpaired (CSdg+, $n = 15$), CS+ with averted gaze unpaired (CSag+, $n = 15$), CS− with direct gaze (CSdg−, $n = 30$) and CS− with averted gaze (CSag−, $n = 30$). The four faces used were randomly ascribed to one of the conditions [(i) CSdg+/CSdg+UCS, (ii) CSag+/CSag+UCS, (iii) CSdg− and (iv) CSag−] for each subject. Each face appeared for 990 ms with a jittered intra-stimulus time of between 10 800 and 14 400 ms between each presentation. Shock was given at the end of the face presentation. The faces were randomly presented in the centre of the screen or 5 mm to the left or to the right. In order to assure attention to the task, subjects were instructed to indicate as quickly as possible with one of three keys on a response box in which position the face appeared. Thus, the faces were presented with a short duration and a simultaneous cognitive task was applied to the task. The total time for the conditioning was 24 min during the one scanning session.

The shocks were delivered with a Digitimer DS7A electrical stimulator (Digitimer Ltd., Welwyn Garden City, UK), and shocks consisted of electrical pulses of up to 20 mA and 1 ms duration through a silver chloride electrode. Stimulation parameters were individually adjusted prior to the experiment to achieve maximum tolerable pain. The stimulation intensity was set individually to 80% on a visual analogue scale, where 0 meant no painful sensation and 100 meant most intense painful sensation imaginable. The stimulation had also to be tolerable—thus if a subject perceived that highest tolerable sensation was below 80% on the VAS that intensity was used in the experiment.

Before conditioning subjects were instructed to indicate how sympathetic each face was experienced on a visual-analogue scale from 0 to 100, where 0 meant that the subject did not perceive the face as sympathetic and 100 meant that he perceived the face as the most sympathetic person he could imagine. Subjects were instructed to complete the same rating directly after the conditioning, ratings that defined subjective value of the faces. Subjects were never told about the contingency nor were they told that this was a conditioning experiment. After the experiment subjects were debriefed about what they had learned about the relation between shock and the different face stimuli. We expected the experimental pairing to lead to a decrease in value for the CS+ relative to the CS− where a larger effect involves a more negative value for the CS+ and a more positive value for the CS−. We assessed differences between acquired ratings for CS+ and CS− stimuli with non-parametric tests since subjective rating scales were used and only one rating was obtained for each condition and time-point.

During stimulus presentation, we acquired standard skin conductance response (SCR) at a sampling rate of 1000 Hz from electrodes on the middle and ring finger of the left hand using an AT64 SCR apparatus (Autogenic Systems, Wood Dale, IL, USA). The SCR data was down-sampled to 100 Hz and mean filtered. We defined the SCR as the highest increase in a time period in-between 1 and 5 s after the start of stimulus presentation minus the baseline skin conductance just preceding each measurement period (Buchel *et al.*, 1998; Kalisch *et al.*, 2006). To normalize the data, we divided the SCR responses for each individual with their average SCR after the UCS. In this way, it is possible to interpret the normalized SCR directly in relation to the response to USC. The effect of conditioning has often been assessed by measuring the first and second interval response (Buchel *et al.*, 1998; LaBar *et al.*, 1998; Morris *et al.*, 2001; Critchley *et al.*, 2002; Tabbert *et al.*, 2005, 2006). Here, we grouped together the first half of the stimuli (first 7 CS+ and first 15 CS−) and defined these responses as the first block, while the remaining (last 8 CS+ and last 15 CS−) were defined as the second block. The level of conditioning was assessed as SCR for (CS+ vs CS−)_{second} vs (CS+ vs CS−)_{first} for each subject. We assessed awareness of contingency between faces and shock after the study by showing all four faces together and asking for an explicit judgement of which of the faces was/were associated with shock. Subjects that could identify both CS+’s correctly were defined as belonging to the high-aware group and the others defined as belonging to the low-aware group.

Q-learning

Q-learning is a simple, but plausible, model of how the brain updates its affective knowledge about the expected value of new stimuli, on the basis of their previous association with rewards and punishments. The model implements the Rescola–Wagner (RW) rule (Rescorla and Wagner, 1972) under an assumption that learning is directly proportional to a prediction error that reflects a difference between actual and expected affective outcomes. This class of model is known to be relevant for learning about rewards (O’Doherty *et al.*, 2003, 2006), and punishments (Seymour *et al.*, 2004, 2005). We used Q-learning to model how subjects change their evaluation of facial stimuli depending on their associated reinforcement pattern. The parameters of the model were optimized so that the first and last Q-values fit the group-average pre- and post-learning ratings for each face stimuli.

Distinct face-stimuli (CSdg+, CSag+, CSdg− and CSag−) were modelled separately. Q-values were set at zero before learning, and after every trial $t > 0$ the value of each stimulus was updated according to the rule $Q(t+1) = Q(t) + \alpha * \delta(t)$. The prediction error was defined as $\delta(t) = R(t) - Q(t)$, with $R(t)$ as outcome (shock or no-shock). The learning model was fitted with a single set of parameters (α and R s) across all subjects. The constant α (learning rate) was set to 0.1, a value suggested as a plausible learning rate in previous

conditioning tasks (O'Doherty *et al.*, 2003, 2006). The values for R was adjusted to minimize the distance between the change in Q -values and change in ratings, comparing pre- and post-learning (the average change in rating of the faces was -0.9 for CSdg+, -15.1 for CSag+, 7.7 for CSdg- and 11.3 for CSag-). R was estimated to be -26 for shock and $+10$ for no-shock.

