

Electrophysiological Correlates Of Emotion-Induced Recognition Bias

Sabine Windmann

and

Marta Kutas

Department of Cognitive Science

University of California, San Diego

Corresponding author:

Dr. Sabine Windmann
Cognitive Science
University of California San Diego
9500 Gilman Drive
La Jolla, CA 92093-0515
Fax (858) 534 1128
Email swindman@cogsci.ucsd.edu

Acknowledgments:

This work was supported by grants MH52893, HD22614, and AG08313 to M.K. and a post-doctoral scholarship of the DAAD (Bonn, Germany) to S.W. within the "Hochschulsonderprogramm III von Bund und Ländern". Thanks to Jeff Elman (UCSD) for the HAL analysis, and to Yixiu Li (University of Colorado) for the LSA. Many thanks also to Tom Urbach for his support in technical and other matters.

Abstract

The question of how emotions influence recognition memory is of interest not only within basic cognitive neuroscience but from clinical and forensic perspectives as well. Emotional stimuli can induce a "recognition bias" such that individuals are more likely to respond "old" to a negative item than to an emotionally neutral item, whether the item is actually old or new. We investigated this bias using event-related brain potential (ERP) measures by comparing the processing of words given "old" responses with accurate detection of old/new differences. ERPs to correctly recognized old (hits) and new words (correct rejections) were influenced by emotional valence relatively late (~450 ms+), similarly, and to the same extent; i.e., regardless of emotional valence, the ERP associated with hits was characterized by a widespread positivity between 300 and 700 ms relative to that for correct rejections. By contrast, the analysis of ERPs to old and new items that were judged "old" (hits and false alarms, respectively) revealed a differential effect of valence by 300 ms: Neutral items showed a large old/new difference over prefrontal sites whereas negative items did not. These results are the first clear demonstration of response bias effects on ERPs linked to recognition memory. They are consistent with the idea that frontal cortex areas may be responsible for relaxing the retrieval criterion for negative stimuli so as to ensure that emotional events are not as easily "missed" or forgotten as neutral events.

Keywords: Affect, Emotion, Event-Related Potentials, Evoked Potentials, False Memories, Recognition, Prefrontal, Bias, Unconscious.

The question of how emotions affect information processing is an important one, not only from the perspective of basic cognitive neuroscience, but also for its clinical (Bremner, Staib, Kaloupek, Southwick, Soufer, & Charney, 1999; Drevets, 1998, Gorman, Kent, Sullivan, & Coplan, 2000; Reiman, 1997) and forensic implications (Christianson, 1992a; Kiehl, Hare, McDonald, & Brink, 1999; Raine, Meloy, Bihrl, Stoddard, LaCasse, & Buchsbaum, 1998). Scientists have studied the multiple ways by which emotional affect can enhance, impair, distort, or otherwise influence memory performance for decades (for an overview see Christianson, 1992b). For example, memory for emotional (relative to neutral) events have been described as more focused, more vivid, more distinct, and more robust to forgetting (Kleinsmith & Kaplan, 1963; Ochsner, 2000). Some of these phenomena have been attributed to the positive effects of emotional arousal on memory consolidation via a cooperation between the amygdala and the medial temporal lobes, mediated by adrenergic and glucocorticoid neuromodulation (Cahill & McGaugh, 1998; McGaugh, 2000; LaBar & Phelps, 1998).

Neuropsychological and imaging studies in humans have also implicated ventromedial/medial prefrontal regions in emotional learning and memory (Bechara, Damasio, Damasio & Lee, 1999; Bremner et al., 1999; Paradiso et al., 1999). These regions seem to provide a crucial interface between the evolutionarily old, preconscious stimulus-evaluation systems within the limbic system, and the more flexible, higher-order control systems within the dorsolateral prefrontal cortex required for decision making, reversal learning, and goal-directed behavior (Bechara et al., 1999; Bechara, Damasio & Damasio, 2000; Dias, Robbins, & Roberts, 1996; Rolls, Hornak, Wade & McGrath, 1994). The prefrontal cortex is also important for the retrieval of episodic memories (Buckner, 1996; Tomita, Ohbayashi, Nakahara, Hasegawa, & Miyashita, 1999), as well as for the suppression of currently irrelevant memories (Schnider, Treyer, & Buck, 2000). Thus, it may play an executive, modulatory role in the retrieval or active reorganization of

emotional memories in addition to and independent of the arousal-related effects of emotions on memory consolidation that have been linked primarily to the amygdala (c.f. Bechara et al., 2000).

In laboratory studies, memory performance has often been found to be greater for emotionally-arousing than neutral stimuli (Bradley, Greenwald, Petry, & Lang, 1992; Palomba, Angrilli, & Mini, 1997; Ochsner, 2000), even in amnesic patients (Hamann, Cahill, McGaugh, & Squire, 1997). This advantage, however, does not come without costs, as Windmann & Krüger (1998) noted: Negatively-charged, potentially threatening stimuli are accompanied not only by more *correct* recall and recognition than neutral stimuli, as may be mediated by enhanced attention, but also by a higher probability of *incorrect* recall and recognition (Cross, 1999; Ehlers, Margraf, Davies, & Roth, 1988; Leiphart, Rosenfeld, & Gabrieli, 1993¹; Windmann & Krüger, 1998). Specifically, subjects seem to adopt a different guessing criterion, i.e., a different *response bias*, to negative than to neutral items, under conditions in which they are not explicitly instructed to *attend* to the emotional dimension, nor have any reason to expect that doing so would improve their performance. Here, we refer to this phenomenon as a *recognition bias induced by negative emotional valence*: Subjects are more likely to think that an item is "old" when it is negative as opposed to neutral, whether the item is actually old or new. At present, the mechanism and functional relevance of this phenomenon are completely unknown. At the first blush, this appears to be a cognitive error – a “memory illusion” as it were, much like the ‘false memories’ that emerge from highly interrelated true memories (Cross, 1999; Roediger, McDermott, & Robinson, 1998; Nessler, Mecklinger, & Penney, 2000). Insofar as negative items are more interrelated, e.g. form a more coherent category than neutral items, and/or encourage more categorical, gist-based thinking (c.f., Heuer & Reisberg, 1992; LaBar & Phelps, 1998; Maratos, Allan, & Rugg, 2000), some of the theoretical proposals offered to account for false

memories (Roediger et al., 1998; Miller & Wolford, 1999; Schacter, Norman, & Koutstaal, 1998) also may account for emotion-induced recognition bias (Maratos et al., 2000). Still, this bias may reflect an adaptive cognitive function – an automatic or otherwise elementary mechanism built-in to ensure that events/things with a potentially high survival value are not "missed", even when the focus of attention is directed elsewhere (Windmann & Krüger, 1998).

The present study is aimed at elucidating the mechanisms whereby emotional valence induces a recognition bias. Specifically, we used ERPs to find out how and when (i.e., at what approximate stage of processing) negative emotional connotation influences decisions about whether an item is old or new as ERPs are sensitive to both recognition memory functions and the processing of emotional information, as reviewed below.

ERPs and Recognition Memory. One of the better established findings in the literature on the electrophysiology of recognition memory is the *ERP old/new effect* (for reviews, see Allan, Wilding & Rugg, 1998; Johnson, 1995; Rugg, 1995). It refers to the finding that items that were presented previously (i.e., old items) during a study phase elicit more positive ERPs in a subsequent recognition memory test than (new) items that were not presented during study. The old/new effect typically occurs between 300 and 1000 ms post stimulus onset, thereby overlapping both the N400 and the P3 or "Late Positive Complex" (LPC). With word-like stimuli, the old/new effect is usually largest parietally between 400 and 600 ms post stimulus onset with a slight left hemisphere predominance (Allan et al., 1998; Donaldson & Rugg, 1999). Since this left-parietal old/new effect is not seen for words that are incorrectly recognized, and is reduced or even absent in amnesics, it has been linked to successful item retrieval mechanisms mediated by the medial temporal lobes (Allan et al., 1998; Johnson, 1995). More sustained

¹ For the Leiphart et al. (1993) study, this has to be inferred from the reported hit and correct rejection

old/new effects are seen frontally, especially at right hemisphere sites. They are observed even for new stimuli that are falsely classified as “old” (Walla, Endl, Lindinger, Deecke, & Lang, 2000). They have tentatively been associated with the active maintenance of subjectively retrieved item representations by prefrontal cortical regions for further action planning and decision making, as required e.g. for source judgments (e.g., Allan et al., 1998).

Many researchers have suggested that the early (300-500 ms) posterior old/new effects are more closely related to implicit memory, priming effects, and stimulus familiarity/fluency due to repetition (Johnson, 1995; Paller, Kutas, & McIsaac, 1995; Rugg, 1995; Rugg, Mark, Walla, Schloerscheidt, Birch & Allan, 1998), whereas the later portions (modulating the LPC) are more closely related to conscious, intentional, and episodic recollection processes (Allan et al., 1998; Düzel, Yonelinas, Mangun, Heinze, & Tulving, 1997; Rugg et al., 1998). Evidence in support of this distinction can be found e.g. in Rugg et al. (1998), who found that old items produced more positivity parietally than new items between 300 and 500 ms (N400), regardless of recognition accuracy, and independent of a levels-of-processing manipulation, in both an explicit and an implicit test of memory, and that the LPC amplitude (500-800 ms) varied with the levels-of-processing manipulation, being larger for recognized words studied deeply than either for recognized words studied shallowly or for unrecognized words (see also Paller & Kutas, 1992).

ERPs and Emotion Perception. The perception and experience of emotional cues in pictorial (Johnston, Miller, & Burleson, 1986; Kayser, Tenke, Nordby, Hammerborg, Hugdahl, & Erdmann, 1997; Palomba et al., 1997; Schupp et al., 2000) and verbal stimuli (Stormark, Nordby, & Hugdahl, 1995; Naumann, Bartussek, Diedrich, & Laufer, 1992; Naumann, Maier, Diedrich, Becker, & Bartussek, 1997) have also been reported to yield larger P3/LPC amplitudes relative to

rates. Ehlers et al (1988) report effects of emotional word valence on the response bias measure β .

emotionally neutral stimuli. As the effect of positive emotional valence on ERPs is qualitatively similar to that of negative emotional valence, albeit somewhat smaller, both have been interpreted in terms of the affective stimuli's enhanced motivational significance and arousal value (Johnson et al., 1986; Kayser et al., 1987; Leiphart et al., 1993; Schupp et al., 2000).

A few studies wherein emotional affect was incidental to the primary task, i.e. secondary to same-different judgments, letter counting, or yes/no recognitions, have failed to find effects of emotional arousal on P3 amplitude (Carretié, Iglesias, García, & Ballestros, 1997; Leiphart et al., 1993; Naumann et al., 1997). However, others employing subliminal stimulation (Bernat, Bunce, & Shevrin, 2000) or visual masking techniques (Zimmer & Schmitt, 1987) to prevent attentional processing found that emotional valence modulated ERPs as early as 100 to 400 ms post word onset (see also Carretié et al., 1997; Wong, Bernat, Bunce, & Shevrin, 1997). Thus, it may be that only conscious processing of affect influences P3/LPC amplitude, while unconscious processing has an earlier and perhaps more frontally-distributed influence (Bernat et al., 2000; Zimmer & Schmitt, 1986). This would be consistent with the dramatic cell responses seen in the medial prefrontal cortex of humans to complex aversive pictures within 200 ms of their presentation (Kawasaki et al., 2000).

