

## 1 ABSTRACT

2 The Implicit Association Test (IAT) is a reaction time based categorization task that measures the  
3 differential associative strength between bipolar targets and evaluative attribute concepts as an approach  
4 to indexing implicit beliefs or biases. An open question exists as to what exactly the IAT measures, and  
5 here EEG (Electroencephalography) has been used to investigate the time course of ERPs (Event-related  
6 Potential) indices and implicated brain regions in the IAT. IAT-EEG research identifies a number of early  
7 (250-450 ms) negative ERPs indexing early-(pre-response) processing stages of the IAT. ERP activity in  
8 this time range is known to index processes related to cognitive control and semantic processing. A central  
9 focus of these efforts has been to use IAT-ERPs to delineate the implicit and explicit factors contributing  
10 to measured IAT effects. Increasing evidence indicates that cognitive control (and related top-down  
11 modulation of attention/perceptual processing) may be components in the effective measurement of  
12 IAT effects, as factors such as physical setting or task instruction can change an IAT measurement.  
13 In this study we further implicate the role of proactive cognitive control and top-down modulation of  
14 attention/perceptual processing in the IAT-EEG. We find statistically significant relationships between  
15 D-score (a reaction-time based measure of the IAT-effect) and early ERP-time windows, indicating where  
16 more rapid word categorisations driving the IAT effect are present, they are in at least partly explainable  
17 by neural activity not significantly correlated with the IAT measurement itself. Using LORETA, we  
18 identify a number of brain regions driving these ERP-IAT relationships notably involving left-temporal,  
19 insular, cingulate, medial frontal and parietal cortex in time regions corresponding to the N2- and P3-  
20 related activity. The identified brain regions involved with reduced reaction times on congruent blocks  
21 coincide with those of previous studies.

22 **Keywords:** Event-related potentials, EEG, implicit association test, LORETA, brain regions, inhibition, word association, N200

## 1 INTRODUCTION

23 The implicit-association test (IAT) is a measure of implicit bias based on the principle that if  
24 a congruent association between two concepts (e.g. target and stereotypical attribute) is readily  
25 accepted as accurate by a decision maker (e.g. disease → negative), then reaction time (RT) to  
26 categorizing such associations as equivalent is very rapid. In contrast, if an incongruent association  
27 between two concepts (e.g. target and counter-stereotypical attribute) is not readily accepted as  
28 accurate (e.g. disease → positive), then RT is comparatively slower due to inhibitory processes  
29 required to override an automatic tendency to associate congruent concepts. Response bias toward  
30 concept-pairings (e.g. fast responding to congruent; slow responding to incongruent) is not only  
31 influenced by knowledge of concrete characteristics (e.g. perceptual, functional) of bipolar concepts,  
32 but also by how we encode emotional valence in these concept associations though this is not always  
33 apparent in explicit self-report attitude measures (**Greenwald and Banaji, 1995**).

34 The IAT effect or measure of implicit bias is based on the standardised difference (D) between the  
35 mean RT to congruent and to incongruent pairings. A positive D-score indicates that individuals  
36 are either slow to respond to incongruent pairings, fast to respond to congruent pairings or a  
37 combination of both (**Forbes et al. (2012)**). A decision maker's D-score can be used to measure  
38 a range of implicit beliefs reflecting social norms (**Greenwald et al. (1998)**, **Fazio and Olson**  
39 **(2007)**), and these measures have proven effective in predicting later decision-making (**Glasman**  
40 **and Albarracín, 2006**).

41 Opponents of the IAT argue that issues like the low degree to which implicit IAT measures fail to  
42 corroborate with explicit measures such as questionnaires, warrants strong consideration to what  
43 exactly the IAT is measuring (**De Houwer et al., 2009**). This line of evidence has been used  
44 to establish that the IAT reflects automatic beliefs through activation of stereotyped associations  
45 which are often dissociated from self-reported explicit beliefs (**Greenwald et al., 1998**), especially  
46 for socially sensitive topics due to factors such as social desirability (**Hofmann et al., 2005**).

47 Despite reported dissociation between implicit and explicit beliefs, IAT measures show moderate  
48 correlation with explicit measures (**Hofmann et al.**, 2005) and are known to be sensitive to a  
49 number of external influences (**Boysen et al.**, 2006). Such studies ultimately indicate that the  
50 IAT captures meaningful information but its use must be considered with care.

51 An implicit measure of personal connectedness to nature based on latency to bipolar mappings  
52 of targets ('Me', 'Other') and attributes ('Nature', 'Built') is known as the n-IAT. Mean RT to  
53 congruent (e.g. 'Nature-me'/'Built-Other') and incongruent (e.g. 'Built-Me'/'Nature-Other') map-  
54 pings are associated with emotional concern (e.g. anxiety) about the environment (**Schultz et al.**,  
55 2004). **Bruni and Schultz** (2010) reported strong associations in the n-IAT with natural relative  
56 to built environments among a sample of environmentalists. Despite observing similar high scores  
57 on self-reported measures of concern for the environment, significant correlations with explicit mea-  
58 sures were restricted to a participant pool of college students and not environmental activists or  
59 children whom had higher n-IAT scores. **Bruni et al.** (2011) show the n-IAT is robust to framing  
60 effects and valence of the stimuli. In the n-IAT there are 4 categories of words used ('Me', 'Other',  
61 'Nature', 'Built'). In compatible (congruent) trial blocks a participant is instructed to indicate  
62 by button press to which of the two category pairings ('Nature-Me' or 'Built-Other') the stimulus  
63 (word) belongs. In the incompatible (incongruent) trial blocks these category pairing are 'Built-Me'  
64 or 'Nature-Other'. Task switching is understood to exist in both of these block types as partici-  
65 pants must switch between classifying stimuli as attributes ('Nature', 'Built') or self-referential  
66 target categories ('Me', 'Other').

67 The extent to which the IAT effect is caused by involuntary processes independent of the goal to  
68 inhibit pre-potent IAT responding remains unclear (**De Houwer et al.**, 2009). Although a great  
69 deal of research has looked at faking IATs by manipulating response times (**Verschuere et al.**,  
70 2009), there has also been some success in faking IATs merely by being instructed to respond in a  
71 certain way. **McDaniel et al.** (2009) instructed participants to respond as extravertly as possible  
72 on an IAT which measured personality types and found that participants could successfully fake  
73 their level of extraversion. **van Nunspeet et al.** (2014) highlight a related finding where in an IAT  
74 to measure bias towards muslim woman, framing the IAT task as a measure of competence (the  
75 participants ability to process new information) vs. morality (a test to measuring their 'values')  
76 resulted in reduced negative-bias IAT scores when the task was framed in terms of morality. Here  
77 they show how ERPs associated with early perceptual processing, selective attention and social  
78 categorization (namely N1, P150, and N450) are sensitive to this framing effect, further indicating  
79 the role of motivational states in modulating aspects of perceptual attention and conflict monitoring.  
80 In the next subsection we outline previous studies using EEG measures to study the IAT and in  
81 the following subsection we posit our research aims with respect to gaps in the existing studies.

## 1.1 PREVIOUS IAT-EEG WORK

82 Studies that have examined ERPs in IAT tasks (hereafter referred to as IAT-ERPs), have implicated  
83 the late positive potential (LPP) as a component of interest (e.g. **Hurtado et al.** (2009)) or other  
84 later occurring ( $> 300ms$ ) ERP components (e.g. **O'Toole and Barnes-Holmes** (2009)). Those  
85 focusing on both early and late activity, such as the study by **Williams and Themanson** (2011)  
86 investigating the IAT effect in a group bias IAT (gay-straight) reported no differences across IAT  
87 conditions for early components (N1, P2) but found later component differences (N400, LPP) for  
88 concept pairings. This is suggestive that early perceptual and attentional processes might not be  
89 associated with the IAT measurement in their study, but later semantic categorization processes  
90 are responding to congruent/incongruent concept pairings.

91 While many of the previous IAT-EEG studies examine ERP phenomena in response to the IAT  
92 task stimuli, others have explored ERP measures taken from separate (but related) tasks on the  
93 same participants in order to understand the time-course of neural processing of stimuli involved  
94 with implicit bias. **He et al.** (2009) demonstrated a relationship between the IAT effect on a

95 racial IAT and early P2 and N2 categorization of faces (e.g. White, Black, Asian). Here they  
96 found correlations between ERP amplitudes in a racial face categorisation task with an IAT-based  
97 measure of implicit racial bias in a group of non-muslim university students. Additionally, a later  
98 positive component (LPC) was observed for extended same-different race faces. **Ibanez et al.**  
99 (2010) has also shown that early ERP components of race-face processing (e.g. N170 component)  
100 can be modulated by the valence of evaluative attributes used in the IAT such as positive or  
101 negative valence words, and also by the social face categories such as in-group or out-group. This  
102 is suggestive of early integration of contextual information related to racial attitude during face  
103 processing in the IAT.

104 A recent study by **Forbes et al.** (2012) investigating EEG correlates of the IAT effect in an  
105 attempt to examine causal factors, showed more positive ERPs at frontal and occipital regions  
106 at automatic processing speeds, as well as occipital regions at controlled processing speeds, when  
107 responding to congruent versus incongruent pairings. Here they investigate ERP timings as de-  
108 terminants of automaticity in order to gain insight into the timing at which implicit and explicit  
109 processes unfold, as these may be less susceptible (in short duration processes) to control. More-  
110 over, they find higher D-scores (or bias) were identified by greater coherence between frontal and  
111 occipital regions in time periods as early as 92ms with no significant difference present between con-  
112 gruent/incongruent conditions. These findings the authors consider could be indicative of top-down  
113 modulation of attention and perceptual processing. When taken in tandem with lesion study data  
114 they indicate the potential for the facilitated performance seen on stereotypic-congruent blocks to  
115 be associated with more efficient neural processing.

116 A number of other ERP components have been observed in the IAT-EEG such as the P3 which is  
117 associated with a range of cognitive processes one of which is attentional focus on novel, salient or  
118 unexpected to-be attended items or on distractors (unattended items) which produce an orienting  
119 response (**Polich**, 2007). The P3 has also been shown to index explicit attention toward self-referent  
120 material (**Tacikowski and Nowicka**, 2010), and also to direct implicit attention in an IAT toward  
121 self-positivity biased words (**Chen et al.**, 2014).

122 **Williams and Themanson** (2011) demonstrate effects surrounding an N400 ERP where larger  
123 amplitudes are present in incompatible trials compared to congruent trials 'suggesting greater se-  
124 mantic congruency in the compatible condition of the IAT'. They note N400 amplitude for both  
125 congruent and incongruent blocks at FCz is correlated with IAT incongruent-congruent reaction  
126 times with no apparent statistically strong relationships present with respect to reaction times  
127 in either block. The N400 ERP (as being distinct from the error-related N450 (**Folstein and**  
128 **Van Petten**, 2008)) is sensitive to semantic anomalies and violations with structures in the 'im-  
129 mediate vicinity of the auditory cortex' (with a left-hemispheric dominance) associated with the  
130 processing of semantically anomalous sentences (**Van Petten and Luka**, 2006). While the N400  
131 was initially thought to reflect linguistic anomalies and violations, further study has identified its  
132 role in semantic priming (**Deacon et al.**, 2000) and expectancy (**Curran et al.**, 1993). There  
133 is evidence that it too does not reflect a purely automatic process (**Holcomb**, 1988) involving  
134 attentional related factors. **Lau et al.** (2008) identify a dominant (left-hemispheric) pattern across  
135 a range of studies utilising EEG and non-EEG imagining modalities investigating the N400, and  
136 indicate the posterior middle temporal cortex as being the only area to show consistent effects  
137 across studies.

