

## **Audio-visual Stroop Matching Task with First and Second Language**

### **Colour Words and Colour Associates**

Iva Šaban and James R. Schmidt

LEAD-CNRS UMR 5022, Université de Bourgogne, Dijon

#### **Author notes**

This work was supported by the French “Investissements d’Avenir” program, project ISITE-BFC (contract ANR15-IDEX-0003) to James R. Schmidt. Replication package is available on the Open Science Framework

([https://osf.io/48q2p/?view\\_only=3c4cdec3f832446984291fc5f22f6392](https://osf.io/48q2p/?view_only=3c4cdec3f832446984291fc5f22f6392)).

1 **Abstract**

2 In the audio-visual Stroop matching task, participants compare one Stroop stimulus dimension  
3 (e.g., the colour of a written word) to a second stimulus (e.g., a spoken word) and indicate  
4 whether these two stimuli match or mismatch. Slower responses on certain trials can be due to  
5 conflict which occurs between colour representations (*semantic conflict*) or due to conflict  
6 between responses evoked by task comparisons (*response conflict*). The contribution of these  
7 conflicts has been investigated with colour word distracters. This is the first study which  
8 explores how two types of first and second language words affect audio-visual matching.  
9 Native French speakers performed a bilingual Stroop matching task with intermixed French  
10 (L1) and English (L2) colour words (Experiment 1) and colour associates (Experiment 2)  
11 presented in congruent and incongruent colours simultaneously with spoken French colour  
12 words. Participants were instructed to indicate whether the spoken word “matches” or  
13 “mismatches” the font colour, while ignoring written word meaning. Interestingly, the results  
14 were similar for the critical “mismatch” trials for both French and English words. The  
15 responses were the fastest on trials in which task comparisons activate fewer response  
16 alternatives, supporting the assumption of the response conflict account.

17

18 **Keywords:** audio-visual matching, between-language interference, within-language  
19 interference, semantic conflict, response conflict

20

## 21 **Cognitive control measured by the Stroop task and corresponding conflict effects**

22           People make everyday decisions about allocating cognitive control in order to pursue  
23 their goals (e.g., what to pay attention to, what to stop themselves from doing). For instance,  
24 when confronted with multiple sources of information, our cognitive system adapts our  
25 attentional resources away from distracting (i.e., non-goal relevant) stimuli and/or toward the  
26 goal-relevant stimuli and the action we are supposed to make. The Stroop task is one  
27 particularly useful tool in assessing the ability of the cognitive control system to control  
28 selective attention. In the Stroop task, participants are instructed to name the ink colour of the  
29 written word while ignoring its meaning. The standard finding of slower and less accurate  
30 responding on incongruent (e.g., “red” in green) relative to congruent (e.g., “red” in red) trials  
31 is known as the *congruency* or *Stroop effect* (Stroop, 1935; for a review, see MacLeod, 1991).  
32 Among other things, the Stroop effect indicates that control over selective attention is not  
33 absolute: the distracting word influences colour naming, indicating that it is not ignored  
34 entirely.

35           One other question of interest concerns the source of this congruency effect.  
36 According to *response conflict accounts*, word reading and colour naming compete for a  
37 single response channel (Goldfarb & Henik, 2007; Morton, 1969; Posner & Snyder, 1975).  
38 The word reading response becomes available prior to a colour naming response, because it is  
39 a faster and more automatized process than colour naming (for the automaticity of reading  
40 debate, see Augustinova & Ferrand, 2014; Besner et al., 1997). Thus, word reading disrupts  
41 colour naming but not vice versa. Alternatively, *semantic (or stimulus) conflict accounts*  
42 assume that the conflict occurs in an earlier phase of processing (Luo, 1999; Seymour, 1977;  
43 Simon & Berbaum, 1988). When the ink colour and word meaning are incongruent (e.g.,  
44 “red” in green), two distinct semantic representations (“red” and “green”) are simultaneously  
45 activated. This semantic conflict takes time to resolve, presumably before response selection.

46 Various authors have discussed the relative contribution of semantic and response conflict in  
47 explaining the source of congruency. Nowadays, the current consensus is that both effects  
48 contribute to the standard Stroop effect (Ferrand & Augustinova, 2014). The presence of  
49 semantic and response conflict indicates that the distracting word slipped through the  
50 attentional filter, either at an early semantic processing phase, or later response selection  
51 phase. Most models (Glaser & Glaser, 1989) assume that semantic processing occurs earlier  
52 in the stimulus processing, with the response being selected at a later stage.

### 53 **Stroop matching task**

54 In a Stroop task, a to-be-ignored written word stimulus and the oral response (e.g.,  
55 colour naming and word reading) are compatible, which has been suggested as an inherent  
56 limitation of the Stroop task (Treisman & Fearnley, 1969). That is, a response in the form of a  
57 spoken word is required in both colour naming and word reading tasks. This might produce a  
58 congruency effect only when the irrelevant stimulus attribute (e.g., word) belongs to the same  
59 class as the response. This limitation has inspired a novel variant of the Stroop task, named  
60 the *Stroop matching task*, in which responses are neither words nor colours.

