



25 **Abstract**

26 Predictive coding models can simulate known perceptual or neuronal  
27 phenomena, but there have been fewer attempts to identify a reliable neural  
28 signature of predictive coding for complex stimuli. In a pair of studies, we test  
29 whether the N300 component of the event-related potential, occurring 250-350  
30 ms post-stimulus-onset, has the response properties expected for such a  
31 signature of perceptual hypothesis testing at the level of whole objects and  
32 scenes. We show that N300 amplitudes are smaller to representative (“good  
33 exemplars”) compared to less representative (“bad exemplars”) items from  
34 natural scene categories. Integrating these results with patterns observed for  
35 objects, we establish that, across a variety of visual stimuli, the N300 is  
36 responsive to statistical regularity, or the degree to which the input is “expected”  
37 (either explicitly or implicitly) based on prior knowledge, with statistically regular  
38 images evoking a reduced response. Moreover, we show that the measure  
39 exhibits context-dependency; that is, we find the N300 sensitivity to category  
40 representativeness when stimuli are congruent with, but not when they are  
41 incongruent with, a category pre-cue. Thus, we argue that the N300 is the best  
42 candidate to date for an index of perceptual hypotheses testing for complex  
43 visual objects and scenes.

44  
45  
46  
47

## 48 **Introduction**

49

50 The stars in the night sky are not arranged in the shape of a great bear and there  
51 is no rabbit on the moon; it is our prior knowledge of these shapes that invokes  
52 such descriptions. Increasingly, it is clear that perception does not depend on the  
53 sensory stimulus alone but is also dynamically influenced by our prior knowledge  
54 (Smith and Loschky 2019; Gordon et al. 2017; Caddigan et al. 2017; Lupyan  
55 2017; Vo and Wolfe 2013; Voss et al. 2012; Summerfield et al. 2006). Indeed,  
56 many models of perception include some form of perceptual hypothesis testing  
57 (PHT), in which perception, a hard inverse problem, is conceived of as a process  
58 of generating a hypothesis on the basis of both sensory input and prior  
59 knowledge and the current context (Clark 2013; Gregory 1980; Hochberg 1981;  
60 Huang and Rao 2011; Rock 1983; Helmholtz 1925). Recently, one class of PHT  
61 models has garnered increased interest: predictive coding models (Rao and  
62 Ballard 1999; Friston 2005; Spratling 2010), which posit that each area of, for  
63 example, visual cortex learns statistical regularities from the world that it then  
64 uses, jointly with the input from the preceding area, to make predictions about the  
65 stimulus. In particular, the prediction and incoming sensory signal are proposed  
66 to undergo an iterative matching process at each stage of the processing  
67 hierarchy. Most of these models are hierarchical in nature, with the prediction  
68 feeding back on the preceding area. The mismatch (“prediction error”), if any,  
69 between the prediction and the incoming sensory signal is then propagated to  
70 higher layers in the processing hierarchy, revising the weights of the hypotheses,  
71 until the feedback matches the incoming signal and the error is zero (Rao and

72 Ballard 1999; Friston 2005; Lange et al. 2018). These predictive coding models  
73 have risen to prominence in recent years, in part because they represent an  
74 efficient coding scheme for the complexity of the visual world and, perhaps more  
75 importantly, because they posit a role for the abundant feedback connections  
76 known to exist between visual areas.

77

78 The bulk of support for predictive coding models has come from the models'  
79 ability to simulate known perceptual or neuronal phenomena (reviewed in  
80 Spratling 2016). The empirical data used for such models have primarily come  
81 from experiments manipulating basic features of simple stimuli, such as  
82 variations in grating orientation or color (Kok et al. 2017; Marzecová et al. 2017,  
83 2018; Rungratsameetaweemana et al. 2018; Smout et al. 2019, 2020). However,  
84 it should also be possible to find signatures of predictive coding at higher levels  
85 of visual analysis. Such a signature would be observed to a variety of types of  
86 complex visual stimuli (objects, faces, natural scenes) across most or all viewing  
87 conditions. More importantly, it should be responsive to statistical regularity, or  
88 the degree to which features in the input are “expected” (either explicitly or  
89 implicitly) by the system based on prior knowledge. We learn regularities of  
90 object and natural scene features by being exposed to prototypical objects and  
91 natural environments over our lifetime. This prior knowledge facilitates our  
92 processing when the regularities in the incoming sensory stream meet our  
93 expectations (Caddigan et al. 2017). Thus, a good measure of predictive coding  
94 would index when stimuli deviate from the regularities we expect to see. In

95 particular, the measured response should increase with increasing irregularity, in  
96 keeping with the increased iterations, or inference-based error, proposed to  
97 occur when an item does not match the prediction. Importantly, the measure  
98 should also show context-dependency, as statistical regularities need to be  
99 sensitive to the immediate context in order to be of use to the system.

100

101 Using complex visual objects, Schendan and colleagues (Schendan and Kutas  
102 2002, 2003, 2007) have shown that the N300 component of the event-related  
103 potential (ERP) can be interpreted as an index of object model selection  
104 processes, a framework that fits within PHT (Schendan and Ganis 2012;  
105 Schendan 2019). Here we build on these findings, addressing the question of  
106 whether the N300 is also sensitive to statistical regularity for complex visual  
107 stimuli other than objects -- in particular, for good and bad examples of visual  
108 scenes. Moreover, critically, we ask whether the N300 is sensitive to in-the-  
109 moment expectations for visual information, as established by, in the present  
110 work, verbal cues. Taken together, this kind of evidence would support the  
111 characterization of the N300 more broadly as a signature of predictive coding  
112 mechanisms, operating in occipitotemporal visual cortex at the scale of whole  
113 objects and scenes.

114

115 *The N300*

116 The N300 is a negative going component with a frontal scalp distribution that  
117 peaks around 300 ms after the onset of a visual stimulus. It has been shown to

118 be sensitive to global perceptual properties of visual input (Mcpherson and  
119 Holcomb 1999; Schendan and Kutas 2002, 2003) but not to manipulations limited  
120 to low level visual features (e.g., color, or small-scale line segments; Schendan  
121 and Kutas 2007) that are known to be processed in early visual cortex.  
122 Components that precede the N300 in time have instead been linked to  
123 processing of and expectations for such low-level features. For example, a  
124 component known as the visual mismatch negativity (vMMN) occurs between  
125 100-160 ms in target-oddball paradigms, where it is larger for the visual oddball  
126 stimuli. The vMMN has sometimes been associated with predictive coding  
127 (Stefanics et al. 2014; Oxner et al. 2019). However, given its sensitivity to the  
128 current experimental context – and, importantly, not to statistical regularities built  
129 up over a lifetime – as well as its source location to occipital cortex (Susac et al.  
130 2014; File et al. 2017), the vMMN would be classified as indexing early stage  
131 PHT processing. In contrast, the N300 is a “late” visual component, with likely  
132 generators in occipitotemporal cortex (Schendan, 2019; Sehatpour et al., 2006).  
133 It immediately precedes access to multimodal semantic memory (reflected in the  
134 N400, which is observed later in time than the N300 when both are present;  
135 Kutas and Federmeier 2011). The N300 is therefore well positioned to capture  
136 the iterative, knowledge- and context-sensitive process of visual processing of  
137 the global features of stimuli, as proposed by predictive coding models, and thus  
138 seems promising as a candidate index of intermediate to late stage PHT  
139 processing.  
140

141 Importantly, as hypothesized by predictive coding models, the amplitude of the  
142 N300 increases for less “expected” (i.e., less statistically regular) stimuli. The  
143 N300 is larger to pictorial stimuli that lack a global structure as compared to when  
144 the global structure of the object is clearly discernible (Schendan and Kutas  
145 2003). The N300 is also sensitive to repetition, with a reduced amplitude for  
146 repeated presentations; importantly, however, N300 repetition effects (but not  
147 those on earlier components) depend on knowledge, as they are larger when the  
148 visual stimulus is meaningful (Voss and Paller 2007; Schendan and Maher 2009;  
149 Voss et al. 2010). Similarly, and critically, N300 amplitudes are sensitive to a  
150 variety of factors that reflect the degree to which an object fits with prior  
151 experience. For example, N300 amplitudes are sensitive to the canonical view of  
152 an object; an open umbrella oriented horizontally (non-canonical) elicits a larger  
153 N300 amplitude than an open umbrella oriented vertically (Schendan and Kutas  
154 2003; Vo and Wolfe 2013). Amplitude modulations are also linked to factors such  
155 as object category membership, presence of category-diagnostic object features,  
156 and (rated) match to object knowledge (Gratton et al., 2009; Schendan, 2019;  
157 Schendan & Maher, 2009). This pattern of data suggests that the N300 may be a  
158 good marker for not only the global structure of an object but the degree to which  
159 the input matches learned statistical regularities more generally, with larger N300  
160 amplitudes for stimuli that do not match predictions based on learned regularities  
161 and hence require further processing.

162

163 Thus far, empirical data have largely linked the N300 to object processing,  
164 sometimes in the context of a scene (Mudrik et al. 2010; Vo and Wolfe 2013;  
165 Lauer et al. 2020), but still ostensibly elicited by an object. Indeed, Schendan  
166 (2019) has specifically linked the N300 to object model selection processes, in  
167 which an input is matched to possible known objects. This model selection  
168 process includes PHT computations. Here, however, we hypothesize that the  
169 N300 may reflect a more general signature of hierarchical inference within higher  
170 level visual processing. If so, it should be elicited by other meaningful visual  
171 stimuli, such as natural scenes. Scenes differ from individual objects in a few  
172 ways. Scenes often contain multiple objects rather than prominent objects that  
173 overshadow their backgrounds. Moreover, the spatial layout of the environment  
174 is much more critical for understanding a photograph of a scene than a  
175 photograph of an object. Finally, it is clear that the human visual system sees  
176 objects and scenes as importantly different as they have sub-systems dedicated  
177 to processing them (Epstein & Kanwisher, 1998). Thus, if the N300 reflects, not  
178 a specific facet of object processing but, more generally, the computations  
179 associated with PHT in higher level vision, then it should also be sensitive to  
180 statistical regularity and prediction during scene processing.