These optimal parameters (α and R) were then used to generate a series of Q -values for each face trial (Csdg+ etc) for each individual (see above), which were then used as Q -learning regressors (QR) in SPM to model value learning for each face-stimulus. The regressors differed depending on the individual pattern of reinforcement received by each subject. Therefore, the CS- regressors (which were never paired with the shocks) were the same across all stimuli (CSdg- and CSag-) and subjects, but the regressors differed for CS+ (which were paired with the shocks on half of the trials in a random way) both across stimuli (Csdg+ and CSag+) and subject. Examples of the individual regressors are shown in Supplementary Figure 4.

fMRI scanning and imaging analysis

The imaging data (T_2^* -weighted echo planar images-EPI) measuring blood-oxygen level dependency (BOLD) contrast were acquired using a 1.5-Tesla Siemens Sonata system. We used a 30° tilted orbitofrontal sequence (Deichmann *et al.*, 2002) with a flip angle of 90° covering the whole brain in 44 planes. The TR was set to 3.96 s (90 ms per slice) and TE to 50 ms in a single session of 24 min. Images were processed using SPM5 (www.fil.ion.ucl.ac.uk/spm; Ashburner *et al.*, 2004). Scans were realigned, normalized and spatially smoothed by an 8 mm full-width half-maximum Gaussian kernel. A high-pass filter (with a cut-off at 120 s) was applied to the time series. The data was then analysed in an event-related fashion.

We modelled the conditions and regressors for each subject in a first-level analysis in SPM and the resulting t -maps were then taken to a second-level analysis. The contrasts of interest for the group average were assessed using a random effect general linear model to allow for statistical inference across the population. Activations are described in the regions involved in pain processing (Peyron *et al.*, 2000) and FG at an uncorrected P -value <0.001 .

Although, the regressors for all six conditions were modelled, the regressors for CS+USC were not used in the critical analysis in order not to contaminate our results with shock events. We first determined regions showing an increase in activity with a value decrease for CS+ but not for CS-, as in $[QR_{(CSdg+, CSag+, CSdg-, CSag-)}] = [-1 -1 -1 -1]$. Although we are contrasting the learning effects in CS+ vs CS-, the regressors have the same sign (-1) since they already have different directions, and we are interested in areas that will increase as the value-learning regressors for CS+ decrease specifically (Table 1). We also performed analyses of simple main effects

Table 1 Reinforcement learning of face-values in amygdala and FG

Changes in value with Q-learning model			
Region	Co-ordinate	Z-value	ROI Corr P-value
Main effect of conditioning*			
Right FG	50 -52 -20	3.75	<0.05
Left FG	-44 -54 -20	2.55	N.S.
Left amy	-24 4 -16	Set-level	<0.05
	-22 -4 -24		
Correlation with face rating†			
Right FG	36 -56 -22	3.52	<0.05

Amy, amygdala; FG, fusiform gyrus; ROI Corr, Region of interest corrected.

*Activations in the amygdala- and FG-ROIs that correlate with the output of a reinforcement learning model.

†Regions where the reinforcement learning model correlates with the individual learning of face values.

All activations are ROI corrected.

$[QR_{(CSdg+, CSag+, CSdg-, CSag-)}] = [-1 0 -1 0]$ and $[QR_{(CSdg+, CSag+, CSdg-, CSag-)}] = [0 -1 0 -1]$, and interactions $[QR_{(CSdg+, CSag+, CSdg-, CSag-)}] = [-1 1 -1 1]$ and $[QR_{(CSdg+, CSag+, CSdg-, CSag-)}] = [1 -1 1 -1]$ to analyse whether there was any difference in the value learning for the directed vs the averted gaze. Given that we were specifically interested in how the amygdala and the FG were involved in affective processing of faces, we defined ROI's in amygdala (radius = 8 mm; centre coordinate n : $[x y z] = [\pm 22 -2 -21]$) and FG (radius = 10 mm; centre coordinate: $[x y z] = [\pm 44 -53 -19]$) based on two previous studies where the specific effect of threat on faces processing was investigated (Vuilleumier *et al.*, 2001, 2004), and used those ROIs for correcting the P-values from our results (small volume correction - svc). To test whether individual subjective change of values correlated with Q -value learning $[QR_{(CSdg+, CSag+, CSdg-, CSag-)}] = [-1 -1 -1 -1]$ in the FG and the amygdala we performed simple second-level regression analyses.

RESULTS

Behavioural results

Effect of conditioning on rating of faces. A prior hypothesis of a change in face evaluation after exposure was confirmed (i.e. how sympathetic the face was perceived depending on whether a shock was associated to the face stimuli or not; Figure 1A and B). The value for CS+ faces decreased (mean change in value for CS+ = -7.98 , s.d. = 13.80), while the value for CS- faces increased (mean change in value for CS- = 9.50 , s.d. = 15.58) post-, as compared with pre-conditioning (Wilcoxon Signed Ranks Test, $Z = -3.486$, $P < 0.001$; Figure 1B). Thus, while the CS+ faces were rated as less sympathetic, the CS- faces were rated as more sympathetic after the conditioning procedure. A majority of the subjects ($n = 20$) were not able