61 In the Stroop matching task, participants are instructed to make matching/mismatching  
62 judgements on two simultaneously presented stimuli (Treisman & Fearnley, 1969). That is,  
63 participants are asked to indicate whether two stimulus dimensions “match” or “mismatch”  
64 (e.g., two colour words or a word and colour). Most importantly, this task permits a test of the  
65 contribution of two contrasting potential sources of conflict: *semantic* and *response conflict*.  
66 For instance, in the *meaning decision task* of Dyer (1973), participants were asked to compare  
67 a colour word to a colour patch and to ignore the print colour of the word.  
68 Matching/mismatching judgements were slower when the colour word was printed in an  
69 incongruent colour. However, responses are slower to “match” trials when the word

70 mismatches the colour (e.g., “red” in blue) than when the word and colour match (e.g., “red”  
71 in red). This is because the incongruent colour activates a semantic representation (i.e., blue)  
72 that competes with the representations activated by the other stimuli (i.e., red). According to  
73 this perspective, then, semantic conflict interferes with the matching/mismatching response  
74 (Dyer, 1973; Flowers, 1975). This finding challenges the assumptions of certain response  
75 conflict accounts because the supposedly slower colour naming response (i.e., “blue”)  
76 influenced responding more than the faster word meaning response (i.e., “red”).

77         Similar findings were observed with the *visual decision task* in which participants are  
78 asked to decide whether two stimuli have the same ink colour (Egeth et al., 1969; Virzi &  
79 Egeth, 1985). For instance, on a trial with the word “red” printed in blue and a blue patch, the  
80 required response is “match”. Interestingly, the conflicting verbal information provided by the  
81 word (i.e., “red”) did not produce interference, seemingly indicating that the word meaning is  
82 not fast enough to compete with the semantic unit (“blue”) accessed by the word’s ink colour  
83 (Egeth et al., 1969; Treisman & Fearnley, 1969). This finding again contradicts the  
84 assumptions of the response conflict account, since word reading, although faster than colour  
85 naming, produced no interference with responding. However, when the colour names were  
86 replaced with the words “SAME” and “DIFF”, interference reappeared. That is, two  
87 simultaneously presented words “DIFF” printed in the same colour (e.g., red) resulted in  
88 interference, because the correct response for the colours (i.e., “matching” or “SAME”)  
89 competes with the response suggested by the distracters (i.e., “mismatching” or “DIFF”). This  
90 indicates that participants had difficulties to ignore the written words and respond to the ink  
91 colour exclusively, as assumed by the response conflict account (Egeth et al., 1969).

92         The *meaning decision* and *visual decision* tasks have been integrated within a single  
93 matching procedure to directly test whether interference is due to semantic or response  
94 conflict. Luo (1999) replicated both the interference in the meaning decision task and the

95 absence of interference in the visual decision task. Luo argued that only the meaning decision  
96 task required participants to access the semantic system. In this task, when a Stroop stimulus  
97 “red” printed in blue is presented with a red patch (i.e., “matching” response is required), the  
98 ink colour and the colour patch activate two competing semantic representations (e.g., “blue”  
99 and “red”). According to Luo (1999), this generates a semantic conflict. In contrast, these  
100 findings are difficult to explain by the response conflict account because it did not matter  
101 whether the response was “matching” or “mismatching” since the response latencies were  
102 faster for related ink colours than for unrelated ink colours.

103           However, Goldfarb and Henik (2006) pointed out that Luo’s (1999) analysis on the  
104 *meaning decision task* only distinguished between a “mismatching” condition in which  
105 coloured patches appeared together with either an incongruent colour word (e.g., “red” in blue  
106 paired with a blue rectangle) or a congruent colour word (e.g., “red” in red paired with a blue  
107 rectangle). Goldfarb and Henik suggested that the congruency of the colour word stimuli  
108 could play a role in producing a conflict. For both “matching” and “mismatching” responses,  
109 Stroop stimuli could be either congruent or incongruent. Thus, in addition to the four  
110 conditions contrasted by Luo (1999), Goldfarb and Henik (2006) introduced a condition in  
111 which both dimensions of the incongruent Stroop stimuli mismatch with the colour of the  
112 patch (e.g., “red” in blue with a green patch). They observed that “matching” responses were  
113 faster when Stroop stimuli were congruent (e.g., “red” in red with a red patch) than when they  
114 were incongruent (e.g., “red” in green with a red patch). The “mismatching” responses were  
115 the slowest when the word and ink colour were congruent (e.g., “red” in red with a green  
116 patch). Delays were similar when the ink colour and patch colour matched (e.g., “red” in  
117 green with a green patch) and when they mismatched (e.g., “red” in blue with a green patch).  
118 To sum up, response latencies to incongruent trials were slower during “matching” responses  
119 and faster during “mismatching” responses. According to Goldfarb and Henik, participants

120 erroneously made an irrelevant match between the word and its ink colour. That is, seeing  
121 congruent and incongruent Stroop stimuli leads to a covert “matching” and “mismatching”  
122 response, respectively, which can either facilitate or interfere with the actual response  
123 required. Thus, they suggested that the results are clearly in line with the response conflict  
124 account.

125           In a related matching task variant, Bornstein (2015) asked participants to make an  
126 audio-visual matching judgement based on the task-relevant auditory (i.e., spoken colour  
127 word) and visual stimuli (i.e., ink colour of a written word). On each trial, participants were  
128 instructed to indicate whether the colour of a written word (while ignoring its meaning)  
129 corresponds to a simultaneously presented spoken word. Bornstein (2015) compared the  
130 interference produced by congruent and incongruent written stimuli on matching spoken word  
131 and font colour. Bornstein observed that incongruent distracters (e.g., “red” in blue while  
132 hearing “blue”) interfered more than congruent distracters (e.g., “blue” in blue while hearing  
133 “blue”) with “matching” responses, similarly to Goldfarb and Henik (2006). Furthermore,  
134 written words that were congruent with either task-relevant dimension (i.e., ink colour or  
135 spoken word) interfered with “mismatching” responses relative to trials in which the word  
136 mismatched both (e.g., “green” in red while hearing “blue”).