181

182 In fact, scrambled scenes (created by recombining parts of the scene image into  
183 a random jigsaw) have been found to elicit larger N300 amplitudes compared to  
184 intact and identified scenes (Pietrowsky et al. 1996). Because the scrambled  
185 scenes were degraded, however, it is not clear whether these effects simply



186 reflect the disruption to the global structure of the image or a deviation from  
187 statistical regularity more generally. Here we use intact scenes that are either  
188 highly representative of their category (e.g., good exemplars of that category) or  
189 less representative of their category (bad exemplars). Importantly, all the images  
190 are good photographs of real world scenes (i.e., they are not degraded); they are  
191 statistically regular or irregular by virtue of how representative they are of their  
192 category. A highly representative exemplar of its category, by definition, contains  
193 better information about its category and thus serves as a better initial prediction  
194 (i.e., has high statistical regularity). We ask whether such statistically regular and  
195 irregular stimuli elicit differential N300s, as would be hypothesized if this  
196 component is indexing hierarchical inference or predictive coding beyond objects.

197

### 198 *Good and bad scenes*

199 We have previously found that good scene exemplars are more readily detected  
200 than bad exemplars (Caddigan et al. 2010, 2017); that is, participants are better  
201 at discriminating briefly presented and masked intact photographs from fully  
202 phase-scrambled versions when those images are good exemplars of their  
203 category (i.e., beaches, forests, mountains, city streets, highways, and offices).  
204 Good and bad exemplar status was determined with a separate rating task in  
205 which participants rated on a 1-5 scale how representative the image was of its  
206 category. We took the 60 highest and 60 lowest rated images from each  
207 category, and verified that participants were significantly faster and more  
208 accurate at categorizing the good scene exemplars than the bad, indicating that

209 our manipulation captured the degree to which the image exemplified the  
210 category (Torralbo et al. 2013). Importantly, again, there were no artificially  
211 introduced objects in any of the bad exemplars nor were they impoverished or  
212 degraded in any way. Instead, their good and bad status derived entirely from  
213 how representative they were of the category being depicted. Note that,  
214 although category was relevant to the choice of stimuli and whether they were  
215 designated good or bad, in Caddigan et al.' experiments it was completely  
216 irrelevant to the intact/scrambled judgement being made (was the stimuli an  
217 intact photo or noise?). Nonetheless, participants had significantly higher  
218 sensitivity ( $d'$ ) for good than bad exemplars (Caddigan et al. 2010, 2017),  
219 suggesting that with the very brief (34–78 m) masked exposures good exemplars  
220 perceptually cohere into a intact photograph sooner than bad exemplars.

221

222 Relatedly, the categories of those same good exemplars are better decoded,  
223 using fMRI multi-voxel pattern analysis, than are the categories of the bad  
224 exemplars in a number of visual areas, including V1 and the parahippocampal  
225 place area (PPA; Torralbo et al. 2013). Interestingly, the BOLD signal for those  
226 same bad exemplars is larger than that for good exemplars in the PPA (Torralbo  
227 et al. 2013), in keeping with predictions from hierarchical predictive coding (i.e.,  
228 increased activity for the less statistically regular images). The poorer detection  
229 with brief presentations, weaker representations in the brain, and greater activity  
230 evoked by bad than good scene exemplars make these stimuli good candidates  
231 for eliciting a neural signature of hierarchical predictive coding.

232

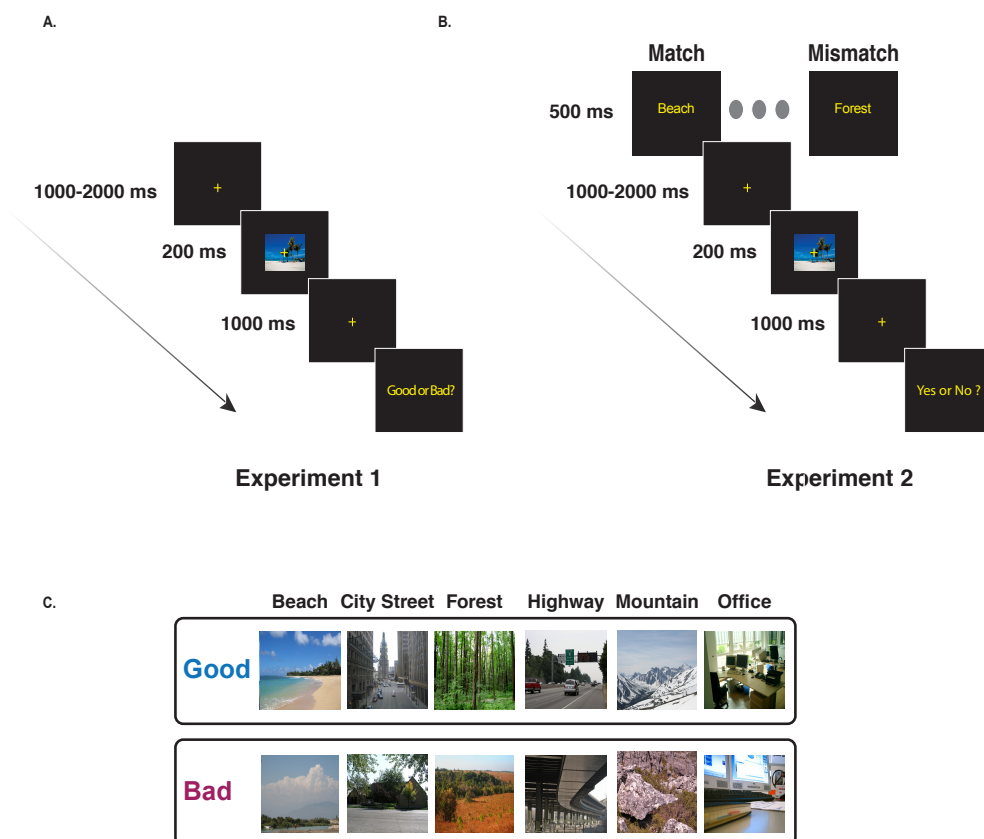
233 *Design of the current experiments*

234 In **Experiment 1**, we recorded scalp EEG while participants viewed good and  
235 bad scene exemplars and made a good/bad judgment. If the N300 serves as an  
236 index of matching incoming stimuli to learned statistical regularities, then N300  
237 amplitude should be smaller for good exemplars of natural scenes than the bad  
238 exemplars. In this first experiment participants viewed the stimuli without any  
239 forewarning of what to expect (category and good/bad status were fully  
240 randomized; see **Figure 1A**), and all the stimuli were unique images with no  
241 repeats in the experiment. If we observe an effect of statistical regularity, then the  
242 particular regularity brought online must stem from the current input, as there  
243 was no confound of repetition priming or episodic memory.

244

245 However, an effective prediction process must also be sensitive to context. Thus,  
246 in **Experiment 2** we then manipulated the expectations of the participants at the  
247 beginning of each trial by presenting a word cue (e.g., 'Beach') that either  
248 matched the upcoming scene's category (on 75% of trials) or mismatched the  
249 upcoming image category (e.g., preceding a forest with the 'Beach' cue; see  
250 **Figure 1B**). If the N300 reflects a PHT process then it should also be sensitive to  
251 the particular template (i.e., statistical regularity) activated by the cue. In  
252 particular, we would predict that a cue with a 75% validity would activate the  
253 statistical regularities associated with the cued category. For images that come

254 from the cued category, then, we should observe smaller N300s for good than  
255 bad exemplars, as in Experiment 1, since good exemplars are a better match to  
256 the statistical regularities of their category. However, in contrast, when the input  
257 image does *not* come from the cued category (i.e., for mismatches), we would  
258 predict a reduction or even elimination of the good/bad N300 effect, since neither  
259 the good nor bad exemplar would fit well with the cued statistical regularity. For  
260 example, good beach exemplars should not systematically provide a better  
261 match to the statistical regularities of a forest than a bad beach does. Experiment  
262 2, then, provides a critical test of the idea that the N300 reflects the process of  
263 matching input to the currently activated template – i.e., the prediction.



264

265 **Figure 1.** Schematic of one trial in each of the experiments. **A.** In **Experiment 1**,  
 266 a fixation cross was shown in the center of screen for a randomly chosen interval  
 267 between 1000-2000 ms. A good or bad exemplar image from one of the six  
 268 categories was then presented for 200ms, followed by a fixation cross. After a  
 269 delay of 1000ms, the subjects respond to the question "Good or Bad?" with a  
 270 button press and the next trial begins. **B.** In **Experiment 2**, the trial sequence is  
 271 similar to **Experiment 1** with the following differences. At the start of each trial a  
 272 word cue (e.g., "Beach") from one of six categories (beaches, city streets,  
 273 forests, highways, mountains, and offices) is shown. At the end of the trial the

274 subjects make a delayed response, with a button press, to the question "Yes or  
275 No?" ("Yes" if the image matches the cue and "No" otherwise) and the next trial  
276 begins. Cue validity was kept high (75%) to promote prediction; on 25% of the  
277 trials, there is a mismatch between the word cue and the image category. **C.** A  
278 sample of good and bad exemplars from each category used in our study.

279

## 280 **Materials and Methods**

281

### 282 Participants

283 The data for **Experiment 1** came from 20 right-handed college-age subjects  
284 (mean age = 24.36 years, range = 18 to 33 years, 12 women), and the data for  
285 **Experiment 2** from a separate set of 20 right-handed subjects (mean age =  
286 22.44; range 18-30 years; 14 women). In both experiments, participants gave  
287 written, informed consent and were compensated for their participation in the  
288 study with course credit or cash. The study was approved by the Institutional  
289 Review Board of the University of Illinois at Urbana-Champaign. All participants  
290 were right-handed, as assessed by the Edinburgh Inventory (Oldfield 1971) and  
291 none had a history of neurological disease, psychiatric disorders, or brain  
292 damage.