138 A common finding among these studies is that both early and late time regions of the EEG  
139 signal following stimulus presentation, demonstrate effects related to IAT congruency condition  
140 and D-score. Earlier effects typically reflect neural mechanisms at work outside of a post-perceptual  
141 processing time region, namely one that occurs later within a time window following a response  
142 (**van Nunspeet et al.**, 2014; **Guex et al.**, 2011). There is no clear consensus on what ERPs and  
143 related morphologies should be found when examining a new IAT task. For instance, **Fleischhauer**  
144 **et al.** (2014) do not find evidence of significant effects in the N2 (as expected by the authors) or

145 N400 time-range but do find relationships for P1/P3b amplitudes relating to early facilitation of  
146 relevant visual input and efficiency of stimulus categorization.

147 An open question remains from the literature as to what extent the IAT-effect can be modulated  
148 by external and other top-down related factors.

## 1.2 IAT & COGNITIVE CONTROL

149 **Hilgard et al.** (2014) investigated the relationship between the medial frontal negativity (MFN)  
150 during an IAT task as it has been identified as linked to proactive cognitive control. This is distinct  
151 from a neurocognitive process of reactive control due to switching in incongruent and congruent  
152 block trials in the IAT task indicated by ‘a positive voltage deflection over frontocentral scalp  
153 locations’ (D-pos). They posit their analysis and hypothesis in terms of a Dual Mechanisms of  
154 Control (DMC) model (**Braver**, 2012) where proactive control relates to ‘the sustained maintenance  
155 of goal information in working memory that serves to bias information in a goal-congruent manner’  
156 and reactive control ‘a late correction mechanism for dealing with cognitive and behavioural conflict  
157 as it arises’.

158 A common pattern of negatives have been consistently identified in IAT tasks in time regions  
159 related to the N2, N450 and related ERPs. Broadly, these ERP components are understood to  
160 be typically implicated in conflict monitoring processes including proactive and reactive cognitive  
161 control. Such negativities are often referred to as medial frontal negativities (MFNs). Identifying  
162 the onset/offset latencies of IAT-EEG sensitive ERP components like these is made difficult by  
163 their overlapping nature (variable latency), IAT-task parameters/stimulus affecting ERP waveform  
164 characteristics, the multifaceted nature of ACC-generated signals related to cognitive control and  
165 external (e.g. environmental) effects driving top-down attentional modulation.

166 An Error-related Negativity (ERN) typically follows post-error conflict detection between in-  
167 compatible responses and the N2 has been found to reflect this monitoring and conflict detection  
168 function (**Yeung et al.**, 2004). N2 amplitude is modulated by the amount of conflict present  
169 between possible choices prior to response selection and performance. **Chee et al.** (2000) used  
170 functional magnetic resonance imagery (fMRI) on participants completing an IAT and found that  
171 the left dorsolateral prefrontal cortex (dPFC) and the anterior cingulate cortex (ACC) mediate  
172 response inhibition in the incongruent condition. The ACC is particularly sensitive to response  
173 conflict in the IAT, and therefore N2 involvement in IAT performance at least ostensibly reflects  
174 conflict detection/cognitive control processes. Numerous other studies have implicated the ACC as  
175 being involved in the generation of a broad range of conflict-monitoring related ERP components  
176 (**Bekker et al.**, 2005). **Clayson and Larson** (2013) demonstrated the N2 shows reliable conflict  
177 adaptation and these conflict adaptation indices were stable in a 2-week test-retest. **Larson et al.**  
178 (2014) highlight with regard to cognitive control theory and goal-directed behaviour the N2 ‘rep-  
179 represents an empirical marker of both a control mechanisms to handle conflict’ and relevant to our  
180 study ‘a reflection of the level of cognitive control implemented during the ... task’. They highlight  
181 other issues which can confound interpretations of N2 amplitude such as in flanker trials where N2  
182 amplitude being ‘sensitive not only to the degree of conflict for a given stimulus but also to the  
183 extent to which task-irrelevant information is processed’. An important distinction between these  
184 studies (e.g. Eriksen flanker task) is that in IAT-EEG reactive control, conflict related changes  
185 arise as a result of task switching within condition blocks. Here, a participant is required to change  
186 between categorisation of stimuli as evaluative (Nature, Built) or target categories (Me, other) with  
187 no knowledge of the upcoming trial type **Hilgard et al.** (2014).

188 Given the overlapping time regions of these early negative components in existing IAT studies,  
189 one of the aims of our study introduced later was to disentangle ERP activity in these time win-  
190 dows using LORETA source analysis to explain patterns of correlated ERP activity with respect

191 to implicated cortical generators of the N2/MFN and N400-related ERPs (cingulate cortex and  
192 temporal lobe structures).

193 This evidence would suggest ERP activity manifesting negatively in the 250ms-450ms time range  
194 indexes a range of distinct neural processes related to cognitive control. Moreover, as activity in  
195 this time range is understood to be involved in proactive control processes, we posit the relation-  
196 ships observed to D-score (without an apparent explanation based purely on reaction time) may  
197 be indicating participant variability with regard to enhanced motivational/attentional aspects to  
198 perform the task 'as quickly as possible', thus engendering the measurement of an IAT-effect. **Jodo**  
199 **and Kayama** (1992) show that the N2 amplitude is enhanced by reaction time constraints in a  
200 go/no-go task varying in amplitude depending on the neuronal activity required for response in-  
201 hibition, indicating increased amplitudes are related to a 'greater effort' needing to be employed  
202 in tasks where response inhibition is constrained by reaction-time constraints. Such reaction-time  
203 constraints are integral to the measurement of the IAT effect where faster responding is presumed  
204 to be less susceptible to being faked.

205 Given evidence that groups typically have positive IAT scores on the n-IAT, we suspected those  
206 participants with lower D-scores (less standardised difference in reactions times) might be engag-  
207 ing in the a with different (less) motivational effort and consequently not engendering conditions  
208 needed to capture reaction time effects underlying implicit associations. Other authors highlight  
209 issues suggesting IAT effect measurement is potentially related to the 'degree of involvement of  
210 the participants' (**Vargo and Petroczi**, 2013). **Agosta et al.** (2013) suggest a neutral D-score  
211 window for scores between -.2, .2 where results are 'inconclusive' i.e. have low accuracy.

### 1.3 AIMS AND RESEARCH QUESTIONS

212 In summarising previous related work on IAT-EEG issues we see that an open question exists as  
213 to what exactly the IAT measures given its susceptibility to be sometimes difficult to relate to  
214 explicit measures of attitude. In the study reported in this paper, we highlight a potential issue  
215 here, namely that successful measurement of an IAT effect likely involves factors of participant  
216 motivation to engage in the task such that some participants might be more likely than others to  
217 produce an IAT effect.

218 Accordingly, in this study we examine how ERP measures in the IAT might offer insight into  
219 the neural mechanisms underlying the more rapid associations that drive IAT effects. Of primary  
220 interest to our work was examining how ERP measures underlying both congruent and incongruent  
221 block types could offer evidence on the neural mechanisms underlying more rapid associations  
222 driving the IAT effect. Our hypothesis is that such a shared relationship would exist, further  
223 implicating proactive cognitive control and top-down modulation of attention/sensory processing  
224 (involved with potential motivational factors) in biasing the IAT measurement.

225 In summary, the aim of our study was 1) to examine how ERP measures in the IAT might offer  
226 insight on the neural mechanisms underlying the more rapid associations that drive the IAT effect  
227 and 2) given such relationships, to investigate potential sources of correlated brain-activity using  
228 LORETA as previous IAT research has implicated a range of early negativities overlapping in time  
229 and scalp topography serving arguably different processes in the IAT.

## 2 MATERIAL & METHODS

### 2.1 PARTICIPANTS

230 30 participants aged between 18 and 45 years were recruited through advertisement using Dublin  
231 City University staff and undergraduate/postgraduate email lists. In total, 8 participants were

232 excluded due to noisy EEG, high error response rates, or to ensure a counter-balanced block  
233 design for the order of congruent/incongruent IAT blocks (i.e. same number of congruent first and  
234 incongruent first block orderings). The 22 remaining participants were predominantly right handed.

## 2.2 APPROVAL BY UNIVERSITY ETHICS COMMITTEE

235 This study was carried out in accordance with the Declaration of Helsinki - Ethical Principles  
236 for Medical Research Involving Human Subjects, and also Dublin City University's guidelines on  
237 Best Practice in Research Ethics with informed written consent from participants. The study  
238 was approved by Dublin City University's Research Ethics Committee (DCU REC/2013/205). All  
239 participants gave written informed consent.

## 2.3 IMPLICIT ASSOCIATION TASK

240 Participants in this study completed a modified version of the IAT outlined by **Bruni and Schultz**  
241 (2010). The purpose of this was to measure the strength of the association between the 'Me'  
242 target category and two evaluative attribute categories ('Nature', 'Built') relative to the associa-  
243 tive strength of the 'Other' target and attribute categories. The experiment consisted of 7 blocks  
244 with 2 blocks (of 32/48 trials respectively) measuring congruent association RTs (e.g. 'Nature-  
245 Me', 'Built-Other'), and another 2 blocks (of 32/48 trials respectively) measuring incongruent  
246 association RTs (e.g. 'Nature-Other', 'Built-Me'). The remaining blocks were practice blocks.  
247 Congruent/Incongruent ordering was counter-balanced across participants. In total, 80 trials were  
248 collected for each of the congruent and incongruent mappings respectively.

249 Participants were required to sort stimuli into category pairings of 'Me-Nature', 'Other-Built' in  
250 the congruent case and 'Me-Built', 'Other-Nature' in the incongruent case. Each participants' name  
251 was used in conjunction with the 'Me' category and a random list of other names for the 'Other'  
252 category. 'Tree', 'Mountain', 'Butterfly' and 'Flower' were used as stimuli for the 'Nature' category.  
253 'Boat', 'Car', 'Chair', 'Truck' were used as stimuli for the 'Built' category. Category pairings (e.g.  
254 'Me-Nature, Other-Built') appeared in the upper left and right corners of the screen respectively.  
255 Participants indicated the category to which stimuli belong by pressing keys '1' (pairing appearing  
256 on left side of screen) or '2' (pairing appearing on right side of screen) on the keyboard with their  
257 dominant hand.

258 Stimuli were preceded with a 1 second central fixation located centrally on-screen. A word stim-  
259 ulus was presented on screen (centrally) until a response was given. A feedback screen appeared  
260 post-response based on correct/incorrect categorization. Participants were not required to cor-  
261 rect response errors. In Figure 2 we show an example of this trial structure within a compatible  
262 (congruent)-mapping block.

263 D-score was calculated for each participant as the difference in mean reaction time between  
264 trials from pooled incongruent and pooled congruent blocks (incongruent-congruent) divided by  
265 the pooled standard deviation of trials from both block types (**Greenwald et al.**, 2003).