Design and Hypotheses

This study was designed to investigate the effects of negatively-charged emotional valence as compared to neutral valence on accurate old/new word recognition judgments and on the bias to recognize a word as "old". We expected that subjects would be more likely to say that words of negative emotional value were "old" than neutral words, thereby leading to not only enhanced hit rates but also enhanced false alarm rates. In other words, we expected to find that a word's emotional valence would influence measures of response bias, but not measures of accurate

old/new discrimination. We recorded ERPs to elucidate the processes underlying this emotion-induced recognition bias.

(I) ERP correlates of emotion on *correct old/new discrimination* were defined as valence-related modulations of old/new effects in ERPs to correctly recognized words. This analysis determined whether negative emotional valence had any influence the ERP difference between correctly recognized old items (hits) and correctly recognized new items (correct rejections). This is the standard comparison in ERP studies of recognition memory and reflects processes involved in accurate detection of the old/new difference. If emotional valence alters an individual's ability to distinguish old from new items or the way s/he performs this discrimination, this should appear as a significant interaction between valence and old/new effects. If, on the other hand, emotion has no significant effect on the old/new discrimination, then this comparison should yield comparable ERP old/new effects for negative and neutral words.

(II) ERP correlates of the *valence-induced recognition bias* were defined as the effects of emotional valence on ERPs to items given "old" responses, whether or not they are actually old (i.e., hits and false alarms). If negative emotional valence influences a subject's decision to respond "old", e.g., if it biases them to say "old" or affects their criteria for this decision differently than neutral valence, then traces of this influence should appear in the ERPs associated with "old" responses as either significant Valence effects or as Valence x Old/New interaction effects. Main effects of Valence would indicate that negative valence influenced the decision to respond "old" similarly for old (hits) and new items (false alarms). Valence x Old/New interaction effects would indicate that valence effects were asymmetric for old items that were correctly recognized as "old" (hits) as compared to new items that were incorrectly recognized as "old" (false alarms).

This analysis is based on the logic that investigating the processes involved in the bias to

respond “old” requires a comparison of the ERPs associated with “old” responses in a high-bias versus a low-bias condition. In the present experiment, we expected negative words to represent the high-bias condition and neutral words to represent the low-bias condition. However, this comparison would yield sufficiently process-pure effects of emotion-induced recognition bias only if (1) old/new discrimination performance for the two bias conditions is comparable, and (2) the ERP effects of emotional valence *per se* (i.e., independent of the valence-related response bias shift) are controlled. Any effects of emotional valence in analysis II cannot be interpreted unambiguously in terms of the emotion-induced recognition bias unless we can somehow show that they are *specific* to items classified as “old” as a result of guessing, and not as a function of correct recognition or of valence *per se*. Fortunately, we can estimate the effects of emotional valence on correct recognition as well as the effects of valence *per se* from analysis I. This analysis will reflect the effects of valence on correctly recognized items, independent of the response given. By contrast, analysis II specifically shows the effects of emotional valence on items considered “old”, whether they were in fact old or new. Thus, if analysis I yields no effects of emotion on ERP measures of correct old/new recognition, then any effects of valence that do emerge in analysis II can only be due to the emotion-induced bias to respond “old”.

In line with our expectations for the behavioral data, we hypothesized that negative emotional valence would affect ERPs associated with "old" responses due to a valence-induced shift in the response bias, but would not significantly affect ERPs associated with correct old/new recognition. A critical aspect of these comparisons will be the timing of any ERP valence effects. Those occurring before 500 ms will be within a time range typically affected by unconscious memory and priming processes whereas those occurring after 500 ms (e.g., during LPC) will be within a time range typically viewed as more sensitive to conscious and attentionally-controlled processes. Thus, the timing of the experimental effects will enable us to draw inferences about

the stage(s) of processing at which recognition processes are influenced by emotional valence.

Given findings on the modulatory role of prefrontal areas on memory for emotional events, we adopted the working hypothesis that valence-memory interactions would be more evident in ERPs over frontal than posterior sites. Furthermore, given the evidence for greater right than left hemisphere sensitivity to negative, withdrawal-related emotions (Davidson, 1998; Kayser et al., 1997; Windmann, Daum & Güntürkün, 2000), we expected the ERP effects to be asymmetric.

To our knowledge, there are no previous studies on the effects of emotional valence on brain correlates of the bias to recognize a word as "old". Hence, to pinpoint these effects in time, we performed quasi-continuous F-tests analyzing experimental effects on ERP amplitudes in consecutive 50 ms windows across the recording epoch (1500 ms). Our experimental hypotheses, however, will be tested using ANOVAs of ERP amplitudes measured in early (300-500 ms) and late (500-700 ms) time windows as usually defined in the literature on recognition memory.

RESULTS

Behavioral Data. As expected, both hit and false alarm rates were elevated for negative relative to emotionally-neutral words (see Figure 1). That is, while negative old words were correctly recognized more often than neutral old words, about the same proportion of negative new words relative to neutral ones was also more often falsely recognized as "old". Accordingly, old/new discrimination accuracy (Pr) for negative and neutral words did not differ ($F(1,16)=0.46$) whereas the bias to respond "old" (Br) was significantly higher for the negative words ($F(1,16)=5.30, p<.05$), reflecting the expected *emotion-induced recognition bias*. This difference in bias (negative minus neutral) was not significantly correlated with the overall old/new

discrimination accuracy (i.e., Pr collapsed across negative and neutral items). The overall Br (collapsed across negative and neutral items) also was not significantly correlated with overall Pr , thus supporting the assumption of statistical independence between the two measures (see Snodgrass & Corwin, 1988). All Pearson correlation coefficients were below .20.

An ANOVA of the Reaction Times (RTs) with three repeated factors of Valence (negative/neutral), Response Type (“old” versus “new”) and Response Correctness (correct/incorrect) revealed that “old” responses were somewhat, although not significantly, faster than “new” responses ($F(1,16)=2.98, p<.11$). Correct responses were overall significantly faster than incorrect responses ($F(1,16)=6.03, p<.05$); this effect was accompanied by a significant Response Type by Correctness interaction ($F(1,16)=11.46, p<.005$). Most importantly, the Valence by Response Type interaction was significant ($F(1,16)=23.02, p<.001$). Post hoc tests were performed separately for negative and neutral words to further examine the nature of these interactions. For negative words, there was a significant main effect of Response Type ($F(1,16)=8.52, p<.01$), indicating faster “old” than “new” responses, and a significant main effect of Response Correctness ($F(1,16)=4.73, p<.05$), indicating faster correct than incorrect responses. Furthermore, there was a marginal Response Type by Response Correctness interaction ($F(1,16)=4.27, p<.055$) reflecting disproportionately shorter RTs to correct “old” responses (hits; see Figure 1). For neutral items, there was only a significant Response Type by Response Correctness interaction ($F(1,16)=5.55, p<.05$), reflecting shorter RTs to hits relative to the other responses (see Figure 1). In summary, when words were negative, subjects made significantly faster “old” than “new” responses, whether or not they were correct, whereas for neutral words, “old” responses were faster *only* when they were correct (hits).

--- **Please insert Figure 1** ---

Effects of Valence on ERP Correlates of Old/New Discrimination. Figure 2 shows the grand average ERPs (N=17) for correctly recognized old words superimposed with those to correctly recognized new words for negative and emotionally neutral words. The outcome of corresponding F-tests are provided in Table 1 (A).

--- **Please insert Table 1 & Figure 2** ---

The first train of significant results occurs between 150 and 700 ms post word onset, subsuming both the N400 (between 300 and 450 ms) and the peak of the Late Positive Complex (between 500 and 700 ms). At practically all sites, the ERPs to old items (hits) were more positive than those to new items (correct rejections). No reliable effects of Valence emerged prior to 450 ms; between 450 and 700 ms, however, there was a train of significant Valence effects at a subset of sites as indicated by significant Valence x Site interactions. There were no significant Old/New x Valence interactions within this interval. Figure 4 shows the mean amplitudes in the early (300-500 ms) and late (500-700 ms) time windows typical of ERP research on recognition memory.

--- **Please insert Figures 3 & 4** ---

For the *early time-window (300-500 ms)*, the ANOVA revealed a significant Old/New main effect ($F(1,16)=53.32, p<.0001$) reflecting greater positivity for old than new words (see Figures 3A and 4A). A significant Old/New x Anteriority interaction was also observed ($F(1,16)=4.92, p<.05$), reflecting larger old/new differences over anterior than posterior sites. No other effects were even marginally significant. All effects including Valence were associated with $p>.20$.

In the *late time-window (500-700 ms)*, the Old/New effect continued to be significant ($F(1,16)=13.57, p<.003$), while the Old/New x Anteriority interaction effect remained

marginally significant ($F(1,16)=3.52, p<.08$). In addition, there was a significant Valence x Anteriority interaction ($F(1,16)=22.73, p<.0003$) reflecting greater positivity for negative than neutral words at posterior but not at anterior sites (see Figure 3A right panel).

For the *very late time-window* (900-1200 ms) the quasi-continuous F-tests (Table 1A) also revealed a train of significant Old/New x Site interactions and some isolated effects of Valence and Valence x Site interactions (see Figure 2). An ANOVA on mean ERP amplitudes in this epoch yielded neither any significant main effects nor any effects involving Valence.

Effects of Valence on the decision to respond "old". Figure 5 shows the grand average ERPs to negative (old and new) and neutral (old and new) words given an "old" response (i.e., hits and false alarms). The continuous F-test (Table 1B) indicated a train of significant Valence x Site and Old/New x Site effects between 300 and 500 ms, and some less reliable Old/New x Valence x Site interactions extending up to 550 ms post stimulus onset. There were no significant Old/New main effects until 450 ms, after which there were six consecutive F-tests showing significant Old/New effects, accompanied by some less consistent Valence x Site interaction effects.

--- **Please insert Figure 5** ---

Analysis of the ERP amplitudes in the *early time-window* (300-500 ms) revealed a significant Old/New x Anteriority interaction ($F(1,16)=6.78, p<.02$) as ERPs to old items were more positive than new items anteriorly but slightly more negative posteriorly. The Valence by Anteriority interaction was also significant ($F(1,16)=9.10, p<.009$), reflecting a different pattern of valence effects across the scalp: ERPs to negative words were more positive than those to neutral ones frontally whereas the opposite tendency held over posterior sites (Figure 3B left). Most importantly, there was a significant Old/New x Valence x Anteriority interaction

($F(1,16)=6.52, p<.025$), reflecting a large ERP difference for old versus new items over frontal sites for emotionally neutral words that was virtually absent for negative words (see Figure 4B). The difference was due mainly due to the effect of negative valence on the ERPs to incorrectly recognized new words (i.e., false alarms; see Figure 3B left).²

In addition, a marginal Valence by Hemisphere interaction ($F(1,16)=4.39, p<.055$), indicated that ERPs to negative words were more positive ($\sim.19 \mu\text{V}$) than those to neutral words over the left hemisphere, but less positive ($\sim-.20 \mu\text{V}$) over the right hemisphere. The marginally significant Old/New by Hemisphere interaction ($F(1,16)=3.75, p<.08$) reflected the tendency for the old/new difference to be larger over the left than the right hemisphere (Figure 3A left).

