

**Learning new words via feedback – association between feedback-locked ERPs and recall
performance**

Christine Albrecht^{1*}, Ruben van de Vijver², and Christian Bellebaum¹

¹ Institute of Experimental Psychology, Heinrich Heine University Düsseldorf, Germany

² Institute of Linguistics and Information Science, Heinrich Heine University Düsseldorf,
Düsseldorf, Germany

* corresponding author

Christine Albrecht

Institute of Experimental Psychology

Heinrich Heine University Düsseldorf

Universitätsstraße 1

building 23.03

room number 00.89

40225 Düsseldorf

Germany

christine.albrecht@hhu.de

Abstract

Feedback learning is thought to involve the dopamine system and its projection sites in the basal ganglia and anterior cingulate cortex (ACC), regions associated with procedural learning. Under certain conditions such as when feedback is delayed, activation shifts to the medial temporal lobe (MTL), which is associated with declarative learning. Feedback elicits the feedback-related negativity (FRN), an event-related potential component which originates in the ACC and is pronounced for immediate feedback. As a potential correlate of MTL activity, the N170 has been shown to be related to delayed feedback processing. In the current study, we investigated whether N170 amplitude predicts memory performance in a test for declarative memory (free recall), especially if feedback is delayed. To test this, we adapted a paradigm by Arbel et al. (2017) in which participants learned associations between non-objects and non-words with either immediate or delayed feedback, and added a subsequent free recall test. We indeed found that N170, but not FRN amplitude depended on later free recall performance, with smaller amplitudes for remembered words. This effect was further modulated by feedback valence, feedback delay and laterality, with particularly low left-hemispheric amplitudes for delayed correct feedback. This finding shows that the N170 reflects an important process especially during delayed feedback processing that is related to expectations and their violation, but is distinct from the process reflected by the FRN.

Keywords. Feedback Learning, Memory, N170, FRN, Feedback Timing

Introduction

Learning from feedback is crucial in everyday life, whether in simple tasks such as finding the right knob to turn on the stove, or in much more complex tasks, such as when a gymnast learns to perform a flip. Research points to an involvement of the dopamine system and in particular the basal ganglia in feedback learning: Evidence from functional neuroimaging studies in humans suggests, for example, that the basal ganglia represent a reward prediction error during feedback learning (Dobryakova & Tricomi, 2013; Foerde & Shohamy, 2011b; Lighthall et al., 2018; O'Doherty et al., 2004), which is coded by single dopamine neurons in the midbrain (Hollerman & Schultz, 1998; Schultz et al., 1997; Zaghoul et al., 2009). Accordingly, patient studies have shown that participants with Parkinson's disease (PD) display altered or impaired feedback learning (Foerde & Shohamy, 2011b; Frank et al., 2004; see Foerde & Shohamy, 2011a). PD involves a degeneration of dopamine neurons and severely impairs basal ganglia functionality.

To study the temporal dynamics of feedback processing in humans, electroencephalography (EEG) has frequently been used. In particular, an event-related-potential (ERP) component between 200 ms and 350 ms following feedback has been linked to feedback processing. As it occurs as a pronounced negativity in the signal it has been referred to as feedback-(related) negativity (FN or FRN; Hajcak et al., 2006; Miltner et al., 1997; Yeung et al., 2005). It is usually more pronounced for negative feedback (Gehring & Willoughby, 2002; Miltner et al., 1997; Nieuwenhuis et al., 2004), and the term FRN has by some researchers been used to refer to the difference wave, subtracting the ERP signal for positive feedback from that for negative feedback (Hajcak et al., 2007; Holroyd et al., 2009).

It has recently been pointed out, however, that the signal can be better described as reward positivity, which emerges when the signal subtraction is performed in the opposite direction (see Proudfit, 2015). In one of the studies in our lab we have distinguished between the FRN_{diff} , derived from the difference wave, and the FRN_{peak} , referring to the negativity in the original waveform (Peterburs et al., 2016). In this study we will use the term FRN in the latter sense, i.e. for peak amplitudes in the original ERPs instead of difference waves.

The FRN is believed to originate from the anterior cingulate cortex (ACC; Gehring & Willoughby, 2002; Nieuwenhuis et al., 2004), which in turn receives input from the dopaminergic system (see Chau et al., 2018). With functional magnetic resonance imaging (fMRI) measures it has been shown that indeed ACC activity is sensitive to feedback valence with higher activation for negative feedback (Holroyd et al., 2004). When fMRI and EEG measures were acquired simultaneously, feedback-locked ERPs, especially for positive feedback, were found to be associated with basal ganglia activity (Becker et al., 2014; Holroyd et al., 2004). This further strengthens the link between the FRN and the dopaminergic and reward system of the brain. Corroborating this link, the ERP amplitude in the FRN time window has been shown to reflect a reward prediction error (i.e. the degree to which an outcome is better or worse than expected; Burnside et al., 2019; Fischer & Ullsperger, 2013).

The degree of involvement of the basal ganglia during learning seems to affect the representation of what has been learned. Material learned via the basal ganglia tends to be acquired procedurally (see Gasbarri et al., 2014; Shohamy et al., 2008; Yin & Knowlton, 2006), meaning that the information cannot be used flexibly and is thus difficult to transfer to new situations (Myers et al., 2003; see Squire, 2004). Accordingly, FRN amplitude has been linked to error correction and adjustment of behavior acquired by means of reinforcement learning (Cohen

& Ranganath, 2007; de Bruijn et al., 2020), but not adjustment of behavior based on explicit rule-based knowledge (Chase et al., 2011). This supports the notion that the FRN reflects procedural rather than declarative learning.

While these studies suggest a key role of the basal ganglia and the procedural learning system in feedback learning, evidence suggests that under certain conditions also the medial temporal lobe (MTL), including the hippocampus, and thus the declarative memory system (Squire & Zola-Morgan, 2015), can be involved (Dickerson & Delgado, 2015; Foerde et al., 2013; Foerde & Shohamy, 2011a). During feedback learning, the relative involvement of declarative and non-declarative processes (and the respective brain regions) seems to be affected by the temporal proximity between feedback and the event it is related to: when feedback is delayed, even only by a few seconds, MTL (hippocampal) activity increases while striatal activity decreases, indicating a shift from procedural to declarative learning (Foerde & Shohamy, 2011a). Patient studies revealed a double dissociation with respect to the neural structures underlying learning from immediate and delayed feedback: PD patients are impaired when learning from immediate feedback, but not when feedback is delayed by 6 s (Foerde et al., 2013; Foerde & Shohamy, 2011a; Weismüller et al., 2018); amnesic patients with suspected MTL damage show the reverse pattern (Foerde et al., 2013).

Recent studies have shown that the differences in feedback processing depending on feedback delay are also evident in feedback-locked ERP components. In accordance with its link to basal ganglia function, the FRN is more strongly involved in immediate than delayed feedback processing, as the difference between positive and negative feedback in the FRN time window is more pronounced for immediate feedback (Arbel et al., 2017; Peterburs et al., 2016; Weinberg et al., 2012; Weismüller & Bellebaum, 2016). An ERP component possibly reflecting more

let participants learn associations between novel objects and novel words. Tasks involving deterministic feedback could be more suitable for declarative learning than tasks with probabilistic feedback. Despite the comparable learning performance, we expect differences in memory characteristics between material learned with different feedback timings: As learning with delayed feedback is based on the MTL, the acquired knowledge should be more declarative and flexible (Fera et al., 2014; Shohamy & Wagner, 2008). When memory is probed, this may be reflected in the performance in a specific type of test. MTL activity during encoding in declarative learning tasks has been associated with performance in a free recall test (Danckert & Craik, 2013; Leshikar et al., 2017; Staresina & Davachi, 2006), but not with performance in a recognition test as, for example, used in the deterministic feedback learning study involving words by Arbel et al. (2017). In free recall memory tests, participants are asked to freely reproduce specific information (like a previously studied word), while in recognition memory tests participants are asked to identify the correct answer from two or more possible alternatives (like selecting the studied word from two presented words).

In the present study, we attempt to establish a link between feedback delay effects on neural feedback processing and on subsequent memory by adapting the paradigm used by Arbel et al. (2017). As mentioned above, participants in this task were asked to learn associations between novel objects and novel words with either immediate or delayed feedback, and EEG recordings during learning allowed to examine feedback-locked ERPs. We extended the study by Arbel et al. (2017) by adding a free recall test after learning in order to assess a more declarative, MTL-based representation of the learned material. With this we aimed to investigate feedback-learning-related brain activity as a function of subsequent memory performance, with a potential role of feedback timing.

Due to its stronger involvement in more declarative types of learning and its potential link to MTL processing (Arbel et al., 2017; Baker & Holroyd, 2013), we expected the N170 amplitude, rather than the FRN amplitude, to be related to free recall memory performance. This effect should be stronger for delayed than immediate feedback.

Method

Participants

Thirty-two participants took part in the experiment. Two participants were excluded due to artefacts in the EEG data or missing data. The final analysis thus included 30 participants, 16 of which identified as women, 14 as men. The participants' age ranged between 19 and 30 years ($M = 23.4$, $SD = 2.6$). All of the 30 participants entering the analysis were right-handed, reported no previous neurological or psychiatric illnesses, did not take regular medication affecting the central nervous system, and spoke German as a native language. Participants took part voluntarily and were reimbursed with 20€ or course credit for participation. The study was approved by the ethics committee of the Faculty of Mathematics and Natural Sciences at Heinrich Heine University Düsseldorf, Germany.

Stimuli

Participants completed an experimental task in which they learned associations between 56 novel objects, referred to as non-objects, and 56 corresponding novel words, referred to as non-words. 56 additional non-words were used as distractors. As in Arbel et al. (2017), the non-objects were adopted from Kroll and Potter (1984).

The non-words used by Arbel et al. (2017) were designed for English-speaking participants. Since our sample consisted of German native speakers, we created new non-words by retrieving the 180 most frequent (sorted by Mannheim Frequency) German monosyllabic

nouns consisting of at least 3 letters from the Celex Online Database (Baayen et al., 1995). When retrieving the words from the database, we ignored duplicate values (some words occurred twice in the list because of other grammatical forms and upper- and lower case) and given names (e.g. the names of cities), so we selected 180 distinct, regular nouns in their basic grammatical form. We then used the software Wuggy (Keuleers & Brysbaert, 2010) to create a maximum of three subsyllabic (syllable fragments were exchanged) non-words for each of the previously found words. The software could not create non-words for 39 of the selected words, reducing the number of words to 141. As the software could not create three non-words for two of the remaining words, 421 non-words were created for the 141 real words in total. In a subsequent step, we deleted non-words that were created more than once: one word was excluded if two identical non-words were created for the same real word or if the non-word was identical to the real word (not considering upper and lower case); both words were excluded if two identical non-words were created for two different real words. As non-words created from one real word were fairly similar, we deleted both double values to reduce similarity (a non-word created for multiple real words would be similar to other non-words created for all these real words). The remaining 374 non-words (based on 139 real words) were then checked for their existence in the German language: the non-word was deleted if there was an entry in the Duden online dictionary (Dudenredaktion, 2020) for the non-word, again not considering upper and lower case, or if the non-word corresponded to a declination or conjugation of an existing word listed in the Duden (e.g. the non-word “lies” was removed because it is a conjugation of the German word “lesen”).

After this step, 297 non-words remained (based on 135 real words), which were rated for *ease of pronunciation*, *similarity to specific German words* and *similarity to specific English words* on 7-point Likert scales by 19 independent raters (14 identified as women, 5 as men) aged

between 18 and 45 years ($M = 24.9$, $SD = 8.1$). We excluded all non-words that were rated as similar to a specific German (3 non-words) or English word (7 non-words) with an average of 5.5 or higher. We then selected the non-word with the highest scores for *ease of pronunciation* for each real word, leaving 135 non-words. Of those, we selected the 112 non-words for which *ease of pronunciation* was rated highest. The *ease of pronunciation* of the remaining 112 non-words was rated between 4.6 and 6.4, with an average of 5.5 ($SD = 0.5$). Please find a list of all non-words in Appendix A.