## 2.4 EEG RECORDING AND ANALYSIS PRELIMINARIES

266 EEG was recorded using a 32-channel ActiCHamp recording system with a 10-20 ActiCap. A  
267 virtual ground was used as an online reference and later re-referenced offline to a digitally linked-  
268 mastoids reference (TP9 + TP10). Prior to this, signals were filtered with an FIR sinc hamming  
269 window filter to between 0.1 Hz and 30Hz. ICA (Independent Component Analysis) was used to  
270 remove artifacts such as eye-blinks and eye-blink related components in particular those described  
271 by **Plöchl et al.** (2012) as CRD (corneo-retinal dipoles), eyelid and eyelid-CRD artefacts. These  
272 were identified from scalp topography and amplitude characteristics and similarly confirmed using  
273 EOG (Electrooculogram) channels VEOG and HEOG. ICA weights were trained on 1Hz to 30Hz

274 filtered data and then applied to the 0.1Hz-30Hz band-passed signals. Analysis revealed strong pre-  
275 stimulus activity related to block conditions. Subsequently, ERP averages were generated on epochs  
276 extracted from signals band-passed to between 4Hz and 30Hz (we explain this in the next section).  
277 Epochs were extracted from -200ms to +1,000ms with respect to the onset of word stimuli to be  
278 categorised for compatible and incompatible blocks. Trials (post ICA clean-up) which exceeded  
279 70mV or contained other noise-like artifacts were discarded. This resulted in a maximum of 7.5%  
280 trial loss across participants with one subject exceeding this at near 20%.

281 EEG recording was carried out in an electrically shielded environment. Participants were seated  
282 approximately 70cm from the screen and reported no issues reading word associations.

283 EEGLAB 13.32 was used for EEG pre-processing and clean-up. Scalp and statistical scalp plots  
284 including grand average ERP plots were generated in IPython. SPSS 21 was used for conduct-  
285 ing repeated-measures ANOVA. sLORETA 20081104 was used for source localization measures  
286 (Pascual-Marqui, 2002).

287 *2.4.1 Baselineing* A baseline of  $-200ms$  to 0 pre-stimulus was initially planned for stimulus-  
288 locked epoch extraction. However, upon analysis of the EEG it was found that a CNV (Contingent  
289 Negative Variation)-like component was present surrounding stimulus onset for many participants,  
290 and upon further inspection differentiated between congruent and incongruent conditions.

291 Baselineing serves to remove noise sources like inter-subject differences and slow-drifts, thus allow-  
292 ing inter-subject measures to be comparable as ERP amplitudes align relative to a zero measure  
293 (baseline) across electrodes and participants. An expectation here is that pre-stimulus baseline ac-  
294 tivity is not systematically affected with respect to factors or conditions being measured. Herein  
295 the issue exists with the IAT experimental structure, that is, a participant is aware of the upcoming  
296 condition type and thus may, through the recruitment of different cognitive preparatory mecha-  
297 nisms for that stimulus type (congruent/incongruent), introduce it into the baseline period activity  
298 which could systemically affect the correct baselineing of later ERP components. This is particularly  
299 relevant as pre-stimulus activity diminishes during the epoch window. One such ERP component  
300 typically seen, a CNV-E (Contingent Negative Variation), is present following a warning stimulus  
301 (S1) such as a fixation cross indicating upcoming stimulus (S2) and results in an expectant pattern  
302 of activity locked to S2.

303 In our study, when examining ERP activity in these early time regions in individual participants'  
304 data plots (without baselineing across a range of incrementally high-pass filtered signals) we found  
305 there was a general trend of pre-stimulus activity extending into early periods of the ERP waveform,  
306 overlapping notably with the P1. Other IAT-ERP studies might not have considered or encountered  
307 such issues with some studies not citing whether a baseline was used (Hurtado et al., 2009;  
308 Barnes-Holmes et al., 2004; Egenolf et al., 2013), others where a prestimulus baseline was  
309 used (O'Toole and Barnes-Holmes, 2009; Williams and Themanson, 2011; Hilgard et al.,  
310 2014) and others where a prestimulus baseline was used but measures were taken to lessen the  
311 impact of pre-stimulus activity such as post-movement ERP activity related to previous trials  
312 Forbes et al. (2012).

313 We found ERP average waveforms from participants without pre-stimulus baselineing indicate  
314 these differences in some instances do not degrade until  $\sim 150ms$  and would suggest these differences, if  
315 included in the baseline measurement, could systematically affect later ERP components, resulting  
316 in these component time-windows containing ostensible effects.

317 Time-frequency decomposition of epochs and related ITC (Inter-Trial Coherence) (Makeig et al.,  
318 2002) revealed that pre-stimulus activity is comprised of contributions across a wide range of  
319 frequencies. Inspection of stimulus-locked ICA components revealed these patterns are not well  
320 captured by a single set of ICs that are not entangled with other post-stimulus trial-locked ERP-  
321 related activity. Ultimately we felt this precluded us from meaningfully interpreting earlier ERP

322 component time-windows that overlap with this potential systematic bias. Later ERP components  
323 are subsequently increasingly affected if a pre-stimulus baseline is used if differentiating CNV ac-  
324 tivity stemming from condition type (congruent vs. incongruent) is present during this baselining  
325 period.

326 One strategy to reduce confounding systematic differences in this late CNV component is not  
327 to allow the participant to be aware of the upcoming stimulus type/task, that is, by not having  
328 blocks with consistent conditions allowing for different neural preparatory mechanisms to affect  
329 pre-stimulus time regions where baselines are typically extracted from. This strategy would deviate  
330 somewhat from the typical IAT task structure as it would require adding another dimension of task  
331 switching in the IAT (compared to just between attribute and target categories). Furthermore, in  
332 this instance at the time of stimulus onset, a participant would need to be aware of the condition  
333 type, thus further introducing deviations of the IAT experiment structure. Merely changing the  
334 corner labels to inform the participant would not likely be perceptible until foveation, further  
335 introducing confounds related to required eye-movements and very likely degrading time-locking  
336 characteristics of the ERP components being studied. Other strategies include varying the S1-S2  
337 difference timings to mitigate consistent pre-stimulus-locked activity but this process may merely  
338 serve to obscure the level to which preparatory-related EEG signals and other time-locked within-  
339 block ERP activity might be affecting baselines.

340 A primary reason for baselining is to remove slow drifts present in the EEG, which when removed  
341 by high-pass filtering can result in obscured ERP amplitude/latency characteristics (**Rousselet,**  
342 2012), particularly so when the ERP is generated as a result of lower frequency band activity.  
343 An issue with this is comparability with other studies as some ERP components may have lost  
344 contributing stimulus-locked and phase-coherent related activity in lower frequencies. However,  
345 doing this allows us to overcome some of the problems for which we use baselining in the first  
346 place, that is to remove slow drifts and other noise sources which complicate comparison of ERP  
347 amplitude activity across participants.

348 The restrictions imposed by the IAT experimental design give rise to a number of confounds  
349 when adopting a typical ERP processing strategy. This does not reflect a fundamental flaw in the  
350 IAT task itself, but rather a characteristic of it that does not fit the typical ERP processing. In  
351 this study we high-pass filter EEG signals in order to overcome these limitations at some cost to  
352 the comparability of amplitude characteristics to other EEG studies. By not doing so, however, we  
353 introduce systematic confounds across the analysis time-window. Examination of the impact of this  
354 high-pass filtering on ERP waveforms would indicate it is largely non-detrimental to activity in early  
355 ERP-component time-windows (N1, P2, N2, P3) but does largely affect (attenuate) stereotyped  
356 late P3b activity, which is notably comprised of lower EEG frequencies in the delta band (0-4Hz)  
357 (see **Demiralp et al.** (2001b)). We focus our analysis on time regions where activity related to  
358 N1, P2, N2 and early P3 contributions are present.

359 **Bidet-Caulet et al.** (2012) outline similar issues encountered pre-stimulus with regard to the  
360 baselining and CNV activity, and they use an approach of high-pass filtering at 4Hz in order to  
361 effectively analyze early ERP stimulus-locked ERP components.

362 ERP plots (including a range of other scalp plots and graphs) used in this study are given in  
363 the supplementary data for this paper using a variety of frequency and baselining methods to  
364 further highlight this problem and how our solution results in earlier ERP waveform characteristics  
365 being largely retained, both in amplitude and timing. Time-frequency wavelet analysis too indicates  
366 there are different frequency and spatial topographies for the ERP components of interest and their  
367 correlative relationship to D-score, indicating high-pass filtering artefacts do not contribute to this  
368 result.

369 *2.4.2 Electrode Reference Choice* Other ERP studies investigating the IAT have typically used  
370 an averaged (linked) mastoid reference (TP9 + TP10). EEG reference choice is known to affect the



371 spatial, temporal and polarity characteristics of ERP waveforms and hence the chosen reference site  
372 should be carefully considered not only to allow comparability of results to other studies but also  
373 such that it is not affected by activity related to the factors being investigated in the experiment.

374 The spatial dispersion of statistical activity seen in statistical scalp plots in our study suggests  
375 that the linked-mastoid reference choice may not be entirely optimal and should at least war-  
376 rant consideration as these electrode sites are located near to temporal-lobe regions implicated in  
377 language-processing. While there is generally high agreement in our study for the locus points of  
378 statistical activity between linked-mastoid and common average reference schemes, differences are  
379 evident notably in terms of higher spatial dispersion of statistical activity for the linked-mastoid  
380 reference to a common average reference scheme. In the supplementary appendix to this paper, we  
381 provide grand average ERP waveforms using a common average reference scheme to highlight the  
382 potential contribution of activity at the TP9 and TP10 reference sites. Similar issues surrounding  
383 EEG referencing schemes are also considered by Dien (1998), Hagemann et al. (2001) and Luck  
384 (2005).

## 2.5 ANALYSIS OF NEURAL DATA

385 *2.5.1 ERP time-windows and Channel Selection* : ERP time-windows were determined by  
386 inspecting grand-averaged ERP plots across participants irrespective of condition type. ERP time-  
387 windows were selected so as to include within the window a primary, and any secondary, troughs  
388 or peaks characteristic of ERP activity of that type. An important point to note is that peaks in  
389 ERP averages are not the same as ERP components, as ERP components contributing to averaged  
390 activity can have varying latencies and overlap. In this work we refer to ERP time-windows as time  
391 periods known to contain stereotyped underlying ERP activity. A further discussion of this can be  
392 found in Luck (2005).

393 There are topographic variations of ERP activity in the IAT literature implicating a number of  
394 fundamental ERP components active in time regions corresponding to the P1, N1, P2, N2 and  
395 P3. To our knowledge, as we are first to investigate ERPs in the nature-IAT, we did not preselect  
396 explicit channel-ERP mappings in our study. Instead, we identified these channels and time regions  
397 from visual inspection of ERP time-topographies on grand-averaged epochs - averaging across  
398 participants and conditions. With respect to time regions and channels, the literature identifies a  
399 variety of stereotyped ERP morphologies that can be present in the IAT. Importantly here, there  
400 are variations in expected ERP channel  $\times$  time morphologies determined by the IAT task itself  
401 and the stimulus content used (pictures vs. words) introducing uncertainty with regards to what  
402 the expected ERP patterns will be in an untested IAT.