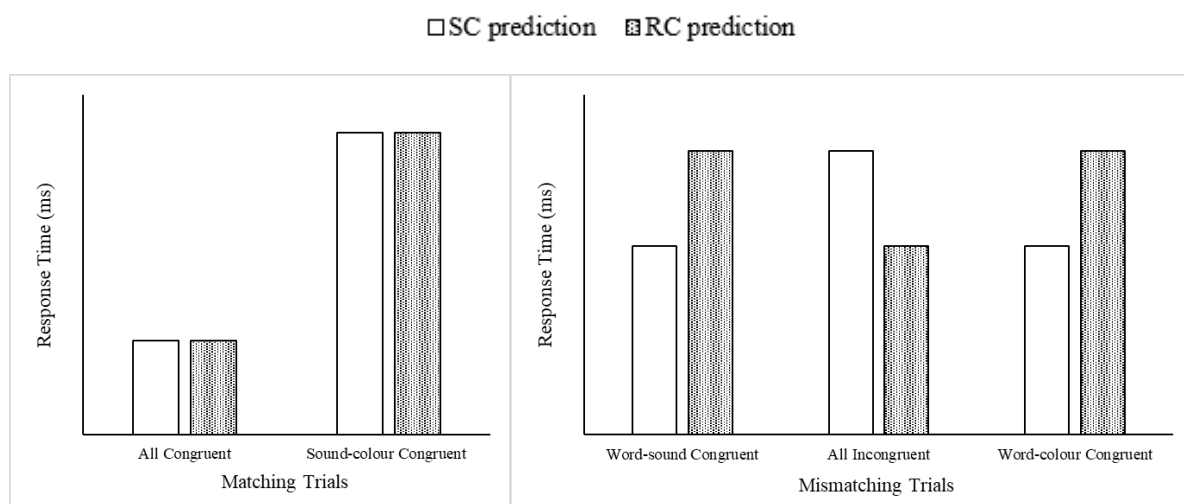
137           Both the semantic and response conflict accounts assume the same outcome for  
138 “matching” responses with faster responses on congruent (i.e., *All congruent*) relative to  
139 incongruent colour words (i.e., *Sound-colour congruent*). According to the semantic conflict  
140 account, this is due to the fact that for congruent colour words, all three task dimensions refer  
141 to the same colour (i.e., blue). The response conflict account explains this difference in  
142 response speed by the three stimulus comparisons, which all suggest the same response  
143 alternative (i.e., “match”). Critically, the assumptions of these two accounts differ for  
144 “mismatching” trials. According to the semantic conflict account, *All incongruent* trials, in

145 which a written colour word is incongruent (e.g., “green” in red, hear “blue”) with the  
 146 remaining two colour dimensions, should produce the largest interference. Three different  
 147 semantic representations (i.e., blue, red, and green) are simultaneously activated, thus slowing  
 148 down responding. In contrast, the response conflict account suggests that incongruent colour  
 149 word distracters should facilitate responding when both dimensions (e.g., green and red) are  
 150 incompatible with a spoken word (e.g., blue). This is because all three comparisons (i.e.,  
 151 written vs. spoken word, written word vs. colour, and spoken word vs. colour) provide  
 152 evidence toward the same response alternative (i.e., “mismatching”), resulting in faster  
 153 response latencies (Bornstein, 2015; Caldas et al., 2012; Goldfarb & Henik, 2006). The shared  
 154 prediction of semantic and response conflict accounts for “matching” trials and contrasting  
 155 predictions for “mismatching” trials are visualised in Figure 1.

### 156 **Figure 1**

157 *Prediction of semantic and response conflict accounts for “matching” and “mismatching”*  
 158 *trials*

159



160

### 161 **Colour Associates**

162

All previously described Stroop matching task studies made use of colour words.

163

However, similar studies have not been conducted with another common word type with a



164 strong colour dimension, namely, colour associates, which could help further evaluate conflict  
165 effects in the Stroop matching task. Colour associates are words that are closely related to  
166 colour words (e.g., “sky” with blue) and their semantic representations (Tanaka & Presnell,  
167 1999). Colour associates do produce interference with colour naming in the Stroop task.  
168 Similar to colour words, colour associates can be congruent (e.g., “sky” in blue) or  
169 incongruent (e.g., “sky” in red) with the ink colour. When contrasting the response latencies  
170 of these two types of trials, a congruency occurs, with slower and less accurate responses on  
171 incongruent relative to congruent colour associates (Glaser & Glaser, 1989; Klein, 1964;  
172 Risko et al., 2006; Schmidt & Cheesman, 2005).