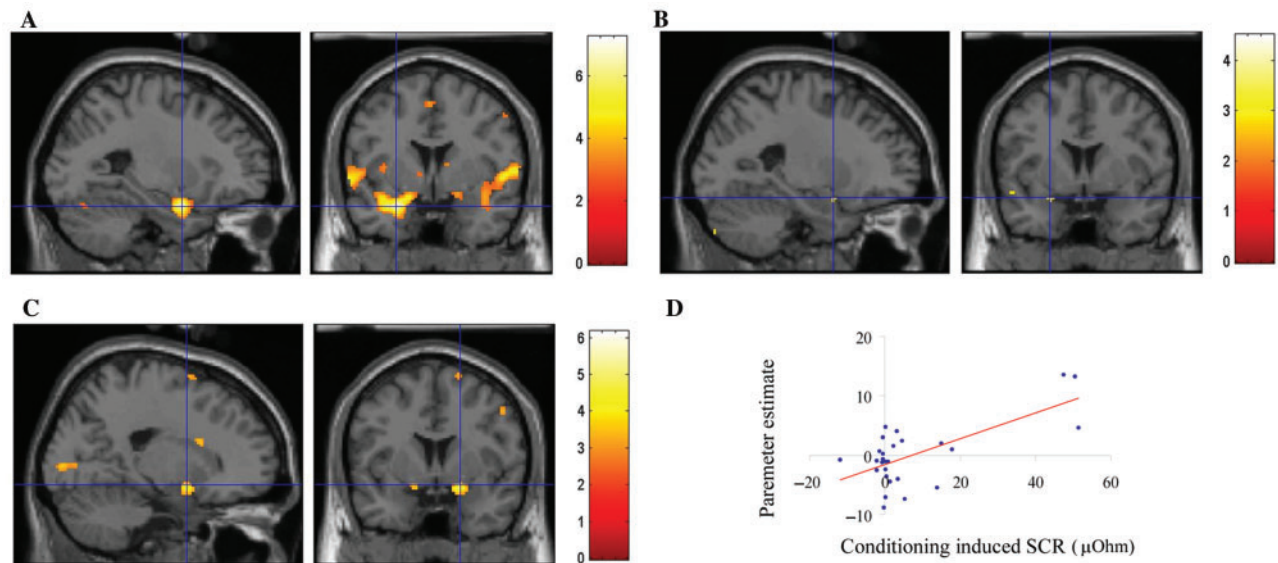


Fig. 2 Amygdala involvement in pain processing, conditioning response and reinforcement learning of affective values. (A) Main effect of pain showing amygdala activation bilaterally ($[x\ y\ z] = [18\ -2\ -12]$; Z -value = 5.05; P -value <0.001 uncorrected; $[x\ y\ z] = [-26\ 2\ -22]$; Z -value = 3.79; P -value <0.001 uncorrected). (B) The reinforcement Q-learning model for change in face value correlated with activity in a left amygdala ROI volume ($[x\ y\ z] = [-22\ -4\ -24]$ and $[-24\ 4\ -16]$; svc set level P -value <0.05). All the activations are thresholded on $P < 0.005$. (C) Conditioning induced learning effect, i.e. $(CS+ \text{ vs } CS-)$ _{second} vs $(CS+ \text{ vs } CS-)$ _{first} in SCR correlated significantly with the same learning effects in the fMRI data bilaterally in the dorsal amygdala ($[x\ y\ z] = [20\ 4\ -12]$; Z -value = 3.86; $[x\ y\ z] = [-16\ 6\ -12]$; Z -value = 3.42; P -value <0.001 uncorrected). Note the similarity between SCR-dependent learning effects and learning of experiential values in amygdala. (D) Scatterplot of the right amygdala activation presented above.

to accurately identify which faces was associated with the shock in a post-experimental structured interview, but still showed the same conditioning effect in how they rated faces (Wilcoxon Signed Ranks Test, $Z = -2.539$, $P < 0.005$). *Post hoc* we tested for whether the conditioning effect was different between the two gaze contexts. No significant effects were observed when the conditioning effect ($CS+ \text{ vs } CS-$) for direct gaze was compared with the conditioning effect for averted faces (Wilcoxon signed ranks test; $Z = -1.177$, $P = 0.239$, two-tailed).

Effect of conditioning on SCR. A three-by-two repeated measure ANOVA for SCR with factors of conditioning ($CS+$ or $CS-$), gaze (direct or averted) and time (first or second block) showed a time by conditioning interaction ($F = 4.681$, P -value < 0.05). The interaction effect was attributed to a larger mean SCR for $CS+ \text{ vs } CS-$ in the second block vs the first block mirroring the induced changes in affective rating (Figure 1C) ($CSag+$ _{second} = 0.2501 μ Ohm s.d. = 0.2068; $CSdg+$ _{second} = 0.2526 s.d. = 0.2198; $CSag-$ _{second} = 0.2159 s.d. = 0.1446; $CSdg-$ _{second} = 0.2030 s.d. = 0.1379; $CSag+$ _{first} = 0.2880 s.d. = 0.2395; $CSdg+$ _{first} = 0.3011 s.d. = 0.1628; $CSag-$ _{first} = 0.3166 s.d. = 0.2249; $CSdg-$ _{first} = 0.2884 s.d. = 0.1813). No significant main effects of gaze were observed. The conditioning effect on SCR (total change of normalized SCR in second vs first block of acquisition for $CS+ \text{ vs } CS-$ stimuli) showed a significant correlation with the change in face values (Spearman's rho, Correlation coefficient = 0.464, $P < 0.05$; Figure 1D). Thus, subjects that

showed highest conditioning effect as indexed by the magnitude of acquired autonomic response also expressed the highest effect in explicit rating of faces.