403 In our study, the N100 was identified as being present in the 110ms-150ms time window, the  
404 P200 in the 160ms-230ms time window, a pattern of fronto-central tending negativity hereafter  
405 referred to as N200 in the 250ms-310ms time window and a frontal P300-like component in the  
406 330ms-450ms time-window.

407 From the existing IAT (and EEG) literature we know a broad range of ERP components are to  
408 be expected such as the P1, N2, P2, P3, and N400 in the IAT-EEG. From N200 studies, we know  
409 variations of this component can manifest with anterior (Fz), central (Cz) and posterior (Pz) scalp  
410 distributions. Similarly, N400 ERP effects are described occurring in overlapping time periods on  
411 these electrode locations. Given our focus investigating early negative ERPs (N2, N400, MFN)  
412 electrode sites Fz, Cz and Pz were chosen as *regions of interest* (ROIs) with regard to the IAT and  
413 key electrode sites for comparisons.

414 *2.5.2 Repeated-Measures ANOVA* : Repeated-measures ANOVAs were used to identify signifi-  
415 cant neural activity during ERP time regions. Channels for each repeated-measures ANOVA were  
416 identified from grand-average ERP plots without differentiating trials based on D-score type or

417 condition (congruent/incongruent), selecting those that displayed stereotyped ERP activity of the  
418 N1, P2, N2 and P3.

419 Repeated-measures ANOVA models were used for each identified ERP time frame examining  
420 electrode site  $\times$  Condition (Congruent/Incongruent) as within-participant factors and a between-  
421 participant factor of 'D-score range' identifying high, medium and low D-scorers (a 7/8/7 split, 22  
422 in total). Greenhouse-Geisser corrected p-values and statistics are reported.

423 *2.5.3 Repeated-Measures ANOVA Post-hoc Analysis* : Correlation based measures are used as  
424 part of our post-hoc RM-ANOVA analysis given the presence of between-subject effects of D-score  
425 magnitude. These are presented both in terms of contrast, explaining significant effects found in  
426 our ANOVAs and in parallel as measures to capture a type of statistical relationship not readily  
427 captured by repeated-measures ANOVA analysis.

428 Correlations are examined using EEG time-window average amplitudes. In Table 2, we show  
429 Pearson-r correlational coefficients for behavioural measures and ERP time-window activity across  
430 selected electrode sites Fz, Cz and Pz, and for electrodes of peak correlation.

431 *2.5.4 The LORETA Approach* : eLORETA is used alongside correlation analysis with D-score,  
432 to identify potential functionally and spatially distinct brain regions that are active in ERP time  
433 ranges. Given the complexity of the resulting relationships, either temporal or spatial in nature,  
434 which are introduced by utilizing reference channels that are not electrically silent (i.e. located  
435 near to language areas), scalp plots of ERP or statistical activity can be misleading as activity at a  
436 particular site might be indicative of two or more channels (and/or ERP components) interacting  
437 in a complex way.

438 In this study, LORETA is used to identify, within the precision of LORETA's localization error,  
439 brain regions and structures involved with early ERP component activity which gives a better sense  
440 of cortical regions that are involved. Both approaches are carried out here as they are considered  
441 complimentary in understanding brain activity driving early IAT-ERP effects.

442 Reported LORETA correlation p-values are adjusted for multiple comparison and presented in  
443 the format [r=.51, p=.005].

## 2.6 CONVENTIONS USED IN THE ANALYSIS DESCRIPTION

444 Further references to congruent and incongruent EEG and reaction times will be described in a  
445 format of *measure-type(measure-src)*:  $RT(C)$  = congruent reaction time,  $RT(I)$  = incongruent  
446 reaction time,  $RT(I - C) = RT(I) - RT(C)$ ,  $E(C)$  = congruent EEG amplitude measure,  $E(I) =$   
447 incongruent amplitude measure and  $E(I - C) = E(I) - E(C)$ .

448 Significant trends are reported for  $alpha < 0.05$  and weakly significant trends for  $alpha < 0.10$ .

449 Statistics for both multivariate and univariate are reported inside square brackets e.g. [ $r(21) =$   
450  $0.8, p = 0.001$ ].

451 *2.6.1 Other methods* : There is evidence for the presence of non-linear relationships surrounding  
452 ERP measures with regard to IAT-effect in our experiment as has been found in other studies  
453 **Williams and Themanson** (2011). Although we do not explore these relationships in the paper,  
454 we include them in the supplementary materials.

### 3 RESULTS

#### 3.1 BEHAVIOURAL IAT ANALYSIS

455 Analysing the behavioural RT data for participants between congruent ( $M = 731.73$  ms,  $s.e. =$   
456  $296.13$ ) and incongruent ( $M = 822.96$  ms,  $s.e. = 338.2$ ) conditions, there was a significant difference  
457 found in reaction times. Reaction times for each condition for each subject submitted to a Wilcoxon  
458 signed-rank test revealed significant differences in reaction time [ $Z=19$ ,  $p = 0.000483$ ]. This confirms  
459 our group shows a pro-nature bias.

460 In Figure 1 we can see that a significant correlation exists between a participant's D-score and  
461 reaction time in congruent (Pearson- $r$   $p = 0.01023$ ) conditions compared to incongruent (Pearson- $r$   
462  $p = 0.74158$ ) conditions. This indicates our measured IAT-effect is being driven by reduced reaction  
463 times in congruent blocks without corresponding related increases in incongruent block reaction  
464 times.

#### 3.2 NEURAL IAT ANALYSIS (ANOVA)

465 *3.2.1 Repeated Measures ANOVA analysis* Amplitude averages across participants for ERP time-  
466 windows were submitted to a repeated measures ANOVA with congruency conditions and channels  
467 as within-subject factors, and D-score range as a between-subject factor. D-score ranges were  
468 acquired by using a 7/8/7 split (by D-score) of available participants. Effects with a significance of  
469  $alpha < 0.10$  are reported.

##### 470 *N100*

471  
472 The N100 was examined across electrode sites Fz, Cz, Pz, F3, F4, C3, C4, P3, P4, CP1, CP2,  
473 FC1 and FC2. A significant main effect was found for channels [ $F(2.914, 55.357) = 37.682, \eta^2 =$   
474  $0.665, p < 0.001$ ].

##### 475 *P200*

476  
477 The P200 was examined across electrode sites Fz, Cz, F3, F4, C3, C4, FC1 and FC2. A significant  
478 main effect for channels was found [ $F(2.562, 48.683) = 35.202, \eta^2 = 0.478, p < 0.001$ ].

##### 479 *N200*

480  
481 The N200 was examined across electrode sites Fz, Cz, Pz, F3, F4, C3, C4, P3, P4, CP1, CP2, FC1  
482 and FC2. Main effects were found for channels [ $F(2.588, 49.171) = 24.279, \eta^2 = 0.561, p < 0.001$ ],  
483 conditions [ $F(1, 19) = 3.252, \eta^2 = 0.146, p = 0.087$ ] and D-score range [ $F(2, 19) = 4.866, \eta^2 =$   
484  $0.339, p = 0.02$ ]. A weakly significant interaction effect for condition  $\times$  D-score range was found  
485 [ $F(2, ) = 1.34, \eta^2 = 0.124, p = 0.079$ ].

##### 486 *P300*

487  
488 The P300 was examined across electrode sites Fz, F3, F4, FC1, FC2, Pz, P3, P4, C3, C4, Cz,  
489 CP1, CP2, CP5 and CP6. Main effects were found for channels [ $F(2.580, 49.023) = 15.586, \eta^2 =$   
490  $0.451, p < 0.001$ ] and D-score range [ $F(2, 19) = 7.529, \eta^2 = 0.442, p = 0.004$ ]. No main effect was  
491 found for congruency condition.

### 3.3 NEURAL CORRELATES OF D-SCORE

#### 3.3.1 N100, P200

492  
493  
494 As neither the N100 or P200 time windows emerged with significant effects (i.e.  $p < .10$ ) we do  
495 not report them further in this study.

#### 3.3.2 N200

496  
497 Repeated-measures ANOVA revealed a number of significant effects for the N200 for between-  
498 subjects (i.e. D-score is predictive of ERP amplitudes) and of within subject-effects, such that N2  
499 amplitudes congruent ( $M = -1.983$  mV,  $s.e. = .274$ ) were enhanced (more negative) compared to  
500 incongruent ( $M = -1.798$  mV,  $s.e. = .25$ ) conditions. There was an effect for between-subjects for  
501 D-score range indicated mean amplitudes more negative for high D-scores ( $M = -3.037$  mV,  $s.e. =$   
502  $.456$ ) compared to low D-scores ( $M = -1.339$  mV,  $s.e. = .456$ ) and different from medium D-scores  
504 ( $M = -1.295$  mV,  $s.e. = .426$ ).

505 Condition  $\times$  D-score emerged as a significant interaction where medium D-scores displayed greater  
506 mean amplitude differences between congruent ( $M = -1.537$  mV,  $s.e. = .486$ ) and incongruent con-  
507 ditions ( $M = -1.053$  mV,  $s.e. = .414$ ) compared to differences between conditions for high D-score  
508 congruent ( $M = -2.978$  mV,  $s.e. = .486$ ) and incongruent ( $M = -3.095$  mV,  $s.e. = .443$ ), and low D-score  
509 congruent ( $M = -1.433$  mV,  $s.e. = .486$ ) and incongruent ( $M = -1.245$  mV,  $s.e. = .443$ ) conditions.

510 Significant linear relationships were present for the N200 time-window examining Pearson-r cor-  
511 relation between D-score congruent [CP6,  $r = -.54$ ,  $p = .009$ ] and incongruent conditions [C4,  $r =$   
512  $-.54$ ,  $p = .009$ ] (Table 2).

513 Examining Table 2 we see these linear relationships are primarily constrained with respect to  
514 D-score  $\times$  amplitude with no significant (univariate) correlations present at electrode sites (matched  
515 for the electrode site with the most significant correlation) comparing other behavioural measures.  
516 Here we can see stronger patterns of correlation across electrode sites for congruent reactions times  
517 to neural measures than incongruent reactions times. Similarly, we see increased correlations for the  
518 standardised reaction differences between congruent blocks (D-score) compared to non standardised  
519 differences (i.e.  $rt(I-C)$ ).

520 LORETA analysis shown in Figure 3 (*a* and *b*) reveals characteristic shared activations between  
521 congruent and incongruent conditions in similar brain structures with these outlined in Table 1 and  
522 Figure 4. Broadly, most significant correlations with D-score were found in areas extending from  
523 anterior, inferior, and insular regions of the left temporal lobe (BA13) and postcentral gyrus (BA  
524 43) Table 1 (BA42, BA13, BA43 and BA22).

#### 3.3.3 P300

525  
526 Repeated-measures ANOVA for the P300 time-window revealed a significant relationships for  
527 between-subjects effect for D-score range such that mean ERP amplitudes were larger in high  
528 ( $M = .763$  mV,  $s.e. = .140$ ) compared to medium ( $M = .119$  mV,  $s.e. = .131$ ) and low ( $M = .087$  mV,  
529  $s.e. = .130$ ) groups.

531 LORETA revealed significant patterns of activation revealed as correlated with D-score for congru-  
532 ent and incongruent conditions. Figure 3(*c* and *d*) show the respective congruent and incongruent  
533 correlated activations.