An ANOVA of the ERPs in the *late time-window* (500-700 ms) revealed a significant Old/New main effect ($F(1,16)=8.52, p<.015$), reflecting larger positivity to old than new words (see Figure 3B). The Valence by Anteriority interaction was marginal ($F(1,16)=3.86, p<.07$). Unlike in the early time-window, this interaction now results from more *positive* ERPs to unpleasant items relative to neutral items mainly at posterior sites. The valence-induced positivity was also larger over the left than right hemisphere, as indicated by a significant Valence by Hemisphere interaction ($F(1,16)=7.42, p<.02$). The Old/New by Valence by Anteriority interaction that had been significant in the early time-window was marginally significant in this window ($F(1,16)=3.04, p=.10$) reflecting the tendency for larger old/new effects for neutral than negative items at anterior sites. Figure 4B shows the grand average ERPs elicited by old and new words that subjects considered "old", separately for the negative and neutral words at five left hemisphere medial sites. At prefrontal/frontal sites the ERPs to old and

² For $1 > Pr > 0$, false alarms (FA) are most indicative of all response types of the bias to guess "old", while correct rejections (CR) are least indicative. Thus, when FA and CR are compared, the emotion-induced recognition bias should lead to a larger ERP difference for negative than neutral items. We did find a significant three-way Response x Valence x Anteriority interaction in the early time-window ($F(1,16)=6.46, p<.025$), indicating a larger anterior FA/CR difference for ERPs to negative items relative to neutral items, as

new words clearly differ when they are emotionally neutral but not when they are negative. The old/new difference for neutral items begins around 200 ms at the ventral prefrontal sites³ just as for correctly recognized items (see Figure 4A). This difference peaks between 300 and 500 ms with a mean amplitude of 2.25 μV ($F(1,16)=6.29, p<.025$). At more posterior sites, the old/new difference appears increasingly later and weaker. At medial frontal sites (LMFR in Figure 4B), it does not start before 400 ms poststimulus, and at medial central and occipital sites (LMCE and LMOC in Figure 4B), there is no difference. Note that the opposite is true for ERPs associated with correct responses (Figure 4A) where posterior old/new differences for neutral and negative words become maximal after 500 ms post stimulus onset.

DISCUSSION

We examined recognition memory processes for emotionally neutral and negative words using behavioral speed and accuracy measures and scalp-recorded electrical brain activity. Specifically, we compared the effects of emotional valence on the ERPs for correct old/new word discriminations to those associated with the decision to respond "old" regardless of accuracy. We replicated the *emotion-induced recognition bias effect*: Words with a negative connotation were classified as "old" more often and more quickly than emotionally neutral words, whether or not they were actually old. As indicated by old/new discrimination performance, however, negative

expected. In the late time-window, this effect was marginally significant ($F(1,16)=3.79, p<.07$).

words were not recognized more accurately than neutral words. In fact, the valence-induced recognition bias and old/new discrimination performance were not correlated with each other. Similar bias effects were seen in the responses of individuals with poor and high recognition memory, suggesting that they are unlikely to reflect any controlled, attention-based processes.

The ERP analyses support this interpretation. In short, ERPs associated with correct recognitions showed typical old/new effects, essentially unaffected by emotional valence until quite late. In contrast, valence affected the ERPs associated with "old" responses much earlier, in a latency range typically more sensitive to automatic, unconscious memory processes than to controlled, conscious ones. In this time window, only neutral (and not negative) words showed an ERP old/new difference at frontal/prefrontal sites. We elaborate on these findings below.

ERPs to words correctly identified as "old" were more positive than those to words correctly identified as "new" from 150 to 700 ms post word onset. This is the typical old/new effect observed in ERP studies of recognition memory (Allan et al., 1998; Johnson, 1995; Rugg, 1995). It was broadly distributed with a frontal maximum between 300 and 450 ms. Taken at face value, this pattern is consistent with the proposal that intentional item retrieval is initiated by the prefrontal cortex (Buckner, 1996; Tomita et al., 1999). The results of two ERP studies in which relatively process-pure reflections of recollection were obtained (Paller & Kutas, 1992; Allan, Doyle, & Rugg, 1996) suggest that the old/new divergence starts at ~250-300 ms post-stimulus at frontal and prefrontal sites, and influences the amplitude of the subsequent posterior LPC *only*

³ There was an old/new difference in the N1 region of neutral items that seemed to be due to differences in prestimulus noise and the potential built up prior to stimulus onset given that stimuli occurred at a fixed rate. ERPs to hits are more positive than those to all other response types prior to stimulus onset (see Figure 4). This early difference could spuriously enhance the later old/new effects making it difficult to pinpoint its onset. This difference is attributable to three subjects. We thus repeated all relevant analyses i) excluding the data of these 3 subjects, and ii) using a 100 ms pre-stimulus baseline in all subjects. Both these analyses eliminated the early differences while leaving the relevant effects between 300-500 ms and 500-700 ms intact. For the analysis with the three subjects excluded, the Old/New x Valence x Anteriority interaction in the analysis of "old" responses was significant ($F(1,13)=6.29, p<.03$) in the early and the late time-windows ($F(1,13)=5.43, p<.04$). For the 100 ms baseline analysis, it was significant in the early time-window ($F(1,16)=4.84, p<.05$).

if studied items are consciously discriminated from new ones (c.f. Allan et al., 1998). More importantly for present purposes, the first effects of emotional valence on old/new effects for correctly recognized words did not appear before 450 ms post stimulus onset. Around this time, negative words elicited more positivity over posterior sites than neutral words, in line with previous findings on the effects of emotion on the LPC (Johnston et al, 1986; Naumann et al., 1992; Palomba et al., 1997, Schupp et al., 2000). This effect of emotion was slightly more pronounced in the ERPs associated with correct “old” (hits) than correct “new” responses (correct rejections; see right side of Figure 3A). However, old/new effects continued to be significant in this region, suggesting that processing emotional valence did not disrupt or otherwise influence successful old/new discrimination processes. The only evidence of interactions between emotional valence and the old/new status of the items appeared quite late (between 950 and 1100 ms post stimulus; see Table 1 and Figure 2, especially at the prefrontal sites), by which time most of the recognition decisions had already been rendered (as indicated by average reaction times). The relative lateness of this interaction suggests that it may be part of a post-retrieval verification process (c.f. Donaldson & Rugg, 1999; Maratos et al., 2000).

This conclusion is only partly consistent with the results of a similar study by Maratos et al. (2000) who found reduced old/new effects for correctly recognized negative compared to neutral words not only in a late frontal slow wave (1100-1400 ms) but also earlier in the region of the LPC (500-800 ms). This reduction is not surprising, however, given that in their study (unlike ours), old/new recognition accuracy, not just bias, was affected by valence: it was poorer for negative than neutral items. Both effects may be due at least in part to the greater semantic interrelatedness ("cohesiveness") among the negative than the neutral items (see below). By contrast, in our data, where accuracy was unaffected by valence, ERP old/new effects associated with correct recognition were also unaffected by valence for almost a second after stimulus onset.

This suggests that the emotional dimension was considered relatively late when subjects successfully discriminated between old and new items.

A very different picture emerged when we examined the effects of emotional valence on ERPs to words that subjects considered "old" – the very effects that reflect the neural processes leading subjects to classify negative words as "old" more often than neutral words⁴. In this analysis, ERPs were affected by emotional valence as early as 300 ms poststimulus. At posterior sites, ERPs showed some sensitivity to a word's emotional valence (greater negativity for negative words). At frontal sites, emotional valence interacted with the old/new status of the items: ERPs to neutral words exhibited a marked old/new difference (greater positivity for old words) over prefrontal/frontal sites, broadly consistent with the results of Walla et al. (2000), while the ERPs to negatively-charged words did not show any old/new effects over frontal sites. This difference between negative and neutral items cannot be attributed to differences in old/new discriminability, nor to ERP effects of emotional valence *per se*, because the effects of these effects are negligible in this latency range as the analysis of the correct responses showed. Instead, the interaction was mainly due to a larger positivity to new items considered "old" – i.e., to unstudied negative items that elicited false alarm responses, the response type that is most indicative of the tendency to guess "old" when recollection fails. Thus, this finding can only reflect emotion-related influences on the bias to respond "old".

Between 500 and 700 ms, we observed reliable old/new ERP effects for both negative and neutral items together with some interactions involving emotional valence. These effects were similar to those seen in the analysis of correct responses in this latency range. Surprisingly, the valence effects were somewhat more pronounced in the left than right hemisphere in both analyses, perhaps because the materials were verbal and presented in an overlearned visual

⁴ Maratos et al. did not perform this analysis as they did not look at ERPs associated with false alarms.

format (c.f. Windmann et al., 2000, and Phelps, LaBar, & Spencer, 1997).

In sum, it seems that the enhanced bias to classify items as "old" when they are emotionally negative as opposed to neutral was associated with relatively early (300-500 ms) ERP effects. It is during this same latency range that ERPs typically show a sensitivity to both unconscious memory processes (Rugg et al., 1998; Paller et al., 1995), and unconscious or incidental processing of emotional valence (Bernat et al., 2000; Carretié et al., 1997; Zimmer & Schmitt, 1987). Conscious recollection (Rugg et al., 1998; Paller & Kutas, 1992; Allan et al., 1996) and focused processing of emotional valence (Naumann et al., 1997), by contrast, usually modulate later (~500-700 ms) portions of the ERP, especially over posterior sites. Thus, our results suggest that emotional valence biased participants' recognition memory for words primarily at unconscious, automatic rather than at conscious, strategic levels of processing. This is consistent with the view that negative stimulus valence can "deceive" or "misdirect" information processing at preattentive stages (Windmann & Krüger, 1998; Windmann et al., 2000). As this bias is also associated with faster reaction times, it may actually serve an adaptive function, prompting the cognitive system to assign greater significance and a higher priority to the processing of a potentially threatening stimulus compared to a neutral one.

The prefrontal locus of the bias-related ERP effect fits with this hypothesis. Prefrontal areas are known to be crucially involved in the regulation of emotional information processing (Bremner et al., 1999; Paradiso et al., 1999; Rolls et al., 1994) as well as in monitoring and "criterion setting" functions during recollection (Schacter et al., 1998; Swick & Knight, 1999). These areas may automatically switch to a different processing mode whenever limbic regions signal the presence of potential threat (LeDoux, 2000; Windmann, 1998). Cells in the medial prefrontal cortex are informed about the aversive nature of complex pictures by ~150 ms after stimulus onset, mediated perhaps by dopamine (Kawasaki et al. 2000). Within a memory task,

such alarm signals might encourage orbitofrontal regions to relax their tendency to inhibit currently irrelevant memories (Schnider et al., 2000), or to set a more liberal threshold for verifying the anticipated retrieval results offered by memory-related structures in the medial temporal lobes (Swick & Knight, 1999). By allowing emotional stimuli to engage this sort of mechanism, the brain can ensure that biologically significant events are not "missed" or forgotten as readily as are emotionally neutral events.

More generally, prefrontal cortical responses in emotional contexts as discussed here might reflect the active withdrawal/removal of inhibition over impulsive cognitive, behavioral and physiological fight-or-flight reactions that are normally under top-down control. Indeed, it is of some interest to find out whether such "disruptions" of controlled cognitive processes by fearful stimuli are stronger, more enduring, and/or more generalizable across stimuli of differing emotional valence in various clinical populations. Of particular interest are patients with anxiety disorders (Reimann, 1997; Windmann, 1998), post-traumatic stress disorder (Bremner et al., 1999), and depression (Drevets, 1998), as it has been suggested that these individuals show information processing biases (e.g. Beck & Clark, 1997) and disinhibition of anxiety, presumably due to prefrontal dysfunction (c.f. Bremner et al., 1999; Davidson, 1998; Gorman et al., 2000; Reiman, 1997; Windmann, 1998). Similarly, we might expect that individuals with psychopathy (e.g. Kiehl et al., 1999) whose information processing is often described as "cold" and less empathetic than normal, will show weaker or perhaps no effects of negative emotional valence on prefrontal functioning in various cognitive tasks.