For each participant, 56 pairs of non-words were randomly created from the 112 non-words, and each pair was randomly assigned to one non-object. One of the words of each pair was randomly defined as the correct name of the non-word during the first learning block of our experimental learning task, the other word of the pair was defined as distractor (see below). The non-object-non-word-combinations were randomly sorted into four sets (14 combinations per set).

Experimental Task

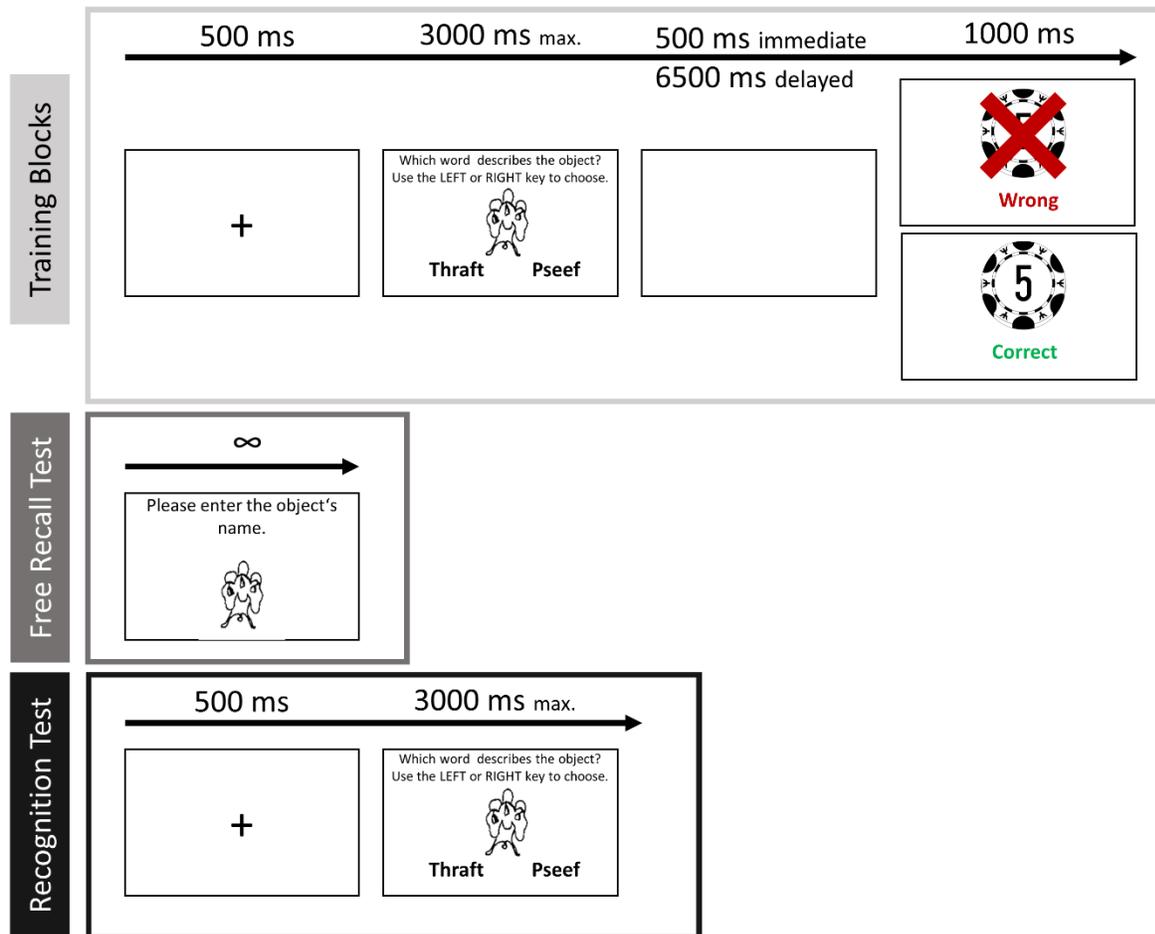
The experiment consisted of four sequences. In each sequence, one of the four sets of combinations was used. Each sequence began with a training phase that consisted of 5 blocks in each of which every combination of non-object and non-words was presented once in an experimental trial (see below for the sequence and timing of events in one trial). The participants were asked to try to learn which non-word was associated with the displayed non-object. To achieve this, they were instructed to select one of the two non-words by pressing a corresponding key, whereupon they received positive or negative feedback. In two sequences delayed feedback was used and in two other sequences immediate feedback was used, in alternating order. To avoid sequence effects, half of the participants started with immediate feedback and the other

half with delayed feedback. While this paradigm was adopted from Arbel et al. (2017), we attempted to improve motivation by adding gamification elements. This included telling participants that they would receive casino chips for each correct answer and that they could try to beat a (fictitious) high score. Participants received five chips for every correct answer during the training phase.

Trials in the training phase began with a fixation cross, shown for 500 ms. Subsequently, one non-object was presented in the center of the screen with two non-words below it. To remind participants of the task, the display included an instruction: “Which word describes the object? Use the RIGHT or LEFT key to choose.” The display ended when the participant pressed a button or after 3000 ms, and a white screen was shown for either 500 ms (sequences with immediate feedback) or 6500 ms (sequences with delayed feedback) before feedback was presented. Feedback could be negative (“Wrong” in red font; a crossed-out casino-chip was displayed), positive (“Correct” in green font; a casino-chip was displayed) or indicating a late response, when participants had failed to answer within 3000 ms (“Too slow!” in black font; a snail was displayed). For an illustration of the training trials, see Figure 1.

Figure 1

Trial structure during the three parts of a sequence



Note. The experiment consisted of four sequences that each contained five training blocks, followed by a free recall and a recognition test. The upper row shows the sequence of events in one training trial, the middle and bottom rows show the procedure for the recall and the recognition tests (refer to the text for details). If the participants did not answer in a 3000 ms time window in the trials of the training blocks and the recognition test, they were informed that they answered too slowly. In the free recall test, there was no time limit.

As in Arbel et al. (2017), all participants completed five training blocks in each sequence. In the first block, feedback was counterbalanced, so that there was an equal number of positive and negative feedback. That means that the assignment of the non-object to its associated name

(non-word) was pseudo-randomly determined in the first training block. If participants answered too slowly, it was randomly determined which word was correct. In the subsequent blocks the feedback corresponded to the feedback given in the first block. Feedback was deterministic in the sense that the same choices for a given non-object were always followed by the same feedback. Each block consisted of 14 trials, so that all 14 combinations of a set appeared once. In addition to Arbel et al.'s setting, the current score was presented to the participants at the end of each block and a progress bar was displayed.

After the training phase, a free recall test was conducted. Participants were instructed that the non-objects that were presented during the training sequences would be displayed and they were asked to recall the associated correct non-word and enter it using the keyboard. Participants were told that they would receive 20 chips for each correct answer and that no feedback would be given during the free recall test, but that their total score would be displayed at the end of the test. Trials in the free recall test consisted of a display of a non-object, above which the question appeared "Do you remember the word that describes this object?". The participant's answer was displayed below the object. There was no time limit. Participants ended their entry by pressing the Enter key. Each non-object of the current set appeared once in the free recall test, resulting in 14 trials.

We then conducted a recognition test, which was identical to the one used in Arbel et al. (2017). As in the training phase, participants were shown an object and two possible names, from which they were to select the correct one. For each non-object the same two non-words were presented as during the training trials. Participants were informed that they would not receive feedback in this task. As part of our gamification, participants were told that they would receive

10 chips for each correct answer and that they would get to know their total score at the end of the test. Each combination appeared once, resulting in 14 recognition trials.

As mentioned above, participants completed four sequences consisting of training phase (5 learning blocks), free recall test and recognition test. Participants had the opportunity to take a break at any time between blocks and tests and between sequences, as they paced the instructions themselves. At the end of the experiment, participants were shown a leader board in which their score was compared to nine fictitious other scores.

EEG Recording

We applied passive scalp electrodes according to the international 10-20 system. Electrodes were attached to the scalp sites F7, F3, Fz, F4, F8, FT7, FC3, FCz, FC4, FT8, T7, C3, Cz, C4, T8, CP3, CPz, CP4, P7, P3, Pz, P4, P8, PO7, PO3, POz, PO4, PO8. We also used four eye electrodes, two for vertical electro-oculography (EOG; FP2 and an electrode below the right eye) and two for horizontal EOG (F9 and F10) for 13 of the participants. During data acquisition, we decided to add the electrodes TP7 and TP8 and reduce the eye electrodes to two (FP2 and F10, respectively) so that P7 and P8 could be recreated (as a mean of TP7 and PO7 or TP8 and PO8, respectively) if the signal was too noisy. However, only the 28 electrodes that were used on all participants were used in the analysis. We used a BrainAmp Standard amplifier (Brain Products, Munich, Germany) and the software BrainVision Recorder (version 1.20.0506, Brain Products, Munich, Germany) to record EEG data during the experiment at a 1000 Hz sampling rate. Electrodes were online-referenced to the average of two mastoid electrodes. All impedances were kept below 5 k Ω .

Procedure

Upon arrival at the laboratory, participants were informed about the experimental procedure and gave written consent to participate in the study. They also completed a demographic questionnaire. Afterwards, the electrodes for the EEG measurement were attached before we began the computer experiment. The experimental stimuli were presented on a 1920 * 1080 px desktop monitor. In total, the computer experiment took between 45 minutes and 1 hour. Together with the preparation of the EEG recording the procedure lasted about 2 hours. The experiment was controlled by Presentation Software (Version 20.0, Neurobehavioral Systems, Albany, CA, USA). After completing the experiment, participants received compensation for their participation in the form of money or course credit.

Data analyses

Behavioral data

We calculated the percentage of correct answers for the free recall test, the recognition test and during the training blocks as dependent variables. This was done separately for each participant for immediate and delayed feedback, and thus averaged across the two sequences with identical feedback timing. For the training blocks, we also calculated values for every block separately. We used IBM SPSS Statistics 25 (IBM Corporation, Armonk, New York, USA) for statistical analyses. A paired-samples t-Test was used to compare performance for associations learned with delayed feedback with performance for associations learned with immediate feedback (factor Feedback Timing) in the free recall test. For the replication of results obtained by Arbel et al. (2017), we additionally computed a paired-samples t-Test to compare performance for the different feedback timings in the recognition test. Since both test measures reflect memory performance and are thus not independent, we did not directly compare the performance in the free recall with performance in the recognition test. Again for the replication

of Arbel et al., we also calculated a 2x5 repeated measures ANOVA to compare performance depending on Feedback Timing across Training Blocks.

EEG Data

EEG preprocessing. BrainVision Analyzer software (version 2.1; Brain Products, Munich, Germany) was used for EEG data preprocessing. In a first step, EEG data were rereferenced to the average signal across all scalp electrodes (only those 28 that were used for all participants). This is a standard procedure for analyzing the N170 component, as this component is measured close to the mastoids and a mastoid reference might cancel out possible N170 effects (Rellecke et al., 2013; Wang et al., 2019; for similar procedures on corresponding tasks, see Arbel et al., 2017; Hölting & Mecklinger, 2020). At the same time, average references are also not uncommon in studies assessing the FRN (e.g. Becker et al., 2014; Chase et al., 2011; Fischer & Ullsperger, 2013). We then applied a 20 Hz low-pass and a 0.5 Hz high-pass filter. Blink artefacts were removed as follows: an independent component analysis was performed on the filtered EEG data. Of the resulting components, one displaying a frontocentral maximum and corresponding to the blinks recorded in the vertical EOG was flagged and removed in an independent component analysis back-transformation. Epochs of 800 ms (200 ms before to 600 ms after feedback onset) were created for each of the conditions positive immediate feedback, positive delayed feedback, negative immediate feedback and negative delayed feedback. The segments were then baseline-corrected, the 200 ms before event serving as baseline. An automatic artifact rejection was applied: all segments were removed that contained either voltage steps above 50 $\mu\text{V}/\text{ms}$, an amplitude difference of more than 100 μV between any data points, or any data point with an amplitude higher than 100 μV or lower than -100 μV .