293

### 294 Materials and Procedures

295 ERP-eliciting stimuli were pictures of natural scenes from six categories:  
296 beaches, forests, mountains, city streets, highways and offices (Figure 1C). In a  
297 previous study, these images were collected from the internet and rated for their

298 representativeness of the named category on Amazon Mechanical Turk, with  
299 participants answering, e.g., for beaches, “How representative is this image of a  
300 BEACH?” for each image, with the interpretation of the term representativeness  
301 left to the participants (Torralbo et al. 2013). In a separate experiment,  
302 participants were significantly faster and more accurate at categorizing the good  
303 exemplars than the bad, further confirming that our manipulation captured the  
304 degree to which the image exemplified the category. The 60 top rated images  
305 were used as good exemplars for each category, and the 60 lowest rated images  
306 were used as bad exemplars for each category (for details on the choice of good  
307 and bad exemplars see (Torralbo et al. 2013). Images were resized to 340 x 255  
308 pixels and presented on a black background with a fixation cross at the center.  
309 The images were randomly presented at one of three locations: the center of the  
310 scene, or with nearest edge 2 degrees to the left or right of fixation, with a total of  
311 120 good images and 120 bad images presented at each location. Here, we  
312 report only results for centrally-presented images<sup>1</sup>. The stimuli were all unique  
313 images with no repeats in the presentation sequence.  
314

---

<sup>1</sup> The laterally presented scenes were included to separately answer questions about hemispheric biases in scene processing that are outside the scope of this manuscript. Because ERP waveforms for laterally presented stimuli have important morphological differences compared to those from centrally presented stimuli, the data from the two presentation conditions cannot be combined.

315 In **Experiment 1**, participants were instructed at the beginning of the study that  
316 they would be seeing good and bad exemplars of six scene categories and that  
317 their task at the end of each trial was to indicate via button press whether the  
318 image was a good or a bad exemplar of its category. Participants first practiced  
319 with 9 trials to acclimatize to the task environment, and these images were not  
320 repeated in the main experiment. Then, they completed 3 blocks each consisting  
321 of an equal number of trials, for a total of 240 centrally presented trials (trials  
322 were also presented to the left and right visual fields in each block). The trial  
323 counts for centrally presented stimuli, for each category (good and bad  
324 combined) are as follows: beaches = 39; cities = 41; forests = 38; highways = 42;  
325 mountains = 36; offices = 44.

326

327 Participants were seated at a distance of 100 cm from the screen, and the  
328 images subtended a visual angle of  $7.65^\circ \times 5.73^\circ$  (width x height). Subjects were  
329 instructed to maintain fixation on the central fixation cross and to try to minimize  
330 saccades and eye blinks during stimulus presentation. As depicted in **Figure 1A**,  
331 each trial began with a fixation cross presented on a blank screen for a duration  
332 jittered between 1000-2000 seconds (to reduce the impact of slow, anticipatory  
333 components on the ERP signal). The scene image, either a good exemplar or a  
334 bad exemplar from one of the six categories, was presented for a duration of 200  
335 ms, followed by a fixation cross on a blank screen for 500 ms. At the end of the  
336 trial a prompt with "Good or Bad?" was displayed on the screen, and participants  
337 pressed one of two response buttons, held in each hand (counterbalanced



338 across participants), to indicate their judgment. The experiment lasted for  
339 approximately one hour and fifteen minutes. Subjects were given two five-minute  
340 breaks at roughly 25 minutes and 60 minutes from the start of the experiment.

341

342 **Experiment 2** was identical to **Experiment 1**, except that each trial began with a  
343 word cue, presented for 500 ms (**Figure 1B**), which corresponded to one of the  
344 six scene categories used in the experiment: Beach, City Street, Forest,  
345 Highway, Mountain, and Office. For each category, we ensured that five trials of  
346 each type (good and bad exemplars) were mismatched. There were thus 75%  
347 matched trials (15 trials each of good and bad within each of the six scene  
348 categories) and 25% mismatched trials, for a total of 180 (90 good, 90 bad)  
349 matched trials and 60 mismatched trials (30 good, 30 bad). Overall cue validity  
350 was kept high to promote the use of the cue to form expectations about what kind  
351 of image would appear next, while still ensuring that we would nevertheless have  
352 a sufficient number of mismatch trials to obtain a stable ERP to that condition as  
353 well. Instead of making a good or bad judgment, at the end of each trial  
354 participants were prompted to respond “yes” or “no,” with a button press, to the  
355 question of whether or not they thought that the picture had matched the cue.  
356 Hand used to respond “yes” or “no” was counterbalanced.

357

### 358 ERP Setup and Analysis

359 EEG was recorded from 26 channels of passive electrodes that were  
360 equidistantly arranged on the scalp, referenced online to the left mastoid and re-

361 referenced offline to the average of the left and right mastoids. Additional  
362 electrodes placed on the outer cantus of each eye and on the orbital ridge below  
363 the left eye were used to monitor saccadic eye movements and blinks.  
364 Impedances were kept below 5 K $\Omega$  for scalp channels and 10 K $\Omega$  for eye  
365 channels. The signal was bandpass filtered online (0.02 Hz - 100 Hz) and  
366 sampled at 250 Hz. Trials with artifacts due to horizontal eye movements or  
367 signal drift were rejected using fixed thresholds calibrated for individual subjects.  
368 Trials with blinks were either rejected, or, for subjects with higher numbers of  
369 blink artifacts (12 in **Experiment 1** and 8 in **Experiment 2**), were corrected using  
370 a blink correction algorithm (Dale 1994). We confirmed that the analytical results  
371 were unchanged if blinks were rejected instead of corrected. On average,  
372 in **Experiment 1**, 6.83% of good exemplar trials and 9.04% of bad exemplar  
373 trials were rejected due to artifacts, and no condition had fewer than 63 trials per  
374 subject in the analysis. The average number of retained trials was, for good  
375 exemplars, 112 (range 81 to 119) and, for bad exemplars, 109 (range 63 to  
376 120). In Experiment 2, in the match condition, 10.8% of good exemplar trials and  
377 11.09% of bad exemplar trials were rejected due to artifacts and no condition had  
378 fewer than 56 trials per subject in the analysis (retained good exemplar trials:  
379 mean 80 (63-90); retained bad exemplar trials: mean 80 (56-90)). In the  
380 mismatch condition, 10.38% of good exemplar trials and 13.89% of bad exemplar  
381 trials were rejected due to artifacts (retained good exemplar trials: mean 27 (19-  
382 30); retained bad exemplar trials: mean 26 (19-30)).  
383

384 ERPs were epoched for a time period spanning 100 ms before stimulus onset to  
385 920 ms after stimulus onset, with the 100 ms prestimulus interval used as the  
386 baseline. This processed signal was then averaged for each condition within  
387 each subject. A digital bandpass filter (0.2 Hz - 30 Hz) was applied before  
388 measurements were taken from the ERPs. Based on prior work showing that the  
389 N300 is frontally distributed and occurs between 250 ms to 350 ms (Federmeier  
390 and Kutas 2001; Schendan and Kutas 2002, 2003), we measured N300 mean  
391 amplitudes in this time window across the 11 frontal electrode sites: MiPf  
392 (equivalent to Fpz on the 10-20 system), LLPf, RLPf, LMPf, RMPf, LDFr, RDFr,  
393 LMPf, RMPf, LLFr, and RLFr (first letter: R=right, L=left, Mi=midline; second  
394 letter: L=lateral, M=medial, D=dorsal; Pf = prefrontal and Fr= frontal); on the 10-  
395 20 system, this array spans from Fpz to just anterior of Cz and from mastoid to  
396 mastoid laterally, with equidistant coverage. Statistics were computed using R (R  
397 Core Team 2020). Specifically, we used the functions `t.test`, to compute t-tests,  
398 and `ttestbf` (from the package: `BayesFactor`) to compute Bayes Factors. The t-  
399 test and Bayes factor calculations compared the measured condition difference  
400 to 0. For within-subject calculations of confidence intervals, we used the function  
401 `summarySEwithin()` that is based on (Morey 2008). The function `anovaBF` (from  
402 the package: `BayesFactor`) was used to compute Bayes factors for interactions.  
403  
404 For completeness, we also analyzed two ERP components in the time-window  
405 after the N300: the N400 and the Late Positive Complex (LPC). Prior work  
406 examining the N400 to pictures has shown a frontal distribution (Ganis et al.

407 1996), and thus we again used the 11 frontal electrode sites, but now in the time-  
408 window 350-500 ms. For the LPC we chose posterior sites in the time-window of  
409 500-800 ms based on prior work characterizing the distribution and timing of the  
410 LPC (Finnigan et al. 2002).

411

## 412 **Results**

### 413 **Experiment 1**

#### 414 **Behavior**

415 To motivate participants to attend to the scenes, we asked participants to make a  
416 delayed response on each trial, judging whether the exemplar was a good or bad  
417 exemplar of the scene category to which it was presumed to belong. Participants  
418 labeled most good exemplars as “good” (mean = 86.2%, std. dev = 13.9%) and  
419 labeled bad exemplars as “bad” about half the time (mean = 56.2%, std. dev =  
420 15.6%). All trials (irrespective of the choice of the participants) were used for the  
421 planned ERP analyses, but, as described below, we also confirmed that the  
422 results hold when conditionalized on subjects’ responses.

423

#### 424 **ERPs**

425 Grand-averaged ERPs at eight representative sites are plotted in **Figure 2**.  
426 Responses to good and bad exemplars can be seen to diverge beginning around  
427 250 ms after stimulus onset, with greater negativity for bad exemplars than for  
428 good exemplars. The polarity, timing, and frontal scalp distribution of this initial

429 effect is consistent with prior work describing the N300 (Mcpherson & Holcomb,  
430 1999; Schendan & Kutas, 2002, 2003, 2007); see **Supplementary Materials** for  
431 a formal distributional analysis.

#### 432 N300

433 To characterize the good/bad effect on the N300, mean amplitudes were  
434 measured from all 11 frontal electrode sites between 250 and 350 ms. Bad  
435 exemplars elicited significantly larger (more negative) N300 responses (mean = -  
436 6.4  $\mu$ V) than did good exemplars (mean = -5.3  $\mu$ V);  $t(19)=-5.4$  and Bayes Factor  
437 = 747.7 (**Table 1**; for a full distributional analysis see **Supplementary Materials**).  
438 In other words, we see the predicted differential response to statistically irregular  
439 exemplars (bad exemplars) as compared to the statistically regular exemplars  
440 (good exemplars). The larger amplitude for the bad exemplars, as compared to  
441 the good exemplars aligns with PHT predictions that would posit greater  
442 inference error, and, hence, greater iterative processing for the bad exemplars as  
443 compared to the good exemplars. These results also confirm that the N300  
444 indexes a match to statistical regularities of natural scenes and thus extend the  
445 validity of the N300 to not only objects, or objects in scene contexts, but more  
446 broadly to complex natural scenes.