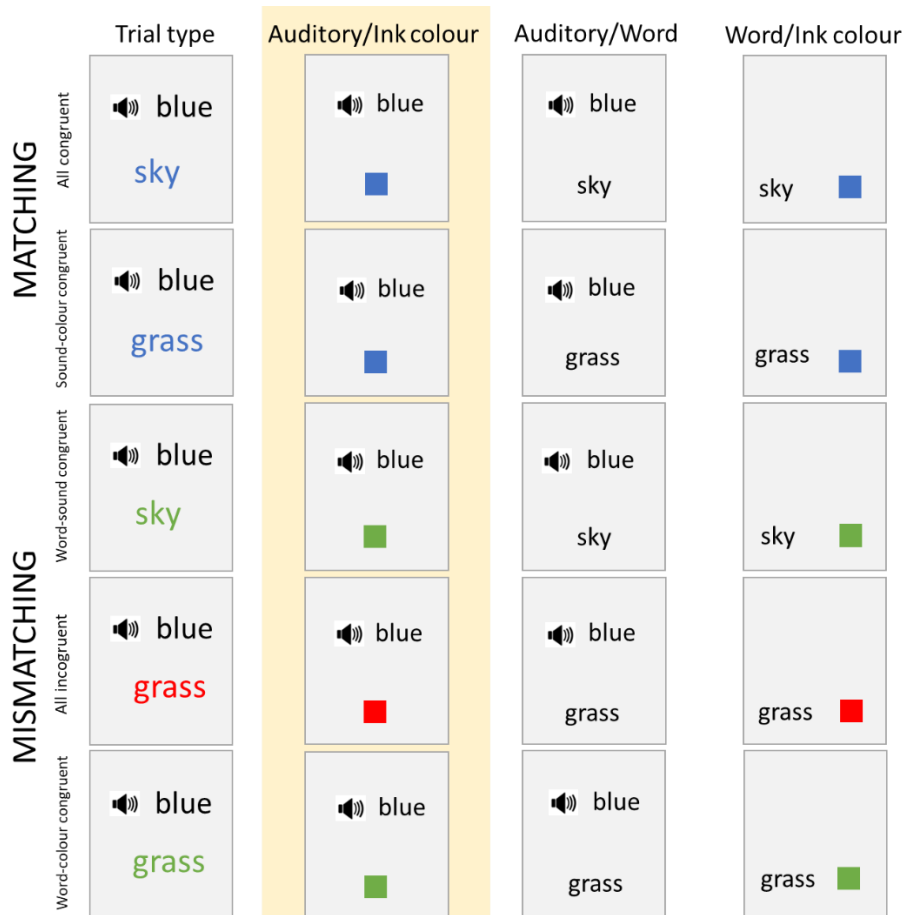
173         This difference in performance might be due to early semantic processes (Glaser &  
174 Glaser, 1989). When a colour word distracter is printed in an incongruent colour (e.g., “sky”  
175 in red), two competing colour representations (i.e., red and blue) are simultaneously activated,  
176 thus producing semantic conflict. According to this perspective, colour associate congruency  
177 effects arise from early, semantic processes. Another account suggests that colour associates  
178 might directly produce the colour response linked to the colour associate. That is, when the  
179 word “sky” is printed in red, both the responses linked to the colour blue (i.e., the colour  
180 associated with “sky”) and the response linked to the colour red (i.e., which is associated to  
181 the ink colour) will be activated. Thus, according to this perspective, incongruent colour  
182 associates produce response competition, resulting in response conflict exclusively, rather than  
183 semantic conflict (Klein, 1964). Third, Sharma and McKenna (1998) suggested that  
184 interference should occur only when vocal responses are required and should be eliminated  
185 with manual responses, though subsequent research clearly indicates the presence of conflict  
186 effects in keypress tasks (e.g., Schmidt & Cheesman, 2005).

187         One reason why colour associates might be especially interesting in the context of the  
188 matching task relates to a peculiarity of the matching task. For “matching” trials, both the

189 semantic and response conflict accounts make identical predictions. For “mismatching” trials,  
190 the two accounts make exactly opposite predictions. Specifically, the semantic conflict  
191 account suggests that *All incongruent* trials should be slower than the two other types of  
192 “mismatching” trial types, whereas the response conflict account suggests that *All*  
193 *incongruent* trials should be faster than the two other types of “mismatching” trial types.  
194 Therefore, if *both* semantic and response conflict occur, the larger of the two effects will  
195 “mask” the other. In particular, evidence of a response conflict effect *could* indicate that only  
196 response conflict occurs in the matching task but could also indicate that response conflict is  
197 merely larger than semantic conflict. Thus, if the response conflict effect can be eliminated,  
198 then we might expect that the “true” effect of semantic conflict would be revealed. Although  
199 some competing accounts of colour associates conflict effects exist (as discussed above), we  
200 hypothesized that colour associates would produce only semantic conflict. Some evidence  
201 suggests this to be the case in standard Stroop studies (e.g., Schmidt & Cheesman, 2005). All  
202 task comparisons (one relevant and two irrelevant) for each colour associate trials are  
203 visualised in Figure 2.

204 **Figure 2**

205 *Types of trials and example stimuli with relevant (highlighted column) and irrelevant task*  
 206 *comparisons*



207

208 **Bilingualism**

209 The Stroop effect has been frequently investigated in bilingual people (Altarriba &  
 210 Mathis, 1997; Dyer, 1971; MacLeod, 1991; Mägiste, 1982; Preston & Lambert, 1969;  
 211 Tzelgov et al., 1990). These previous studies showed that congruency can be observed with  
 212 both first language (L1) and second language (L2) words. However, the interference is  
 213 generally larger for L1 words than for L2 words. This could be explained by the nature of L2  
 214 connections. For instance, there has been debate about whether L2 words 1) have strong direct  
 215 connections to semantic representations but weak connections to the L1 lexicon, 2) are

216 strongly connected to the L1 lexicon but not semantics, or 3) have both semantic and lexical  
217 connections (Altarriba & Mathis, 1997; Kroll & Stewart, 1994; Schmidt et al., 2018). Thus, it  
218 is unclear whether L2 words would lead to semantic conflict, response conflict, or a  
219 combination of both. Specifically, L2 words would not be expected to generate semantic  
220 conflict if they have no (or very weak) connections to semantics. If the exact reverse is true  
221 and L2 words function as semantic associates to their L1 translations, then only semantic  
222 conflict might be expected, as discussed in the previous section on colour associates.