On average, 1.0% ($SD = 2.6\%$) of segments were removed in the correct immediate condition, resulting in an average of 121.1 segments (93 - 155, $SD = 19.5$). In the correct delayed condition, an average of 1.6% ($SD = 3.2\%$) of segments were removed, resulting in 115.2 segments on average (83 - 154, $SD = 19.8$). In the incorrect immediate condition, 1.5% of the segments were removed ($SD = 4.5\%$), resulting in an average of 69.7 segments (41 - 101, $SD = 17.7$). And finally in the incorrect delayed condition, an average of 2.0% of segments were removed ($SD = 4.0\%$), resulting in an average of 73.2 segments (30 - 110, $SD = 20.7$). As participant's performance improved during training (see behavioral results) this resulted in more trials with positive compared to negative feedback. As the last step, averages for all feedback conditions were created.

Components. For extracting the amplitudes for each of the ERP component of interest, we used MATLAB, version R2017b (Mathworks, Natick, Massachusetts, USA). The focus was on two ERP components, namely the FRN and the N170. Please note that the statistical analysis was based on single-trial data. Thus, amplitude values for the FRN and the N170 had to be determined for every single trial. To determine these values, however, we also took the subjects' averages for each condition into account. For the FRN, we first pooled the ERP signal over the frontocentral electrodes Fz, FCz and Cz, where the FRN is most pronounced. We then identified the latency of the maximum negative peak between 200 and 400 ms in the averages for each condition and participant, and the latency of the preceding maximum positive peak between 100 ms and the negative peak. This information about the peak latency was then used to extract single-trial amplitude values: We determined the amplitude values at the latency of the negative peak and at the latency of the positive peak for each segment and exported the difference between these values for further analyses (negative peak - positive peak), as it reflects the single-

trial correlate of the common peak-to-peak measure that is frequently used in average-based analyses (Holroyd et al., 2003; Hölftje & Mecklinger, 2020; Peterburs et al., 2016). A similar procedure was used for the N170 (Hölftje & Mecklinger, 2020), but at different electrode sites: Over the left-hemisphere we pooled the signal over electrodes P7 and PO7 and over the right-hemisphere we pooled over P8 and PO8. The latency of the maximum negative peak was identified in a window between 80 and 220 ms in each subject's average of each condition and separately over the left and right hemisphere, as well as the latency of the preceding positive peak in a window between 30 ms and the latency of the negative peak for each participant and condition. Again, these latencies were used to export single-trial amplitude values. For each trial, the difference between the value at the negative peak and the value at the positive peak was calculated, separately for the signal over the left and the right hemisphere.

Before the statistical analysis we excluded all trials with amplitude values that differed by more than two standard deviations from the mean amplitudes per participant, separately in each condition (Feedback Valence and Feedback Timing) and, for the N170, over the left and right hemisphere. On average, 4.5% of the trials were excluded ($SD = 1.9\%$) with a maximum of 11.1% in one participant.

Statistical analyses. We performed linear mixed effect (LME) analyses to examine the effects of Feedback Timing, Feedback Valence and Memory Performance in Free Recall on both the FRN and N170, using the lme4 package (version 1.1-27.1) in R (version 3.5.3). For the FRN analysis, we defined single-trial peak-to-peak amplitude values as dependent variable. As fixed effects, we included the three within-subject factors Feedback Timing (immediate and delayed, coded as -1 and 1, respectively, in the model), Feedback Valence (negative and positive, coded as -1 and 1) and Free Recall Memory (not-remembered and remembered, coded as -1 and 1). We

included random slopes and intercepts per participant. According to best practice, models should include all within-subjects effects and interactions as random slopes and intercepts per participant, but at the same time, singular fit and convergence failure should be avoided (Meteyard & Davies, 2020), which can be a problem for a large number of random factors. To determine the model with the highest possible complexity, we used an iteration process: We started with the lowest possible complexity, so only random intercepts for each participant, and then added random effects (interaction and main effects) gradually, testing in each case whether their addition resulted in a singular fit or convergence failure. This resulted in a model that used random slopes per participant for Feedback Timing, Feedback Valence and the interaction between Feedback Timing and Feedback Valence.

We used the same model for the analysis of N170 single-trial peak-to-peak amplitudes as dependent variable, with the additional fixed factor Laterality (left and right hemisphere, coded as -1 and 1, respectively). Random slopes were again determined iteratively, resulting in a model including random slopes per participant for Feedback Timing, Feedback Valence, Laterality and all two-way interactions between these factors. Note that no model including random effects for the main or interaction effects of Free Recall Performance could be considered a valid model for FRN or N170; this is probably due to the already high variability in Free Recall Performance between participants in the data.

We decided to not include the Training Block as factor in the EEG analyses, because there were generally very few trials with negative feedback in the later blocks, in particular trials with negative feedback for which the non-word associated to the non-object was correctly recalled (see Table S3).

Outlier detection. For each of the two models, we performed Cook's distance outlier detection using the R-package stats (version 3.6.2). Cook's distance calculates how much each position of a given level (in our case, we calculated Cook's distance on the subject level) affects the values of the model. If the influence estimate is larger than a cut-off value of $4 / (n - \text{factors} - 1)$, the data point is excluded from further analysis. For 30 participants and 3 or 4 factors for the two analyses, respectively, our cutoff values were .15 for the FRN and .16 for the N170.

We excluded one participant from the FRN analysis, resulting in 29 participants (19 – 30 years old, $M = 23.4$, $SD = 2.7$), of which 16 identified as women, 13 as men. For the N170 analysis, no participant was excluded.

Interactions. Significant interactions were resolved in a step-wise manner: Conditional slopes were calculated (the slope of a specific effect when one predictor was held constant) with -1 or 1 as constants, according to the coding of the respective variable, to test for significant lower-level interaction effects. If these were found, this procedure was repeated for this lower-level interaction until all factors were resolved.

Results

Behavioral Data

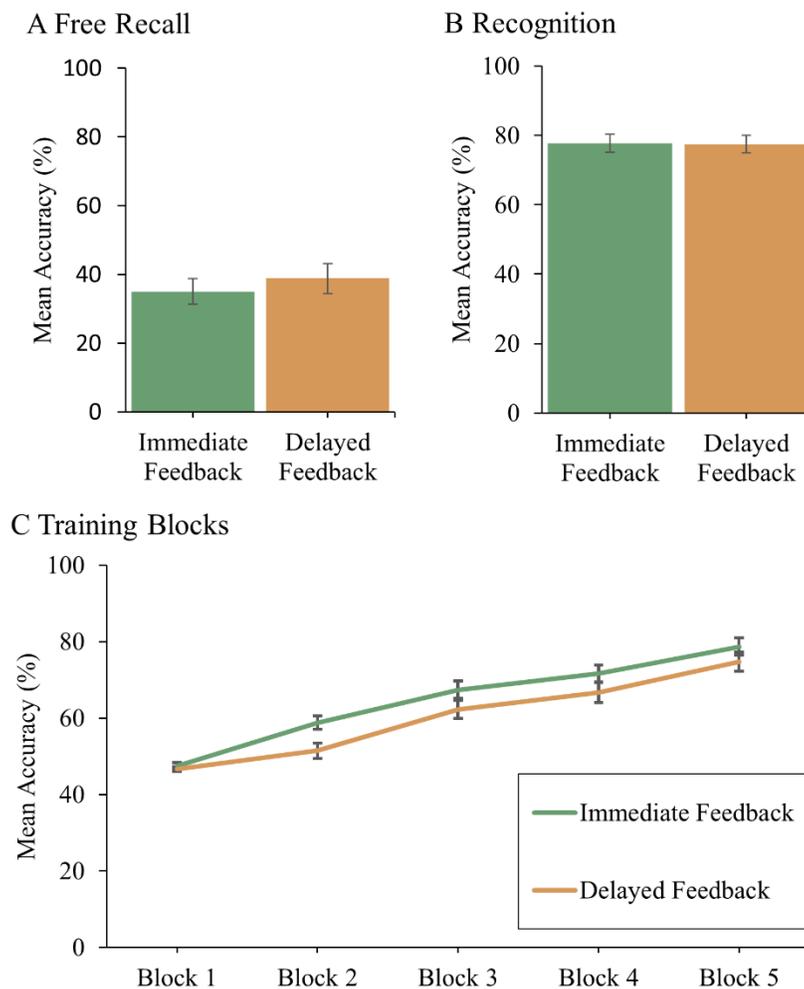
We found no significant difference in memory performance between associations learned with immediate and delayed feedback, neither in free recall performance, $t(29) = -1.14$, $p = .265$, $d = 0.21$, nor in recognition performance, $t(29) = 0.11$, $p = .910$, $d = 0.02$. Descriptive results for these two measures are displayed in Figure 2a and b.

For the analysis of performance during the training blocks we found a significant main effect of Feedback Timing, $F(1,29) = 11.42$, $p = .002$, $\eta_p^2 = .28$ with better performance with immediate ($M = 64.81\%$, $SD = 7.86\%$) than delayed feedback ($M = 60.36\%$, $SD = 8.13\%$). We also found a

significant effect of block, $F(2.56, 74.16) = 88.74$, $p < .001$, $\eta_p^2 = .28$. Performance accuracy increased significantly with every block ($p \leq .035$ for all pairwise comparisons). There was no interaction effect ($p = .251$). Descriptive results for performance in the training blocks can be seen in Figure 2c. Note that accuracy in the first blocks is slightly lower than 50% because of too slow answers which were counted as incorrect.

Figure 2

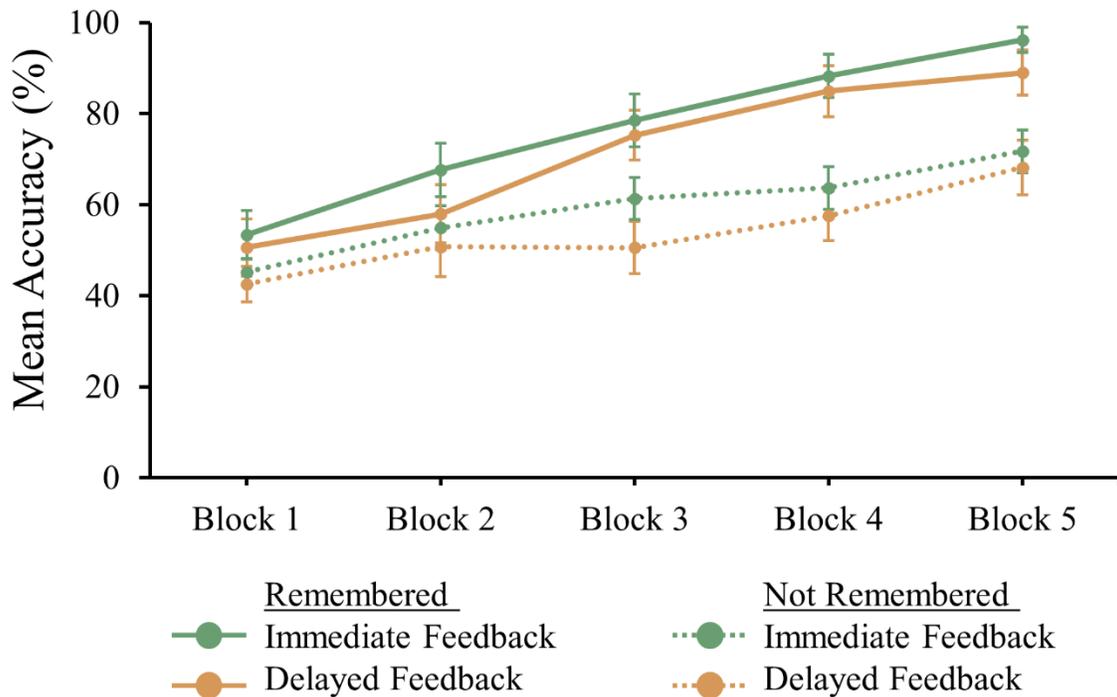
Means and standard errors of Memory Performance in Free Recall, Recognition and Training Blocks



To investigate the relationship between performance in the training blocks and subsequent memory performance, we additionally investigated potential differences in choice accuracy during the training blocks between later recalled and not recalled associations. In addition to the aforementioned effects (see LME analysis on training block performance), we found a main effect of Free Recall Performance, $F(1,32.72) = 107.80, p < .001, b = 8.73$, indicating that accuracy was significantly higher already in the training blocks for associations that were later remembered in the free recall test. Furthermore, we found a two-way-interaction between Free Recall Performance and Block, $F(1,524.00) = 37.90, p < .001$. While a significant effect of Free Recall Performance emerged for both early, $F(1,61.61) = 31.57, p < .001, b = 5.55$, and late blocks, $F(1,61.61) = 145.66, p < .001, b = 11.91$, the effect for late blocks was much larger ($b = 11.91$ as opposed to $b = 5.55$). We found no other interaction effects involving Free Recall Performance (all $p \geq .925$). For a display of the descriptive data (means and standard errors) underlying this analysis see Figure 3.