An important issue with regards to our findings relates to the distinction between emotional valence and arousal. Empirical research has shown that negative emotional valence is positively correlated with arousal (e.g. Bradley et al., 1992). The positive effects of affect on memory consolidation are usually attributed to emotional arousal, and not to emotional valence (Cahill &

McGaugh, 1998; Cross, 1999; Bradley et al., 1992; Phelps et al., 1997; McGaugh, 2000). However, whether this is also true for the effects of emotion on memory retrieval processes is less clear. We have referred to emotional valence rather than to emotional arousal throughout this report because we are interested primarily in emotion-related information processing patterns, not in processes associated with emotional *experiences*. We purposely used words that are only mildly negative in connotation, rated 3.32 on a 7-point scale, and thus not physiologically arousing. Our participants were exposed to these words on a computer screen in a completely safe and neutral context for almost an hour. Moreover, they were asked to focus only on whether the words were old or new, so their attention was not explicitly drawn to the emotional meaning of the stimuli. Using similar procedures, Phelps et al. (1997) did not observe any enhanced arousal in their subjects as indicated by skin conductance responses (SCR) – in fact, neutral words elicited significantly *larger* SCRs than did emotional words. All in all, we believe that our stimuli probably did not induce any significant physiological arousal in our subjects. Hence, we feel safe in interpreting the observed effects in terms of emotional valence rather than arousal. At the same time, we note that our negative words do differ from the neutral ones in their arousal value in a purely *informational* (i.e., *semantic*) sense insofar as they refer to fight-or-flight related concepts. In that sense, then our results suggest that operating on these words (concepts) in the context of a recognition task is sufficient to activate brain mechanisms that are typically involved in the control of emotional affect, even when these processes are not accompanied by any significant subjective feelings.

We conclude with a discussion of an alternative account for the effects of emotion reported herein that makes recourse to explanations commonly offered for “false” memories. Presented with a list of study words like “attack, ocean, teeth, bite, fish, fin”, subjects often falsely and confidently remember having seen the word “shark”. Apparently, the likelihood of falsely

classifying a new item as “old” in a memory test increases dramatically when this new item is strongly (semantically, associatively, thematically) related to actually studied items (e.g. Nessler et al., 2000; Roediger et al., 1998). Several mechanisms including semantic and associative priming, feature overlap, semantic categorization, source confusion, among others have been proposed to account for this phenomenon; some of which might alter response bias (Miller & Wolford, 1999). Thus, if the negative words in our study are more interrelated than the neutral ones, then it could be argued that the effects we attributed to negative valence are instead due to one of these factors (Cross, 1999; Maratos et al., 2000). As these processes (e.g., priming) not only affect memory but perceptual performance as well, it is a potential confound in all studies including emotional stimuli, regardless of the experimental task used.

However, we believe this not to be a major concern in our study. First, we made every effort to equate the negative and neutral lists for interrelatedness. We included as many sets of semantically related words in the neutral list (e.g., formulate, paraphrase, interpret, verbalize, discuss, describe, articulate, explicate, elucidate, delineate, outline, illustrate, illuminate, clarify, inform, reveal) as in the negative list. Analyses in two publicly available databases indicated that we had succeeded in this attempt. Second, even if we had been unsuccessful in equating the lists, the effects we attribute to emotion cannot easily be explained in terms of either the controlled or the automatic processes typically invoked to account for false memories. Controlled effects would probably have affected ERPs later, i.e., 500 ms or beyond (see Rugg et al. 1998; Paller et al. 1995; Düzel et al. 1997). More automatic semantic priming or categorization processes are also unlikely explanations, as these generally *reduce* the N400 amplitude (Gunter, Jackson & Mulder, 1998; Nessler et al., 2000; Schwartz, Kutas, Butters, Paulsen, & Salmon, 1996), whereas we found that negative words had slightly *larger* N400 amplitudes relative to neutral words, especially over right posterior sites (see Figure 3B left panel).

Finally, it is important to note that we are not claiming that the general pattern of ERP effects we observed are unique to response biases induced by negative emotions. We believe that other variables that may alter an individual's bias to respond "old" are likely to yield a similar pattern of ERP effects, albeit with somewhat different scalp distributions if they are less indicative of prefrontally controlled top-down processes than a recognition task involving emotional stimuli.

METHODS

Participants. Twenty-one subjects were paid ~ \$18 for their participation. Four subjects' data were not analyzed due to excessive eye movements, antidepressant medication, psychiatric diagnosis, or low trial counts. The final sample thus consisted of 17 right-handed, native English speakers (mean age 21, range 18 – 31 yrs; 5 men) with normal or corrected-to-normal vision.

Stimuli. Word lists are shown in the Appendix. 158 verbs with a negative connotation and 158 emotionally-neutral (~90%) or slightly positive (~10%) verbs were chosen, matched for frequency (Kucera & Francis, 1967), word length, and abstractness (using the MRC database, see Wilson, 1988). Since positive and negative items were found to behave similarly relative to neutral items (e.g., Naumann et al., 1992; Palomba et al., 1997; Schupp et al., 2000), if anything, including a few positive items in the neutral list worked against rather than for our hypothesis. After the experiment, a subsample of 11 subjects rated ~50% of the words on a 7-point scale (0=not at all negative; 6=extremely negative). These subjects rated the negative words ($Mean=3.32$, $SD=1.15$) as significantly more negative than the neutral words ($Mean=0.57$, $SD=0.57$; $t(10)=27.75$, $p<.00001$).

We matched the neutral and negative lists on degree of semantic interrelatedness by choosing

related words from the MS-WORD Thesaurus and the Edinburgh Association Thesaurus (<http://www.itd.clrc.ac.uk/Activity/Psych>). We estimated the degree of item-interrelatedness on the two lists from co-occurrence measures in Hyperspace Analogue of Language (HAL) based on a corpus of ~300 million words (Burgess & Lund, 1997) and semantic similarity in the Encyclopedia corpus of ~60,000 words of the Latent Semantic Analysis (LSA; Landauer, Foltz, & Laham, 1998; <http://lsa.colorado.edu>). HAL yielded a total cumulative co-occurrence of all words with every other of 13,254 for the negative and of 14,507 for the neutral list (frequencies collapsed across 2-, 3-, and 4-word windows starting from the critical word moving either forwards or backwards). LSA yielded an average semantic similarity estimate of .049, ($SD=0.078$) for the negative, and of .048 ($SD=0.074$) for the neutral list (collapsed across pairwise comparisons of each word with every other word). Hence, both analyses showed that the negative and neutral lists had about about the same degree of interrelatedness.

Seventy negative and 70 neutral words were assigned to lists A and B, respectively. Participants saw either list A or list B at study (balanced across subjects), and all these words at test in a quasi-randomized order. Thus, 70 neutral and 70 negative stimuli were presented at study, and these words plus 88 new words of each valence type were presented at test.

Procedures. Subjects sat in a comfortable chair in a light and sound attenuated chamber facing a 21" monitor ~1.5 m away. A yellow frame (6 cm x 16 cm) in the center of the screen throughout recording helped subjects maintain fixation. Words were presented in the middle of the frame, in Univers20 font, yellow on a black background, for 400 ms with an interstimulus interval of 2200 ms ($SOA=2600$ ms). Subjects were asked to memorize the words for a subsequent memory test.

After study, subjects performed a lexical decision task (on different stimuli) for ~30 minutes, followed by the recognition memory test wherein they indicated whether each word (400 ms

duration) was old or new via button presses by the left and right hand, respectively (balanced across subjects), guessing as needed. Each word appeared 1600 ms after a response was given.

ERP Recordings. The electroencephalogram (EEG) and electrooculogram (EOG) were recorded using tin electrodes, 26 of which were embedded in an elastic cap (see Figure 6). Two additional electrodes (LVPf and RVPf) were attached at left and right "ventromedial" PFC sites (5% of the nasion-inion distance up from the nasion, and 10% of the interaural distance laterally). EEG recordings were referenced to the left mastoid, and re-referenced offline to the average of the left and right mastoids. Vertical eye movements and blinks were recorded with an electrode below the right eye, vertically aligned with and referenced to the right ventral prefrontal (RVPf) electrode. Horizontal eye movements were recorded with electrodes placed at the outer canthi of both eyes.

Signals were amplified (Nicolet SM2000) with bandpass filter of .016 to 100 Hz at 12dB/octave, and digitized at 250 Hz. The recording epoch was 2040 ms (500 ms prestimulus). All trials were scanned offline for artifacts and contaminated trials (~16%) were excluded from further analyses. Blinks were corrected using an adaptive spatial filter developed by A. Dale.

After artifact rejection, average bin trial counts ranged from 10 to 60: means were 37 (hit negative), 32 (false alarm negative), 32 (hit neutral), 26 (false alarm neutral), 37 (correct rejection negative) and 44 (correct rejection neutral). We determined that our results did not depend on low trial counts by repeating all relevant analyses in the thirteen subjects who had at least 17 trials in each bin, and by examining a trial-weighted grand average. In these analyses, the most important effects were even slightly stronger. ERPs were digitally filtered with a bandpass of .2 to 20 Hz.

Data analysis. Data were analyzed with repeated measures ANOVAs. Old/new discrimination

accuracy $Pr (=Hit-FA)$ and the Response Bias $Br (=FA/(1-Pr))$ were computed according to two-high-threshold theory (Snodgrass & Corwin, 1988), where $Hit=probability\ of\ "old"\ response\ to\ an\ old\ item$, and $FA=probability\ of\ an\ "old"\ response\ to\ a\ new\ item$. Mean ERP amplitudes were taken and collapsed across electrode sites to constitute the Hemisphere (left/right) and Anteriority (frontal/posterior) factors as depicted in Figure 6.