Neuroimaging

Face and UCS responses. The unconditioned stimuli, i.e. the painful shock, extracted from the contrast [($Csdg+UCS$ plus $CSag+UCS$) minus ($Csdg+$ plus $CSag+$)], induced increased activation in pain-related regions including thalamus bilaterally, caudal ACC, posterior and anterior insula bilaterally, secondary somatosensory cortex bilaterally and several brainstem regions (Supplementary Table 1 and Figure 1). Bilateral activations were evident in amygdala (Figure 2A) and in functionally defined bilateral fusiform gyri (Supplementary Figure 3). To show areas involved in processing faces, we compared all face events that were not coupled with shocks (both $CS+$ and $CS-$) against baseline. This showed significant activation in the FG that extended throughout extrastriate cortex into primary visual cortex (Supplementary Figure 2).

The categorical analysis showed no significant conditioning by time interaction [($CS+ \text{ minus } CS-$)_{second} vs ($CS+ \text{ minus } CS-$)_{first}] in the amygdala. However, conditioning related effects in bilateral amygdala [($CS+ \text{ minus } CS-$)_{second} vs ($CS+ \text{ minus } CS-$)_{first}] were significantly correlated with the conditioning specific SCR as expressed across time [($CS+ \text{ minus } CS-$)_{second} vs ($CS+ \text{ minus } CS-$)_{first}] (Figure 2C and D). Thus, subjects showing the largest

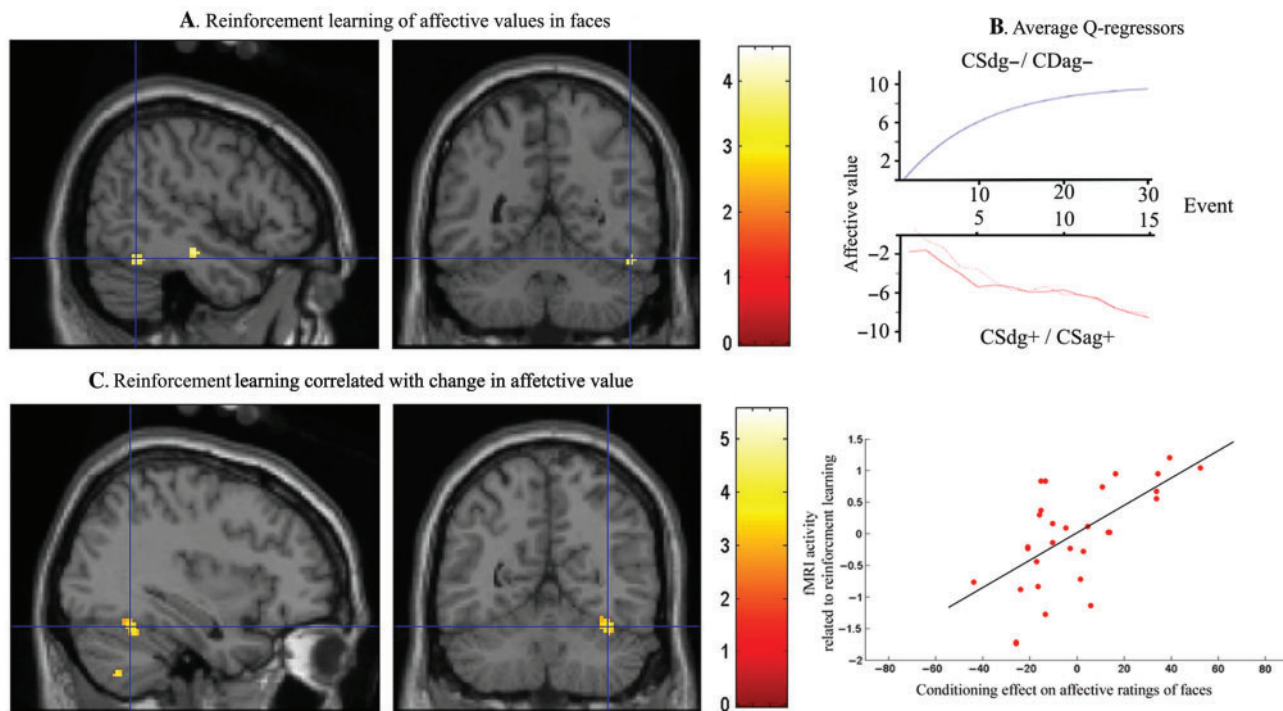


Fig. 3 FG involvement in reinforcement learning of affective values. (A) The reinforcement Q-learning model for change in face values showed a significant correlation with activity in right FG ($[x\ y\ z] = [50\ -52\ -20]$; Z -value = 3.75; svc P -value <0.05). (B) FG activity correlated negatively with the Q-regressors for reinforcement learning of affective values (the average model is displayed in the figure). Thus, FG showed more activity the more values decreased for CS+ but not for CS-. (C) An adjacent region of FG showed a correlation between the Q-learning model and individual indices of learning for face-values ($[x\ y\ z] = [36\ -56\ -22]$; Z -value = 3.52; svc P -value <0.05). Thus, subjects that showed greatest change in face evaluation also showed the largest reinforcement learning effect in fusiform area within the Q-learning model. The plot shows the correlation described in (C) in the maximally activated voxel. All the activations are thresholded on $P < 0.005$.

conditioning effect as measured by acquired change in SCR also showed the largest conditioning effect in the amygdala.

Brain changes and Q-learning of affective ratings. Activity in left amygdala (Table 1, Figure 2B) and right FG (Table 1, Figure 3A) correlated with the output of the Q-learning model (Figure 3B) that expressed change in face values across all subjects. In other words, activity increased as a function of decreasing subjective value for CS+ relative to CS- faces, indicating that the less sympathetic the subjects perceived the CS+ faces, according to the reinforcement model, the more these regions activated. When the reinforcement learning model was correlated with how actual face evaluations changed on a subject by subject basis (i.e. second-level regression), we observed a significant effect in right FG (Table 1; Figure 3C). Thus, subjects that showed the most extensive change in how they evaluated faces also showed the largest effect of reinforcement learning in FG as determined by Q-learning.