Figure 3

Means and standard errors for performance during training blocks, as a function of Block, Feedback Timing and Free Recall Performance

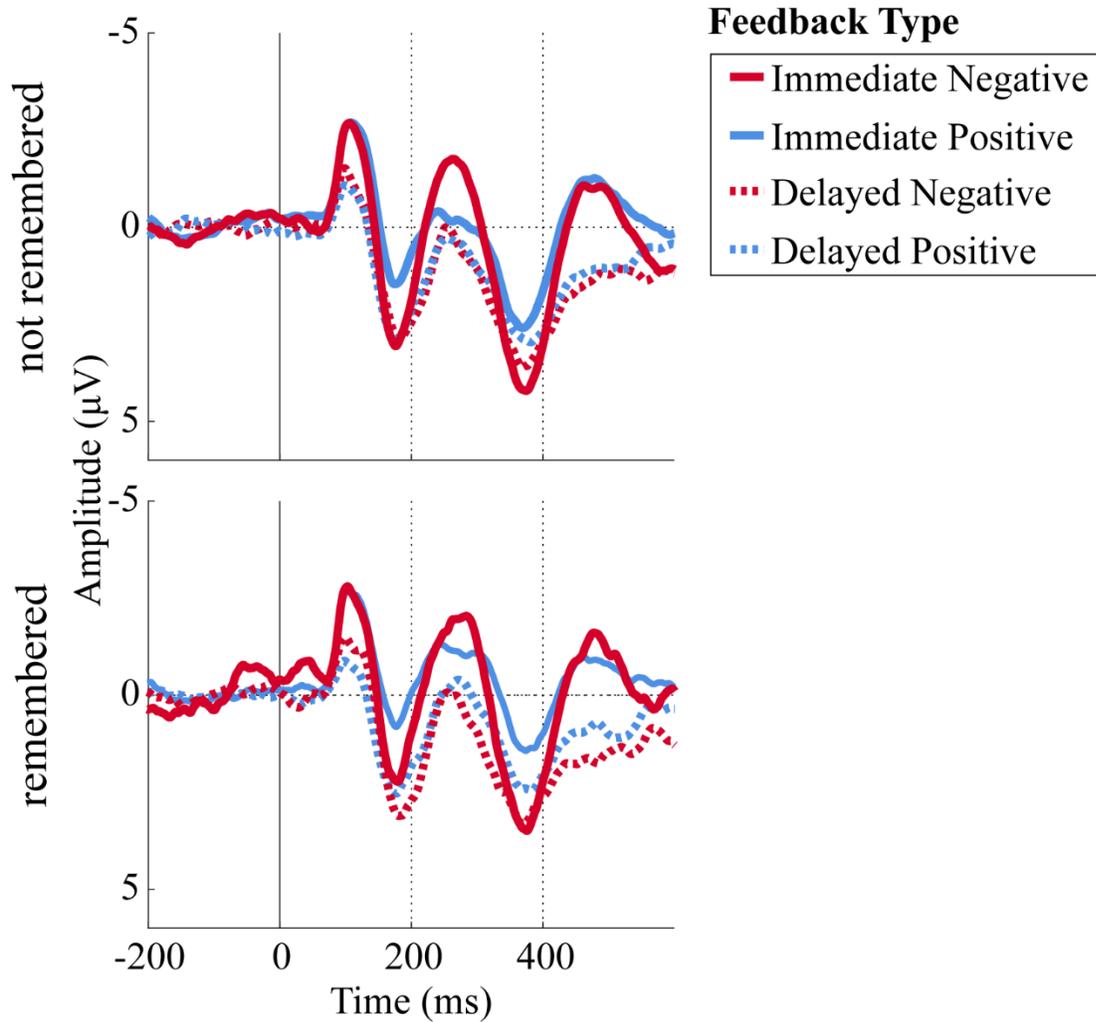


EEG Data

FRN. Please find a display of the feedback-locked Grand Average ERPs depending on Feedback Timing, Feedback Valence and Memory Performance, pooled over electrode sites Fz, FCz and Cz, in Figure 3. These ERPs show the grand average FRN in the different conditions, means and standard errors of its amplitude are shown in Figure 4. All statistical parameters can be found in Table S1 in the supplementary material.

Figure 3

Grand Average ERPs pooled over Fz, FCz and Cz



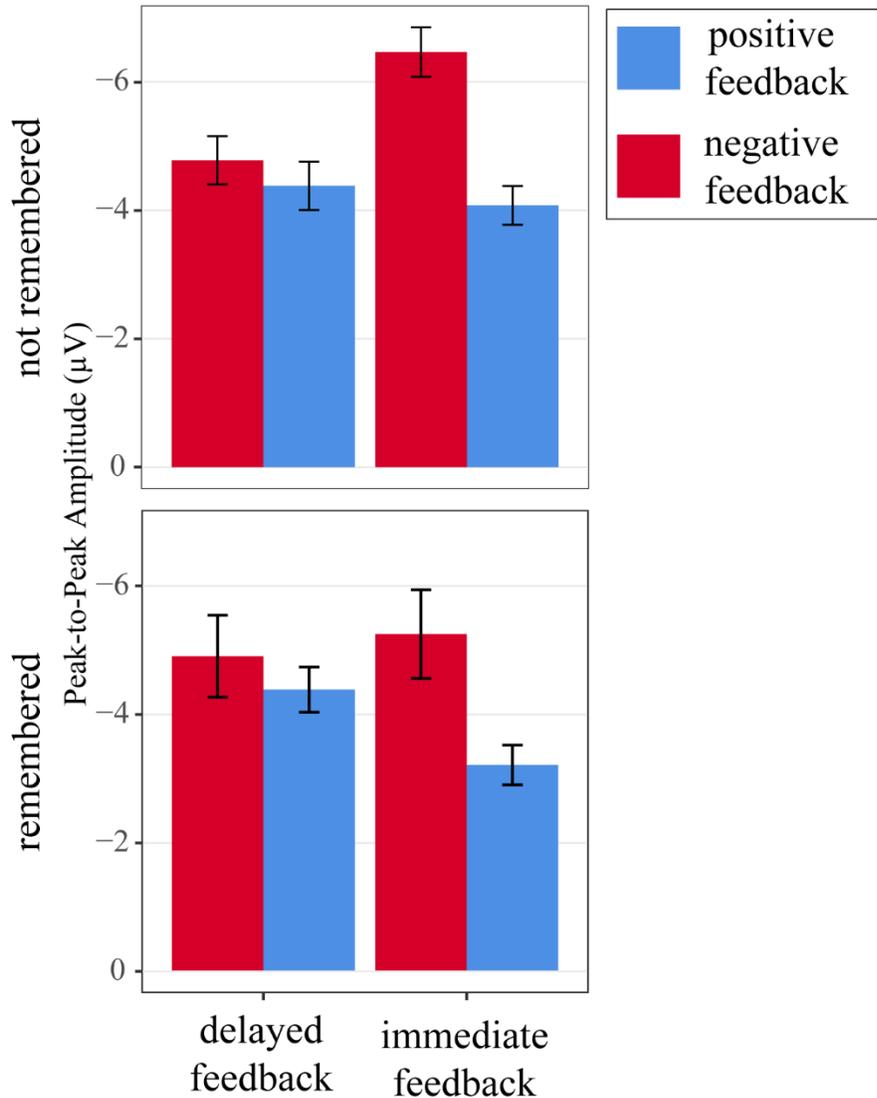
Notes. The dotted vertical lines represent the start and end point for the maximum negative peak detection for the FRN amplitude.

The model revealed a main effect of Feedback Valence, $F(1,33.60) = 43.87, p < .001, b = 0.70$, with larger amplitudes for negative than positive feedback. No other main effects emerged (all $p \geq .537$). The analysis also revealed an interaction effect between Feedback Timing and Feedback Valence, $F(1,31.90) = 10.47, p = .003$. Resolving this interaction by Feedback Timing, we found a main effect of Feedback Valence for immediate feedback, $F(1,32.40) = 34.12, p < .001, b = 1.14$, with larger amplitudes for negative than positive feedback. No such effect

emerged for delayed feedback ($p = .078$). The main effect of the factor Free Recall Performance as well as all other interaction effects were not significant (all $p \geq .135$).

Figure 4

Means and standard errors of the FRN amplitude depending on Feedback Valence, Feedback Timing and Free Recall Performance.

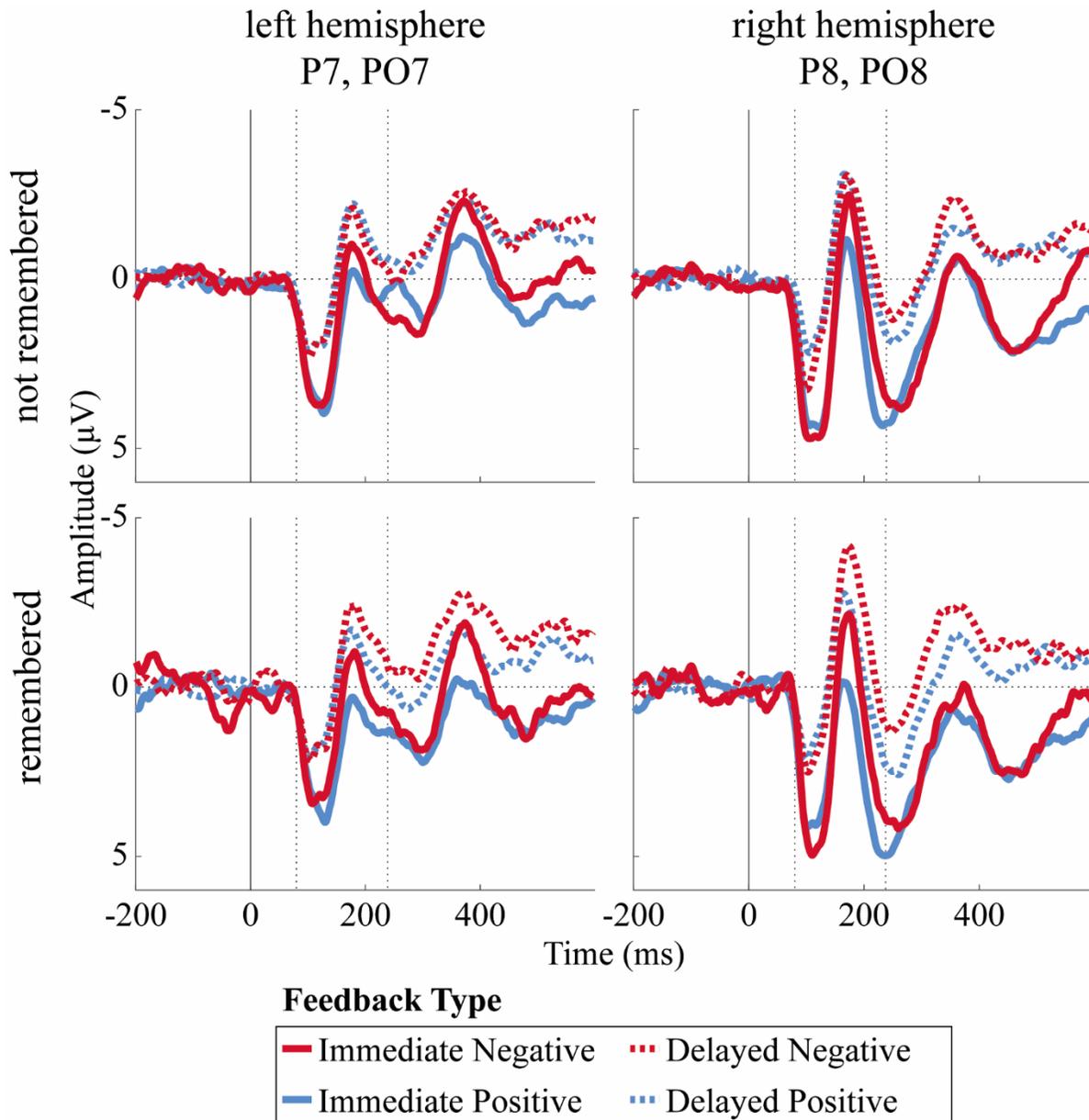


N170. A display of the Grand Average ERPs depending on Feedback Timing, Feedback Valence, Memory Performance and Laterality, pooled over P7 and PO7 and P8 and PO8,

respectively, can be found in Figure 5. These ERPs show the grand average N170 in the different conditions, means and standard errors of its amplitude are shown in Figure 6. Statistical parameters are documented in Table S2 in the supplementary material.