447 The above analysis was computed on all trials, to avoid confounding N300  
448 response patterns with the outcome of late stage decision making processes.  
449 However, for completeness, we also analyzed the results conditionalized on  
450 participants' responses (i.e., including only good trials judged as good and bad

451 trials judged as bad). This yielded the same effect pattern (Bayes factor for  
452 good/bad difference = 5.4;  $t = -2.89$ ,  $p = 0.0094$ ). For details see **Supplementary**  
453 **Materials**. We also analyzed the bad exemplar trials, as about half of them were  
454 judged as good, and did not see an N300 effect based on participants'  
455 judgements of only the bad exemplars (see **Supplementary Materials**).

#### 456 Post N300 Components

457 Although the N300 was the component of primary interest, to more completely  
458 characterize the brain's response to the scenes, we also examined good/bad  
459 differences in later time windows encompassing the N400 (350-500 ms) and Late  
460 Positive Complex (LPC) (500-800 ms). The details of the analyses and results  
461 are provided in the **Supplementary Materials** and summarized here. N400  
462 responses, which index multimodal semantic processing, were larger for bad (-  
463 3.3  $\mu\text{V}$ ) than for good exemplars (-2.2  $\mu\text{V}$ ), suggesting that items that better fit  
464 their category allow facilitated semantic access. We note however, that given the  
465 similar scalp distribution of the N300 and the N400 to picture stimuli (Ganis et al.  
466 1996), it is difficult to tell where the boundary of the two components might be  
467 and thus how much the N400 pattern might be influenced by the preceding N300.  
468 LPC responses were larger -- more positive -- to good (4.5  $\mu\text{V}$ ) than to bad (3.3  
469  $\mu\text{V}$ ) exemplars. The LPC amplitude is known to positively correlate with  
470 confidence in decision making (Finnigan et al. 2002; Schendan and Maher 2009).  
471 Larger LPC responses to good items, therefore, is consistent with the behavioral

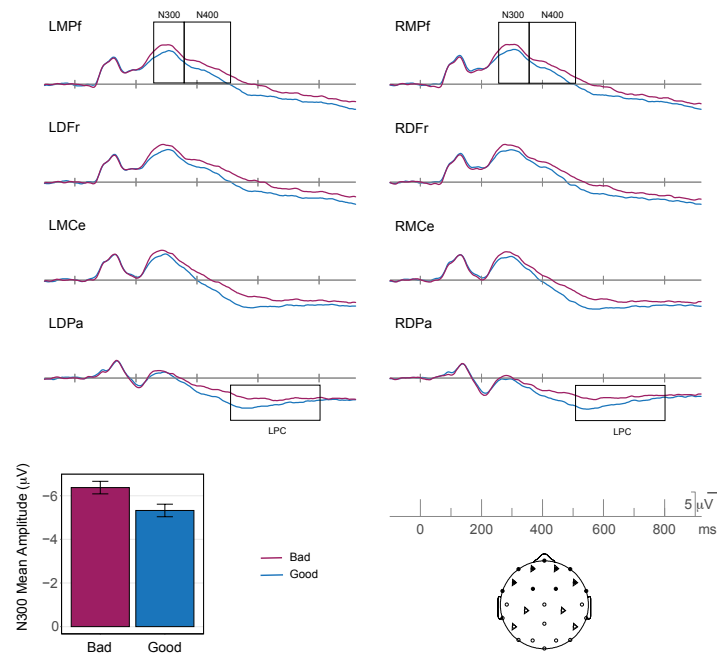
472 pattern in which good exemplars were classified more consistently than bad

473 exemplars.

474 **Figure 2.**

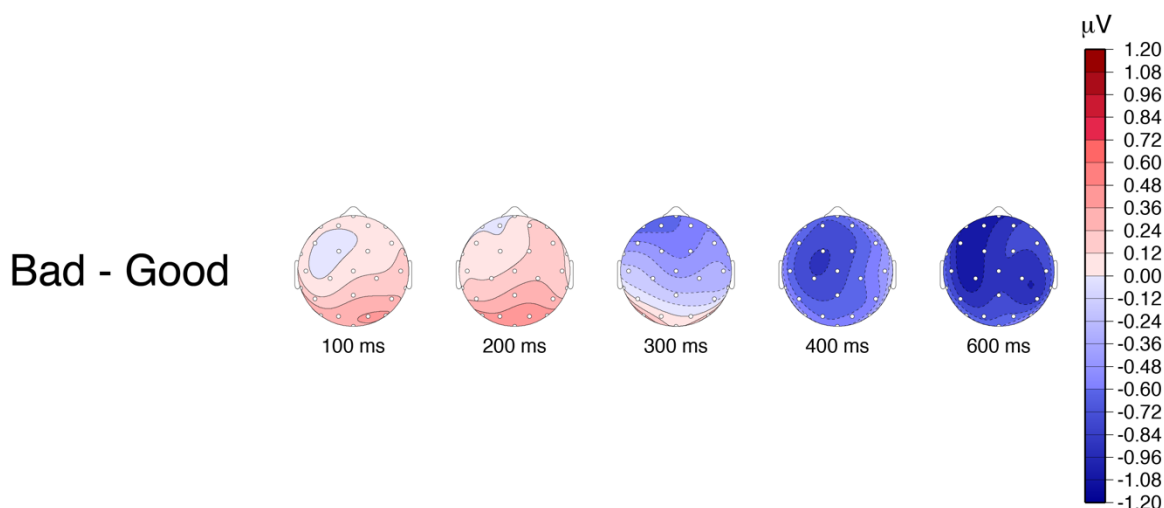
475

476 **A**



477

478 **B**



479

480 **Figure 2 A.** Grand average ERP waveforms for good (blue) and bad (maroon)

481 exemplars in **Experiment 1** are shown at 8 representative electrode sites

482 distributed over the head. Plotted channel locations are marked as triangles on

483 the schematic of the scalp (LMCe and RMCe are just posterior of and lateral to

484 Cz on the 10-20 system). Negative voltage is plotted upwards. The waveforms

485 differ over frontal sites beginning in the N300 time-window (250-350 ms), with

486 greater negativity for bad exemplars as compared to good exemplars. The bar

487 plot gives mean amplitude over the 11 frontal electrode sites (darkened electrode

488 sites on the schematic of the scalp) used for the primary statistical analyses. The

489 error bars plotted are within-subject confidence intervals. N=20. **B.** Topographic

490 plots of the difference waves for the main effect of representativeness (Bad –

491 Good). In the N300 time-window we see a frontal distribution, whereas in the

492 N400 time-window we see a centro-parietal distribution, with a slightly left

493 laterality.

494



495

496

497

498

499

500 **Table 1. Experiment 1**, mean amplitudes in the N300 time-window (250-350 ms)  
501 over 11 frontal electrode sites (see **Figure 2**), along with t-test and Bayes factor  
502 values. The N300 response to bad exemplars is more negative (larger) than that  
503 to good exemplars. The t- test and Bayes factor calculations compared the within  
504 subject Good/Bad difference to 0.

Condition	N	Mean ( $\mu\text{V}$ )	Mean Bad/Good Difference ( $\mu\text{V}$ )	Bad/Good Difference 95% C.I.	t(19)	p	Bayes Factor
Bad	20	-6.4 $\pm$ 0.61	-1.05	-1.46 to -0.64	-5.4	3.3E-05	747.7
Good	20	-5.3 $\pm$ 0.61					

505 Note:  $\pm$  values reflect the normed standard deviation within subjects. C.I. =  
506 confidence interval.

507

## 508 **Experiment 2**

509

510 As mentioned in the introduction, a predictive coding signal should be sensitive to  
511 context. In particular, if the context predicts a specific stimulus category then  
512 initial predictions should reflect the statistical regularities associated with the  
513 predicted category. The good/bad difference observed in **Experiment 1** was  
514 elicited without any expectation regarding the specific category to be presented  
515 (i.e., category and good/bad status were completely randomized). Thus, the  
516 particular template or statistical regularity with which the image was compared  
517 must have been initially elicited by the input itself. This is also the case in almost  
518 all previous work examining the N300 to objects. In **Experiment 2**, therefore, we  
519 set out to examine whether the N300 is sensitive to expectations induced in the  
520 moment by context. We preceded each image with a word cue that either  
521 matched or mismatched the upcoming category. If the N300 difference observed  
522 in **Experiment 1** reflects the matching of incoming stimuli to learned statistical  
523 regularities, we should be able to modulate that difference by activating either the  
524 appropriate (match cue) or inappropriate (mismatch cue) statistical regularity. In  
525 particular, since neither a good nor a bad exemplar of, e.g., a beach, should be a  
526 better match to an inappropriate category (e.g., a forest), we should find that the  
527 N300 good/bad difference is reduced or eliminated when the cue mismatches the  
528 current category.

529

## 530 **Behavior**

531 On each trial, participants were asked to respond if the stimulus matched the  
532 verbal cue (“Yes” or “No”) via a button press. In the match condition, participants  
533 responded “Yes” to good exemplars (mean = 98.7%, std. dev = 2.4%) more often  
534 than to bad exemplars (mean= 67.9% and std. dev = 14.6%). In the mismatch  
535 condition, wherein the exemplars did not fit the cued category, participants  
536 responded “No” to good exemplars (mean = 95.9%, std. dev = 4.6%) more often  
537 than to bad exemplars (mean = 94.0% and std. dev = 5.5%). All trials were used  
538 for the ERP analyses.

539

#### 540 **ERPs**

541 Scenes elicited an N300 response (**Figure 3**) with similar timing, polarity and  
542 scalp distribution to that observed in **Experiment 1**; see the **Supplementary**  
543 **Materials** for a formal distributional analysis. Analyses of N300 mean amplitudes  
544 were conducted using the same time window (250-350 ms) and frontal electrode  
545 sites as in **Experiment 1**, here comparing good and bad exemplars under the  
546 two cueing conditions: match and mismatch.