--- **Please insert Figure 6** ---

REFERENCES

- Allan, K., Wilding, E.L., & Rugg, M.D. (1998). Electrophysiological evidence for dissociable processes contributing to recollection. *Acta Psychologica*, 98, 231-252.
- Allan, K., Doyle, M.C., & Rugg, M.D. (1996). An event-related study of word-stem cued recall. *Cognitive Brain Research*, 4, 251-262.
- Bechara, A., Damasio, H., Damasio, A. R. (2000). Emotion, decision making, and the orbitofrontal cortex. *Cerebral Cortex*, 10, 295-307.
- Bechara, A., Damasio, H., Damasio, A.R., & Lee, G.P. (1999). Different contributions of the human amygdala and ventromedial prefrontal cortex to decision-making. *Journal of Neuroscience*, 19, 5473-5481.
- Beck A.T. & Clark, D.A. (1997). An information processing model of anxiety: Automatic and strategic processes. *Behaviour Research and Therapy*, 35, 49-58.
- Bernat, E., Bunce, S., & Shevrin, H. (2000). Event-related potentials differentiate positive and negative mood adjectives during both supraliminal and subliminal visual processing. *Manuscript accepted for publication by International Journal of Psychophysiology*.
- Bradley, M.M., Greenwald, M.K., Petry, M., & Lang, P. (1992). Remembering pictures: Pleasure and arousal in memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 18, 379-390.
- Bremner, J.D., Staib, L.H., Kaloupek, D., Southwick, S.M., Soufer, R., & Charney, D.S. (1999). Neural correlates of exposure to traumatic pictures and sound in vietnam combat veterans with and without posttraumatic stress disorder: A positron emission tomography. *Biological Psychiatry*, 45, 806-816.
- Buckner, R.K. (1996). Beyond HERA: Contributions of specific prefrontal brain areas to long-

- term memory retrieval. *Psychonomic Bulletin & Review*, 3, 149-158.
- Burgess, C., & Lund, K. (1997). Modelling parsing constraints with high-dimensional context space. *Language and Cognitive Processes*, 12, 177-210.
- Cahill, L. & McGaugh, J.L. (1998). Mechanisms of emotional arousal and lasting declarative memory. *Trends in Neuroscience*, 21, 294-299.
- Carretié, L., Iglesias, J., García, M. & Ballestros, M. (1997). N300, P300 and the emotional processing of visual stimuli. *Electroencephalography and Clinical Neurophysiology*, 103, 298-303.
- Christianson, S.-A. (1992a). Emotional stress and eye-witness memory: A critical review. *Psychological Bulletin*, 112, 284-309.
- Christianson, S.-A. (1992b, Ed.). *The Handbook of Emotion And Memory: Research and Theory*. Hillsdale, N.J.: Lawrence Erlbaum.
- Cross, V.L. (1999). *Effects of Semantic Arousal on Memory: Encoding, Retrieval, and Errors* [Dissertation, University of California Davis]. Ann Arbor, MI: UMI Microform.
- Davidson, R.J. (1998). Affective style and affective disorders: Perspectives from affective neuroscience. *Cognition & Emotion*, 12, 307-330.
- Dias, R., Robbins, T.W., & Roberts, A.C. (1996). Dissociation in prefrontal cortex of affective and attentional shifts. *Nature*, 380, 69-72.
- Donaldson, D.I. & Rugg, M.D. (1999). Event-related potential studies of associative recognition and recall: Electrophysiological evidence for context dependent retrieval processes. *Cognitive Brain Research*, 8, 1-16.
- Drevets, W.C. (1998). Functional neuroimaging studies of depression: The anatomy of melancholia. *Annual Review of Medicine*, 49, 341-361.
- Düzel, E., Yonelinas, A.P., Mangun, G.R., Heinze, H.-J., & Tulving, E. (1997). Event-related

- brain potential correlates of two states of conscious awareness in memory. *Proceedings of the National Academy of Sciences (USA)*, 94, 5973-5978.
- Ehlers, A., Margraf, J., Davies, S., & Roth, W.T. (1988). Selective processing of threat cues in subjects with panic attacks. *Cognition and Emotion*, 2, 201-219.
- Gorman, J.M., Kent, J.M., Sullivan, G.M., & Coplan, J.D. (2000). Neuroanatomical hypothesis of panic disorder, revised. *American Journal of Psychiatry*, 157, 493-505.
- Gunter, T.C., Jackson, J.L., & Mulder, G. (1998). Priming and aging: An electrophysiological investigation of N400 and recall. *Brain & Language*, 65, 333-355.
- Hamann, S.B., Cahill, L., McGaugh, J.L., & Squire, L.R. (1997). Intact enhancement of declarative memory for emotional material in amnesia. *Learning & Memory*, 4, 301-309.
- Heuer, F. & Reisberg, D. (1992). Emotion, arousal, and memory for detail. In: S.-A. Christianson (Ed). *The Handbook of Emotion and Memory: Research and Theory* (pp. 151-180). Hillsdale, NJ: Lawrence Erlbaum.
- Johnson, R., Jr. (1995). Event-related potential insights into the neurobiology of memory systems. In: F. Boller & J. Grafman (Eds.), *Handbook of Neuropsychology, Vol. 10* (pp. 135-163). Amsterdam: Elsevier.
- Johnston, V.S., Miller, D.R., & Burleson, M.H. (1986). Multiple P3s to emotional stimuli and their theoretical significance. *Psychophysiology*, 23, 684-694.
- Kawasaki, H., Adolphs, R., Kaufman, O., Damasio, H., Bakken, H., Howard M. III, & Tori, T. (2000). Responses to emotional visual stimuli recorded in human prefrontal cortex. *CNS Annual Meeting Program; Journal of Cognitive Neuroscience (Supplement)*, p. 56.
- Kayser, J., Tenke, C., Nordby, H., Hammerborg, D., Hugdahl, K., & Erdmann, G. (1997). Event-related potential (ERP) asymmetries to emotional stimuli in a visual half-field paradigm. *Psychophysiology* 34, 414-426.

- Kiehl, K.A., Hare, R.D., McDonald, J.J., & Brink, J. (1999). Semantic and affective processing in psychopaths: An event-related potential study. *Psychophysiology*, *36*, 765-774.
- Kleinsmith, L.J., & Kaplan, S. (1963). Paired-associate learning as a function of arousal and interpolated interval. *Journal of Experimental Psychology*, *65*, 190-193.
- Kucera, H. & Francis, N. (1967). *Computational Analysis of Present-Day American English*. Providence: Brown University Press.
- LaBar, K.S. & Phelps, E.A. (1998). Arousal-mediated memory consolidation: Role of the medial temporal lobe in humans. *Psychological Science*, *9*, 490-494.
- Landauer, T.K., Foltz, P.W., & Laham, D. (1998). Introduction to Latent Semantic Analysis. *Discourse Processes*, *25*, 259-284.
- LeDoux, J. (2000). Emotion circuits in the brain. *Annual Review of Neuroscience*, *23*, 155-185.
- Leiphart, J., Rosenfeld, P., & Gabrieli, J.D. (1993). Event-related potential correlates of implicit priming and explicit memory tasks. *International Journal of Psychophysiology*, *15*, 197-206.
- Maratos, E., Allan, K., & Rugg, M.D. (2000). Recognition memory for emotionally negative and neutral words: an ERP study. *Neuropsychologia*, *38*, 1452-1465.
- McGaugh, J. (2000). Memory — A century of consolidation. *Science*, *287*, 248-251.
- Miller, M.B. & Wolford, G. L. (1999). Theoretical commentary: The role of criterion shift in false memory. *Psychological Review*, *106*, 398-405.
- Naumann, E., Bartussek, D., Diedrich, O., & Laufer, M.E. (1992). Assessing cognitive and affective information processing functions of the brain by means of the late positive complex of the event-related potential. *Journal of Psychophysiology*, *6*, 285-298.
- Naumann, E., Maier, S., Diedrich, O., Becker, G. & Bartussek, D. (1997). Structural, semantic, and emotion-focused processing of neutral and negative nouns: Event-related potential correlates. *Journal of Psychophysiology*, *11*, 158-172.

- Nessler, D., Mecklinger, A., & Penney, T.B. (2000). Event-related potentials and illusory memories: The effects of differential encoding. *Cognitive Brain Research (in press)*.
- Ochsner, K.N. (2000). Are affective events richly recollected or simply familiar? The experience and process of recognizing feelings past. *Journal of Experimental Psychology: General, 129*, 242-261.
- Paller, K.A. & Kutas, M. (1992). Brain potentials during memory retrieval provide neurophysiological support for the distinction between conscious recollection and priming. *Journal of Cognitive Neuroscience, 4*, 375-391.
- Paller, K.A., Kutas, M., & McIsaac, H.K. (1995). Monitoring conscious recollection via the electrical activity of the brain. *Psychological Science, 6*, 107-111.
- Palomba, D., Angrilli, A., & Mini, A. (1997). Visual evoked potentials, heart rate responses and memory to emotional pictorial stimuli. *International Journal of Psychophysiology, 27*, 55-67.
- Paradiso, S., Johnson, D.L., Andreasen, N.C., O'Leary, D.S., Watkins, G.L., Ponto, L.L. & Hichwa, R.D. (1999). Cerebral blood flow changes associated with attribution of emotional valence to pleasant, unpleasant, and neutral visual stimuli in a PET study of normal subjects. *American Journal of Psychiatry, 156*, 1618-1629.
- Phelps, E.A., LaBar, K.S., & Spencer, D.D. (1997). Memory for emotional words following unilateral temporal lobectomy. *Brain and Cognition, 35*, 85-109.
- Raine, A., Meloy, J.R., Bihrlé, S., Stoddard, J., LaCasse, L., & Buchsbaum, MS. (1998). Reduced prefrontal and increased subcortical brain functioning assessed using positron emission tomography in predatory and affective murderers. *Behavioral Sciences and the Law, 16*, 319-32.
- Reiman, EM. (1997). The application of positron emission tomography to the study of normal and pathologic emotions. *Journal of Clinical Psychiatry, 58*, 4-12.

- Roediger, H.L., McDermott, K.B., & Robinson, K.J. (1998). The role of associative processes in creating false memories. In: M.A. Conway, S.E. Gathercole, and C. Cornoldi, *Theories of Memory (Vol II)*, pp. 187-245. Sussex, UK: Psychology Press.
- Rolls, E.T., Hornak, J., Wade, D., & McGrath, J.(1994). Emotion-related learning in patients with social and emotional changes associated with frontal lobe damage. *Journal of Neurology, Neurosurgery, & Psychiatry*, *57*, 1518-1524.
- Rugg, M.D. (1995). ERP studies of memory. In: M.D. Rugg and M.G.H. Coles (Eds.), *Electrophysiology of the Mind* (pp. 132-170). Oxford, UK: University Press.
- Rugg, M.D.; Mark, R.E., Walla, P., Schloerscheidt, A., Birch, C.S., & Allan, K. (1998). Dissociation of the neural correlates of implicit and explicit memory. *Nature*, *392*, 595-598.
- Schacter, D.L., Norman, K.A., & Koutstaal, W. (1998) The cognitive neuroscience of constructive memory. *Annual Review of Psychology*, *49*, 289-318.
- Schnider, A., Treyer, V., & Buck, A. (2000). Selection of currently relevant memories by the human posterior medial orbitofrontal cortex. *Journal of Neuroscience*, *20*, 5880-5884.
- Schupp, H.T., Cuthbert, B.N., Bradley, M.M., Cacioppo, J.T., Ito, T., & Lang, P. (2000). Affective picture processing: The late positive potential is modulated by motivational relevance. *Psychophysiology*, *37*, 256-261.
- Schwartz, T.J., Kutas, M., Butters, N. Paulsen, J.S., & Salmon, D.P. (1996). Electrophysiological insights into the nature of the semantic deficit in Alzheimer's disease. *Neuropsychologia*, *34*, 827-841.
- Snodgrass J.G. & Corwin J. (1988). Pragmatics of measuring recognition memory: Applications to dementia and amnesia. *Journal of Experimental Psychology: General*, *117*: 34-50.
- Stormark, K. M., Nordby, H., & Hugdahl, K. (1995). Attentional shifts to emotionally charged cues: Behavioural and ERP data. *Cognition & Emotion*, *9*, 507-523.