Figure 5

Grand Average ERPs pooled over P7 and PO7 and over P8 and PO8.

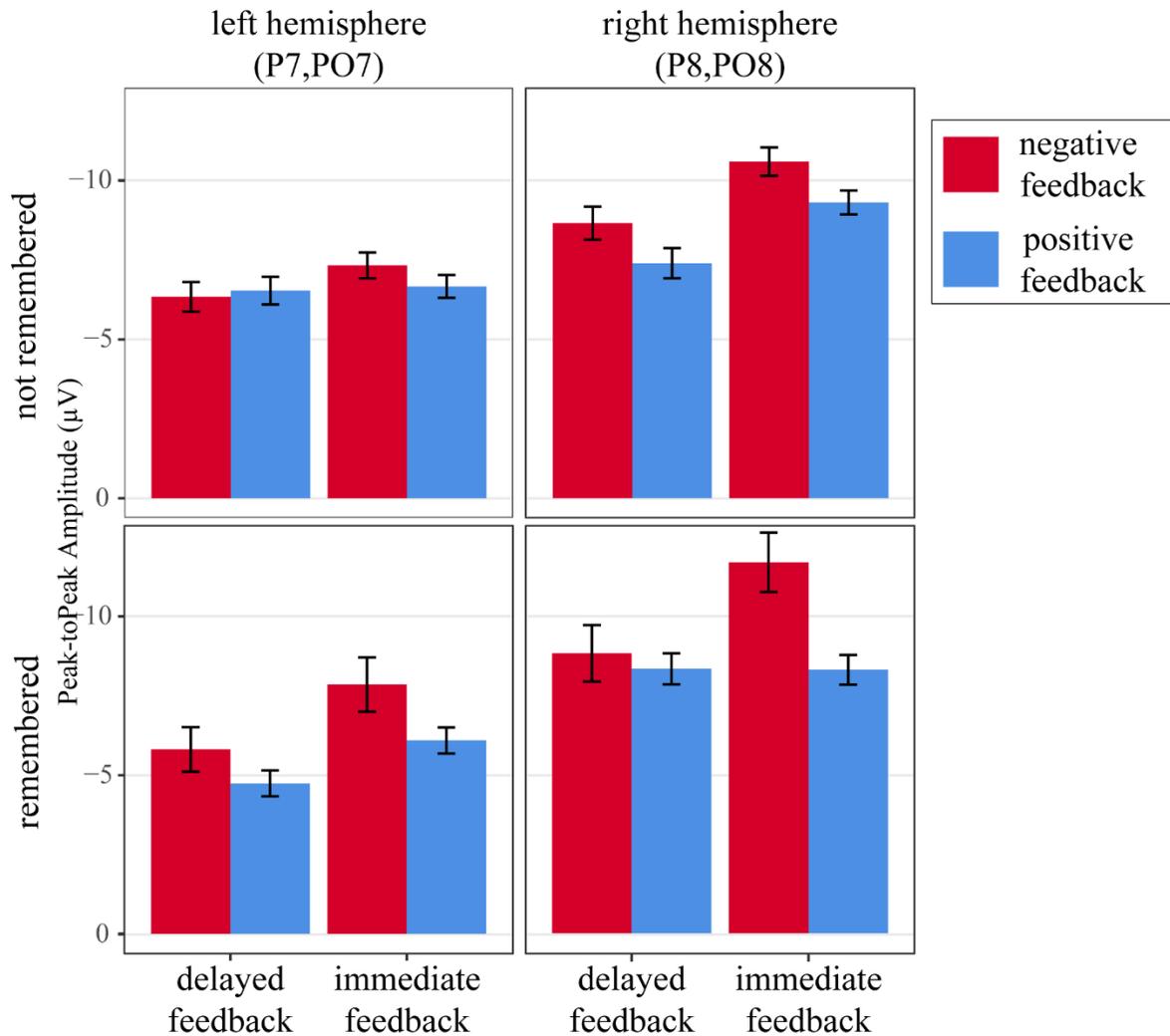


Notes. The dotted vertical lines represent the start and end point for the maximum negative peak detection for the N170.

All main effects were significant: We found a main effect of Feedback Timing, $F(1,28.40) = 5.41$, $p = .027$, $b = 0.69$ (larger amplitudes for immediate compared to delayed feedback), a main effect of Feedback Valence, $F(1,30.70) = 20.26$, $p < .001$, $b = 0.58$ (larger amplitudes for negative compared to positive feedback), a main effect of Laterality, $F(1,28.40) = 18.18$, $p < .001$, $b = -1.30$ (larger amplitudes over the right hemisphere), and a main effect of Free Recall Performance, $F(1,14548.50) = 9.18$, $p = .002$, $b = 0.20$ (higher amplitudes for not-remembered associations).

Figure 6

Means and standard errors of the N170 amplitude depending on Feedback Valence, Feedback Timing, Free Recall Performance and Laterality for the N170.



We also found a significant two-way interaction between Feedback Valence and Free Recall Performance, $F(1,7476.10) = 4.08$, $p = .043$. For positive feedback, there was an effect of Free Recall Performance, $F(1,10668.90) = 21.34$, $p < .001$, $b = 0.33$ with larger amplitudes for not remembered associations. For negative feedback, there was no such effect ($p = .522$).

Finally, a significant four-way interaction between Feedback Timing, Feedback Valence, Laterality and Free Recall Performance emerged, $F(1,13864.20) = 7.66$, $p = .006$. Resolving this interaction by analyzing the N170 for immediate and delayed feedback separately revealed a three-way interaction between Feedback Valence, Laterality and Free Recall

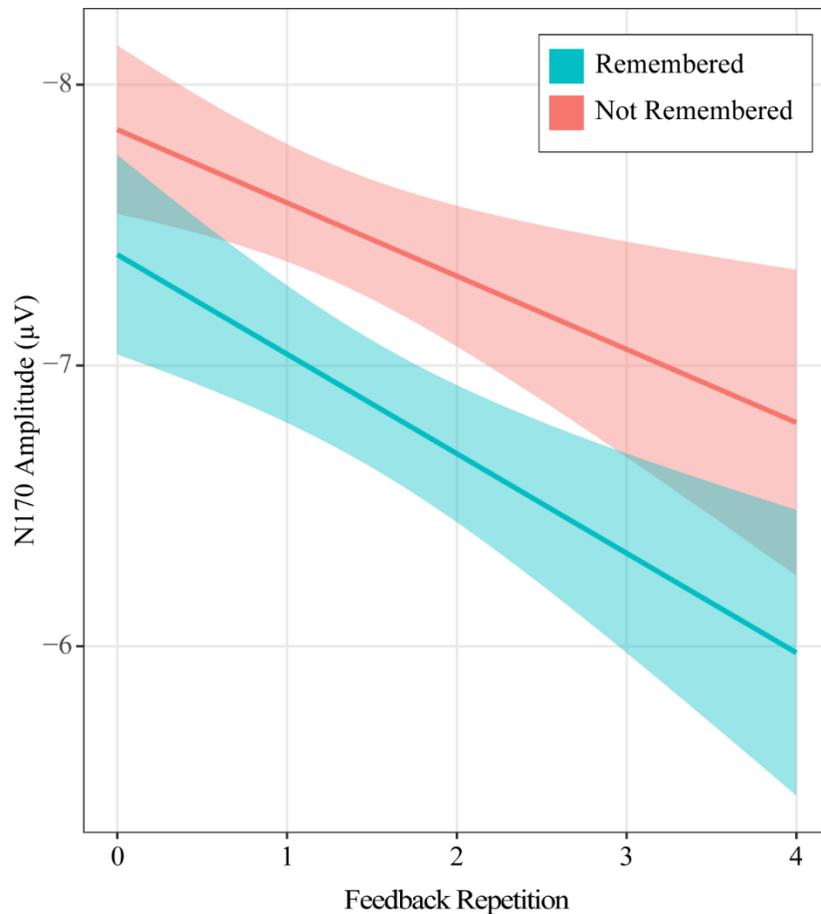
Performance only for delayed feedback, $F(1,8139.60) = 5.08$, $p = .024$, but not for immediate feedback ($p = .109$). Resolving further, a Feedback Valence x Free Recall Performance two-way interaction emerged for left-hemispheric delayed, $F(1,5367.00) = 5.36$, $p = .021$, but not right-hemispheric delayed feedback ($p = .421$). For the left-hemispheric delayed interaction, an effect of Free Recall Performance showed only for positive, $F(1,8774.70) = 11.58$, $p < .001$, $b = .50$, not for negative feedback ($p = .721$). Amplitudes were larger for not-remembered associations.

With a further post-hoc analysis we aimed to explore potential reasons for the selective subsequent memory (recall) effect on the N170 following positive feedback. As became evident in the analyses of the behavioral data reported above, correct responses, and thus positive feedback, became more frequent during the training blocks, and, accordingly, negative feedback became less frequent, especially for remembered associations (see also Table S3 in the supplementary material for the number of trials with positive and negative feedback in the different blocks of trials). As a possible explanation for the reduced N170 amplitudes following positive feedback for later remembered associations, we checked whether the number of times the right non-word had been chosen for an object in the training blocks, i.e. the number of times participants received positive feedback for the same choice, influenced N170 amplitude. We thus conducted an additional analysis only for positive feedback trials, also because the number of trials with negative feedback became too low during the course of the experiment. We expected that predictions would be stronger the more times the correct non-word for the respective non-object had been chosen before, i.e. the more often participants had received positive feedback for their choice. This, in turn, could reduce N170 amplitude if it reflects prediction (violation). We thus included Feedback Repetition as independent variable and N170 amplitude as dependent variable into our model. Moreover, we added Free Recall Performance as independent variable,

as we wanted to explore if a potential effect of Repetition would perhaps be more pronounced for correctly recalled non-word-non-object pairs. Random intercepts and slopes for Feedback Repetition were set per subject (as determined in the iteration process mentioned above). We found a significant main effect of Feedback Repetition on N170 amplitudes for positive feedback, $F(1,30.40) = 15.08, p < .001, b = 0.27$, with smaller amplitudes the higher the Repetition number of Positive Feedback. In accordance with the analysis reported above, we found a main effect of Free Recall Memory, $F(1,9397.40) = 8.57, p = .003, b = 0.23$. Importantly, there was no significant interaction between Feedback Repetition and Free Recall Memory, $F(1,1272.60) = 0.50, p = .482$. As can be seen in Figure 7, the N170 amplitude is smaller for remembered compared to non-remembered associations already when positive feedback is given the first time (0 repetitions).

Figure 7

Means and confidence intervals of the N170 amplitude for positive feedback depending on the Feedback Repetition and Free Recall Performance.