547

#### 548 N300

549 In the match condition, when the scene was congruent with the verbal cue, we  
550 replicated the N300 effect of **Experiment 1** for the good and bad exemplars, with  
551 a frontally distributed negativity that was larger for the bad exemplars than the  
552 good exemplars (**Figure 3, Table 2A, 2B**). Importantly, and as predicted, this  
553 N300 difference between good and bad exemplars was notably reduced –

554 indeed, likely absent altogether (Bayes factor 0.31) – in the mismatch condition  
555 compared to the match condition (Bayes factor for interaction of Good/Bad x  
556 Cueing = 4.0). This is consistent with the idea that the N300 is indexing the fit of  
557 the incoming stimulus to the template activated by the verbal cue. That is, neither  
558 a good or bad exemplar of category A represents a better match to a template for  
559 category B. The same pattern of results is also seen when the analysis is  
560 conditioned on subjects' judgement; i.e., they responded to a cue congruent stimulus  
561 as 'Yes' and cue incongruent stimulus as "No, see **Supplementary Materials, Table**  
562 **S6**. We note that we chose to discuss the interaction in terms of the good/bad  
563 effect being dependent on a matching cue. However, one might also discuss the  
564 interaction in terms of the effect of cueing being different as a function of  
565 good/bad status. Indeed, the good images show a decrease in the N300 when  
566 they are preceded by a match cue than when they are preceded by a  
567 mismatched cue (Bayes Factor = 1.2 for good mismatch - good match;  $t = -$   
568 1.998,  $p = 0.06$ ), consistent with the mismatch cue producing a prediction error.  
569 In contrast, not only is there little evidence for a cueing effect for bad exemplars  
570 (Bayes Factor = 0.68 for bad mismatch - bad match;  $t = 1.59$ ,  $p = 0.13$ ) but the  
571 difference is numerically in the opposite direction (slightly larger for match).

572

573 For completeness, and to compare the N300 in our experiment with its  
574 characterization in the existing literature, we also performed an ANOVA across  
575 multiple factors: Good/Bad x Cueing (Match/Mismatch) x Anteriority x Laterality x

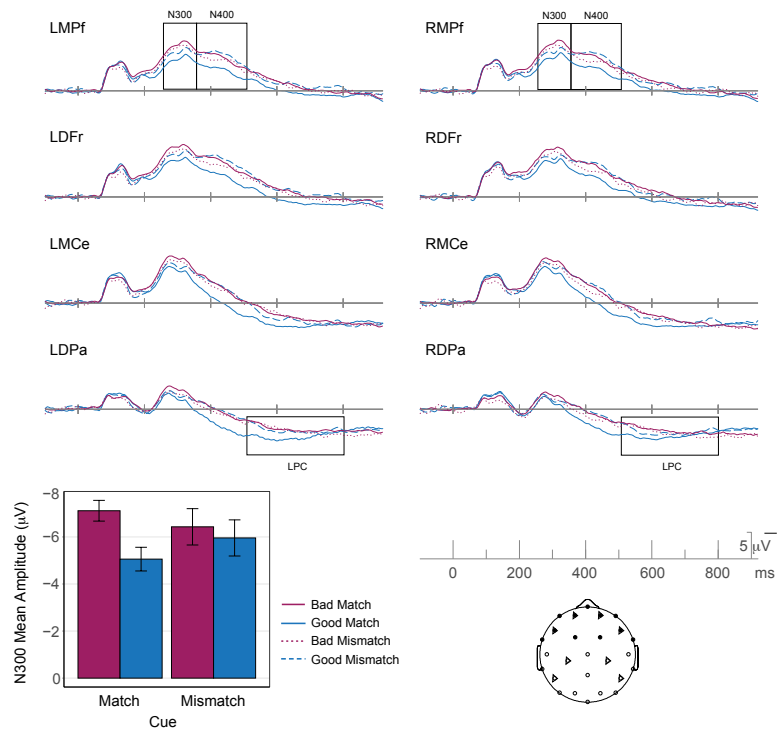
576 Hemisphere. There was a main effect of Good vs. Bad (bad larger than good;  
577  $(F(1,19) = 15.34)$  and an interaction between Good/Bad and Cueing ( $F(1,19)$   
578  $= 5.87$ ), with larger Good/Bad effects when the scene matched the cue. The main  
579 effect of Cueing was not significant ( $F(1,19) = 0$ ). For details on the distributional  
580 analysis see **Supplementary Materials**.

581

582 Finally, to ensure that our results are not due to the differential number of trials in  
583 the match and mismatch condition, we subsampled the trials in the match  
584 condition to be equal to that of the mismatch condition. This subsampling did not  
585 change the results (see **Supplementary Materials, Table S7**).

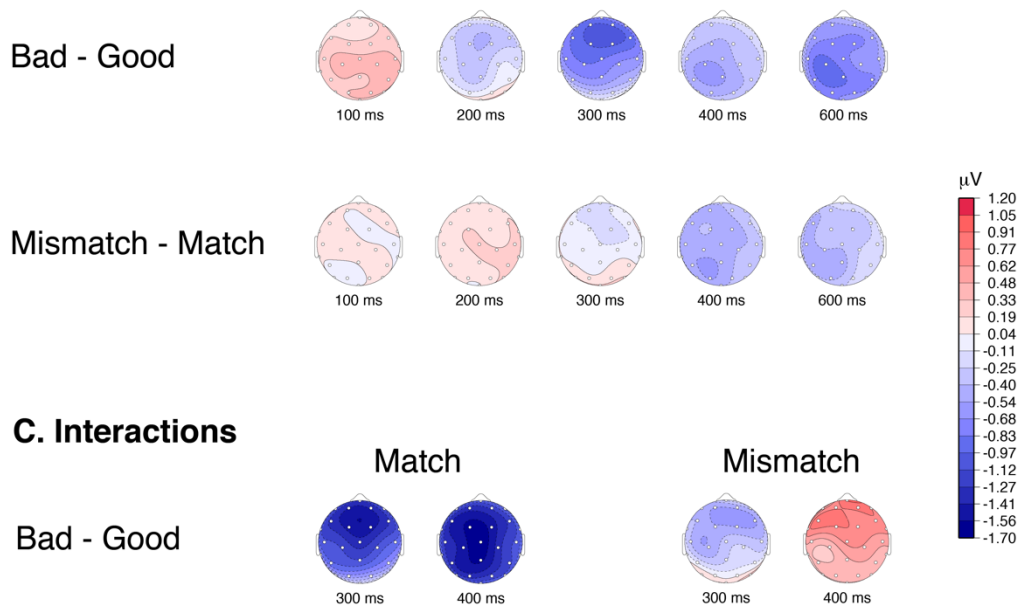
586

587 **A**

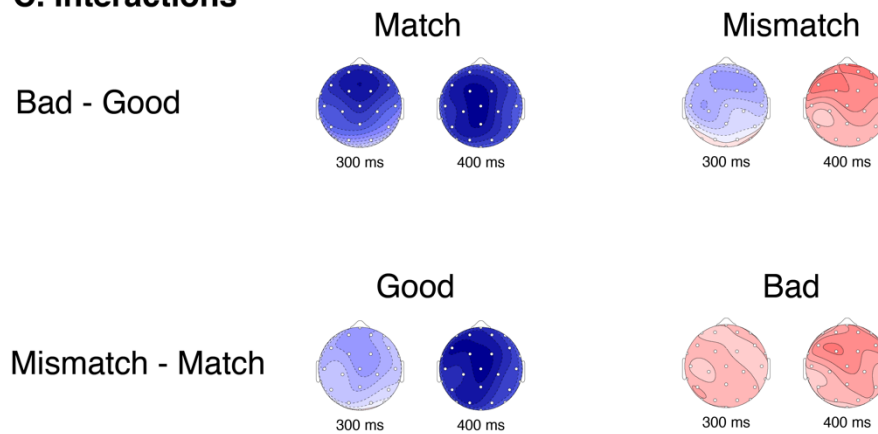


588  
589

## B. Main Effects



## C. Interactions



590

591

592 **Figure 3 A.** Grand average ERP waveforms for the good-match (solid-blue),  
 593 bad-match (solid-maroon), good-mismatch (dashed-blue), and bad-mismatch  
 594 (dotted-maroon) conditions in **Experiment 2** are shown at the same 8  
 595 representative electrode sites. In the match condition, responses to good and  
 596 bad exemplars differ in the N300 time-window (250-350 ms), with greater  
 597 negativity for bad exemplars as compared to good exemplars, over frontal sites  
 598 (darkened electrode sites on the schematic of the scalp). In the mismatch  
 599 condition, the differences between good and bad exemplars on the N300 are  
 600 diminished/eliminated. The bar plot gives the grand average mean of the ERP  
 601 amplitude over the 11 frontal electrode sites (darkened electrode sites on the  
 602 schematic of the scalp) used for the primary statistical analyses (N = 20). The  
 603 plotted error bars are within-subject confidence intervals. **B.** Topographic plots of  
 604 the difference waves for the two main effects of representativeness (Bad – Good)  
 605 and cueing (Mismatch – Match). In the N300 time-window the two main effects  
 606 are qualitatively similar, with both main effects showing a frontal distribution. The  
 607 N300 time-window also shows a quantitatively larger effect for the  
 608 representativeness (Bad – Good) than for the cueing (Mismatch – Match). In the

609 N400 time-window, both effects are centro-parietally distributed with a slight left  
610 laterality. **C.** Topographic plots for the difference in the interactions for Good/Bad  
611 x Cuing are shown in two interpretations: in terms of the Good/Bad effect - (Bad-  
612 Good) x Match and (Bad-Good) x Mismatch; and in terms of the cuing effect -  
613 Good x (Mismatch -Match) and Bad x (Mismatch -Match).

614

615

616

617

618

619

620

621

622

623 **Table 2A.** The grand average mean values, in the N300 time-window (250-350  
624 ms), shown for 11 frontal electrode sites (see **Figure 3**), along with t-test and  
625 Bayes factor values. There is strong evidence (large Bayes factor) for greater  
626 negativity of the N300 for bad exemplars as compared to good exemplars when  
627 the cue matches the stimulus. When there is a mismatch between the cue and  
628 the stimulus there is no evidence (small Bayes factor) for the difference between  
629 good and exemplars in the N300 time-window. The t- test and Bayes factor  
630 calculations compared the within subject Good/Bad difference to 0.