534 More prominent differences emerge, differentiating correlated neural activity of the congruent  
535 and incongruent conditions for regions surrounding: Medial Frontal Gyrus (BA 10) being more  
536 correlated with D-score in congruent [ $r = .605$ ,  $p = .0732$ ] v. incongruent conditions [ $r = .459$ ,  $p =$

537 .333], and for Postcentral Gyrus (BA 3) in congruent [ $r = .0$ ,  $p = .99$ ] v. incongruent conditions [ $r$   
538  $= .627$ ,  $p = .037$ ](Table 2).

539 Non-standardised reaction time differences here (rt(I-C)) seem to show increased patterns of  
540 correlation to neural measures respective to D-score (compared to the case for the N2) suggest-  
541 ing that differences in reaction times (related to response-locked activity) are more likely driving  
542 contributions here for the IAT measurement.

**Table 1.** LORETA-derived regions of peak correlation of D-score across congruent and incongruent conditions. Rows marked with \* are provided to allow comparison of matched MNI (x,y,z) coordinates between respective maxima of peak correlation between congruent and incongruent conditions.

Component	Condition	Area	Brodmann Area	Side	R	P	X	Y	Z
N200	C	Temporal Lobe - STG	42	L	.619	.0530	-55	-30	15
		Temporal Lobe - STG	22	L	.6	.0734	-45	5	-5
		Insular - Sub-Lobar	13	L	.595	.0786	-40	5	-5
		Postcentral Gyrus*	43	L	.598	.075	-65	-20	20
		Insular - Sub-Lobar*	13	L	.550	.144	-45	0	-10
	I	Postcentral Gyrus	43	L	.643	.041	-65	-20	20
		Insular - Sub-Lobar	13	L	.605	.0812	-45	0	-10
		Temporal Lobe - STG*	42	L	.541	.190	-55	-30	15
		Temporal Lobe - STG*	22	L	.521	.230	-45	5	-5
		Insular - Sub-Lobar*	13	L	.560	.149	-40	5	-5
P300	C	Cingulate Gyrus	24	R	.645	.033	10	-20	45
		Insular - Sub-Lobar	13	L	.623	.0548	-45	-25	20
		Medial Frontal Gyrus	10	R	.605	.0732	15	60	5
		Superior Temporal Gyrus*	22	L	.534	0.1842	-50	5	-5
		Cingulate Gyrus*	31	R	.582	.103	20	-25	40
	Postcentral Gyrus*	3	R	0	.99	30	-25	40	
	I	Superior Temporal Gyrus	22	L	.645	.0244	-50	5	-5
		Cingulate Gyrus	31	R	.627	.037	20	-25	40
		Postcentral Gyrus	3	R	.627	.037	30	-25	40
		Cingulate Gyrus*	24	R	.570	.095	10	-20	45
Insular - Sub-Lobar*		13	L	.501	.231	-45	-25	20	
Medial Frontal Gyrus*	10	R	.459	.333	15	60	5		

**Table 2.** Pearson-r correlation coefficients across behavioural and EEG activity measures. Correlated variable pairs presented in the first column as *Behavioural Measure, EEG Measure*. \* indicates univariate p-value < .05 and \*\* indicates univariate p-value < .1. D=D-score, C = Congruent, I = Incongruent, rt() = Reaction Time. Max columns represent the electrode site with smallest p-value for D-score correlated with EEG measure (first 4 rows) where electrode site for each EEG measure type (C+I, I-C, C, I) is maintained across subsequent comparisons as a way to interpret the source of EEG activity driving correlations with D-score at that site.

	N200				P300			
	Fz	Cz	Pz	Max	Fz	Cz	Pz	Max
D,I-C	-0.21	-0.14	-0.08	0.39(P8)	0.35	0.40	0.50**	0.50** (Pz)
D,C+I	-0.46**	-0.47**	-0.42	-0.54** (C4)	0.57**	0.52**	0.34	0.62** (C4)
D,I	-0.48**	-0.49**	-0.44**	-0.54** (C4)	0.61**	0.57**	0.47**	0.64** (C4)
D,C	-0.44**	-0.45**	-0.39	-0.54** (CP6)	0.51**	0.42	0.13	0.54** (F4)
rt(I-C),I-C	-0.17	-0.11	-0.04	0.36(P8)	0.29	0.34	0.37	0.37(Pz)
rt(I-C),C+I	-0.26	-0.32	-0.36	-0.38(C4)	0.51**	0.46**	0.37	0.57** (C4)
rt(I-C),I	-0.29	-0.33	-0.37	-0.39(C4)	0.55**	0.51**	0.44**	0.60** (C4)
rt(I-C),C	-0.25	-0.30	-0.34	-0.41(CP6)	0.48**	0.39	0.22	0.50** (F4)
rt(C),I-C	0.14	0.15	0.10	-0.09(P8)	-0.11	-0.10	-0.17	-0.17(Pz)
rt(C),C+I	0.35	0.36	0.35	0.27(CP6)	-0.32	-0.33	-0.13	-0.28(F4)
rt(C),I	0.35	0.37	0.36	0.40(C4)	-0.30	-0.30	-0.15	-0.27(C4)
rt(C),C	0.33	0.32	0.30	0.23(CP6)	-0.29	-0.30	-0.03	-0.24(F4)
rt(I),I-C	0.05	0.10	0.09	0.11(P8)	0.04	0.08	0.03	0.03(Pz)
rt(I),C+I	0.22	0.20	0.16	0.19(C4)	-0.05	-0.09	0.07	0.01(C4)
rt(I),I	0.21	0.20	0.17	0.19(C4)	-0.01	-0.03	0.08	0.05(C4)
rt(I),C	0.21	0.17	0.13	0.01(CP6)	-0.04	-0.09	0.09	0.02(F4)

## 4 DISCUSSION

543 The aim of the study reported here was to examine how ERP measures in the IAT might offer an  
544 insight into the neural mechanisms underlying the more rapid associations that drive IAT effects.  
545 Of primary interest in our work was examining how ERP measures underlying both congruent and  
546 incongruent block types might offer evidence of the neural mechanisms involved with these rela-  
547 tively more rapid associations. In our results from behavioural measures we find average congruent  
548 reaction times are significantly correlated with participant D-score, while the reaction times in in-  
549 congruent conditions are not. From this we would expected neural activity predictive of D-score to  
550 be present only in congruent blocks. Similarly, we would expect ERP measures for time-windows  
551 during the incongruent blocks to be largely unresponsive of D-score, however, we find this is not  
552 the case.

553 Our hypothesis was positioned such that in a situation where the measurable IAT effect is pri-  
554 marily modulated by reduced congruent reaction times, in the respective incongruent blocks we  
555 should find shared patterns of ERP activity correlated with the size of IAT effect, given the in-  
556 volvement of proactive cognitive control and other top-down control processes. This is related to  
557 the motivational/attention aspects in the IAT affecting the level to which an implicit bias might be  
558 measured. Given evidence that groups typically have positive IAT scores on the n-IAT, we suspected  
559 those participants with lower D-scores (a lower standardised difference in reactions times) might  
560 be engaging in the task differently due to factors like less motivational effort thus not engendering  
561 conditions necessary to capture IAT effects.

562 The aim of searching for such evidence was to disentangle cortical generators involved with the  
563 production of an IAT effect in early time periods of the ERP (early negativities between 250-450ms)  
564 following stimulus presentation that have been previously implicated in other studies to be sensitive  
565 to the IAT effect size, such as cognitive control or error monitoring.

566 In our study we identified an N2-like ERP component in the 250-310 ms range. While we have  
567 labelled activity in this time-region as indicative of an N2 ERP, there is close overlap in time regions  
568 of an N400 described in other related EEG-IAT studies. Importantly, some of these studies iden-  
569 tify correlational relationships between congruent, incongruent and incongruent-congruent activity  
570 ERP measures and D-score (**Williams and Themanson, 2011**). The N400 has been widely used  
571 as a measure of semantic congruency for words (**Kiefer, 2002**) and statements (**Kutas and Hill-  
572 yard, 1980**). **Williams and Themanson (2011)** report a significantly smaller N400 for congruent  
573 conditions compared to incongruent conditions in an IAT suggestive that the N400 is an indicator  
574 of semantic (integration) congruency where greater incongruency results in larger (more negative)  
575 amplitudes. LORETA analysis estimating the source of correlated neural activity and D-score for  
576 both block types in our study implicate a number of left-temporal cortical regions, known gener-  
577 ators in the N400 and more widely understood to be involved with language processing (**Maess  
578 et al., 2006**). **Lau et al. (2008)** identify a dominant (left-hemispheric) pattern across a range of  
579 studies utilising EEG and non-EEG imaging modalities investigating the N400, and indicate the  
580 posterior middle temporal cortex as being the only area to show consistent effects across studies.  
581 It would seem that although no apparent N400 ERP component is present in our averaged wave-  
582 forms, there is evidence overlapping ERP activity from the N400 time-frame might present during  
583 our N200 analysis window.

584 **Forbes et al. (2012)** suggest that a number of brain regions surrounding the left temporal lobe  
585 (as indicated by integrating both EEG source localisation and lesion studies) are implicated as  
586 being important in the production of reduced congruent reaction times in the IAT. Interestingly,  
587 they find patients (lesion) vs. controls show no significant difference on incongruent reaction times  
588 or D-scores but show statistically significant differences where patients were slower to respond in  
589 congruent conditions. Similarly, they identify that volume loss in large regions of the left insula  
590 exhibit robust associations with slower reaction times in congruent blocks. In the context of our  
591 results, these findings support the role of left temporal/insular brain regions as being important in



592 the production of an IAT effect. Another similarity in results is a strong indication that a number  
593 of shared brain structures are recruited across both congruent and incongruent conditions, but  
594 importantly there are differences associated with activations, suggesting different recruitment of  
595 brain regions based on condition.

596 The N2 has been found to reflect a conflict detection function (Yeung et al., 2004) between  
597 possible choices prior to response selection and performance. Hilgard et al. (2014) show that  
598 the medial-frontal negativity (MFN) ERP between 250-450 ms post stimulus at midline regions is  
599 larger for incongruent mappings, compared to congruent mappings, indicating increased proactive  
600 control is required during incongruent blocks. Although we find congruent (instead of incongruent)  
601 mappings in our IAT generated seemingly more negative going waveforms in this time region,  
602 their study highlights how the involvement of reactive control in the IAT (due to task switching)  
603 generates a similar temporal and spatially overlapping positive ERP (D-pos) in this time region  
604 which might be one explanation for the relationship we found. Importantly, as relative differences  
605 in these ERP measures between congruency conditions have been to understand congruency effects  
606 with respect to semantic integration and cognitive control, we find absolute measures here to be  
607 involved as well. There is other evidence that for early ERP negativities being sensitive to task  
608 constraints, for example Jodo and Kayama (1992) show that the N2 amplitude is enhanced by  
609 reaction time constraints on a go/no-go task.

610 One important difference in our study is that we find an N2 component where other studies  
611 have not, in the IAT task. A possible explanation for this is the stimuli used in the n-IAT task;  
612 differences exist when comparing waveforms as a function of content used in the task, i.e. pictures  
613 versus words (Williams and Themanson, 2011). Fleischhauer et al. (2014) do not find evidence  
614 of significant effects in the N2 (as expected). Although the authors here are considering implicit  
615 measures of neuroticism, the stimuli and experimental structure are similar to ours.