Discussion

Humans and other species can learn many different types of associations via feedback. As has recently been shown in humans, also associations between novel words (non-words) and novel objects can be learned in this way. Arbel et al. (2017) found that study participants can acquire non-word-non-object pairs when they choose between two potential non-words for a given non-object and receive feedback for their choice. In the present study we aimed to investigate the relationship between feedback processing and subsequent memory for these associations. Focusing on recall of the associated non-word for each non-object, we hypothesized that the feedback-locked N170 would be related to subsequent memory, as it has, in contrast to the FRN, been linked to processing in the MTL and to declarative memory (Arbel et al., 2017;

Höltje & Mecklinger, 2020). Moreover, we expected that the relationship between the N170 and recall performance would be modulated by the interval between the participant's choice and the feedback, as delayed feedback processing has been shown to recruit the MTL (Foerde et al., 2013; Foerde & Shohamy, 2011b; Lighthall et al., 2018), which, in turn, has been linked to declarative memory and recall (Danckert & Craik, 2013; Leshikar et al., 2017; Staresina & Davachi, 2006).

The results were partially in line with our hypotheses. For the FRN, previous findings were replicated, as FRN amplitude was larger for negative than positive feedback (Gehring & Willoughby, 2002; Miltner et al., 1997; Nieuwenhuis et al., 2004) but only with immediate, not delayed feedback (Arbel et al., 2017; Höltje & Mecklinger, 2020; Peterburs et al., 2016; Weismüller & Bellebaum, 2016). As expected, FRN amplitude did not differ for later remembered and not remembered associations.

Also as expected, we found that the N170 was modulated by later free recall performance. First of all, an interaction between Recall Performance and Feedback Valence was found for the N170: Only for positive feedback there was an N170 amplitude difference between not remembered and remembered associations, with larger amplitudes for the former.

Several aspects of the study design might explain that the effect is found for positive feedback only. First of all, positive feedback was much more frequent than negative feedback throughout the experiment, as participants learned during the training blocks. Furthermore, from the nature of the task, positive feedback might be more relevant for free recall, and thus for declarative learning than for procedural learning. For correct answers in the training blocks as well as the recognition tests it is sufficient to know either which non-word was correct or which non-word was wrong for a given non-object (as non-objects were always paired with the same

two non-words), which is indicative of a number of tasks involving feedback (Foerde & Shohamy, 2011a; Frank et al., 2004; Gluck et al., 2002; Maddox et al., 2003; Peterburs et al., 2016). For the free recall test, participants needed to specifically know the correct non-word and could (and even should) ignore the wrong non-word, which could mean that positive feedback promoted declarative learning involving the MTL.

What remains open is the question which cognitive process the N170 indicates during the evaluation of feedback, especially as N170 amplitudes were reduced for later remembered associations. As a post-hoc analyses revealed, the more confident participants became in their responses (the more often positive feedback was repeated for the same non-word-non-object pair), the smaller were the N170 amplitudes. This may indicate that the N170 amplitude was modulated by expectancy. There is indeed evidence that the neural processing in the MTL is affected by expectancy and prediction. Taking fMRI studies on feedback learning into account, it has been shown that prediction error representations can also be found in the MTL (Dickerson et al., 2011; Lighthall et al., 2018). While the link between MTL processing and the N170 in the context of feedback processing is still unclear, there is recent evidence that the N170 is modulated by predictions and prediction errors. Originally described as an ERP component representing face processing (Rossion, 2014), N170 amplitude has been shown to be modulated by the predictability of faces, but also of other visual stimuli, with lower amplitudes for more strongly predicted stimuli (Baker et al., 2021; Johnston et al., 2017; Robinson et al., 2020). In this context, however, the N170 is seen as an indicator of higher-order visual processing, which is *reduced* if stimuli are predicted (Baker et al., 2021; Johnston et al., 2017; Robinson et al., 2020). Accordingly, for this component sources in fusiform gyrus have been described (Cohen et al., 2000; Gao et al., 2019), which belongs to the ventral path of visual processing (Ungerleider

& Haxby, 1994). In the context of feedback processing, however, *increased* N170 amplitudes have been interpreted as enhanced MTL activity in previous studies (Arbel et al., 2017; Höljtje & Mecklinger, 2020; Kim & Arbel, 2019). The relative contribution of these two processes to the N170 is difficult to determine, especially because also effects of expectations on visual processing are shaped by learning. It seems at least unlikely, however, that the N170 directly reflects a bottom-up visual process, as a modulation of the N170 by feedback timing has been found also for auditory feedback (Kim & Arbel, 2019). Nevertheless, it is conceivable that there are two different, possibly overlapping processes reflected in the N170 window related to visual processing in the fusiform gyrus and to MTL activity.

Importantly, the change in positive feedback expectancy during the course of the experiment cannot explain the N170 amplitude difference between remembered and non-remembered items in the present study. As our post-hoc analysis revealed, the N170 following positive feedback became smaller with each repetition of positive feedback for both, remembered and non-remembered associations. Instead, positive feedback processing already differed between remembered and non-remembered associations on the first encounter of positive feedback during early learning stages. If a non-word is later recalled or not thus seems to partially depend on how positive feedback for an accidentally correct choice in the first trial entailing that non-word is processed. An inherent assumption in feedback learning tasks with two response options and two types of feedback is that, at the beginning of the task, both response options are considered to be equally likely to yield positive (and negative) feedback. While this is probably true on average also for the stimuli used in the present study, it is conceivable that the expectations varied for the different non-object-non-word associations. On some trials in the first block, participants may have just randomly decided which non-word to choose for a given

non-object. On other trials, they may have based their decision on the tendency that they considered one non-word to slightly fit better to the shown non-object. In this case they may have been less surprised by positive feedback, which led to a reduced N170 amplitude. These non-words were then later remembered correctly more often. In the first learning block there may thus have been a bias for some associations and in some participants that may have partially determined learning success. At the same time the expectation of positive feedback grew stronger the more often it was repeated for the same choice. This further reduced N170 amplitude, irrespective of later recall performance.

Another finding in line with our expectations is that feedback delay affected the relationship between N170 amplitude and free recall. More specifically, the interaction between Feedback Valence and Free Recall Performance was especially pronounced at left-hemispheric sites with delayed feedback. As outlined above, the MTL is more strongly involved in delayed feedback processing (Foerde & Shohamy, 2011a; Lighthall et al., 2018), and the N170 has been linked to processing in the MTL (Arbel et al., 2017; Hölting & Mecklinger, 2020). It is thus not surprising that the N170 for delayed feedback is particularly related to learning success. The left lateralization of the effect may, in turn, be caused by the importance of verbal information in our task. It is known that visual word processing involves the left-hemispheric fusiform gyrus (Binder et al., 2006; Cohen et al., 2000). In our task non-words had to be associated with non-objects during learning, and we used the words “Correct” and “Incorrect” as feedback. Visual word processing, as opposed to visual face processing, indeed coincides with a left-hemispheric N170 (Bentin et al., 1999; Mercure et al., 2008; Rossion et al., 2003). While a difference between orthographic and non-orthographic stimuli was found only for left-hemispheric processing, mean amplitudes were comparable between the right and left hemispheric N170

(Bentin et al., 1999). In our findings, overall amplitudes were even larger for the right-hemispheric N170, but only the left-hemispheric N170 was modulated based on Feedback Delay and Free Recall Performance. The sensitivity of the N170 to different visual stimuli needs to be considered in future experiments investigating this component as a candidate for representing MTL activity, using standardized feedback stimuli to control for visual discrepancies.

One aspect of the results of the present study that was unexpected was that N170 amplitudes were reduced for delayed feedback instead of enhanced (as found by Arbel et al., 2017; Höljtje & Mecklinger, 2020; Kim & Arbel, 2019). The relative contributions of the striatum and MTL to feedback learning depend, however, most likely not only on feedback delay, but also on other factors such as feedback contingency. In case of deterministic feedback, as in the present study, feedback learning resembles paired associates learning, which is considered an example of declarative, MTL-based learning (Poldrack et al., 2001; Shohamy et al., 2004). The deterministic nature of the feedback in the present study may thus provide a possible explanation why the N170 was pronounced also for immediate feedback. While Arbel et al. (2017) found a larger N170 for delayed feedback with the same learning task, there were also important differences between the studies. One difference is related to the way in which learning was assessed. In our study memory performance was (also) measured by means of free recall. This may have strengthened the application of declarative learning strategies in the study participants, irrespective of feedback delay, which may, in turn, have increased the amplitude of the N170 also for immediate feedback. Additionally, the N170 is influenced by higher-order visual processing, with larger amplitudes for faces (Rossion, 2014), but also other types of stimuli: for example, there are larger amplitudes for road signs compared to tools and texture (Itier & Taylor, 2004). Arbel et al. (2017) used three Xs or three checks for feedback, Höljtje &

Mecklinger (2020) used indoor and outdoor scenes. The visual display of round casino chips during feedback in the present study might have also increased at least right-lateralised N170 amplitudes. While the described factors provide potential explanations for a pronounced N170 amplitude for immediate feedback, the reason why it is even larger for immediate than delayed feedback may be related to the gamification elements that we added to the task. These might have led to different motivation levels compared to Arbel et al. (2017). Higher motivation could lead to higher functional connectivity between MTL and striatum (Davidow et al., 2016), which we assume would be especially pronounced for immediate feedback, when striatal activity was pronounced.

As a conclusion, we could show for the first time that the feedback-locked N170 (but not the FRN) is modulated by subsequent memory for newly learned (non-)words, as assessed by means of free recall performance. While we also found the link between N170 amplitude and later free recall memory performance to be modulated by feedback timing, the exact neural mechanisms reflected by the N170 in the context of feedback processing still need to be determined. We suspect an influence of prediction with a potential role of the visual properties of the feedback. To differentiate between different factors modulating the N170, the future use of functional imaging methods with high spatial resolution would be a promising approach.

References

- Arbel, Y., Hong, L., Baker, T. E., & Holroyd, C. B. (2017). It's all about timing: An electrophysiological examination of feedback-based learning with immediate and delayed feedback. *Neuropsychologia*, *99*, 179–186.
<https://doi.org/10.1016/j.neuropsychologia.2017.03.003>
- Baayen, R. H., Piepenbrock, R., & Gulikers, L. (1995). *CELEX2 LDC96L14*.
- Baker, K. S., Pegna, A. J., Yamamoto, N., & Johnston, P. (2021). Attention and prediction modulations in expected and unexpected visuospatial trajectories. *PloS One*, *16*(10), e0242753. <https://doi.org/10.1371/journal.pone.0242753>
- Baker, T. E., & Holroyd, C. B. (2013). The topographical N170: Electrophysiological evidence of a neural mechanism for human spatial navigation. *Biological Psychology*, *94*(1), 90–105. <https://doi.org/10.1016/j.biopsycho.2013.05.004>
- Becker, M. P. I., Nitsch, A. M., Miltner, W. H. R., & Straube, T. (2014). A single-trial estimation of the feedback-related negativity and its relation to BOLD responses in a time-estimation task. *The Journal of Neuroscience : The Official Journal of the Society for Neuroscience*, *34*(8), 3005–3012. <https://doi.org/10.1523/JNEUROSCI.3684-13.2014>
- Bentin, S., Mouchetant-Rostaing, Y., Giard, M.-H., Echallier, J.-F., & Pernier, J. (1999). ERP manifestations of processing printed words at different psycholinguistic levels: time course and scalp distribution. *Journal of Cognitive Neuroscience*, *11*(3), 235–260.
- Binder, J. R., Medler, D. A., Westbury, C. F., Liebenthal, E., & Buchanan, L. (2006). Tuning of the human left fusiform gyrus to sublexical orthographic structure. *NeuroImage*, *33*(2), 739–748. <https://doi.org/10.1016/j.neuroimage.2006.06.053>