631

---

Condition	Cue	N	Mean ( $\mu$ V)	Mean Difference ( $\mu$ V)	95% C.I.	t(19)	p	Bayes Factor
-----------	-----	---	--------------------	----------------------------------	----------	-------	---	-----------------

---



Bad	Match	20	-7.1±0.94					
Good	Match	20	-5.1±1.07	-2.06	-2.6 to -1.5	-7.4	5.6E-07	30457
Bad	Mismatch	20	-6.4±1.65					
Good	Mismatch	20	-6.0±1.64	-0.47	-1.7 to 0.73	-0.82	0.42	0.31
Good mismatch – Good match		20		-0.9	-1.84 to 0.04	-1.998	0.06	1.20
Bad mismatch – Bad match		20		0.68	-0.22 to 1.58	1.59	0.13	0.68

632

633 **Table 2B.** The Bayes factor for the main effects and interaction computed using  
 634 Bayesian ANOVA. This shows that there is evidence for the interaction of  
 635 Good/Bad x Cueing in **Experiment 2.**

Effect	Bayes Factor
Good/Bad	118.1
Cueing	0.2
Good/Bad x Cueing	4.0

636

637 Note: ± values reflect the normed standard deviation within subjects.

638

### 639 Post N300 Components

640 Again, for completeness, we also examined effects on the N400 (350-500 ms)  
 641 and Late Positive Complex (LPC) (500-800 ms). These are presented in full in  
 642 the **Supplementary Materials** and summarized here. Given prior work (reviewed

643 in Kutas and Federmeier 2011), we expected the N400 to be particularly  
644 sensitive to the match between the verbal cue and the scene category. Indeed,  
645 overall, N400 responses to good scenes that matched the verbal cue were  
646 facilitated (more positive:  $-3.5 \mu\text{V}$ ) than to good scenes that mismatched their  
647 cues ( $-5.6 \mu\text{V}$ ), consistent with the large literature on N400 semantic priming (see  
648 **Table S4**). Moreover, we replicated the effect in **Experiment 1**: N400 amplitudes  
649 were also larger for bad ( $-5.3 \mu\text{V}$ ) than for good exemplars ( $-3.5 \mu\text{V}$ ) in the match  
650 condition, although, again, we cannot rule out influence from the prior N300  
651 effects on the observed pattern. We see an interaction of Good/Bad x Cuing in  
652 the N400 window ( $F = 13.7$ ;  $p = 0.0015$ ;  $E = 1$ ), with the largest facilitation for good  
653 exemplars in the match condition. LPCs were larger (more positive) for good  
654 exemplars in the match condition ( $2.7 \mu\text{V}$ ) compared to both bad exemplars ( $0.4$   
655  $\mu\text{V}$ ) in the match condition (replicating **Experiment 1**) and to either scene type in  
656 the mismatch condition (Good:  $0.2 \mu\text{V}$ ; Bad:  $0.9 \mu\text{V}$ ), presumably reflecting the  
657 increased ease and confidence of responding to the good match items (see  
658 **Table S5**).

659

## 660 **Discussion**

661 In two experiments, we tested whether the N300 component of the ERP has  
662 response properties expected for an index of hierarchical predictive coding  
663 during late stage visual processing, when global features of the stimulus are  
664 being processed. Across many studies, larger (more negative) N300 responses  
665 have been observed for conditions that might be characterized as statistically

666 irregular (Pietrowsky et al. 1996; Schendan and Kutas 2002, 2003, 2007; Mudrik  
667 et al. 2010; Vo and Wolfe 2013). However, the focus of the literature thus far has  
668 been limited to objects, objects in scenes, or artificially degraded stimuli. If the  
669 N300 more generally reflects predictive hypothesis testing in later visual  
670 processing, then it should be sensitive to statistical regularity outside of the  
671 context of object processing and artificial manipulations of global structure. To  
672 this end, in Experiment 1 we showed that the N300 is sensitive to the difference  
673 between good (statistically regular) and bad (statistically irregular) exemplars of  
674 natural scenes. Because none of the scenes we used were degraded, had any  
675 misplaced elements, or contained objects that were surprising or violated  
676 expectations (e.g., a watermelon instead of the expected basketball; see Mudrik  
677 et al. 2010; Vo & Wolfe, 2013), these results strongly link N300 modulations to  
678 statistical regularity as such.

679

680 Predictive coding posits a larger inference error in processing statistically  
681 irregular items (bad exemplars) as compared to statistically regular items (good  
682 exemplars), and, consistent with this, N300 responses were larger for the  
683 statistically irregular exemplars. Note that the observed pattern cannot be  
684 explained by interstimulus perceptual variance (ISPV; Thierry et al., 2007;  
685 Schendan and Ganis, 2013). The good exemplars we used have more consistent  
686 low-level image statistics, and thus lower ISPV, than the bad exemplars (see  
687 Torralbo et al. 2013). Thus, if the pattern were driven by ISPV, we would have  
688 expected the good exemplars to elicit larger ERP modulations (see Thierry et al.,

689 2007; Schendan and Ganis, 2013). Instead, we found that the good exemplars  
690 have a lower amplitude ERP, consistent with the claim that it is statistical  
691 regularity – and not ISPV – that is responsible for the effect.

692

693 The data from Experiment 1, in combination with prior experiments, show that the  
694 N300 manifests the expected response properties for a general index of  
695 predictive coding mechanisms for late stage visual processing (for studies that  
696 rule out the N300 indexing early visual processing see Schendan and Kutas  
697 2002; Johnson and Olshausen 2003) of complex objects and scenes. Across the  
698 literature, the kinds of stimuli distinguished by the N300 encompass global  
699 structure, canonical viewpoints, probable views of objects in scene contexts, and,  
700 in our own experiment, the category-level representativeness of the stimuli. We  
701 would like to collectively refer to these properties as learned statistical  
702 regularities. We mean statistics in the Bayesian sense: The statistical regularities  
703 reflect the system's prior belief. Although frequency of occurrence may be one  
704 factor that goes into constructing a regularity, the regularities should be more  
705 sophisticated than simple frequency. They should be constructed to maximize  
706 the informativeness of the prediction and minimize, on average, the amount of  
707 updating needed. Thus, canonicity, prototypicality or representativeness will all  
708 be critical determinants of the regularities, as well as frequency or familiarity. A  
709 collection of these regularities can be viewed as a template (see also Johnson  
710 and Olshausen 2003), constructed to reduce, on average, the prediction error.  
711 Thus, we can think of the differences on the N300 component as an indicator of

712 the degree to which an incoming exemplar can be matched with a template, with  
713 greater negativity for a stimulus when it doesn't match a template as compared to  
714 when it does.

715

716 In **Experiment 1**, neither scene category nor exemplar status (good or bad) was  
717 predictable from trial to trial, and thus the statistical regularity driving the  
718 observed effect must have been acquired over the life time (i.e., learning what  
719 does and does not constitute a good exemplar of a category), rather than within  
720 the context of the experiment. However, a key attribute of PHT models, of which  
721 predictive coding is a popular example, is that the hypotheses that are generated  
722 are sensitive to the current context. If the N300 reflects a template matching  
723 process, such that the input is compared against a contextually-relevant learned  
724 statistical regularity, then the N300 sensitivity to statistical regularity should vary  
725 in the moment, as a function of context.

726 In **Experiment 2**, therefore, we set up expectations for a particular category on  
727 each trial using a word cue with high validity, with the aim of pre-activating a  
728 particular scene category template. Critically, however, on 25% of trials the  
729 scene did not match the cued category. We found that the N300 is indeed  
730 sensitive to regularities cued by the current context. When the scenes were  
731 congruent with the cued category, we observed a significant effect of statistical  
732 regularity (good versus bad) in the N300 time-window, replicating the results from  
733 **Experiment 1**. Here the good exemplars provide a better match to the activated

734 template than the bad exemplars, and thus the reduced inference error or  
735 iterative matching is reflected in the amplitude of the N300. In the mismatching  
736 condition, however, the presented stimulus, whether a good or bad exemplar of  
737 its *own* category, does not match the *cued* template (e.g., a “Forest” template  
738 has been cued but a good or bad beach scene was presented). In this case,  
739 notably, we failed to observe a reliable difference between the N300 to good and  
740 bad exemplars. In the language of predictive coding models, similar inference  
741 errors would be generated for both statistically regular (good) and irregular (bad)  
742 exemplars that mismatch the activated template, as they would both violate the  
743 predicted regularities – or, at least, neither good nor bad exemplars of another  
744 category should violate the predicted regularities more than the other. Beyond  
745 the statistical regularities learned over a lifetime, including our increased  
746 familiarity with more prototypical inputs, the N300 shows sensitivity to the specific  
747 expectations the visual system has in the moment, generated from the current  
748 context.

749

750 Others have discussed the use of visual templates in the context of holding  
751 information active in memory to afford optimal performance on, e.g., visual  
752 matching tasks. In the case of sequential match paradigms, it is assumed that  
753 subjects can hold on to a recently seen target object – the “template” in this case  
754 – and then use that information to judge subsequent stimuli. Indeed, in these kind  
755 of paradigms, differences in anterior ERPs (which may be labeled N2s or N300s;

756 see discussion in Schendan 2019) have been observed between the match and  
757 mismatch conditions. Moreover, using a verbal cue for object type (e.g., “dog”  
758 followed by an image), Johnson & Olshausen (2003) observed a significant effect  
759 of cueing on a frontally-distributed negativity between 150 and 300 ms, which  
760 likely is encompassed by what we are calling the N300. Responses were more  
761 positive when the image matched the cue compared to when it did not. They did  
762 not vary the representativeness of their images, but it is reasonable to assume  
763 that they were on average more representative than our bad images, specifically  
764 chosen to be less representative. Thus, our results are in accordance with those  
765 of Johnson & Olshausen (2003), and extend them, not only to natural scenes, but  
766 also by showing that the effect of cuing interacts with sensitivity to statistical  
767 regularity. Thus, Experiment 2 brings together two important facets of visual  
768 processing on a PHT framework. First, is the fact that the visual system builds  
769 templates based on statistical regularities, accumulated over the lifespan, and  
770 routinely uses those templates, elicited by the input itself, to guide its iterative  
771 processing. Second, then, is that fact that context information (such as a verbal  
772 cue) can cause a *particular* template to be activated in advance of the input,  
773 biasing processing toward that template.