616 Caution is warranted in interpreting the P3-like activity identified in our study as it differs  
617 from the more classical and response-locked P3b found in other IAT-EEG studies. The P3 has  
618 been shown to index attention towards self-referent materials in an IAT (Chen et al., 2014)  
619 so our discovery of the involvement of this component in an IAT responding to both congruent  
620 and incongruent self-referent mappings was somewhat expected. The P3b component (the more  
621 common P3 variant identified in existing IAT studies) is notably comprised of lower-frequency  
622 EEG activity (0-4 Hz). As we use high-pass filtering, P3-related ERP activity is highly attenuated  
623 in our ERP measures. Also, the timecourse of this P3-like activity in our study overlaps with  
624 time regions where other studies have found N400 to be present. Importantly, the N400 is not  
625 necessarily characterised by a negative deflection in the ERP waveform, as it is measured relatively  
626 as a more negative going signal with respect to other experimental conditions. The P3a is a frontal-  
627 central tending ERP and is seen in target detection tasks to novel and infrequent stimuli, and it also  
628 reflects attention mechanisms during task processing (Polich, 2007). Wavelet analysis indicates this  
629 component is partly comprised of theta band (4-8 Hz) activity (Demiralp et al., 2001a) and source  
630 localisation analysis reveals a wide range of cortical generators. Our preprocessing strategy - using  
631 high-pass filtering to avoid the use of pre-stimulus baselining confounds - compromises a robust  
632 interpretation of this frontal P3-like activity with respect to existing ERP literature. Although  
633 we find, however, in a variety of alternatively explored filtering strategies this earlier occurring  
634 frontal P3-like activity remains present particularly in conjunction with its posterior counterpart  
635 (P3b). As other IAT-ERP indices have been implicated in this time-region we choose to retain  
636 this time-region in our analysis. There is strong suggestion of overlapping cortical generators of  
637 IAT-sensitive ERPs in both our N2 and P3 analysis time windows as can be seen in our LORETA  
638 analysis particularly for the incongruent block conditions. This would indicate ERP activity might  
639 not only be modulated by differences in amplitude of underlying components but also latency.  
640 We do find, however, significant patterns of correlated activity using LORETA of brain-regions  
641 typically implicated in the generation of the P3 (notably cingulate cortex and medial frontal gyrus)

642 (Volpe et al., 2007). Egenolf et al. (2013) similarly examine covarying relationships of brain-  
643 region activation to the magnitude of behavioural IAT effect and find differences in a time window  
644 of 510ms to 710ms between incongruent and congruent ERP activity that is proportional to the  
645 IAT effect, a time region largely corresponding to the P3 ERP.

646 Given previous results of the n-IAT we would expect averaged D-score across participants to be  
647 more positive (i.e. groups of participants in previous studies tended to show a stereotypic pro-nature  
648 IAT effect like we similarly have). These results might indicate the nIAT-effect cannot be measured  
649 reliably on all participants, due to differences such as motivation during the task and/or different  
650 patterns of top-down control employed in the task. The presence of such an effect is important in  
651 understanding instances where the IAT might be failing to measure an expected bias and would be  
652 a potential source of detrimental noise in measuring relationships to explicit measures.

## 5 CONCLUSIONS

653 The results presented in this paper indicate that EEG is informative in understanding cognitive  
654 processes behind the n-IAT. Our results both confirm patterns of activity seen in other IAT studies  
655 and also extend these by showing novel behavioural-ERP predictive relationships. Importantly, we  
656 identify that N2-/MFN- related amplitudes in our ERP analysis time window show a correlational  
657 relationship with D-score, highlighting the potential involvement of participant motivation via  
658 proactive cognitive control and top-down attention related mechanisms as a source of noise in the  
659 successful measurement of an IAT-effect. This has broad implications for other studies utilising the  
660 n-IAT (and other IATs in general) in that it might offer explanation as to why IAT measures can  
661 often fail to correlate with explicit measures. Such a line of evidence would indicate other secondary  
662 measures (including EEG) alongside the IAT may be useful in measuring these motivational related  
663 factors so as to enable an experimenter to disqualify participants who may be IAT averse.

664 One notable difference between this study and other studies is how the data is preprocessed due  
665 to the presence of pre-stimulus locked ERP activity related to the CNV. Although this potential  
666 problem of the IAT-ERP and its confounds introduced by standard baselining exists, there is little  
667 reported in the IAT literature that it has been at least taken account of.

668 One area of future work to be explored next is to examine predictive relationships and func-  
669 tional/structural brain differences that emerge within, and across, participants for a variety of IAT  
670 tasks.

**DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT**

671 The authors declare that the research was conducted in the absence of any commercial or financial  
672 relationships that could be construed as a potential conflict of interest.

**AUTHOR CONTRIBUTIONS**

673 Graham Healy ideated the project, created the data acquisition protocol, led the data acquisition,  
674 led the data analysis, wrote the initial version of the paper and agrees to be accountable for all  
675 aspects of the work.

676 Lorraine Boran contributed to the design of the work and interpretation of data for the work,  
677 revised the work for intellectual content, contributed to final approval of the version to be submitted  
678 and agrees to be accountable for all aspects of the work.

679 Alan. F Smeaton ideated, planned and oversaw the project, reviewed and refined the paper, and  
680 agrees to be accountable for all aspects of the work.

## ACKNOWLEDGEMENT

681 The authors wish to acknowledge assistance from Miriam Kennedy and Yang Yang in recording of  
682 experimental data, and acknowledge the input from Michael Keane in assisting with recording and  
683 analysis of the experimental data.

684 *Funding:* The work reported here was funded by the International Energy Research Centre /  
685 Enterprise Ireland and by Science Foundation Ireland under grant SFI/12/RC/2289.

**SUPPLEMENTARY DATA**

686 Supplementary Data is provided.

## REFERENCES

- 687 Agosta, S., Pezzoli, P., and Sartori, G. (2013), How to detect deception in everyday life and the  
688 reasons underlying it, *Applied Cognitive Psychology*, 27, 2, 256–262, doi:10.1002/acp.2902
- 689 Barnes-Holmes, D., Stanton, C., Barnes-Holmes, Y., Whelan, R., Stewart, I., Commins, S., et al.  
690 (2004), Interfacing relational frame theory with cognitive neuroscience: Semantic priming, the  
691 implicit association test, and event related potentials, *International Journal of Psychology and*  
692 *Psychological Therapy*, 4, 215–240
- 693 Bekker, E. M., Kenemans, J. L., and Verbaten, M. N. (2005), Source analysis of the N2 in a cued  
694 Go/NoGo task, *Brain Res Cogn Brain Res*, 22, 2, 221–231
- 695 Bidet-Caulet, A., Barbe, P. G., Roux, S., Viswanath, H., Barthelemy, C., Bruneau, N., et al. (2012),  
696 Dynamics of anticipatory mechanisms during predictive context processing, *Eur. J. Neurosci.*,  
697 36, 7, 2996–3004
- 698 Boysen, G. A., Vogel, D. L., and Madon, S. (2006), A public versus private administration of the  
699 Implicit Association Test, *European Journal of Social Psychology*, 36, 6, 845–856, doi:10.1002/  
700 ejsp.318
- 701 Braver, T. S. (2012), The variable nature of cognitive control: a dual mechanisms framework,  
702 *Trends Cogn. Sci. (Regul. Ed.)*, 16, 2, 106–113
- 703 Bruni, C. M., Chance, R. C., Schultz, P. W., and Nolan, J. M. (2011), Natural connections:  
704 Bees sting and snakes bite, but they are still nature, *Environment and Behavior*, doi:10.1177/  
705 0013916511402062
- 706 Bruni, C. M. and Schultz, P. W. (2010), Implicit beliefs about self and nature: Evidence from an  
707 IAT game, *Journal of Environmental Psychology*, 30, 1, 95 – 102, doi:http://dx.doi.org/10.1016/  
708 j.jenvp.2009.10.004
- 709 Chee, M. W., Sriram, N., Soon, C. S., and Lee, K. M. (2000), Dorsolateral prefrontal cortex and  
710 the implicit association of concepts and attributes, *Neuroreport*, 11, 1, 135–140
- 711 Chen, Y., Zhong, Y., Zhou, H., Zhang, S., Tan, Q., and Fan, W. (2014), Evidence for implicit  
712 self-positivity bias: an event-related brain potential study, *Experimental Brain Research*, 232, 3,  
713 985–994, doi:10.1007/s00221-013-3810-z
- 714 Clayson, P. E. and Larson, M. J. (2013), Psychometric properties of conflict monitoring and conflict  
715 adaptation indices: response time and conflict N2 event-related potentials, *Psychophysiology*, 50,  
716 12, 1209–1219
- 717 Curran, T., Tucker, D. M., Kutas, M., and Posner, M. I. (1993), Topography of the N400: brain  
718 electrical activity reflecting semantic expectancy, *Electroencephalography and Clinical Neurophys-*  
719 *iology/Evoked Potentials Section*, 88, 3, 188 – 209, doi:http://dx.doi.org/10.1016/0168-5597(93)  
720 90004-9
- 721 De Houwer, J., Teige-Mocigemba, S., Spruyt, A., and Moors, A. (2009), Implicit measures: A  
722 normative analysis and review, *Psychol Bull*, 135, 3, 347–368
- 723 Deacon, D., Hewitt, S., Yang, C.-M., and Nagata, M. (2000), Event-related potential indices of  
724 semantic priming using masked and unmasked words: evidence that the N400 does not reflect a  
725 post-lexical process, *Cognitive Brain Research*, 9, 2, 137 – 146, doi:http://dx.doi.org/10.1016/  
726 S0926-6410(99)00050-6
- 727 Demiralp, T., Ademoglu, A., Comerchero, M., and Polich, J. (2001a), Wavelet analysis of P3a and  
728 P3b, *Brain Topography*, 13, 4, 251–267, doi:10.1023/A:1011102628306
- 729 Demiralp, T., Ademoglu, A., I Stefanopoulos, Y., Başar-Eroglu, C., and Başar, E. (2001b), Wavelet  
730 analysis of oddball P300, *Int J Psychophysiol*, 39, 2-3, 221–227
- 731 Dien, J. (1998), Issues in the application of the average reference: Review, critiques, and recom-  
732 mendations, *Behavior Research Methods, Instruments, Computers*, 30, 1, 34–43, doi:10.3758/  
733 BF03209414
- 734 Egenolf, Y., Stein, M., Koenig, T., Grosse Holtforth, M., Dierks, T., and Caspar, F. (2013),  
735 Tracking the implicit self using event-related potentials, *Cogn Affect Behav Neurosci*, 13, 4,  
736 885–899