Q-learning effects as a function of gaze direction. We next examined whether Q-learning of face value was influenced by gaze direction. The relevant interaction analysis showed a significant effect in bilateral amygdala (Table 2; Figure 4). This effect was further characterized by analysis of simple effects where Q-learning for averted faces showed a significant effect in the amygdala (Table 2). In contrast,

Table 2 Reinforcement learning of face-values for direct and averted gaze

Changes in value with Q-learning model			
Region	Co-ordinate	Z-value	ROI Corr P-value
Averted gaze*			
Right FG	44 -54 -30	2.91	N.S.
Left amy	-20 -8 -22	3.25	<0.05
Direct gaze†			
Left FG	-44 -54 -20	2.55	N.S.
Left amy	-22 -2 -21	2.68	N.S.
Gaze interaction‡			
Right amy	20 -6 -28	3.51	<0.05
Left amy	-18 4 -24	Set-level	<0.05
	-18 -8 -206		

Amy, amygdala; QR, Q-regressor; dg, direct gaze; ag, averted gaze; +, stimuli paired with UCS; -, stimuli not paired with UCS; ROI Corr, Region of interest corrected. Activity in the amygdala- and FG-ROIs correlate with the reinforcement model of learning face values during conditioning for stimuli displaying. *Averted gaze. †Direct gaze. ‡Their interaction. All activations are ROI corrected.

the simple main effects for direct gaze did not show significant effects in these regions. Thus, the interaction effect in the Q-learning was driven primarily by effects for averted gaze.

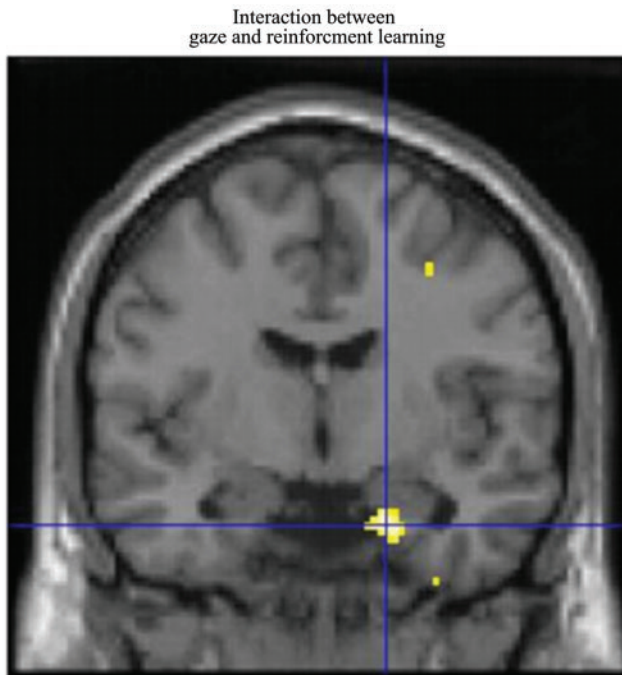


Fig. 4 Difference between direct and averted gaze in neuronal response to reinforcement learning of affective values. An interaction analysis [QR_(CSdg+, CSag+, CSdg-, CSag-)] = [1 -1 1 -1] showing that the Q-reinforcement learning model was significantly more correlated in faces with averted gaze vs direct gaze in the amygdala (right side activation shown here). All the activations are thresholded on $P < 0.005$.

DISCUSSION

Humans update their affective response to individuals associated with, or causal in engendering an aversive event. Consequently, a neutral face may be experienced as highly emotional after it has been associated with a threat stimulus. Although it is known that amygdala is involved in conditioning (Davis and Whalen, 2001; Phelps, 2006), and that face processing in FG is amplified by negative emotional expressions (Vuilleumier *et al.*, 2001, 2004), the computational basis for learning affective values for faces has not been previously described. The present results suggest that there is a direct relation between amygdala and FG activity and how we encode explicit values for faces.

A change in how faces are affectively rated after being paired with shock directly correlated with the change in SCR, a metric usually conceived as a marker of conditioning (Tranel and Damasio, 1989; Tranel, 2000; Cheng *et al.*, 2003; Knight *et al.*, 2005; Carter *et al.*, 2006). Although conditioning is usually considered as an implicit form of learning in humans, in reality there are likely to be both explicit and implicit learning components. We note that not all subjects learnt the contingencies between shock and predictive stimuli (CS+ faces), but subjects that learnt these contingencies tended to have larger conditioning induced effects on affective ratings and SCR (Supplementary Figure 5). This finding is in line with a suggestion of a relation between

awareness of conditioning stimulus contingencies and SCR (Lovibond and Shanks 2002; Tabbert *et al.*, 2006).

A significant conditioning induced change in autonomic response was observed in that the SCR difference between CS+ and CS- increased in the second block of the study (when compared to the first block), although no overall increase of SCR was observed for CS+ vs CS-. We noted a large variability in the first block where some subjects responded strongly to the 'wrong' face, which might have precluded such a finding. Thus, with faces it is difficult to have a neutral pre-learning state and we suggest that an experience-dependent learning effect, i.e. the interaction showing an increase in (CS+ vs CS-) for the second block vs the first block is a better indicator of a learning effect in the present study than an overall differential CS+ vs CS- SCR difference.