- Swick, D. & Knight, R.T. (1999). Contributions of prefrontal cortex to recognition memory: Electrophysiological and behavioral evidence. *Neuropsychology*, *13*, 155-170.
- Tomita, H., Ohbayashi, M., Nakahara, K., Hasegawa, I., & Miyashita, Y. (1999). Top-down signal from PFC in executive control of memory retrieval. *Nature*, *401*, 699-701.
- Walla, P., Endl, W., Lindinger, G., Deecke, L., & Lang, W. (2000). False recognition in a verbal memory task: An ERP study. *Cognitive Brain Research*, *9*, 41-44.
- Wilson, M.D. (1988). The MRC Psycholinguistic Database: Machine Readable Dictionary, Version 2. *Behavioural Research Methods, Instruments and Computers*, *20*, 6-11.
- Windmann S. (1998). Panic disorder from a monistic perspective: Integrating neurobiological and psychological approaches. *Journal of Anxiety Disorders*, *12*, 485-507.
- Windmann, S., Daum, I., & Güntürkün, O. (2000). Dissociation of accuracy and response bias in lexical decision: Do hemispheric differences in the processing of emotional valence require accurate stimulus identification? *Manuscript submitted for publication*.
- Windmann S. & Krüger, T. (1998). Subconscious detection of threat as reflected by an enhanced response bias. *Consciousness and Cognition*, *7*, 603-633.
- Wong, P.S., Bernat, E., Bunce, S., & Shevrin, H. (1997). Brain indices of nonconscious associative learning. *Consciousness and Cognition*, *6*, 519-544.
- Zimmer, K. & Schmitt, R. (1987). Emotionality of words processed at conscious and unconscious level as reflected in event-related potentials (ERPs). In: E. van der Meer and J. Hoffmann (Eds.), *Knowledge aided information processing* (pp. 283-300). Amsterdam, Netherlands: Elsevier.

FIGURE LEGENDS

Figure 1. Behavioral results. *Top:* 'Hit rate' (probability of old items that are correctly classified as "old") and 'false alarm rate' (probability of new items that are incorrectly classified as "old") for negative and neutral words. *Center:* Old/New discrimination accuracy (Pr) and the response bias (Br) for negative and neutral words. *Bottom:* Reaction times associated with correct and incorrect "old" responses (i.e., hits and false alarms (FA)) and with correct and incorrect "new" responses (i.e., correct rejections (CR) and misses).

Figure 2. Grand average ERPs at all recording sites during accurate recognition of emotionally negative and neutral words. ERPs to correctly recognized old items (hits) and correctly recognized new items (correct rejections) are shown. ERPs were digitally filtered with a lowpass of 8 Hz.

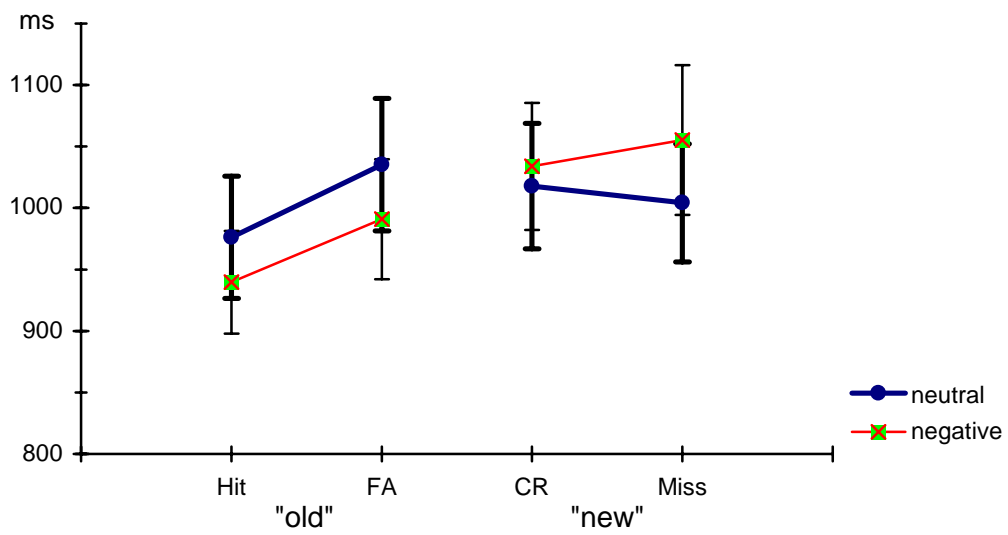
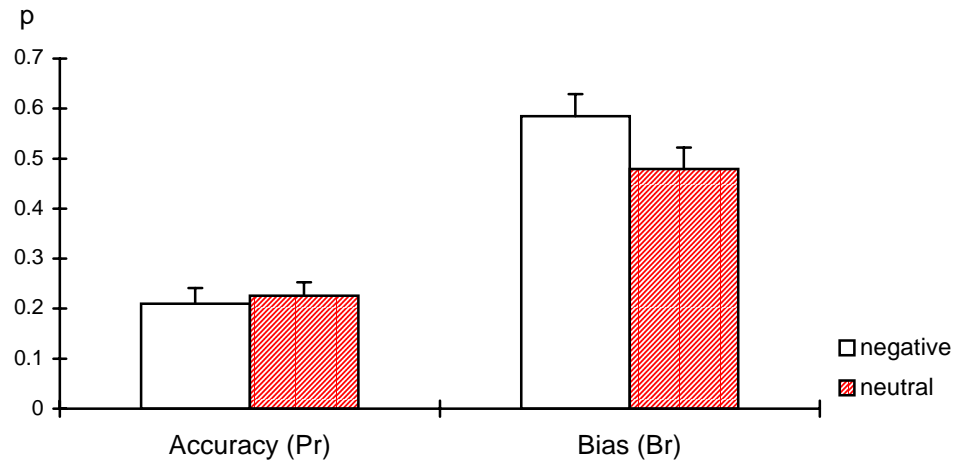
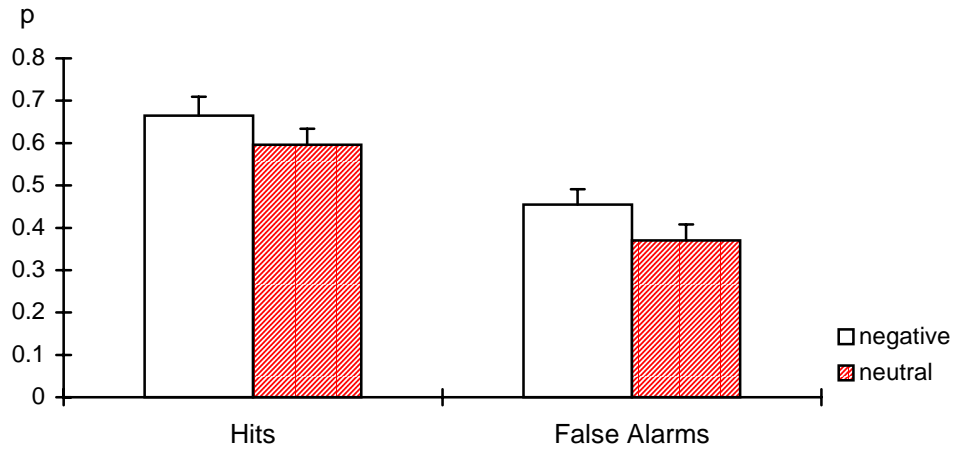
Figure 3. Mean ERP amplitudes measured in the early and late time-windows. Old/new effects for neutral words are compared to old/new effects for negative words. The top panel (**A**) shows these comparisons for correct responses (hits and correct rejections), and the bottom panel (**B**) shows them for words given an "old" response (hits and false alarms).

Figure 4. Subset of ERPs recorded over the left-medial parasagittal midline. **A:** ERPs associated with correct responses to old words (hits) compared to new words (correct rejections), separately for items of negative (left) and neutral valence (right). **B:** ERPs associated with "old" responses which are correct for old words (hits) and incorrect for new words (false alarms), separately for negative (left) and neutral valence (right).

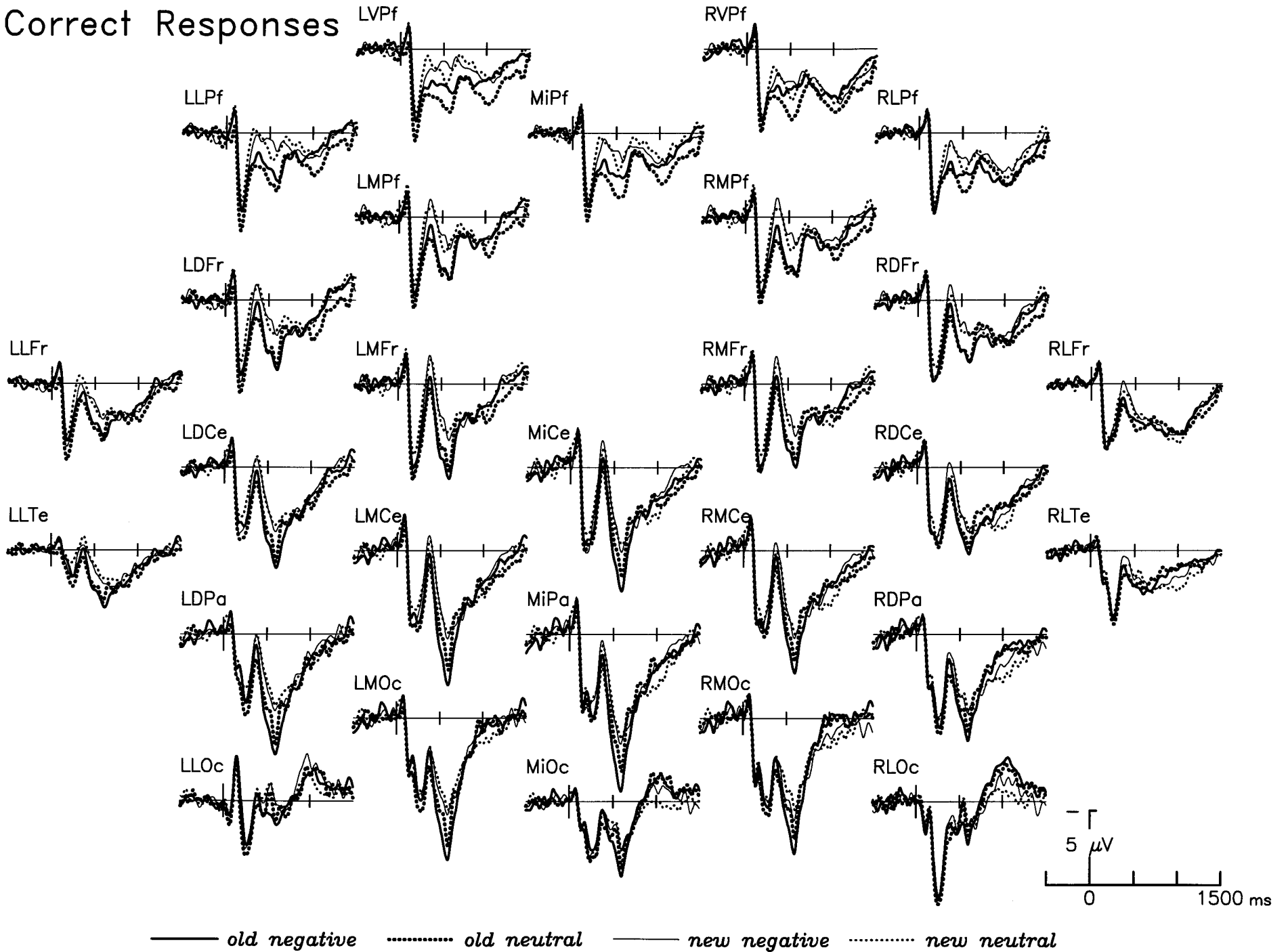
Figure 5. Grand average ERPs associated with the decision to respond "old" to emotionally negative and emotionally neutral words. ERPs to old words (hits) and new words (false alarms) are shown.