- Brackbill, Y., & Kappy, M. S. (1962). Delay of reinforcement and retention. *Journal of Comparative and Physiological Psychology*, *55*, 14–18.
- Brackbill, Y., Wagner, J. E., & Wilson, D. (1964). Feedback delay and the teaching machine. *Psychology in the Schools*, *1*, 148–156.
- Carpenter, S. K., & Vul, E. (2011). Delaying feedback by three seconds benefits retention of face-name pairs: The role of active anticipatory processing. *Memory & Cognition*, *39*(7), 1211–1221. <https://doi.org/10.3758/s13421-011-0092-1>
- Chase, H. W., Swainson, R., Durham, L., Benham, L., & Cools, R. (2011). Feedback-related Negativity Codes Prediction Error but Not Behavioral Adjustment during Probabilistic Reversal Learning. *Journal of Cognitive Neuroscience*, *23*(4), 936–946. <https://doi.org/10.1162/jocn.2010.21456>
- Chau, B. K. H., Jarvis, H., Law, C.-K., & Chong, T. T.-J. (2018). Dopamine and reward: A view from the prefrontal cortex. *Behavioural Pharmacology*, *29*(7), 569–583. <https://doi.org/10.1097/FBP.0000000000000424>
- Cohen, L., Dehaene, S., Naccache, L., Lehéricy, S., Dehaene-Lambertz, G., Hénaff, M. A., & Michel, F. (2000). The visual word form area: Spatial and temporal characterization of an initial stage of reading in normal subjects and posterior split-brain patients. *Brain : A Journal of Neurology*, *123* (Pt 2), 291–307. <https://doi.org/10.1093/brain/123.2.291>
- Cohen, M. X., & Ranganath, C. (2007). Reinforcement learning signals predict future decisions. *The Journal of Neuroscience : The Official Journal of the Society for Neuroscience*, *27*(2), 371–378. <https://doi.org/10.1523/JNEUROSCI.4421-06.2007>
- Danckert, S. L., & Craik, F. I. M. (2013). Does aging affect recall more than recognition memory? *Psychology and Aging*, *28*(4), 902–909. <https://doi.org/10.1037/a0033263>

- Davidow, J. Y., Foerde, K., Galván, A., & Shohamy, D. (2016). An Upside to Reward Sensitivity: The Hippocampus Supports Enhanced Reinforcement Learning in Adolescence. *Neuron*, *92*(1), 93–99. <https://doi.org/10.1016/j.neuron.2016.08.031>
- de Bruijn, E. R. A., Mars, R. B., & Hester, R. (2020). Processing of performance errors predicts memory formation: Enhanced feedback-related negativities for corrected versus repeated errors in an associative learning paradigm. *The European Journal of Neuroscience*, *51*(3), 881–890. <https://doi.org/10.1111/ejn.14566>
- Dickerson, K. C., & Delgado, M. R. (2015). Contributions of the hippocampus to feedback learning. *Cognitive, Affective & Behavioral Neuroscience*, *15*(4), 861–877. <https://doi.org/10.3758/s13415-015-0364-5>
- Dickerson, K. C., Li, J., & Delgado, M. R. (2011). Parallel contributions of distinct human memory systems during probabilistic learning. *NeuroImage*, *55*(1), 266–276. <https://doi.org/10.1016/j.neuroimage.2010.10.080>
- Dobryakova, E., & Tricomi, E. (2013). Basal ganglia engagement during feedback processing after a substantial delay. *Cognitive, Affective & Behavioral Neuroscience*, *13*(4), 725–736. <https://doi.org/10.3758/s13415-013-0182-6>
- Dudenredaktion. (2020). *Duden*. Bibliographisches Institut GmbH. <https://www.duden.de/>
- Fera, F., Passamonti, L., Herzallah, M. M., Myers, C. E., Veltri, P., Morganti, G., Quattrone, A., & Gluck, M. A. (2014). Hippocampal BOLD response during category learning predicts subsequent performance on transfer generalization. *Human Brain Mapping*, *35*(7), 3122–3131. <https://doi.org/10.1002/hbm.22389>

- Fischer, A. G., & Ullsperger, M. (2013). Real and fictive outcomes are processed differently but converge on a common adaptive mechanism. *Neuron*, *79*(6), 1243–1255.
<https://doi.org/10.1016/j.neuron.2013.07.006>
- Foerde, K., Race, E., Verfaellie, M., & Shohamy, D. (2013). A role for the medial temporal lobe in feedback-driven learning: Evidence from amnesia. *The Journal of Neuroscience : The Official Journal of the Society for Neuroscience*, *33*(13), 5698–5704.
<https://doi.org/10.1523/JNEUROSCI.5217-12.2013>
- Foerde, K., & Shohamy, D. (2011a). Feedback timing modulates brain systems for learning in humans. *The Journal of Neuroscience : The Official Journal of the Society for Neuroscience*, *31*(37), 13157–13167. <https://doi.org/10.1523/JNEUROSCI.2701-11.2011>
- Foerde, K., & Shohamy, D. (2011b). The role of the basal ganglia in learning and memory: Insight from Parkinson's disease. *Neurobiology of Learning and Memory*, *96*(4), 624–636. <https://doi.org/10.1016/j.nlm.2011.08.006>
- Frank, M. J., Seeberger, L. C., & O'reilly, R. C. (2004). By carrot or by stick: Cognitive reinforcement learning in parkinsonism. *Science (New York, N.Y.)*, *306*(5703), 1940–1943. <https://doi.org/10.1126/science.1102941>
- Gao, C., Conte, S., Richards, J. E., Xie, W., & Hanayik, T. (2019). The neural sources of N170: Understanding timing of activation in face-selective areas. *Psychophysiology*, *56*(6), e13336. <https://doi.org/10.1111/psyp.13336>
- Gasbarri, A., Pompili, A., Packard, M. G., & Tomaz, C. (2014). Habit learning and memory in mammals: Behavioral and neural characteristics. *Neurobiology of Learning and Memory*, *114*, 198–208. <https://doi.org/10.1016/j.nlm.2014.06.010>

Gehring, W. J., & Willoughby, A. R. (2002). The medial frontal cortex and the rapid processing of monetary gains and losses. *Science (New York, N.Y.)*, *295*(5563), 2279–2282.

<https://doi.org/10.1126/science.1066893>

Gluck, M. A., Shohamy, D., & Myers, C. (2002). How do people solve the “weather prediction” task? Individual variability in strategies for probabilistic category learning. *Learning & Memory (Cold Spring Harbor, N.Y.)*, *9*(6), 408–418. <https://doi.org/10.1101/lm.45202>

Hajcak, G., Moser, J. S., Holroyd, C. B., & Simons, R. F. (2006). The feedback-related negativity reflects the binary evaluation of good versus bad outcomes. *Biological Psychology*, *71*(2), 148–154. <https://doi.org/10.1016/j.biopsycho.2005.04.001>

Hajcak, G., Moser, J. S., Holroyd, C. B., & Simons, R. F. (2007). It’s worse than you thought: The feedback negativity and violations of reward prediction in gambling tasks.

Psychophysiology, *44*(6), 905–912. <https://doi.org/10.1111/j.1469-8986.2007.00567.x>

Hollerman, J. R., & Schultz, W. (1998). Dopamine neurons report an error in the temporal prediction of reward during learning. *Nature Neuroscience*, *1*(4), 304–309.

<https://doi.org/10.1038/1124>

Holroyd, C. B., Krigolson, O. E., Baker, R., Lee, S., & Gibson, J. (2009). When is an error not a prediction error? An electrophysiological investigation. *Cognitive, Affective & Behavioral Neuroscience*, *9*(1), 59–70. <https://doi.org/10.3758/CABN.9.1.59>

Behavioral Neuroscience, *9*(1), 59–70. <https://doi.org/10.3758/CABN.9.1.59>

Holroyd, C. B., Nieuwenhuis, S., Yeung, N., & Cohen, J. D. (2003). Errors in reward prediction are reflected in the event-related brain potential. *Neuroreport*, *14*(18).

https://journals.lww.com/neuroreport/Fulltext/2003/12190/Errors_in_reward_prediction_are_reflected_in_the.37.aspx

- Holroyd, C. B., Nieuwenhuis, S., Yeung, N., Nystrom, L., Mars, R. B., Coles, M. G., & Cohen, J. D. (2004). Dorsal anterior cingulate cortex shows fMRI response to internal and external error signals. *Nature Neuroscience*, 7(5), 497–498.
<https://doi.org/10.1038/nn1238>
- Höltje, G., & Mecklinger, A. (2020). Feedback timing modulates interactions between feedback processing and memory encoding: Evidence from event-related potentials. *Cognitive, Affective, & Behavioral Neuroscience*. Advance online publication.
<https://doi.org/10.3758/s13415-019-00765-5>
- Itier, R. J., & Taylor, M. J. (2004). N170 or N1? Spatiotemporal differences between object and face processing using ERPs. *Cerebral Cortex (New York, N.Y. : 1991)*, 14(2), 132–142.
<https://doi.org/10.1093/cercor/bhg111>
- Johnston, P., Robinson, J., Kokkinakis, A., Ridgeway, S., Simpson, M., Johnson, S., Kaufman, J., & Young, A. W. (2017). Temporal and spatial localization of prediction-error signals in the visual brain. *Biological Psychology*, 125, 45–57.
<https://doi.org/10.1016/j.biopsycho.2017.02.004>
- Keuleers, E., & Brysbaert, M. (2010). Wuggy: A multilingual pseudoword generator. *Brain Research Methods*, 42(3), 627–633.
- Kim, S., & Arbel, Y. (2019). Immediate and delayed auditory feedback in declarative learning: An examination of the feedback related event related potentials. *Neuropsychologia*, 129, 255–262. <https://doi.org/10.1016/j.neuropsychologia.2019.04.001>
- Kroll, J. F., & Potter, M. C. (1984). Recognizing words, pictures, and concepts: A comparison of lexical, object, and reality decisions. *Journal of Verbal Learning and Verbal Behavior*, 23(1), 39–66. [https://doi.org/10.1016/S0022-5371\(84\)90499-7](https://doi.org/10.1016/S0022-5371(84)90499-7)

Leshikar, E. D., Leach, R. C., McCurdy, M. P., Trumbo, M. C., Sklenar, A. M.,

Frankenstein, A. N., & Matzen, L. E. (2017). Transcranial direct current stimulation of dorsolateral prefrontal cortex during encoding improves recall but not recognition memory. *Neuropsychologia*, *106*, 390–397.

<https://doi.org/10.1016/j.neuropsychologia.2017.10.022>

Lieberman, D. A., Vogel, A. C. M., & Nisbet, J. (2008). Why do the effects of delaying reinforcement in animals and delaying feedback in humans differ? A working-memory analysis. *Quarterly Journal of Experimental Psychology (2006)*, *61*(2), 194–202.

<https://doi.org/10.1080/17470210701557506>

Lighthall, N. R., Pearson, J. M., Huettel, S. A., & Cabeza, R. (2018). Feedback-Based Learning in Aging: Contributions and Trajectories of Change in Striatal and Hippocampal Systems. *The Journal of Neuroscience : The Official Journal of the Society for Neuroscience*, *38*(39), 8453–8462. <https://doi.org/10.1523/JNEUROSCI.0769-18.2018>

Maddox, W. T., Ashby, F. G., & Bohil, C. J. (2003). Delayed feedback effects on rule-based and information-integration category learning. *Journal of Experimental Psychology. Learning, Memory, and Cognition*, *29*(4), 650–662. <https://doi.org/10.1037/0278-7393.29.4.650>

Mercure, E., Dick, F., Halit, H., Kaufman, J., & Johnson, M. H. (2008). Differential lateralization for words and faces: Category or psychophysics? *Journal of Cognitive Neuroscience*, *20*(11), 2070–2087. <https://doi.org/10.1162/jocn.2008.20137>

Meteyard, L., & Davies, R. A. (2020). Best practice guidance for linear mixed-effects models in psychological science. *Journal of Memory and Language*, *112*, 104092. <https://doi.org/10.1016/j.jml.2020.104092>