#### 774 **The N300 Indexes Perceptual Hypothesis Testing**

775 We can think of visual identification and categorization as a cascade of  
776 processes, starting with identification of low level visual features, followed by  
777 perceptual grouping of features, and then appreciation of the “whole” visual form

778 of objects and scenes, after which processing moves beyond the visual modality  
779 into multi-modal semantics and decision making. PHT mechanisms can work  
780 within and across each of these stages. In the context of object processing, prior  
781 work on the N300 has posited it as an index of object model selection, an  
782 intermediate stage in the process of object identification and categorization  
783 (Schendan, 2019; Schendan & Kutas, 2002, 2003, 2007). Having extended the  
784 N300 differences to natural scenes, we propose that the N300 reflects PHT  
785 mechanisms in this intermediate stage more broadly, not just object selection.  
786 Similar to other work (Schendan, 2019), we believe that the N300 reflects  
787 processing at the point wherein the input is matched to items in memory with  
788 similar perceptual structures. However, our data show that this process is not  
789 limited to objects and that it makes use of variety of statistical regularities learned  
790 from the world, including those critical for processing both objects and scenes.  
791 The broadened view of the N300 as being reflective of a general visual template  
792 matching process would suggest that its source be occipitotemporal visual areas.  
793 Indeed, the N300 response to objects has been source localized to  
794 occipitotemporal visual areas (Schendan & Lucia, 2010; Sehatpour et al., 2006).  
795 Although the N300 for scenes has not yet been source localized, a high-density  
796 ERP study on scene categorization localized activity in the 200-300 ms time  
797 window to these same occipitotemporal visual areas (Greene and Hansen 2020).  
798 Similarly, Kaiser and colleagues (2019, 2020), using both fMRI and ERPs,  
799 demonstrated a similar sensitivity to intact versus jumbled scenes in the occipital  
800 place area and PPA as they did in the N300 time window. Moreover, our prior



801 fMRI work with good and bad scene exemplars (Torralbo et al. 2013) would  
802 suggest that the N300 for scenes originates in the PPA, a region known to  
803 preferentially process natural scenes (Epstein and Kanwisher 1998). Using the  
804 same good and bad scene exemplars as in our experiments, we found that, in  
805 the PPA, bad exemplars elicited a greater BOLD signal than good exemplars  
806 (Torralbo et al. 2013), mirroring the effect we observed for the N300.  
807 Interestingly, in that same PPA region of interest we observed that good  
808 exemplars were better decoded than bad exemplars; that is, we were better able  
809 to predict the scene category presented on the basis of activity patterns when the  
810 scene was a good exemplar than when it was bad in the same region that  
811 showed greater activity for the bad exemplar (Torralbo et al. 2013). In other  
812 words, it was not the case that reduced activity for good exemplars reflected a  
813 weaker representation but instead likely reflected a more efficient representation,  
814 an interpretation that aligns nicely with our characterization of the N300 effect as  
815 one of visual template matching in occipitotemporal cortex. We suggest that the  
816 N300 may be interpreted as a component that reflects the iterative processing,  
817 as posited by PHT, in occipitotemporal cortical regions, which helps match  
818 previously learned regularities of objects and scenes with the incoming stimulus.

819

820 Although we are arguing that the N300 indexes PHT for late stage visual  
821 processing of complex visual objects and scenes, it is possible that other  
822 components could index PHT at other stages of processing. For example, PHT

823 matching low level sensory features, such as gratings (Kok et al. 2012), to  
824 hypotheses about such low level features should occur at earlier stages in the  
825 processing hierarchy. Earlier visual sensory components can manifest sensitivity  
826 to expected visual features (Boutonnet and Lupyan 2015) or to differences  
827 between well-learned visual categories, such as words vs. objects, and faces vs.  
828 objects (Schendan et al. 1998) – category comparisons that are thus at a much  
829 higher taxonomy than within objects or scenes. Of particular relevance to PHT is  
830 the vMMN which, as overviewed in the introduction, temporally precedes the  
831 N300 and has been observed in experimental contexts wherein a stream of  
832 standard stimuli that share particular low-level visual features (e.g., orientation,  
833 color) is occasionally interrupted by the presentation of a target stimulus that  
834 carries a featural difference (Stefanics et al. 2014; Oxner et al. 2019). Thus, the  
835 vMMN is sensitive to the context of recent exposure to low-level visual  
836 information, possibly reflecting PHT processes at that lower level.

837

838 The N300, instead, does not modulate with low-level differences and manifests  
839 sensitivity to both regularities established through long-term experience and  
840 knowledge-based expectations derived from semantic contextual information. It  
841 may thus index a late stage of visual PHT, at the transition into multimodal,  
842 semantic processing. Immediately after the N300, ERP responses to complex  
843 objects and scenes are characterized by an N400, which we also observe in our  
844 experiment. The N400 is widely accepted as a signature of multi-modal semantic

845 processing, elicited by not only visual words and pictures, but also meaningful  
846 stimuli in other modalities (see review Kutas & Federmeier, 2011), whereas the  
847 N300 seems to be about visual perceptual structure (Schendan, 2019; Schendan  
848 & Kutas, 2002, 2003, 2007). In some cases, it may be difficult to disentangle the  
849 precise contributions of the N300 and N400 to observed effects of object  
850 categorization and match to object knowledge (Gratton et al., 2009; Schendan,  
851 2019; Schendan & Maher, 2009) since the N400 is known to be sensitive to the  
852 fit between, e.g., a picture and its context (Ganis et al. 1996; Federmeier and  
853 Kutas 2002). Importantly, however, this does not impact the critical effect of our  
854 good versus bad scenes, as neither contain contextually inappropriate items, nor,  
855 in **Experiment 1**, did we set up any context prior to an image (i.e., the scene  
856 category is unpredictable).

## 857 **Conclusion**

858 In a set of experiments we have provided support for the hypothesis that the  
859 N300 component is an index of PHT at the level of whole-objects and scenes.  
860 Using statistically regular and irregular exemplars of natural scenes, we showed  
861 that items that are a poorer match to our learned regularities for types of scenes  
862 – and, thus, inputs that should lead to larger inference errors in a predictive  
863 coding framework – indeed evoked a larger N300 amplitude compared to  
864 statistically regular exemplars, even when the upcoming scene category was not  
865 predictable. We further showed, not only that N300 responses to scenes are  
866 modulated by context – such as the scene category predicted by a verbal cue --  
867 but that they behave as expected for a template matching process in which

868 statistically regular images procure their advantage by virtue of matching the  
869 current contextual prediction.

870

871 Our work thus not only extends prior work on the N300 to natural scenes but it  
872 suggests that the N300 reflects a general template/model selection process of  
873 the sort proposed by PHT models, such as predictive coding. We propose that  
874 the N300 indexes visual inference processing in a late visual time-window that  
875 occurs at the boundary between vision and the next stage of multi-modal  
876 semantic processing. Further studies will be needed to explore the full range of  
877 the N300 response. For example, does it require that the object or scene is  
878 attended or might it proceed more automatically? Can it be modulated by  
879 contexts set up in different modalities (e.g., auditory inputs: speech, sounds)?  
880 Regardless, we propose that the N300 can serve as a useful marker of  
881 knowledge guided visual processing of objects and scenes, with templates based  
882 on prior knowledge serving as hypotheses for visual inference as posited by  
883 PHT.

884

885

#### 886 **Funding**

887

888 This work was supported by Office of Naval Research (grant to D.M.B); National  
889 Institutes of Health (R01 AG026308 to K.D.F); and the James S. McDonnell  
890 foundation (grant to K.D.F).

891

#### 892 **Acknowledgments**

893

894 We would like to thank Yanqi Zhang for assistance with running subjects in  
895 **Experiment 1**, and Resh Gupta, and Nirupama Mehrotra for helping with  
896 **Experiment 1** data collection. We also thank Rami Alsaqri, Johan Saelens, Daria  
897 Niescierowicz, and Benjamin D. Schmitt for helping with data collection in  
898 **Experiment 2**.

899

## 900 **References**

- 901 Boutonnet B, Lupyan G. 2015. Words Jump-Start Vision: A Label Advantage in  
902 Object Recognition. *Journal of Neuroscience*. 35:9329–9335.
- 903 Caddigan E, Choo H, Fei-Fei L, Beck DM. 2017. Categorization influences  
904 detection: A perceptual advantage for representative exemplars of natural  
905 scene categories. *Journal of Vision*. 17:21.
- 906 Caddigan E, Walther DB, Fei-Fei L, Beck DM. 2010. Perceptual differences  
907 between natural scene categories. *OPAM 2010 18th Annual Meeting*.  
908 *Visual Cognition*. 18:1498–1502.
- 909 Clark A. 2013. Whatever next? Predictive brains, situated agents, and the future  
910 of cognitive science. *Behavioral and Brain Sciences*. 36:181–204.
- 911 Dale AM. 1994. Source localization and spatial discriminant analysis of event-  
912 related potentials: linear approaches (brain cortical surface).
- 913 Epstein R, Kanwisher N. 1998. A cortical representation of the local visual  
914 environment. *Nature*. 392:598–601.
- 915 Federmeier KD, Kutas M. 2001. Meaning and modality: Influences of context,  
916 semantic memory organization, and perceptual predictability on picture  
917 processing. *Journal of Experimental Psychology: Learning, Memory, and*  
918 *Cognition*. 27:202.
- 919 Federmeier KD, Kutas M. 2002. Picture the difference: Electrophysiological  
920 investigations of picture processing in the two cerebral hemispheres.  
921 *Neuropsychologia*. 40:730–747.
- 922 File D, File B, Bodnár F, Sulykos I, Kecskés-Kovács K, Czigler I. 2017. Visual  
923 mismatch negativity (vMMN) for low- and high-level deviances: A control  
924 study. *Atten Percept Psychophys*. 79:2153–2170.
- 925 Finnigan S, Humphreys MS, Dennis S, Geffen G. 2002. ERP ‘old/new’ effects:  
926 memory strength and decisional factor (s). *Neuropsychologia*. 40:2288–  
927 2304.
- 928 Friston K. 2005. A theory of cortical responses. *Philosophical Transactions of the*  
929 *Royal Society B: Biological Sciences*. 360:815–836.
- 930 Ganis G, Kutas M, Sereno MI. 1996. The Search for “Common Sense”: An  
931 Electrophysiological Study of the Comprehension of Words and Pictures in  
932 Reading. *Journal of Cognitive Neuroscience*. 8:89–106.