223         Another important question in the bilingual Stroop literature concerns the modulation  
224 of Stroop interference by stimulus and response language (i.e., the language of a distracter  
225 and the language of a response, respectively). First, the distracter language can match the  
226 response language. For instance, colour naming of the distracter “red” printed in green  
227 produces within-language (or intralingual) interference when English is a response language  
228 (i.e., a correct response is to say “green”). Second, the distracter language can mismatch the  
229 response language. That is, colour naming of the distracter “rouge” (red in French) printed in  
230 blue produces between-language (or interlingual) interference when English is a response  
231 language (i.e., a correct response is to say “green”).

232         The magnitude of within- and between-language interference has been compared  
233 repeatedly. A standard finding is a larger within-language than between-language interference  
234 effect (Dyer, 1971; Hamers & Lambert, 1972; Kiyak, 1982; MacLeod, 1991; Preston &  
235 Lambert, 1969). For instance, MacLeod (1991) reported that the between-language  
236 interference represents about 75% of within-language interference. However, these findings  
237 mostly originated from the standard visual (MacLeod, 1991) and auditory (Hamers &  
238 Lambert, 1972) Stroop task but have never been confirmed with the Stroop matching task. In  
239 a bilingual Stroop matching task, it might be assumed that distracters that match in language  
240 with a spoken word will produce larger interference relative to those that mismatch. To test

241 this in the present series of studies, we used distracting words from both the first language  
242 (i.e., French) and a second language (i.e., English). However, spoken words were always  
243 French. French distracters are therefore expected to produce larger interference (i.e., within-  
244 language interference) relative to English distracters (e.g., between-language interference).

### 245 **Present Study**

246 In the present series of experiments a bilingual audio-visual Stroop matching task was  
247 designed to further explore the 1) magnitude of interference produced by first (L1) and second  
248 (L2) language colour words and colour associates, and 2) the relative contributions of  
249 semantic and response conflict. In addition to first language colour words, frequently used as  
250 distracters in the literature, we introduced second-language colour words (Experiment 1).  
251 That is, intermixed French (L1) and English (L2) colour words served as distracters, while  
252 participants had to match its ink colour with a spoken French colour word. Thus, this  
253 manipulation allows us to test the consensus of larger within- than between-language  
254 interference. If this is the case, a larger interference effect is expected to occur with French  
255 (L1) than with English (L2) colour word distracters. The design of this study can be found in  
256 the *Audiovisual Stimulus Combination* section. Experiment 2 aims to further expand the  
257 findings by using colour associates instead of colour words. That is, both French and English  
258 colour associates were used as distracters, with participants matching their ink colour with a  
259 spoken French colour word. Note that, in contrast to Experiment 1, a spoken word (e.g.,  
260 “vert”, French for green) does not correspond to a written word (e.g., “herbe”, French for  
261 grass). This manipulation should (according to some views) eliminate response conflict since  
262 “herbe” might be unable to retrieve the response linked to green. Furthermore, this could  
263 reveal the role of the semantic conflict, which is possibly masked by a (larger) response  
264 conflict effect. Apart from that, the question of larger within- relative to between-language

265 interference remains open. That is, French colour-associates are expected to produce more  
266 interference than their English counterparts.

267 The present series of studies also aims to investigate the source of this interference. As  
268 already discussed, the interference could be due to the conflict between semantic  
269 representations (i.e., semantic conflict) or due to the conflict between response alternatives  
270 (i.e., response conflict). Based on the findings of Luo (1999) and Goldfarb & Henik (2003),  
271 these two opposing accounts predict similar outcomes for “matching” responses. That is,  
272 when a correct response is “match”, *Sound-colour congruent* trials will produce slower  
273 responses than *All congruent* distracters. However, semantic- and response-conflict accounts  
274 make different assumptions for “mismatching” responses, based on the congruency between  
275 task dimensions. According to the semantic conflict account, a written distracter should  
276 produce the largest interference by being incongruent with both task dimensions (e.g., on *All*  
277 *incongruent* trials) than by being incongruent with only one of them (e.g., on *Word-sound*  
278 *congruent* and *Word-colour congruent* trials). This is because, on *All incongruent* trials, the  
279 distracting written word is incongruent with both target dimensions, thus producing a delay in  
280 responding. In contrast, the response-conflict account assumes that the smallest interference  
281 will be observed with *All incongruent* trials, when all task comparisons suggest the same,  
282 “mismatching” response. That is, interference will be mostly observed on *Word-sound*  
283 *congruent* and *Word-colour congruent* trials, where one of the irrelevant task comparisons  
284 suggest the same response alternative as the relevant comparison (i.e., “mismatch”), but the  
285 third comparisons suggest the other (incorrect) response alternative (i.e., “match”).

## 286 Experiment 1

287 Experiment 1 contrasts the response latencies on congruent and incongruent French  
288 (L1) and English (L2) colour word distracters, each accompanied by a French spoken word.