The amygdala is a key region in emotional learning (Davis and Whalen, 2001; Phelps, 2006). The magnitude of SCRs, an indicator of conditioning, is correlated with amygdala activity (Tranel and Damasio, 1989; Buchel *et al.*, 1998; Tranel, 2000; Cheng *et al.*, 2003; Knight *et al.*, 2005; Olsson *et al.*, 2005; Carter *et al.*, 2006). Knight *et al.* (2005) suggested that amygdala induced SCR was coupled to the conditioned response and not the CS+ stimuli itself. In line with this suggestion, we show that a conditioning induced SCR correlated with the magnitude of conditioning induced amygdala activity underlining a tight relationship between CR and amygdala activity (Figure 2C). Although some studies have shown that amygdala activity and SCR are expressed more in early phases of conditioning, probably due to habituation in the later phase (Buchel *et al.*, 1998; LaBar *et al.*, 1998; Morris *et al.*, 2001), other studies have indicated an opposite effect (Critchley *et al.*, 2002; Tabbert *et al.*, 2005). Recently, Straube *et al.* (2007) showed a dissociation in early and late amygdala activation during conditioning depending on attentional load. While simple conditioning showed early amygdala activation, conditioning during high attentional load was associated with late amygdala activation. Our late amygdala activation might be explained by the distraction of our use of a speeded spatial discrimination task (as opposed to identity-identification) simultaneous with actual conditioning. Moreover, using two CS+ and two CS- stimuli instead of one of each, the short stimulus presentation, similar face-stimuli, long inter-stimulus interval, low CS-UCS contingencies and sub-optimal presentation environment in the scanner may have influenced the strength of the conditioning procedure.

The main goal of this experiment was to establish whether learnt evaluation of faces could be captured by a computational learning model and whether the output of this model would be reflected in activity in emotional processing brain regions. The amygdala is one such region and a putative functional role is modulation of cortical stimulus processing that has an emotional signature (Dolan and Vuilleumier, 2003). The amygdala is thought to augment

fear face processing in a well-defined area in lateral anterior FG (Vuilleumier and Pourtois, 2006). We hypothesized these regions are also important in updating affective evaluations that predicts the likelihood of an aversive event. In order to model how the conditioning changes the affective evaluation of faces, we used a Q-learning reinforcement model, a simple learning model for both appetitive (O'Doherty *et al.*, 2003, 2006) and aversive reinforcement learning (Seymour *et al.*, 2004, 2005) that incorporates an error-signal that updates prediction regarding the value of the next stimulus. Crucially, the output of our reinforcement learning model correlated with activity in amygdala and anterior FG, regions implicated in emotional perception of faces (Vuilleumier *et al.*, 2001, 2004). This suggests that learning high-level emotional evaluations is supported by both amygdala and fusiform activity.

We modelled each CS paired with shock (CS+UCS) as decreasing the value of a face and every CS that was not paired with a shock (CS−/CS+) as increasing the value of the face based on the measured average ratings for all the subjects. The increase in value for absence of shock is in line with the suggestion that an omitted aversive event is a proxy for reward (Seymour *et al.*, 2005). This was evident in the behavioural results, which showed that CS− faces were rated as more likable after the conditioning procedure. In line with the hypothesis, we showed that the amygdala and the FG significantly correlated with the output of our reinforcement learning model of value changes. Thus, the activity changes in these regions were more related to aversive learning in CS+ faces than for likeability learning in CS− faces. Moreover, this relation was shown to be maximally expressed in a similar region of the right FG for subjects that showed the largest change in affective ratings of face evaluations. We emphasize that involvement of amygdala and FG captures learning-related activity for value, and not perception, after conditioning. Thus, we suggest that amygdala and FG, a circuit known to be involved in the emotional processing of faces, are involved in learning of explicit face values.

The reinforcement-learning effects were more pronounced in the amygdala for the faces displaying averted gaze. No significant difference as a function of exposure to shock was evident between the two gaze conditions for affective ratings or SCR. Direct gaze activates amygdala and FG more than averted gaze in a neutral context (Kawashima *et al.*, 1999; George *et al.*, 2001). It has been suggested that this effect underlines the importance of direct gaze as social and identity signal, while information from averted gaze is more important for guiding spatial attention (Haxby *et al.*, 2002). A logical hypothesis is that a reinforcement learning model should apply more to direct than averted gaze in amygdala and FG. However, the alternative hypothesis that averted gaze should more closely follow the learning of Q-values is more in line with the present data set. Since a face with direct

gaze activates amygdala and FG automatically then any learning effect in these regions is likely to be smaller than for averted faces. These results cannot support a claim that any gaze-stimuli show a stronger evaluative learning effect in general, only that the present Q-learning model is more correlated with activity in amygdala for averted than direct faces. Interestingly, the Q-learning model was differentially expressed for the averted compared with the direct gaze faces in the amygdala in the high-aware subjects compared with the low aware subjects suggesting that the degree of awareness also influence how different gaze faces are conditioned (Supplementary Figure 5). In the FG, there was no tendency for a different effect of the reinforcement learning model depending on gaze. Thus, our Q-learning model seems to apply for both gaze conditions here.

The UCS, i.e. the pain shocks, showed involvement of the amygdala similar to the SCR modulated learning effect. Nociceptive signals may reach amygdala either directly (Bernard *et al.*, 1996) or indirectly (Shi and Davis, 1999) and induce increased activity (Bornhovd *et al.*, 2002), although top-down processes can suppress the amygdala response in pain processing (Petrovic *et al.*, 2004). The amygdala is preferably activated to nociceptive processing when the CS predicts a shock (Carlsson *et al.*, 2006) also indicating its specific involvement in learning and CR. Thus, the close spatial location between the amygdala activation in the UCS, CR-SCR correlation and Q-learning probably indicates a functional relationship. Also, the FG showed increase activation in the UCS condition (Supplementary Figure 3) indicating an involvement of the amygdala–fusiform circuit in the CR, possibly as a teaching signal (prediction error).