- 737 Fazio, R. H. and Olson, M. A. (2007), The SAGE Handbook of Social Psychology: Concise Student  
738 Edition (SAGE Publications Ltd), 123–146, doi:/10.4135/9781848608221
- 739 Fleischhauer, M., Strobel, A., Diers, K., and Enge, S. (2014), Electrophysiological evidence for  
740 early perceptual facilitation and efficient categorization of self-related stimuli during an implicit  
741 association test measuring neuroticism, *Psychophysiology*, 51, 2, 142–151, doi:10.1111/psyp.12162
- 742 Folstein, J. R. and Van Petten, C. (2008), Influence of cognitive control and mismatch on the N2  
743 component of the ERP: a review, *Psychophysiology*, 45, 1, 152–170
- 744 Forbes, C. E., Cameron, K. A., Grafman, J., Barbey, A. K., Solomon, J., Ritter, W., et al.  
745 (2012), Identifying temporal and causal contributions of neural processes underlying the Implicit  
746 Association Test (IAT), *Frontiers in Human Neuroscience*, 6, 320, doi:10.3389/fnhum.2012.00320
- 747 Glasman, L. R. and Albarracín, D. (2006), Forming attitudes that predict future behavior: a  
748 meta-analysis of the attitude-behavior relation, *Psychol Bull*, 132, 5, 778–822
- 749 Greenwald, A. G. and Banaji, M. R. (1995), Implicit Social Cognition: Attitudes, Self-Esteem, and  
750 Stereotypes, *Psychological Review*, 102, 1, 4–27
- 751 Greenwald, A. G., McGhee, D. E., and Schwartz, J. L. (1998), Measuring individual differences  
752 in implicit cognition: the Implicit Association Test, *J Pers Soc Psychol*, 74, 6, 1464–1480,  
753 [PubMed:9654756]
- 754 Greenwald, A. G., Nosek, B. A., and Banaji, M. R. (2003), Understanding and using the implicit  
755 association test: I. An improved scoring algorithm, *J Pers Soc Psychol*, 85, 2, 197–216
- 756 Guex, R., Cerić, F., Hurtado, E., González, R., Navarro, Á., Manes, F., et al. (2011), Performance  
757 Errors of In Group/Out Group Stimuli and Valence Association in the Implicit Association Test:  
758 Brain Bias of Ingroup Favoritism, *Open Neuroscience Journal*, 5, 16–23
- 759 Hagemann, D., Naumann, E., and Thayer, J. F. (2001), The quest for the EEG reference revis-  
760 ited: A glance from brain asymmetry research, *Psychophysiology*, 38, 5, 847–857, doi:10.1111/  
761 1469-8986.3850847
- 762 He, Y., Johnson, M. K., Dovidio, J. F., and McCarthy, G. (2009), The relation between race-related  
763 implicit associations and scalp-recorded neural activity evoked by faces from different races, *Soc*  
764 *Neurosci*, 4, 5, 426–442
- 765 Hilgard, J., Bartholow, B. D., Dickter, C. L., and Blanton, H. (2014), Characterizing switching and  
766 congruency effects in the Implicit Association Test as reactive and proactive cognitive control,  
767 *Soc Cogn Affect Neurosci*, doi:10.1093/scan/nsu060
- 768 Hofmann, W., Gawronski, B., Gschwendner, T., Le, H., and Schmitt, M. (2005), A meta-analysis  
769 on the correlation between the implicit association test and explicit self-report measures, *Pers*  
770 *Soc Psychol Bull*, 31, 10, 1369–1385
- 771 Holcomb, P. J. (1988), Automatic and attentional processing: An event-related brain potential  
772 analysis of semantic priming, *Brain and Language*, 35, 1, 66 – 85, doi:http://dx.doi.org/10.1016/  
773 0093-934X(88)90101-0
- 774 Hurtado, E., Haye, A., Gonzalez, R., Manes, F., and Ibanez, A. (2009), Contextual blending of  
775 ingroup/outgroup face stimuli and word valence: LPP modulation and convergence of measures,  
776 *BMC Neurosci*, 10, 69
- 777 Ibanez, A., Gleichgerrcht, E., Hurtado, E., Gonzalez, R., Haye, A., and Manes, F. F. (2010), Early  
778 Neural Markers of Implicit Attitudes: N170 Modulated by Intergroup and Evaluative Contexts  
779 in IAT, *Front Hum Neurosci*, 4, 188
- 780 Jodo, E. and Kayama, Y. (1992), Relation of a negative ERP component to response inhibition in  
781 a go/no-go task, *Electroencephalography and Clinical Neurophysiology*, 82, 6, 477 – 482, doi:http:  
782 //dx.doi.org/10.1016/0013-4694(92)90054-L
- 783 Kiefer, M. (2002), The N400 is modulated by unconsciously perceived masked words: further  
784 evidence for an automatic spreading activation account of N400 priming effects, *Brain Res Cogn*  
785 *Brain Res*, 13, 1, 27–39
- 786 Kutas, M. and Hillyard, S. A. (1980), Reading senseless sentences: brain potentials reflect semantic  
787 incongruity, *Science*, 207, 4427, 203–205

- 788 Larson, M. J., Clayson, P. E., and Clawson, A. (2014), Making sense of all the conflict: a theoretical  
789 review and critique of conflict-related ERPs, *Int J Psychophysiol*, 93, 3, 283–297
- 790 Lau, E. F., Phillips, C., and Poeppel, D. (2008), A cortical network for semantics: (de)constructing  
791 the N400, *Nat. Rev. Neurosci.*, 9, 12, 920–933
- 792 Luck, S. J. (2005), An Introduction to Event-Related Potentials and their Neural Origins (Chapter  
793 1) (MIT Press, Cambridge)
- 794 Maess, B., Herrmann, C. S., Hahne, A., Nakamura, A., and Friederici, A. D. (2006), Localizing the  
795 distributed language network responsible for the {N400} measured by {MEG} during auditory  
796 sentence processing, *Brain Research*, 1096, 1, 163 – 172, doi:http://dx.doi.org/10.1016/j.brainres.  
797 2006.04.037
- 798 Makeig, S., Westerfield, M., Jung, T. P., Enghoff, S., Townsend, J., Courchesne, E., et al. (2002),  
799 Dynamic brain sources of visual evoked responses, *Science*, 295, 5555, 690–694
- 800 McDaniel, M. J., Beier, M. E., Perkins, A. W., Goggin, S., and Frankel, B. (2009), An assess-  
801 ment of the fakeability of self-report and implicit personality measures, *Journal of Research in*  
802 *Personality*, 43, 4, 682 – 685, doi:http://dx.doi.org/10.1016/j.jrp.2009.01.011
- 803 O’Toole, C. and Barnes-Holmes, D. (2009), Electrophysiological Activity Generated During the  
804 Implicit Association Test: A Study Using Event-Related Potentials, *The Psychological Record*,  
805 59, 4
- 806 Pascual-Marqui, R. (2002), Standardized low-resolution brain electromagnetic tomography  
807 (sloreta): technical details., *Methods and findings in experimental and clinical pharmacology*,  
808 24, 5
- 809 Plöchl, M., Ossandón, J. P., and König, P. (2012), Combining eeg and eye tracking: Identifica-  
810 tion, characterization and correction of eye movement artifacts in electroencephalographic data,  
811 *Frontiers in Human Neuroscience*, 6, 278, doi:10.3389/fnhum.2012.00278
- 812 Polich, J. (2007), Updating P300: an integrative theory of P3a and P3b, *Clin Neurophysiol*, 118,  
813 10, 2128–2148
- 814 Rousselet, G. A. (2012), Does filtering preclude us from studying ERP time-courses?, *Frontiers in*  
815 *Psychology*, 3, 131, doi:10.3389/fpsyg.2012.00131
- 816 Schultz, P., Shriver, C., Tabanico, J. J., and Khazian, A. M. (2004), Implicit connections with  
817 nature, *Journal of Environmental Psychology*, 24, 1, 31 – 42, doi:http://dx.doi.org/10.1016/  
818 S0272-4944(03)00022-7
- 819 Tacikowski, P. and Nowicka, A. (2010), Allocation of attention to self-name and self-face: An erp  
820 study, *Biological Psychology*, 84, 2, 318 – 324, doi:http://dx.doi.org/10.1016/j.biopsycho.2010.  
821 03.009
- 822 van Nunspeet, F., Ellemers, N., Derks, B., and Nieuwenhuis, S. (2014), Moral concerns increase  
823 attention and response monitoring during IAT performance: ERP evidence, *Soc Cogn Affect*  
824 *Neurosci*, 9, 2, 141–149
- 825 Van Petten, C. and Luka, B. J. (2006), Neural localization of semantic context effects in  
826 electromagnetic and hemodynamic studies, *Brain Lang*, 97, 3, 279–293
- 827 Vargo, E. J. and Petroczi, A. (2013), Detecting cocaine use? The autobiographical implicit associ-  
828 ation test (aIAT) produces false positives in a real-world setting, *Subst Abuse Treat Prev Policy*,  
829 8, 22
- 830 Verschuere, B., Prati, V., and Houwer, J. D. (2009), Cheating the lie detector: faking in the  
831 autobiographical Implicit Association Test, *Psychol Sci*, 20, 4, 410–413
- 832 Volpe, U., Mucci, A., Bucci, P., Merlotti, E., Galderisi, S., and Maj, M. (2007), The cortical  
833 generators of P3a and p3b: A LORETA study, *Brain Research Bulletin*, 73, 4 – 6, 220 – 230,  
834 doi:http://dx.doi.org/10.1016/j.brainresbull.2007.03.003
- 835 Williams, J. K. and Themanon, J. R. (2011), Neural correlates of the implicit association test:  
836 evidence for semantic and emotional processing, *Soc Cogn Affect Neurosci*, 6, 4, 468–476
- 837 Yeung, N., Botvinick, M. M., and Cohen, J. D. (2004), The neural basis of error detection: conflict  
838 monitoring and the error-related negativity, *Psychol Rev*, 111, 4, 931–959

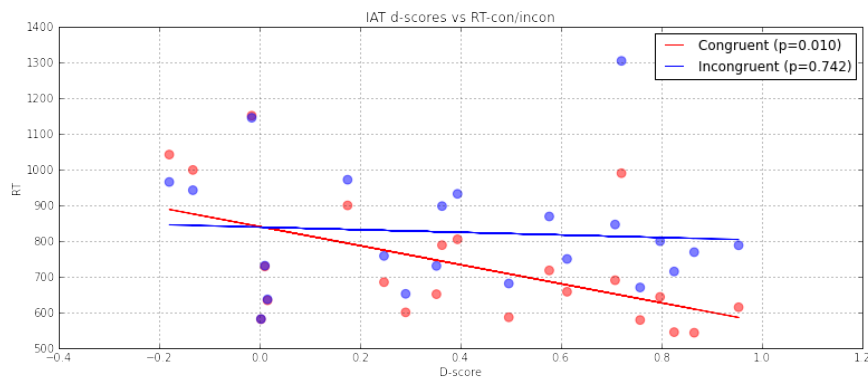


Figure 1: Reaction times across subjects broken down across congruent and incongruent conditions (y-axis) with calculated D-scores (x-axis)

**FIGURES**

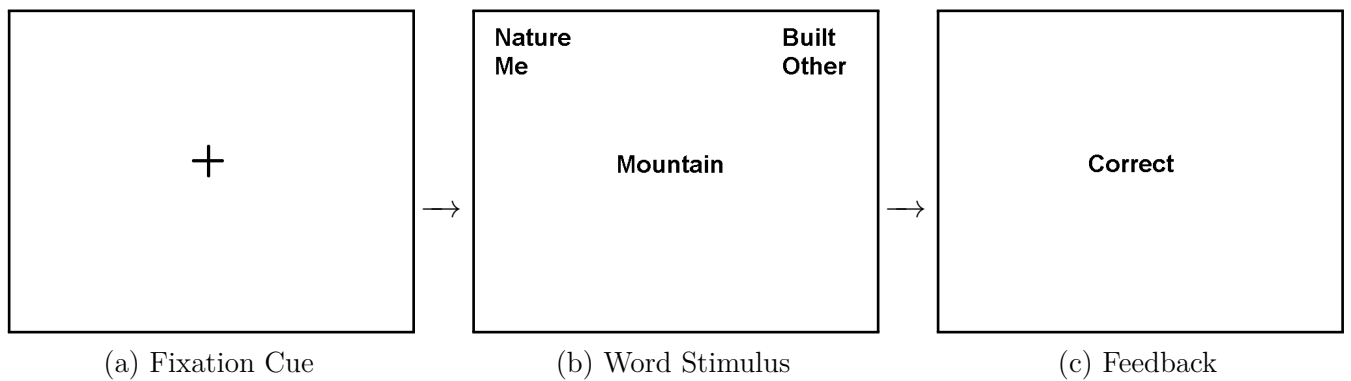


Figure 2: Depiction of trial structure within a congruent block. From left to right: a fixation cross is presented on screen for 1 second, then a word is presented for categorisation and following a key press a feedback screen is presented for 1 second indicating whether the response was correct. 80 were recorded for each congruency condition.