289 Participants were instructed to respond according to whether the ink colour and spoken word  
290 match or mismatch by pressing the corresponding key. The combinations of visual and  
291 auditory stimuli produced five trial types: two “matching” and three “mismatching”,  
292 discussed in detail in the *Audiovisual Stimulus Combination* section. The aim of Experiment 1  
293 was to 1) compare the magnitude of interference produced by first and second language  
294 colour words in the audio-visual Stroop matching task, and 2) investigate whether this  
295 interference is due to semantic or response conflict.

## 296 **Method**

### 297 *Participants*

298 A total of 34 (31 women) [removed for review] undergraduates ( $M_{age} = 19$ ;  $SD = .78$ )  
299 voluntarily participated in the experiment in exchange for course credit. An a priori power  
300 analysis was conducted using G\*Power 3 (Faul et al., 2007) for sample size estimation, based  
301 on data from Goldfarb and Henik (2006),  $N=12$ , which compared response times on matching  
302 and mismatching trials separately. The effect size in Goldfarb and Henik’s (2006) study was  
303  $\eta_p^2 = .57$ , considered to be large. With a significance criterion of  $\alpha = .05$  and power .95, the  
304 minimum sample size needed with this effect size is  $N = 22$  for repeated measures ANOVA.  
305 Preferring more power than minimally necessary, we decided to collect data for at least 30  
306 participants, stopping after a testing week when this number was exceeded (resulting in the  
307 obtained sample size of  $N = 34$ ).

308 All participants had normal or corrected-to-normal visual acuity, normal colour vision  
309 and normal auditory acuity, as assessed via screening questions. Participants gave written  
310 informed consent before the study. All the procedures were conducted in accordance with the  
311 Declaration of Helsinki, although nonbiomedical research in [removed for review] does not  
312 require ethics approval. All participants were native French speakers. A language

313 questionnaire (to be discussed shortly) was used to assess and confirm that participants fit  
314 with these criteria. Average language background scores (mean age and standard errors) are  
315 presented in Table 1 (see Results section).

### 316 *Apparatus*

317         The experiment was conducted in a sound-attenuated room in the laboratory. Stimulus  
318 presentation and response timing were controlled and recorded by Psytoolkit (Stoet, 2010,  
319 2017). The study was conducted using a PC laptop with an AZERTY keyboard and a 15''  
320 monitor. Participants responded with the “D” key when the audio and the ink colour of the  
321 written distracted mismatched (e.g., hear “green” and see “brown” in brown). Participants  
322 responded with the “K” key when the audio and the ink colour matched (e.g., hear “green”  
323 and see “brown” in green). Prior to the Stroop matching portion of the experiment,  
324 participants filled out a short language demographic questionnaire. This questionnaire asked  
325 for gender, age, native language, years of English training in school, a self-rating of English  
326 knowledge ranging from 0 (= almost none) to 5 (= perfect). A subset of questions from the  
327 French version of the Language Experience and Proficiency Questionnaire (LEAP-Q; Marian  
328 et al., 2007) was inserted. In particular, the questions asking participants to list the languages  
329 in order of dominance and acquisition were retained. They were also asked to indicate the  
330 percentage with which they used French and English in the recent period. Also retained from  
331 the LEAP-Q were two boxes, one for French and one for English, asking for the age the  
332 participants began acquiring the language, became fluent in the language, began learning to  
333 read in the language, and became fluent in reading the language. The purpose of this  
334 questionnaire was to assure that participants had the correct language dominance. Finally, as  
335 an addition to these two questionnaires, participants were asked to give the French  
336 translations of the four English words used in the experiment (i.e., “green”, “brown”, “pink”  
337 and “white”).



338 This was followed by the LexTale English vocabulary test (Lemhöfer & Broersma,  
339 2012) with instructions translated into French. This test contains 63 English-looking words (3  
340 practice trials and 60 test trials). 2/3 of the test trials are actual English words (e.g., “moonlit”,  
341 “fluid”), whereas the remaining 1/3 are not (e.g., “plaudate”, “rebondicate”). Participants were  
342 instructed to select the words that they are certain are actual English words. Correct “hits”  
343 were rewarded with one point, and incorrect “false alarms” were penalized by two points.

#### 344 *Materials and design*

345 During the experimental part of the experiment, participants were presented with a set  
346 of French-English translation equivalents (i.e., “green/vert”, “brown/marron”, “rose/pink”,  
347 and “white/blanc”), typed in lowercase Courier New Bold font (size 72). The corresponding  
348 print colours and their RGB codes were green (0, 128, 0), brown (165, 42, 42), hot pink (255,  
349 105, 180), and white (255, 255, 255). These four words were non-cognates, that is, do not  
350 share phonological or orthographic features across languages, unlike several other colour  
351 word pairs (e.g., “blue/bleu” or “red/rouge”). Auditory stimuli consisted of the colour words  
352 (/vert/, /marron/, /rose/, /blanc/, French for green, brown, pink and white, respectively),  
353 spoken by a female speaker.