Figure 6. Locations of the 28 EEG electrodes. LVPf and RVPf were loose electrodes (not embedded in the cap) placed "ventromedial" to LLPf and RLPf. For statistical analyses, mean ERP amplitudes were taken and collapsed across electrode sites to constitute the factors Hemisphere (left/right) and Anteriority (frontal/posterior) as follows. *left frontal:* left ventral prefrontal (LVPf), left lower prefrontal (LLPf), left medial prefrontal (LMPf), left dorsal frontal (LDFr), left lower frontal (LLFr), left medial frontal (LMFr); *left posterior:* left dorsal central (LDCe), left medial central (LMCe), left lower temporal (LLTe), left dorsal parietal (LDPa), left medial occipital (LMOc), left lower occipital (LLOc); and the same on the right side, respectively: *right frontal* (RVPf, RLPf, RMPf, RDFr, RLFr, RMFr), and *right posterior* (RDCe, RMCe, RLTe, RDPa, RMOc, RLOc).

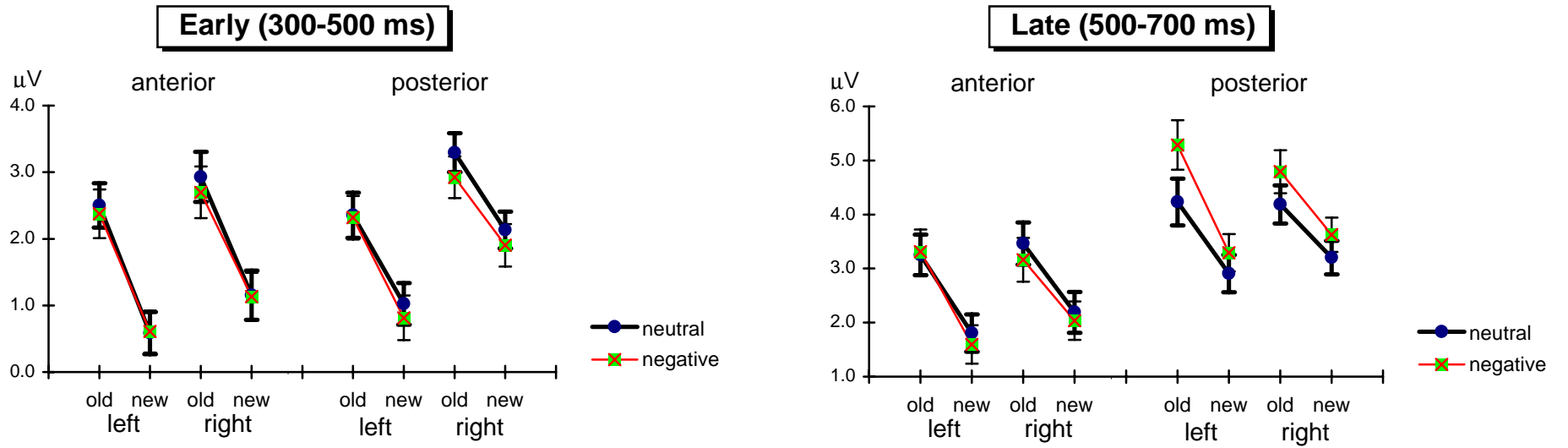
Figure 1



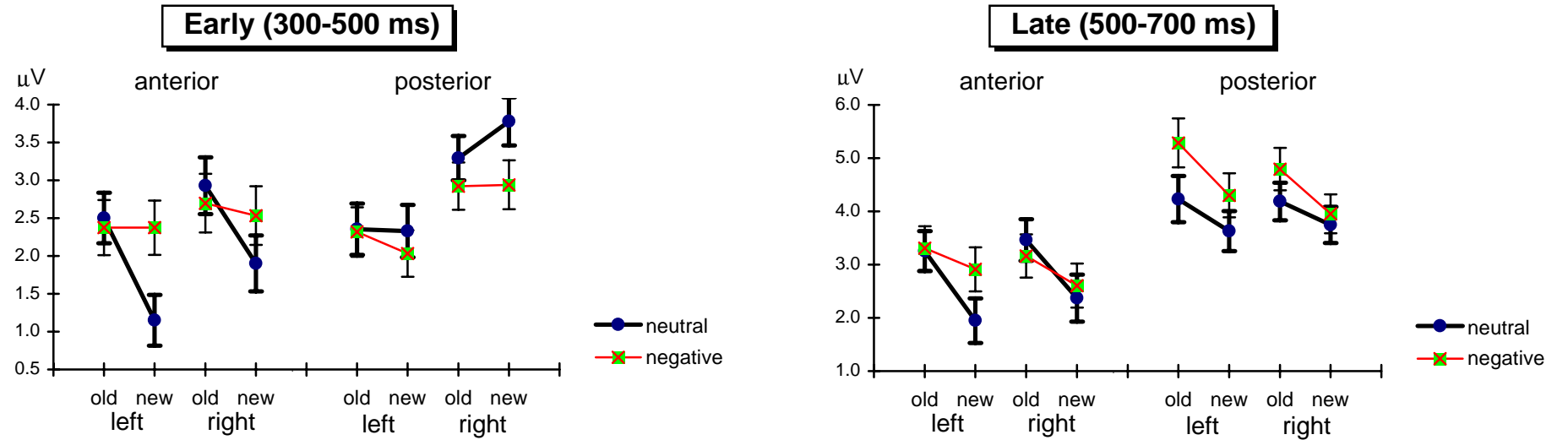
Correct Responses



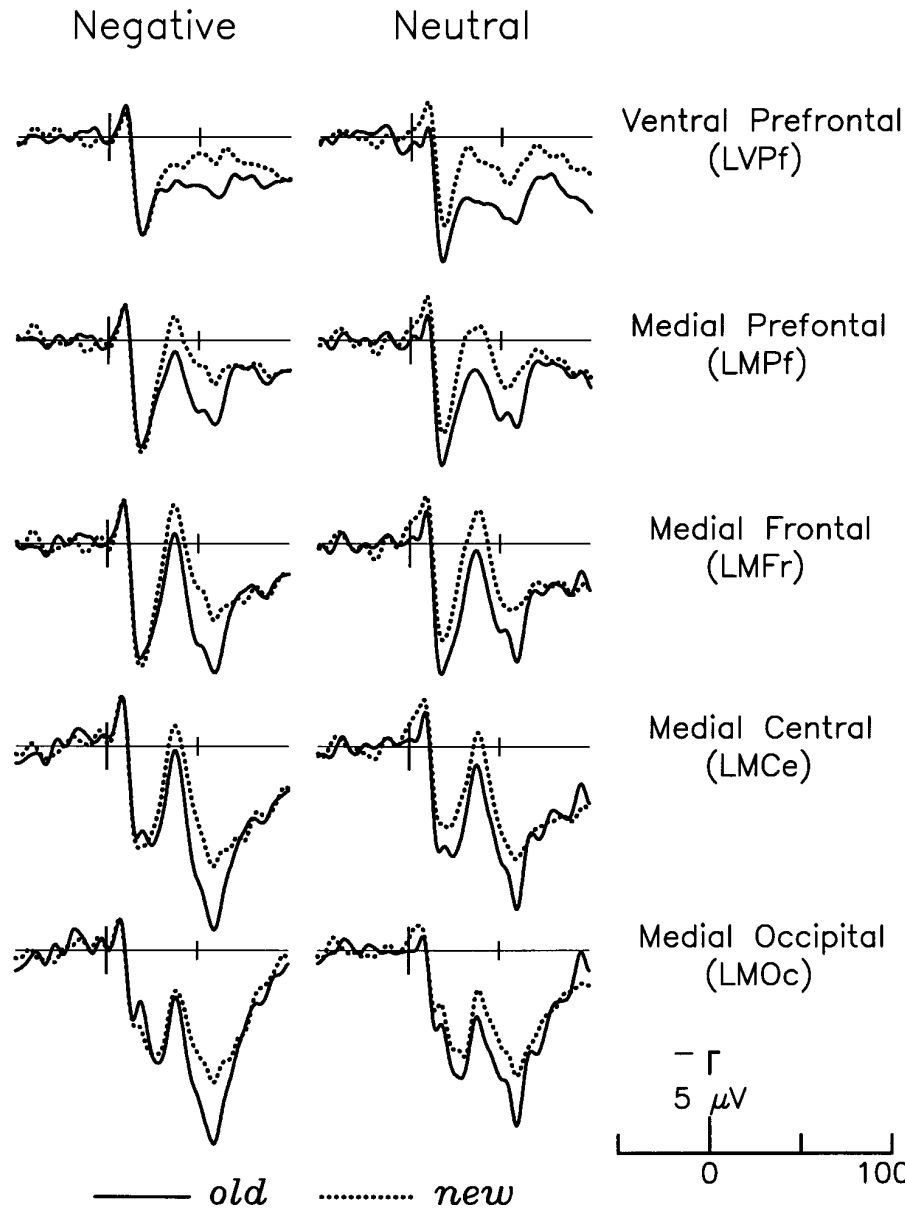
A. Correct Responses To Old And New Items