In conclusion, we show a strong behavioural relation between explicit learning of how sympathetic faces are evaluated and low-level learning as expressed in SCR. This is mirrored by reinforcement learning related changes, as captured in a Q-learning model, in amygdala and FG. These results suggest that the brain updates affective values for faces through implementation of reinforcement learning-like processes implemented in amygdala and FG. It is likely that learning affective values is mediated via an amygdala modulation of FG, reflecting causal interactions between these structures as shown in earlier studies (Vuilleumier *et al.*, 2004). While it has previously been indicated that top-down influences such as subjective preconceptions (Olsson *et al.*, 2005) and instructed learning (Phelps *et al.*, 2001) change low-level conditioning for faces, our study suggests that conditioning is involved in evaluative learning of faces.

SUPPLEMENTARY DATA

Supplementary data are available at SCAN online.

REFERENCES

- Amaral, D.G., Behniea, H., Kelly, J.L. (2003). Topographic organization of projections from the amygdala to the visual cortex in the macaque monkey. *Neuroscience*, 118, 1099–120.
- Ashburner, J., Friston, K., Penny, W. (2004). Imaging neuroscience—Theory and analysis. In: Frackowiak, R.S.J., Friston, K., Frith, C.D. et al. editors. *Human Brain Function*, 2nd edn, San Diego: Academic Press, pp. 599–1104.
- Bernard, J.F., Bester, H., Besson, J.M. (1996). Involvement of the spino-parabrachio-amygdaloid and -hypothalamic pathways in the autonomic and affective emotional aspects of pain. *Progress Brain in Research*, 107, 243–55.
- Bornhovd, K., Quante, M., Glauche, V., Bromm, B., Weiller, C., Buchel, C. (2002). Painful stimuli evoke different stimulus-response functions in the amygdala, prefrontal, insula and somatosensory cortex: a single-trial fMRI study. *Brain*, 125, 1326–36.
- Buchel, C., Morris, J., Dolan, R.J., Friston, K.J. (1998). Brain systems mediating aversive conditioning: an event-related fMRI study. *Neuron*, 20, 947–57.
- Carlsson, K., Andersson, J., Petrovic, P., Petersson, K.M., Ohman, A., Ingvar, M. (2006). Predictability modulates the affective and sensory-discriminative neural processing of pain. *Neuroimage*, 32, 1804–14.
- Carter, R.M., O'Doherty, J.P., Seymour, B., Koch, C., Dolan, R.J. (2006). Contingency awareness in human aversive conditioning involves the middle frontal gyrus. *Neuroimage*, 29, 1007–12.
- Cheng, D.T., Knight, D.C., Smith, C.N., Stein, E.A., Helmstetter, F.J. (2003). Functional MRI of human amygdala activity during Pavlovian fear conditioning: stimulus processing versus response expression. *Behavioural Neuroscience*, 117, 3–10.
- Critchley, H.D., Mathias, C.J., Dolan, R.J. (2002). Fear conditioning in humans: the influence of awareness and autonomic arousal on functional neuroanatomy. *Neuron*, 33, 653–63.
- Davis, M., Whalen, P.J. (2001). The amygdala: vigilance and emotion. *Molecular Psychiatry*, 6, 13–34.
- Deichmann, R., Josephs, O., Hutton, C., Corfield, D.R., Turner, R. (2002). Compensation of susceptibility-induced BOLD sensitivity losses in echo-planar fMRI imaging. *Neuroimage*, 15, 120–35.
- Dolan, R.J., Vuilleumier, P. (2003). Amygdala automaticity in emotional processing. *Annals of the New York Academic Sciences*, 985, 348–55.
- George, N., Driver, J., Dolan, R.J. (2001). Seen gaze-direction modulates fusiform activity and its coupling with other brain areas during face processing. *Neuroimage*, 13, 1102–12.
- Gottfried, J.A., Dolan, R.J. (2004). Human orbitofrontal cortex mediates extinction learning while accessing conditioned representations of value. *Nature Neuroscience*, 7, 1144–52.
- Haxby, J.V., Hoffman, E.A., Gobbini, M.I. (2000). The distributed human neural system for face perception. *Trends in Cognitive Sciences*, 4, 223–33.
- Haxby, J.V., Hoffman, E.A., Gobbini, M.I. (2002). Human neural systems for face recognition and social communication. *Biological Psychiatry*, 51, 59–67.
- Kalisch, R., Korenfeld, E., Stephan, K.E., Weiskopf, N., Seymour, B., Dolan, R.J. (2006). Context-dependent human extinction memory is mediated by a ventromedial prefrontal and hippocampal network. *Journal of Neuroscience*, 26, 9503–11.
- Kanwisher, N., McDermott, J., Chun, M.M. (1997). The fusiform face area: a module in human extrastriate cortex specialized for face perception. *Journal of Neuroscience*, 17, 4302–11.
- Kapp, B.S., Supple, W.F.J., Whalen, P.J. (1994). Effects of electrical stimulation of the amygdaloid central nucleus on neocortical arousal in the rabbit. *Behavioural Neuroscience*, 108, 81–93.
- Kawashima, R., Sugiura, M., Kato, T., et al. (1999). The human amygdala plays an important role in gaze monitoring. A PET study. *Brain*, 122(Pt 4), 779–83.
- Knight, D.C., Nguyen, H.T., Bandettini, P.A. (2005). The role of the human amygdala in the production of conditioned fear responses. *Neuroimage*, 26, 1193–200.
- LaBar, K.S., Gatenby, J.C., Gore, J.C., LeDoux, J.E., Phelps, E.A. (1998). Human amygdala activation during conditioned fear acquisition and extinction: a mixed-trial fMRI study. *Neuron*, 20, 937–45.
- Lovibond, P.F., Shanks, D.R. (2002). The role of awareness in Pavlovian conditioning: empirical evidence and theoretical implications. *Journal of Experimental Psychology Animal Behavior Processes*, 28(1), 3–26.
- Morris, J.S., Buchel, C., Dolan, R.J. (2001). Parallel neural responses in amygdala subregions and sensory cortex during implicit fear conditioning. *Neuroimage*, 13, 1044–52.
- O'Doherty, J.P., Buchanan, T.W., Seymour, B., Dolan, R.J. (2006). Predictive neural coding of reward preference involves dissociable responses in human ventral midbrain and ventral striatum. *Neuron*, 49, 157–66.
- O'Doherty, J.P., Dayan, P., Friston, K., Critchley, H., Dolan, R.J. (2003). Temporal difference models and reward-related learning in the human brain. *Neuron*, 38, 329–37.
- Olsson, A., Ebert, J.P., Banaji, M.R., Phelps, E.A. (2005). The role of social groups in the persistence of learned fear. *Science*, 309, 785–7.
- Petrovic, P., Carlsson, K., Petersson, K.M., Hansson, P., Ingvar, M. (2004). Context-dependent deactivation of the amygdala during pain. *Journal of Cognitive Neuroscience*, 16, 1289–301.
- Peyron, R., Laurent, B., Garcia-Larrea, L. (2000). Functional imaging of brain responses to pain. A review and meta-analysis. *Neurophysiologie Clinique*, 30, 263–88.
- Phelps, E.A. (2006). Emotion and cognition: insights from studies of the human amygdala. *Annual Review of Psychology*, 57, 27–53.
- Phelps, E.A., O'Connor, K.J., Gatenby, J.C., Gore, J.C., Grillon, C., Davis, M. (2001). Activation of the left amygdala to a cognitive representation of fear. *Nature Neuroscience*, 4, 437–41.
- Rescorla, R.A., Wagner, A.R. (1972). A theory of Pavlovian conditioning: variations in the effectiveness of reinforcement and nonreinforcement. In: Black, A.A., Prokasy, W.F., editors. *Classical Conditioning II: Current Research and Theory*. New York: Appleton Century Crofts, pp. 64–99.
- Seymour, B., O'Doherty, J.P., Koltzenburg, M., et al. (2005). Opponent appetitive-aversive neural processes underlie predictive learning of pain relief. *Nature Neuroscience*, 8, 1234–40.
- Seymour, B., O'Doherty, J.P., Dayan, P., et al. (2004). Temporal difference models describe higher-order learning in humans. *Nature*, 429, 664–7.
- Shi, C., Davis, M. (1999). Pain pathways involved in fear conditioning measured with fear-potentiated startle: lesion studies. *Journal of Neuroscience*, 19, 420–30.
- Singer, T., Seymour, B., O'Doherty, J.P., Stephan, K.E., Dolan, R.J., Frith, C.D. (2006). Empathic neural responses are modulated by the perceived fairness of others. *Nature*, 439, 466–9.
- Singer, T., Kiebel, S.J., Winston, J.S., Dolan, R.J., Frith, C.D. (2004). Brain responses to the acquired moral status of faces. *Neuron*, 41(4), 653–62.
- Straube, T., Weiss, T., Mentzel, H.J., Miltner, W.H. (2007). Time course of amygdala activation during aversive conditioning depends on attention. *Neuroimage*, 34, 462–9.
- Sutton, R.S., Barto, A.G. (1998). *Reinforcement Learning - an Introduction*. Cambridge, Massachusetts: MIT Press.
- Tabbert, K., Stark, R., Kirsch, P., Vaitl, D. (2005). Hemodynamic responses of the amygdala, the orbitofrontal cortex and the visual cortex during a fear conditioning paradigm. *International Journal of Psychophysiology*, 57, 15–23.