- 933 Gordon N, Koenig-Robert R, Tsuchiya N, van Boxtel JJ, Hohwy J. 2017. Neural  
934 markers of predictive coding under perceptual uncertainty revealed with  
935 Hierarchical Frequency Tagging. *eLife*. 6:e22749.
- 936 Gratton C, Evans KM, Federmeier KD. 2009. See what I mean? An ERP study of  
937 the effect of background knowledge on novel object processing. *Mem*  
938 *Cognit*. 37:277–291.
- 939 Greene MR, Hansen BC. 2020. Disentangling the Independent Contributions of  
940 Visual and Conceptual Features to the Spatiotemporal Dynamics of Scene  
941 Categorization. *J Neurosci*. 40:5283–5299.
- 942 Gregory RL. 1980. Perceptions as Hypotheses. *Phil Trans R Soc Lond B*.  
943 290:181–197.
- 944 Helmholtz H von. 1925. *Treatise on physiological optics*, Bd. 3 : The perceptions  
945 of vision. English translation of the 3rd edition. ed. The Optical Society of  
946 America.
- 947 Hochberg J. 1981. On cognition in perception: Perceptual coupling and  
948 unconscious inference. *Cognition*. 10:127–134.
- 949 Huang Y, Rao RPN. 2011. Predictive coding. *Wiley Interdisciplinary Reviews:*  
950 *Cognitive Science*. 2:580–593.
- 951 Johnson JS, Olshausen BA. 2003. Timecourse of neural signatures of object  
952 recognition. *J Vis*. 3:499–512.
- 953 Kaiser D, Häberle G, Cichy RM. 2020. Cortical sensitivity to natural scene  
954 structure. *Human Brain Mapping*. 41:1286–1295.
- 955 Kaiser D, Turini J, Cichy RM. 2019. A neural mechanism for contextualizing  
956 fragmented inputs during naturalistic vision. *eLife*. 8:e48182.
- 957 Kok P, Jehee JFM, de Lange FP. 2012. Less is more: expectation sharpens  
958 representations in the primary visual cortex. *Neuron*. 75:265–270.
- 959 Kok P, Mostert P, de Lange FP. 2017. Prior expectations induce prestimulus  
960 sensory templates. *Proceedings of the National Academy of Sciences*.  
961 114:10473–10478.
- 962 Kutas M, Federmeier KD. 2011. Thirty years and counting: Finding meaning in  
963 the N400 component of the event related brain potential (ERP). *Annu Rev*  
964 *Psychol*. 62:621–647.
- 965 Lange FP de, Heilbron M, Kok P. 2018. How Do Expectations Shape  
966 Perception? *Trends in Cognitive Sciences*. 22:764–779.
- 967 Lauer T, Willenbockel V, Maffongelli L, Võ ML-H. 2020. The influence of scene  
968 and object orientation on the scene consistency effect. *Behavioural Brain*  
969 *Research*. 394:112812.
- 970 Lupyan G. 2017. Changing What You See by Changing What You Know: The  
971 Role of Attention. *Front Psychol*. 8.
- 972 Marzecová A, Schettino A, Widmann A, SanMiguel I, Kotz SA, Schröger E. 2018.  
973 Attentional gain is modulated by probabilistic feature expectations in a  
974 spatial cueing task: ERP evidence. *Sci Rep*. 8:54.
- 975 Marzecová A, Widmann A, SanMiguel I, Kotz SA, Schröger E. 2017. Interrelation  
976 of attention and prediction in visual processing: Effects of task-relevance  
977 and stimulus probability. *Biological Psychology*. 125:76–90.

- 978 Mcpherson WB, Holcomb PJ. 1999. An electrophysiological investigation of  
979 semantic priming with pictures of real objects. *Psychophysiology*. 36:53–  
980 65.
- 981 Morey RD. 2008. Confidence intervals from normalized data: A correction to  
982 Cousineau (2005)." 4.2 (2008):61. *Web. Reason*. 61.
- 983 Mudrik L, Lamy D, Deouell LY. 2010. ERP evidence for context congruity effects  
984 during simultaneous object–scene processing. *Neuropsychologia*. 48:507–  
985 517.
- 986 Oldfield RC. 1971. The assessment and analysis of handedness: The Edinburgh  
987 inventory. *Neuropsychologia*. 9:97–113.
- 988 Oxner M, Rosentreter ET, Hayward WG, Corballis PM. 2019. Prediction errors in  
989 surface segmentation are reflected in the visual mismatch negativity,  
990 independently of task and surface features. *Journal of Vision*. 19:9–9.
- 991 Pietrowsky R, Kuhmann W, Krug R, Molle M, Fehm HL, Born J. 1996. Event-  
992 related brain potentials during identification of tachistoscopically presented  
993 pictures. *Brain and Cognition*. 32:416–428.
- 994 R Core Team. 2020. R: A language and environment for statistical computing. R  
995 Foundation for Statistical Computing, Vienna, Austria.
- 996 Rao RPN, Ballard DH. 1999. Predictive coding in the visual cortex: a functional  
997 interpretation of some extra-classical receptive-field effects. *Nat Neurosci*.  
998 2:79–87.
- 999 Rock I. 1983. *The Logic Of Perception*. Cambridge: MIT Press.
- 1000 Rungratsameetaweemana N, Itthipuripat S, Salazar A, Serences JT. 2018.  
1001 Expectations Do Not Alter Early Sensory Processing during Perceptual  
1002 Decision-Making. *J Neurosci*. 38:5632–5648.
- 1003 Schendan HE. 2019. Memory influences visual cognition across multiple  
1004 functional states of interactive cortical dynamics. In: Federmeier KD,  
1005 editor. *Psychology of Learning and Motivation*. Academic Press. p. 303–  
1006 386.
- 1007 Schendan HE, Ganis G. 2013. Face-specificity is robust across diverse stimuli  
1008 and individual people, even when interstimulus variance is zero.  
1009 *Psychophysiology*. 50:287–291.
- 1010 Schendan HE, Ganis G. 2012. Electrophysiological Potentials Reveal Cortical  
1011 Mechanisms for Mental Imagery, Mental Simulation, and Grounded  
1012 (Embodied) Cognition. *Front Psychol*. 3.
- 1013 Schendan HE, Kutas M. 2002. Neurophysiological evidence for two processing  
1014 times for visual object identification. *Neuropsychologia*. 40:931–945.
- 1015 Schendan HE, Kutas M. 2003. Time course of processes and representations  
1016 supporting visual object identification and memory. *Journal of Cognitive  
1017 Neuroscience*. 15:111–135.
- 1018 Schendan HE, Kutas M. 2007. Neurophysiological evidence for the time course  
1019 of activation of global shape, part, and local contour representations  
1020 during visual object categorization and memory. *Journal of Cognitive  
1021 Neuroscience*. 19:734–749.



- 1022 Schendan HE, Lucia LC. 2010. Object-sensitive activity reflects earlier perceptual  
1023 and later cognitive processing of visual objects between 95 and 500ms.  
1024 *Brain Research*. 1329:124–141.
- 1025 Schendan HE, Maher SM. 2009. Object knowledge during entry-level  
1026 categorization is activated and modified by implicit memory after 200 ms.  
1027 *NeuroImage*. 44:1423–1438.
- 1028 Sehatpour P, Molholm S, Javitt DC, Foxe JJ. 2006. Spatiotemporal dynamics of  
1029 human object recognition processing: An integrated high-density electrical  
1030 mapping and functional imaging study of “closure” processes.  
1031 *NeuroImage*. 29:605–618.
- 1032 Smith ME, Loschky LC. 2019. The influence of sequential predictions on scene-  
1033 gist recognition. *Journal of Vision*. 19:14–14.
- 1034 Smout CA, Garrido MI, Mattingley JB. 2020. Global effects of feature-based  
1035 attention depend on surprise. *NeuroImage*. 215:116785.
- 1036 Smout CA, Tang MF, Garrido MI, Mattingley JB. 2019. Attention promotes the  
1037 neural encoding of prediction errors. *PLOS Biology*. 17:e2006812.
- 1038 Spratling MW. 2010. Predictive Coding as a Model of Response Properties in  
1039 Cortical Area V1. *Journal of Neuroscience*. 30:3531–3543.
- 1040 Spratling MW. 2016. Predictive coding as a model of cognition. *Cogn Process*.  
1041 17:279–305.
- 1042 Stefanics G, Kremláček J, Czigler I. 2014. Visual mismatch negativity: a  
1043 predictive coding view. *Front Hum Neurosci*. 8.
- 1044 Summerfield C, Egnér T, Greene M, Koechlin E, Mangels J, Hirsch J. 2006.  
1045 Predictive Codes for Forthcoming Perception in the Frontal Cortex.  
1046 *Science*. 314:1311–1314.
- 1047 Susac A, Heslenfeld DJ, Huonker R, Supek S. 2014. Magnetic Source  
1048 Localization of Early Visual Mismatch Response. *Brain Topogr*. 27:648–  
1049 651.
- 1050 Thierry G, Martin CD, Downing P, Pegna AJ. 2007. Controlling for interstimulus  
1051 perceptual variance abolishes N170 face selectivity. *Nature Neuroscience*.  
1052 10:505–511.
- 1053 Torralbo A, Walther DB, Chai B, Caddigan E, Fei-Fei L, Beck DM. 2013. Good  
1054 Exemplars of Natural Scene Categories Elicit Clearer Patterns than Bad  
1055 Exemplars but Not Greater BOLD Activity. *PLoS ONE*. 8:e58594.
- 1056 Vo ML-H, Wolfe JM. 2013. Differential Electrophysiological Signatures of  
1057 Semantic and Syntactic Scene Processing. *Psychological Science*.  
1058 24:1816–1823.
- 1059 Voss JL, Federmeier KD, Paller KA. 2012. The potato chip really does look like  
1060 Elvis! Neural hallmarks of conceptual processing associated with finding  
1061 novel shapes subjectively meaningful. *Cereb Cortex*. 22:2354–2364.
- 1062 Voss JL, Paller KA. 2007. Neural correlates of conceptual implicit memory and  
1063 their contamination of putative neural correlates of explicit memory. *Learn  
1064 Mem*. 14:259–267.
- 1065 Voss JL, Schendan HE, Paller KA. 2010. Finding meaning in novel geometric  
1066 shapes influences electrophysiological correlates of repetition and



1067 dissociates perceptual and conceptual priming. *NeuroImage*. 49:2879–  
1068 2889.  
1069

1070

1071

1072

1073