354 The manipulation allowed for 2 within-subject factors: Trial Type (“matching”  
355 condition that contained *All congruent* and *Sound-colour congruent* trials vs. “mismatching”  
356 condition that contained *Word-sound congruent*, *All incongruent*, and *Word-colour congruent*  
357 trials) and Language (French vs. English). In each experimental block, there were 25%  
358 matching (6.25% *All congruent*, 18.75% *Sound-colour congruent*) and 75% mismatching  
359 trials (18.5% *Word-sound* and *Word-colour congruent* trials, 37.5% *All incongruent*). This  
360 was because each combination of colour word distracter, print colour, and sound were  
361 presented equally often to avoid contingency biases (i.e., learning of regularities between

362 stimuli; Schmidt et al., 2007; see also, Lorentz et al., 2016).<sup>1</sup> This does mean that  
363 mismatching responses were more frequent than matching responses. However, it is important  
364 to note that all of the key comparisons are within response type. That is, we conducted one  
365 analysis for matching responses and another analysis for mismatching responses, as  
366 previously suggested (Goldfarb & Henik, 2006). This way, even if participants had a learned  
367 strategic tendency to prepare the “mismatching” response, this bias cannot impact “matching”  
368 responses. No systematic biases were produced in our statistical tests, as two trial types were  
369 analysed separately (i.e., none of our comparisons involve comparing a trial with a  
370 “matching” response to a “mismatching” response. In total, there were 3 larger experimental  
371 blocks of 128 trials each (in total 384 trials), presented randomly without replacement. This  
372 main phase of the experiment was preceded by a practice block. The practice block consisted  
373 of 32 trials, with the colour words replaced with the stimulus “xxxx”.

#### 374 *Audiovisual Stimulus Combination*

375 A total of 128 audiovisual stimulus combinations were created from the eight visual  
376 stimuli (“vert”, “marron”, “rose”, “blanc”, “green”, “brown”, “pink”, “white”), four font  
377 colours (green, brown, pink, and white) and four auditory stimuli (“vert”, “marron”, “rose”,  
378 “blanc”). These combinations were grouped into 5 conditions, varying by the congruence or  
379 incongruence between spoken word meaning, font colour, and written word meaning. In two  
380 conditions, the font colour and spoken colour word (task-relevant comparison) were  
381 congruent and thus required a “matching” response. These conditions were: 1) *All congruent*,  
382 and 2) *Sound-colour congruent*. In the other three conditions, the font colour and spoken  
383 colour word were incongruent and thus required a “mismatching” response. These conditions

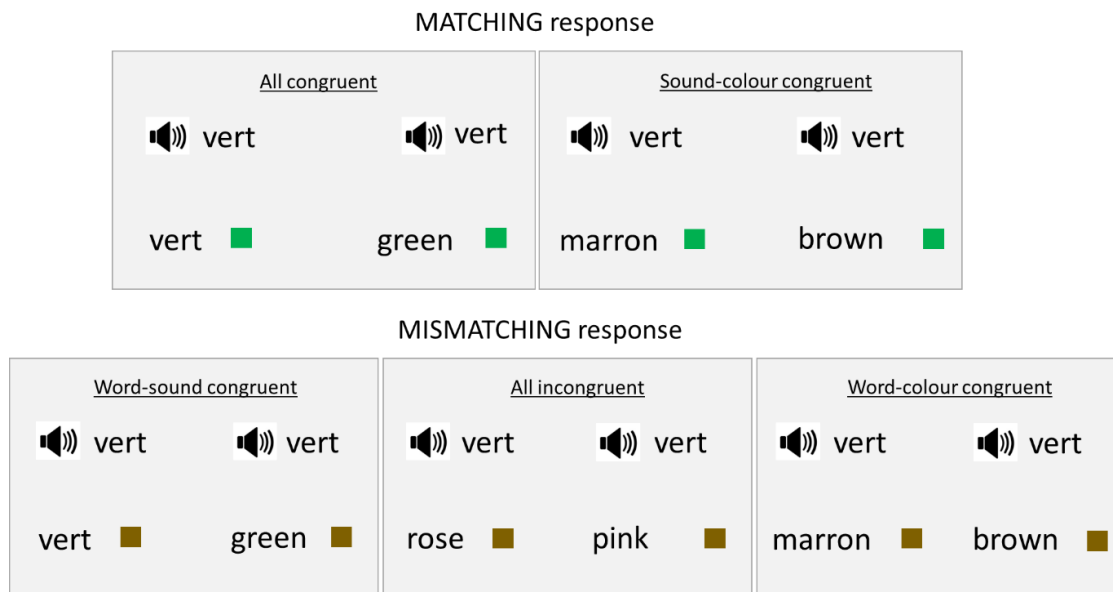
---

<sup>1</sup> In a standard Stroop task, the proportion of congruent trials is often increased, sometimes merely to have the same number of congruent and incongruent trials (e.g., 1:1 congruent:incongruent in a four-choice task) or to increase control demands (e.g., 3:1 congruent:incongruent in Blumenfeld & Marian, 2014). However, this is suboptimal as regularities are introduced between distracting and target stimuli, meaning that congruency effects are confounded by contingency learning effects.

384 were: 3) *All incongruent*, 4) *Word-sound congruent*, and 5) *Word-colour congruent*. All of  
 385 these five conditions applied for both distracter languages. These conditions are presented in  
 386 Figure 3.

387 **Figure 3**

388 *All trial types across two distracter languages (French and English)*



389

390 *Note.* All trial types have two equivalents; one with a French distracter (on the left) and one  
 391 with an English distracter (on the right). Colour patches represent the ink colour in each trial.