- Miltner, W. H., Braun, C. H., & Coles, M. G. (1997). Event-related brain potentials following incorrect feedback in a time-estimation task: Evidence for a “generic” neural system for error detection. *Journal of Cognitive Neuroscience*, *9*(6), 788–798.
<https://doi.org/10.1162/jocn.1997.9.6.788>
- Myers, C. E., Shohamy, D., Gluck, M. A., Grossman, S., Kluger, A., Ferris, S., Golomb, J., Schnirman, G., & Schwartz, R. (2003). Dissociating Hippocampal versus Basal Ganglia Contributions to Learning and Transfer. *Journal of Cognitive Neuroscience*, *15*(2), 185–193. <https://doi.org/10.1162/089892903321208123>
- Nieuwenhuis, S., Holroyd, C. B., Mol, N., & Coles, M. G. (2004). Reinforcement-related brain potentials from medial frontal cortex: Origins and functional significance. *Neuroscience and Biobehavioral Reviews*, *28*(4), 441–448.
<https://doi.org/10.1016/j.neubiorev.2004.05.003>
- O’Doherty, J., Dayan, P., Schultz, J., Deichmann, R., Friston, K., & Dolan, R. J. (2004). Dissociable roles of ventral and dorsal striatum in instrumental conditioning. *Science (New York, N.Y.)*, *304*(5669), 452–454. <https://doi.org/10.1126/science.1094285>
- Peterburs, J., Kobza, S., & Bellebaum, C. (2016). Feedback delay gradually affects amplitude and valence specificity of the feedback-related negativity (FRN). *Psychophysiology*, *53*(2), 209–215. <https://doi.org/10.1111/psyp.12560>
- Poldrack, R. A., Clark, J., Paré-Blagoev, E. J., Shohamy, D., Creso Moyano, J., Myers, C., & Gluck, M. A. (2001). Interactive memory systems in the human brain. *Nature*, *414*(6863), 546–550. <https://doi.org/10.1038/35107080>
- Proudfit, G. H. (2015). The reward positivity: From basic research on reward to a biomarker for depression. *Psychophysiology*, *52*(4), 449–459. <https://doi.org/10.1111/psyp.12370>

- Rellecke, J., Sommer, W., & Schacht, A. (2013). Emotion effects on the n170: A question of reference? *Brain Topography*, 26(1), 62–71. <https://doi.org/10.1007/s10548-012-0261-y>
- Robinson, J. E., Breakspear, M., Young, A. W., & Johnston, P. J. (2020). Dose-dependent modulation of the visually evoked N1/N170 by perceptual surprise: A clear demonstration of prediction-error signalling. *The European Journal of Neuroscience*, 52(11), 4442–4452. <https://doi.org/10.1111/ejn.13920>
- Rossion, B. (2014). Understanding face perception by means of human electrophysiology. *Trends in Cognitive Sciences*, 18(6), 310–318. <https://doi.org/10.1016/j.tics.2014.02.013>
- Rossion, B., Joyce, C. A., Cottrell, G. W., & Tarr, M. J. (2003). Early lateralization and orientation tuning for face, word, and object processing in the visual cortex. *NeuroImage*, 20(3), 1609–1624. <https://doi.org/10.1016/j.neuroimage.2003.07.010>
- Schultz, W., Dayan, P., & Montague, P. R. (1997). A Neural Substrate of Prediction and Reward. *Science (New York, N.Y.)*, 275(5306), 1593–1599. <https://doi.org/10.1126/science.275.5306.1593>
- Shohamy, D., Myers, C. E., Grossman, S., Sage, J., Gluck, M. A., & Poldrack, R. A. (2004). Cortico-striatal contributions to feedback-based learning: Converging data from neuroimaging and neuropsychology. *Brain : A Journal of Neurology*, 127(Pt 4), 851–859. <https://doi.org/10.1093/brain/awh100>
- Shohamy, D., Myers, C. E., Kalanithi, J., & Gluck, M. A. (2008). Basal ganglia and dopamine contributions to probabilistic category learning. *Neuroscience and Biobehavioral Reviews*, 32(2), 219–236. <https://doi.org/10.1016/j.neubiorev.2007.07.008>

- Shohamy, D., & Wagner, A. D. (2008). Integrating memories in the human brain: Hippocampal-midbrain encoding of overlapping events. *Neuron*, *60*(2), 378–389.
<https://doi.org/10.1016/j.neuron.2008.09.023>
- Squire, L. R. (2004). Memory systems of the brain: A brief history and current perspective. *Neurobiology of Learning and Memory*, *82*(3), 171–177.
<https://doi.org/10.1016/j.nlm.2004.06.005>
- Squire, L. R., & Zola-Morgan, A. J. O. (1991). Conscious and unconscious memory systems. *Cold Spring Harbor Perspectives in Biology*, *7*(3), a021667.
<https://doi.org/10.1101/cshperspect.a021667>
- Staresina, B. P., & Davachi, L. (2006). Differential encoding mechanisms for subsequent associative recognition and free recall. *The Journal of Neuroscience : The Official Journal of the Society for Neuroscience*, *26*(36), 9162–9172.
<https://doi.org/10.1523/JNEUROSCI.2877-06.2006>
- Ungerleider, L. G., & Haxby, J. V. (1994). ‘What’ and ‘where’ in the human brain. *Current Opinion in Neurobiology*, *4*(2), 157–165. [https://doi.org/10.1016/0959-4388\(94\)90066-3](https://doi.org/10.1016/0959-4388(94)90066-3)
- Wang, Y., Huang, H., Yang, H., Xu, J., Mo, S., Lai, H., Wu, T., & Zhang, J. (2019). Influence of EEG References on N170 Component in Human Facial Recognition. *Frontiers in Neuroscience*, *13*, 705. <https://doi.org/10.3389/fnins.2019.00705>
- Weinberg, A., Luhmann, C. C., Bress, J. N., & Hajcak, G. (2012). Better late than never? The effect of feedback delay on ERP indices of reward processing. *Cognitive, Affective & Behavioral Neuroscience*, *12*(4), 671–677. <https://doi.org/10.3758/s13415-012-0104-z>

- Weismüller, B., & Bellebaum, C. (2016). Expectancy affects the feedback-related negativity (FRN) for delayed feedback in probabilistic learning. *Psychophysiology*, *53*(11), 1739–1750. <https://doi.org/10.1111/psyp.12738>
- Weismüller, B., Ghio, M., Logmin, K., Hartmann, C., Schnitzler, A., Pollok, B., Südmeyer, M., & Bellebaum, C. (2018). Effects of feedback delay on learning from positive and negative feedback in patients with Parkinson's disease off medication. *Neuropsychologia*, *117*, 46–54. <https://doi.org/10.1016/j.neuropsychologia.2018.05.010>
- Yeung, N., Holroyd, C. B., & Cohen, J. D. (2005). Erp correlates of feedback and reward processing in the presence and absence of response choice. *Cerebral Cortex (New York, N.Y. : 1991)*, *15*(5), 535–544. <https://doi.org/10.1093/cercor/bhh153>
- Yin, H. H., & Knowlton, B. J. (2006). The role of the basal ganglia in habit formation. *Nature Reviews. Neuroscience*, *7*(6), 464–476. <https://doi.org/10.1038/nrn1919>
- Zaghloul, K. A., Blanco, J. A., Weidemann, C. T., McGill, K., Jaggi, J. L., Baltuch, G. H., & Kahana, M. J. (2009). Human substantia nigra neurons encode unexpected financial rewards. *Science (New York, N.Y.)*, *323*(5920), 1496–1499. <https://doi.org/10.1126/science.1167342>

Table S1

Statistical data for FRN ERP LME analysis for Feedback Timing x Feedback Valence x Free

Recall Performance

Effect	Estimate (<i>b</i>)	Std. error	<i>df</i>	<i>t</i> -value	<i>p</i> -value	CI 2.5%	CI 97.5%
(Intercept)	-4.78	0.46	28.29	-10.37	< .001	-5.65	-3.94
Feedback Timing	0.17	0.26	28.76	0.63	0.537	-0.38	0.69
Feedback Valence	0.70	0.11	33.61	6.62	< .001	0.49	0.89
Free Recall Performance	0.03	0.08	7181.4	0.42	0.672	-0.13	0.20
Feedback Timing x Feedback Valence	-0.44	0.14	31.93	-3.24	0.003	-0.70	-0.18
Feedback Timing x Free Recall Performance	-0.12	0.08	6941.75	-1.49	0.135	-0.29	0.04
Feedback Valence x Free Recall Performance	-0.05	0.08	2048.34	-0.65	0.518	-0.20	0.10
Feedback Timing x Feedback Valence x Free Recall Performance	0.05	0.08	3452.33	0.62	0.538	-0.11	0.20

Note. Degrees of Freedom (*df*), *t*- and *p*-values as well as estimates (*b*) based on a restricted maximum likelihood approach, as proposed by Luke (2017) for the Feedback Timing x Feedback Valence x Free Recall Performance LME analysis on the single-trial ERP data for FRN. Satterthwaite approximation was used for the degrees of freedom. Significant values are displayed in bold font.

Table S2*Statistical data for N170 ERP LME analysis for Feedback Timing x Feedback Valence x Free**Recall Performance x Electrode*

Effect	Estimate (b)	Std. error	df	t-value	p-value	CI 2.5%	CI 97.5%
(Intercept)	-7.79	0.66	28.10	-11.78	< .001	-9.02	-6.46
Feedback Timing	0.69	0.30	28.44	2.33	.027	0.10	1.26
Feedback Valence	0.58	0.13	30.71	4.50	< .001	0.34	0.85
Laterality	-1.30	0.31	28.45	-4.26	< .001	-1.80	-0.77
Free Recall Performance	0.20	0.07	14550.00	3.03	.002	0.08	0.33
Feedback Timing x Feedback Valence	-0.14	0.10	31.75	-1.46	.155	-0.35	0.05
Feedback Timing x Laterality	0.20	0.20	29.43	0.97	.339	-0.18	0.57
Feedback Valence x Laterality	0.20	0.11	31.70	1.84	.075	0.00	0.40
Feedback Timing x Free Recall Performance	0.09	0.07	13690.00	1.32	.186	-0.05	0.24
Feedback Valence x Free Recall Performance	0.13	0.06	7476.00	2.02	.043	0.00	0.26
Laterality x Free Recall Performance	0.030	0.07	13850.00	0.45	.653	-0.12	0.14
Feedback Timing x Feedback Valence x Laterality	-0.05	0.06	14680.00	-0.78	.436	-0.17	0.09
Feedback Timing x Feedback Valence x Free Recall Performance	-0.04	0.06	4861.00	-0.57	.570	-0.17	0.09
Feedback Timing x Laterality x Free Recall Performance	0.05	0.06	13290.00	0.70	.486	-0.08	0.18
Feedback Valence x Laterality x Free Recall Performance	-0.02	0.06	6232.00	-0.38	.702	-0.14	0.11
Feedback Timing x Feedback Valence x Laterality x Free Recall Performance	-0.17	0.06	13860.00	-2.77	.006	-0.30	-0.05

Note. Degrees of Freedom (*df*), *t*- and *p*-values as well as estimates (*b*) based on a restricted maximum likelihood approach, as proposed by Luke (2017) for the Feedback Timing x Feedback Valence x Free Recall Performance LME analysis on the single-trial ERP data for FRN. Satterthwaite approximation was used for the degrees of freedom. Significant values are displayed in bold font.

Table S3

Means and standard deviations for the number of trials with positive and negative feedback per condition for all subjects

Feedback Valence	Free Recall Memory	Feedback Timing	Block					
			1	2	3	4	5	
Positive	Not Remembered	Delayed	7.53 (3.36)	8.23 (3.93)	9.00 (4.69)	9.43 (3.90)	10.30 (4.42)	
		Immediate	8.13 (2.84)	9.83 (3.75)	11.00 (3.67)	11.17 (3.80)	12.03 (4.11)	
	Remembered	Delayed	5.53 (3.42)	6.17 (4.13)	8.43 (5.70)	9.03 (5.94)	8.83 (6.40)	
		Immediate	5.17 (3.13)	6.63 (4.28)	7.87 (4.86)	8.70 (5.47)	8.83 (5.70)	
	Negative	Not Remembered	Delayed	8.40 (3.39)	8.30 (4.02)	7.83 (3.71)	7.17 (4.21)	5.30 (3.58)
			Immediate	8.73 (2.68)	8.07 (3.02)	6.93 (3.57)	6.60 (3.62)	5.20 (3.37)
Remembered		Delayed	4.53 (2.92)	3.83 (2.95)	2.27 (1.98)	1.53 (1.81)	0.77 (0.68)	
		Immediate	4.10 (3.16)	2.90 (2.01)	1.83 (1.64)	0.93 (0.91)	0.27 (0.45)	