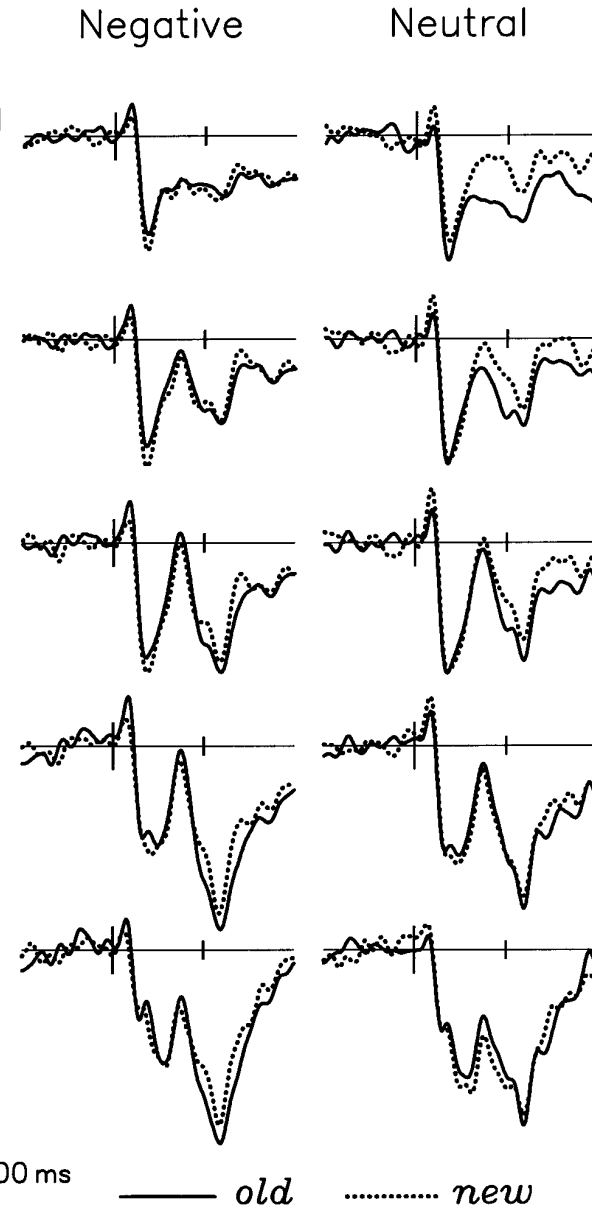
B. "Old" Responses To Old And New Items



A. Correct Responses



B. Old Responses



Old Responses

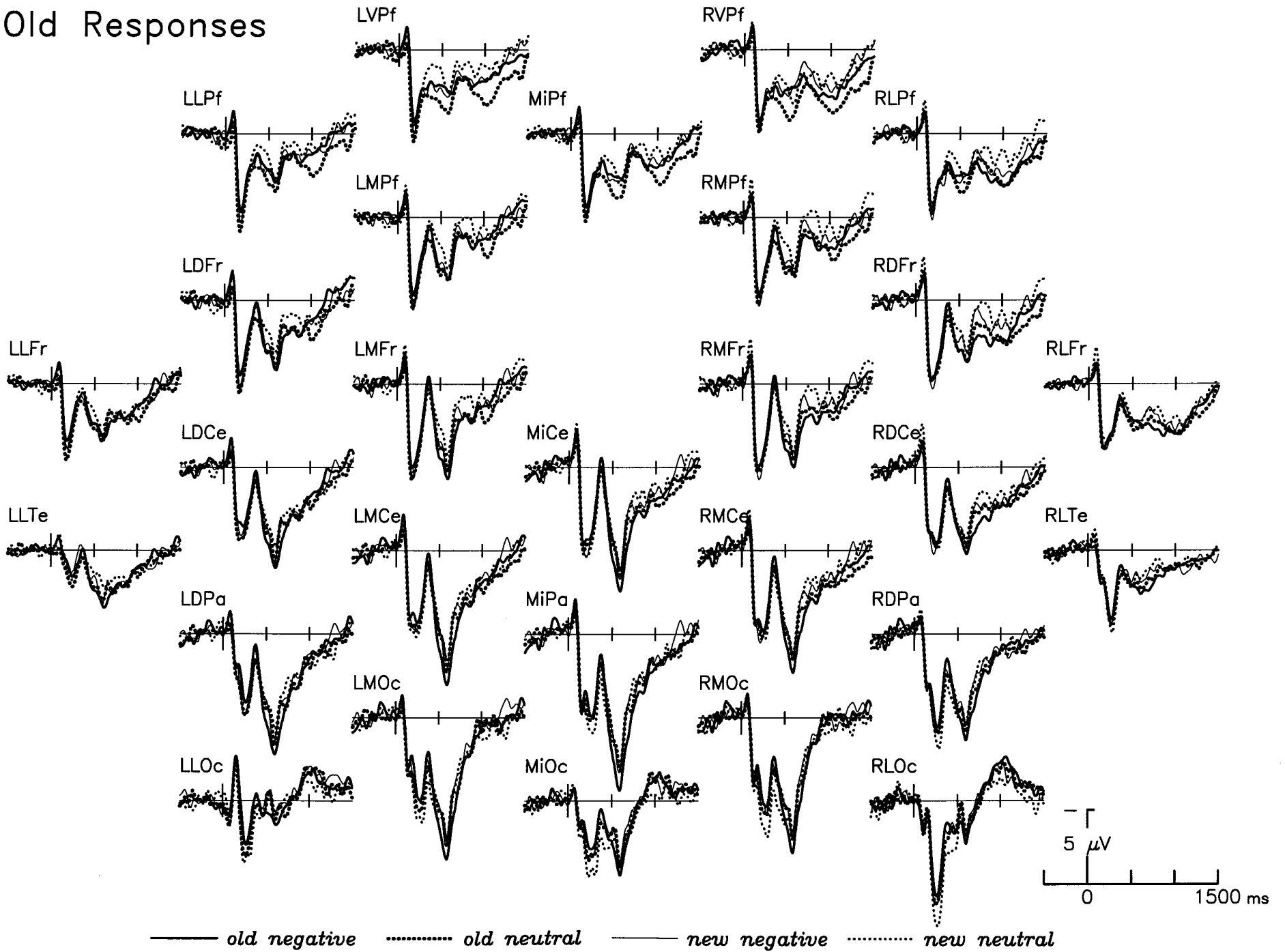
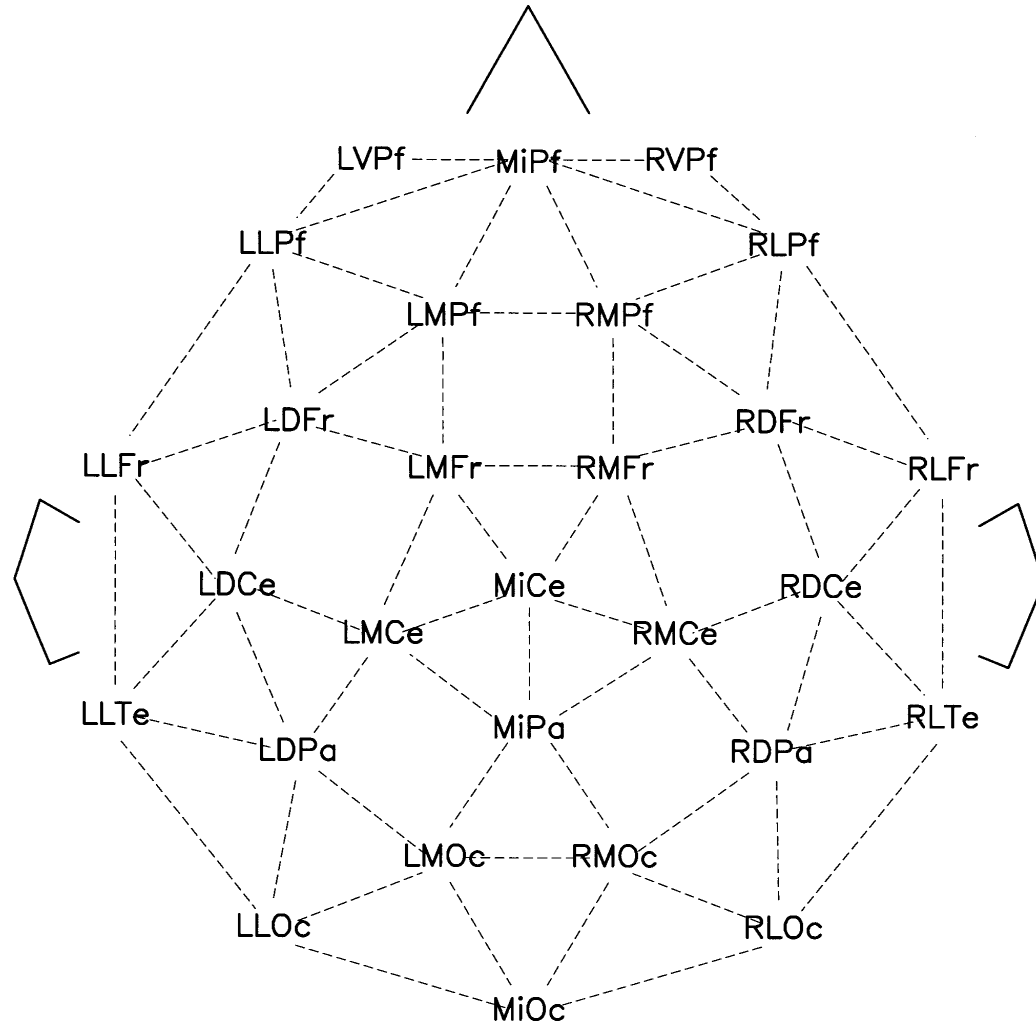


Figure 6



A. Correct Responses	150-450					450-700					900-1200						
Old/new effect				*	*	*	*	*	*	*	*	*					
Valence																*	
Old/New x Valence																	
Old/New x Site				*	*	*	*	*	*				*	*	*	*	*
Valence x Site				*	*	*	*	*	*	*	*	*					*
Old/New x Val x Site													*	*			

B. "Old" Responses	300-450					450-750				
Old/New effect						*	*	*	*	*
Valence										*
Old/New x Valence										
Old/New x Site				*	*	*				*
Valence x Site				*	*	*	*	*	*	
Old/New x Val x Site				*	*	*	*	*	*	*

0 - 50
50 - 100
100 - 150
150 - 200
200 - 250
250 - 300
300 - 350
350 - 400
400 - 450
450 - 500
500 - 550
550 - 600
600 - 650
650 - 700
700 - 750
750 - 800
800 - 850
850 - 900
900 - 950
950 - 1000
1000 - 1050
1050 - 1100
1100 - 1150
1150 - 1200
1200 - 1250
1250 - 1300
1300 - 1350
1350 - 1400
1400 - 1450
1450 - 1500

Table 1: Quasi-continuous F-tests analyzing amplitudes of ERPs associated with correct responses to old and new items, i.e. hits and

correct rejections (**A**); and “old” responses to old and new items, i.e. hits and false alarms (**B**), in 50ms time-steps across the whole recording epoch at all sites. ‘*’

indicates a significant effect at $p < .05$ (Hynh-Feldt corrected). Val=Valence.

Appendix: Word lists

Neutral Stimuli

Both (A+B):	appreciate	protract	estimate	
	immortalize	manifest	liken	
	collect	sketch	confer	
	install	plead	revise	
	gaze	signify	embody	
	designate	convince	marvel	
	reveal	versify	qualify	
	inspire	enunciate	varnish	
	draft	behold	earn	
	illustrate	treat	ravish	
	discuss	festoon	attain	
	describe	cheer	List B:	
	enrapture	inaugurate	commit	
	sponsor	allegorize	introduce	
	accentuate	glaze	brighten	
	compose	contemplate	adjust	
	vaunt	rent	intone	
	elucidate	glorify	modulate	
	honor	spangle	generate	
	renew	hone	compound	
	List A:	prompt	delineate	interpret
		impose	compare	worship
hallow		eternalize	bedeck	
clap		inspirit	clarify	
unravel		applaud	rarefy	
tailor		decorate	tabulate	
hew		induct	negotiate	
edit		adapt	denote	
signalize		restore	visualize	
arise		display	preserve	
prize		adore	refine	
revere		symbolize	doodle	
practice		whittle	transfer	
expose		dip	update	
verbalize		embroider	parse	
persuade		chant	fulfill	
formulate		solve	amaze	
illuminate		explicate	articulate	
impart		hearten	impress	
carve		deploy	melt	
award		conform	bless	
		align	animate	

Negative Stimuli

Both (A+B):	deprive	dispe
	whip	affro
	frustrate	sent
	mortify	pillag
	fluster	gall
	starve	cease
	condemn	tanta
	torment	sland
	antagonize	stunt
	mock	stagg
	offend	List A:
	weaken	scand
	bother	critic
	sicken	banis
	disrupt	harm
	decry	peru
	aggravate	overv
	conquer	damn
	crash	spank
	dishearten	agitat
	worsen	freez
	ruin	upset
avenge	infur	
steal	outla	
excruciate	plun	
ache	ridic	
deject	moan	
curse	revolt	
malign	deport	
alienate	ravage	
toil	expulse	
kick	disqualify	
enrage	hurt	
denounce	madden	
insult	hunt	
punch	sacrifice	
weep	disrepute	
spoil	frighten	
fight	murder	
eliminate	exterminate	
incense	wound	
	intim	