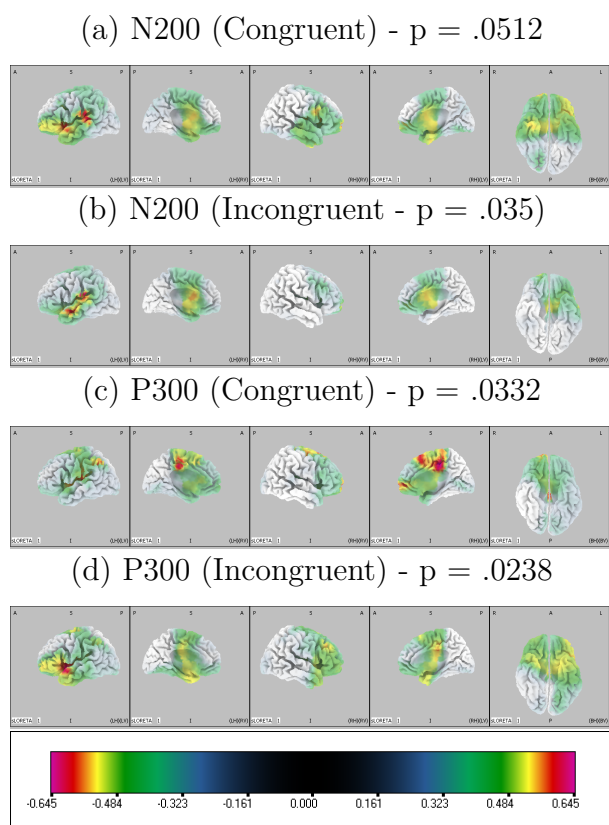


Figure 3: Correlated LORETA voxel activity and D-score. D-score is correlated with congruent and incongruent ERP time-window averages across participants localising activity driving correlated scalp EEG measures. Multiple comparison corrected p-values for peak correlations are presented on top of each condition  $\times$  ERP plot.

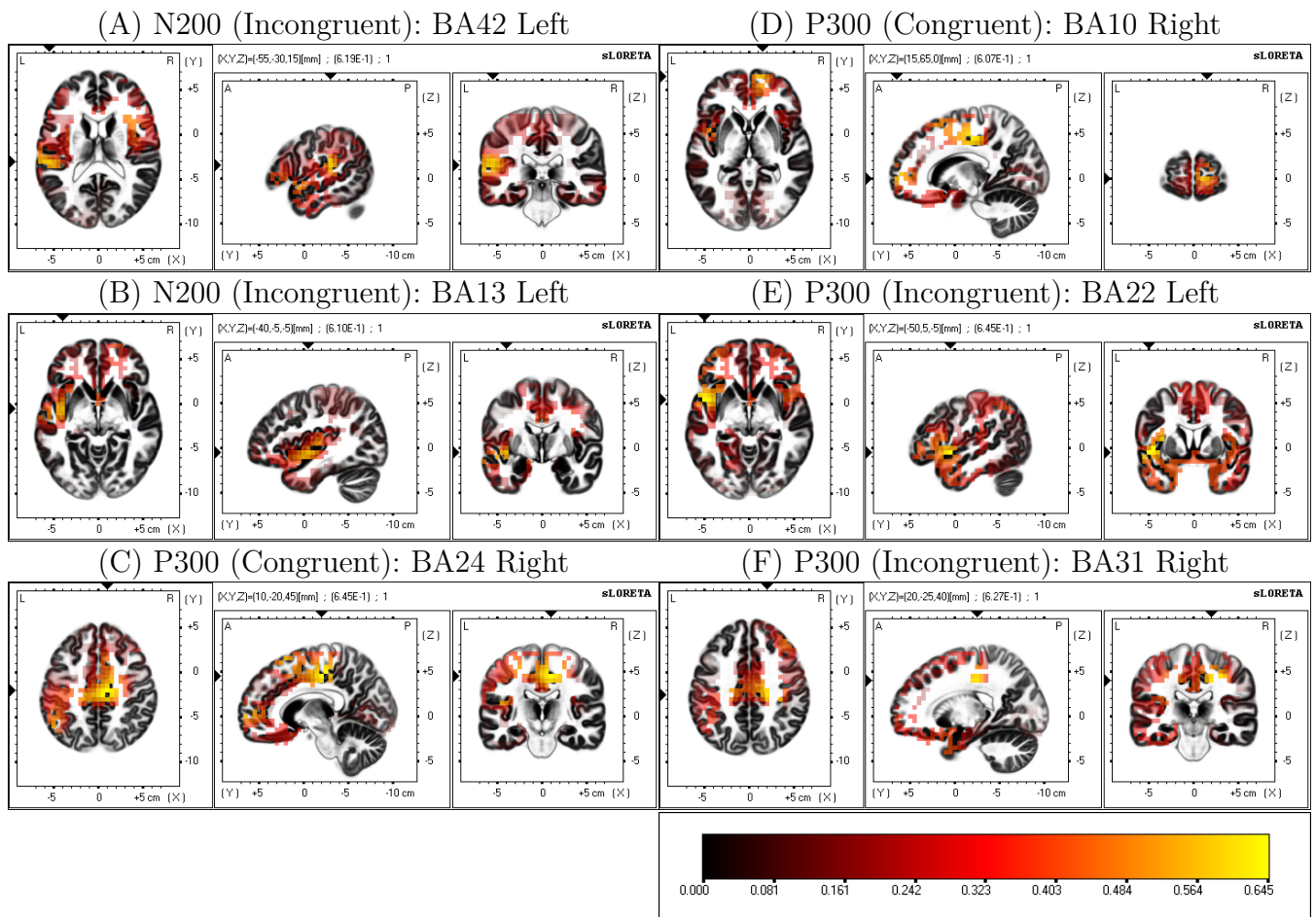
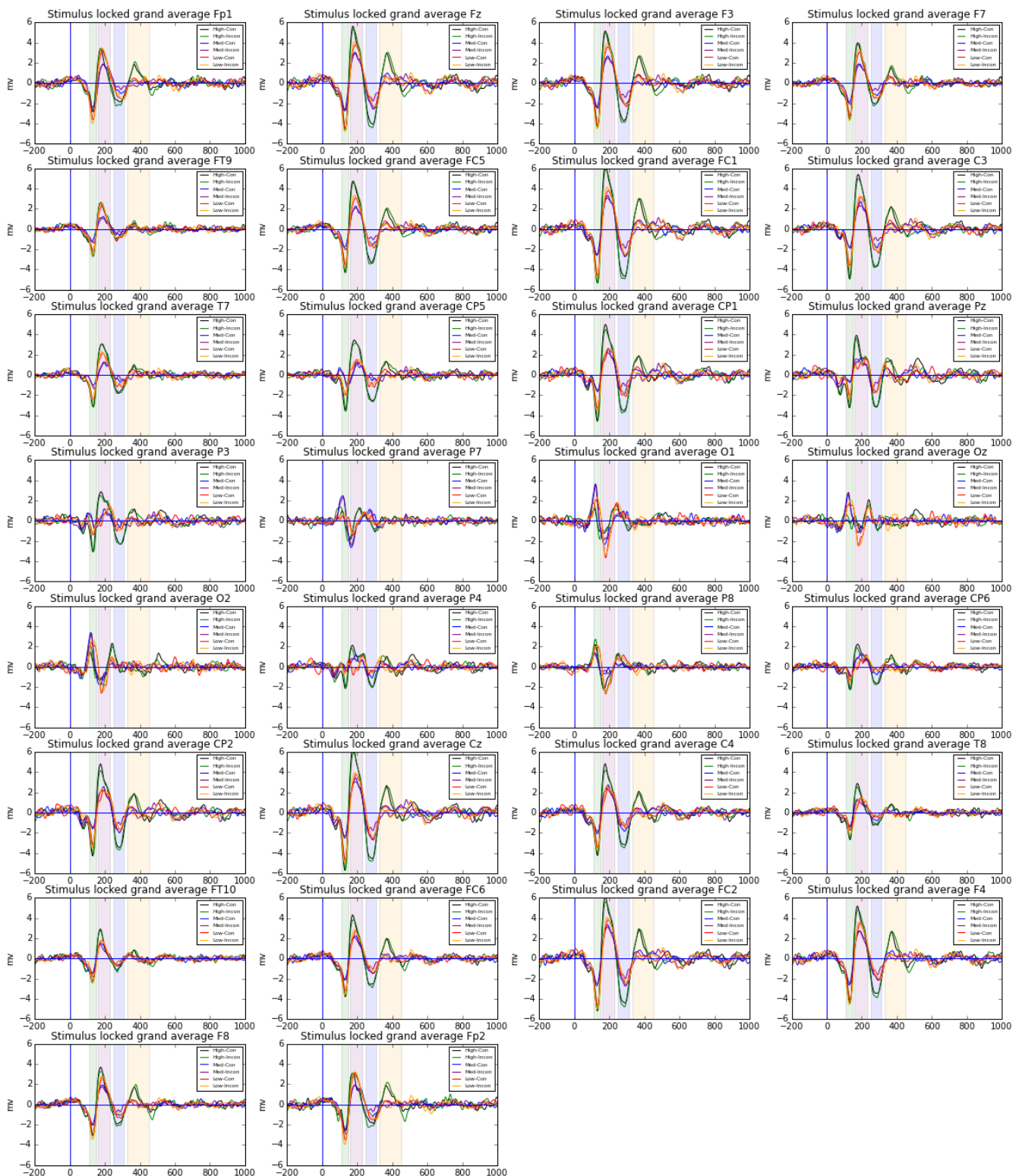


Figure 4: Selected ROIs revealed through LORETA D-score regressions.



**Figure 5.** ERP averages across electrode sites for high, medium and low D-scorers across congruent/incongruent conditions using a linked-mastoids reference. Signals are filtered in the range 4Hz-30Hz.