- Tabbert, K., Stark, R., Kirsch, P., Vaitl, D. (2006). Dissociation of neural responses and skin conductance reactions during fear conditioning with and without awareness of stimulus contingencies. *Neuroimage*, 32, 761–70.
- Tranel, D. (2000). Electrodermal activity in cognitive neuroscience. In: Lane, R.D., Nadel, L., editors. *Cognitive Neuroscience of Emotion*. New York: Oxford University Press, pp. 192–224.
- Tranel, D., Damasio, H. (1989). Intact electrodermal skin conductance responses after bilateral amygdala damage. *Neuropsychologia*, 27, 381–90.
- Vuilleumier, P., Pourtois, G. (2006). Distributed and interactive brain mechanisms during emotion face perception: Evidence from functional neuroimaging. *Neuropsychologia*, 45, 174–94.
- Vuilleumier, P., Armony, J.L., Driver, J., Dolan, R.J. (2001). Effects of attention and emotion on face processing in the human brain: an event-related fMRI study. *Neuron*, 30, 829–41.
- Vuilleumier, P., Richardson, M.P., Armony, J.L., Driver, J., Dolan, R.J. (2004). Distant influences of amygdala lesion on visual cortical activation during emotional face processing. *Nature Neuroscience*, 7, 1271–8.