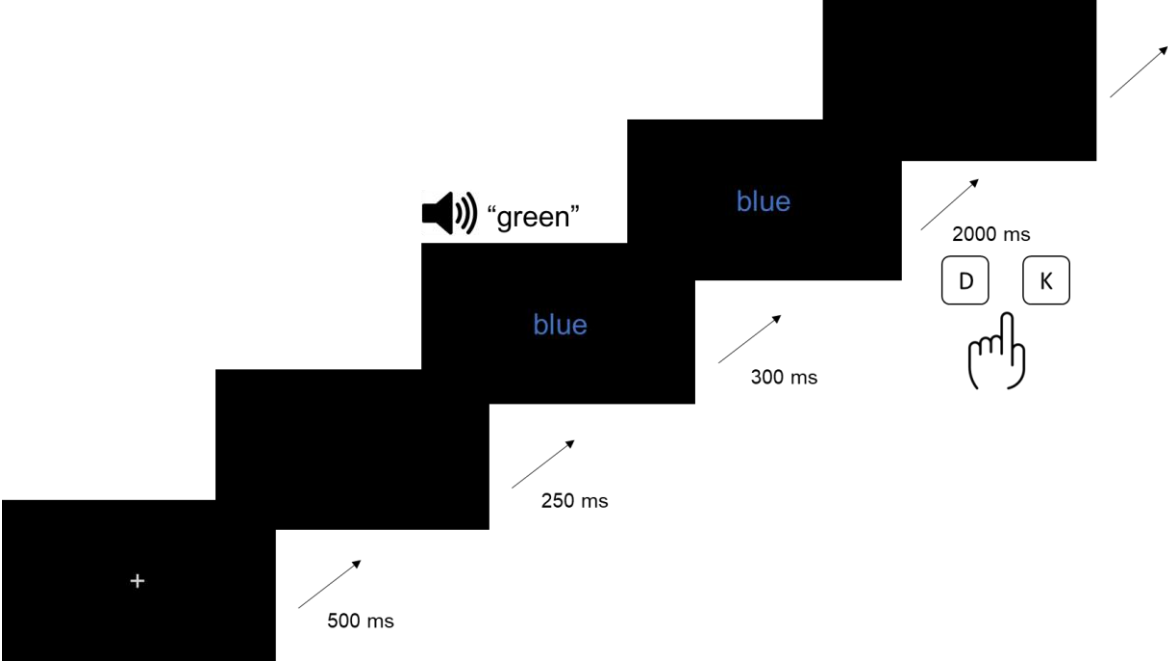
392 ***Procedure***

393 After completing the survey questions, the main experiment began. Stimuli were  
 394 presented on a black (0, 0, 0) screen. On each trial, participants were first presented with a  
 395 fixation “+” in grey (128, 128, 128) for 500 ms. This was followed by blank screen presented  
 396 for 250 ms. Then the coloured distracter appeared on the screen until a response was  
 397 registered or 2000 ms elapsed. The coloured distracter was presented simultaneously with the  
 398 auditory stimulus. Responses could be provided only after 300 ms from the stimulus onset.  
 399 This is due to the programming of the experiment. On each trial, an initial event plays the  
 400 audio and presents the visual stimuli, which is then followed by a second event with only the  
 401 stimulus and where responses are recorded. This was also done because the task required a

402 comparison of the auditory stimulus with the print colour. Thus, a response before the  
 403 auditory stimulus has been played is inevitably an anticipatory response that would be best  
 404 excluded anyway. The next trial began after a 750-ms blank screen. The timeline of each trial  
 405 is visualized in Figure 4. If the participant made an error or failed to respond in time, then the  
 406 message “Erreur” (“Error”) or “Trop lent” (“Too slow”), respectively, appeared in red (255, 0,  
 407 0) for 1000 ms before the next trial. In both experiments, participants were explicitly  
 408 instructed to respond as quickly and as accurately as possible and avoid reading a distracter  
 409 since it represents a task-irrelevant dimension. The “matching” key had to be pressed for trials  
 410 in which the spoken colour word and the font colour matched, and the “mismatching” key for  
 411 trials in which the spoken colour word and the font colour mismatched.

412 **Figure 4**

413 *Timeline of an experimental trial*



414

415 **Results**

416 We used French and English words in this experiment to compare a highly-fluent L1  
 417 with a low-fluency L2. In [removed for review], French is normally the native language and  
 418 English is typically learned later in life and not to a very high level of mastery. To assure that  
 419 this was actually the case for our sample, we first analysed average language metric scores<sup>2</sup>,  
 420 which are presented in Table 1. All participants seemed to sufficiently fit our language  
 421 criteria, as they were native French speakers who acquired the language early in life.  
 422 Importantly, French was ranked as the first language in terms of dominance and order of  
 423 acquisition by all participants. The percentage of French use revealed that participants had  
 424 been using French almost exclusively in their everyday lives. In contrast, English was learned  
 425 much later as a foreign language in primary schools. Participants were only moderately  
 426 proficient in English, as shown by LexTale score and their self-rated English knowledge  
 427 level. Although they studied English for a considerable amount of time (almost 9 years) and  
 428 declared being able to speak and read English fluently (approximately at the age of 15), their  
 429 objective proficiency level is rather low.

430 **Table 1**431 *Mean French and English language scores and standard errors (in brackets)*

	<i>M</i>	<i>SE</i>
LexTale		
Years English	8.94 years	0.332
English level	3 (1-5)	0.158
Score	65.82 (0-100)	1.312

<sup>2</sup> The vast majority (33/34; 1 empty) of participants indicated French as their first language in order of dominance and in order of acquisition. One participant ranked Turkish as the first language in both dominance and acquisition, but further inspection of provided responses revealed that this participant had started acquiring French early enough and thus was therefore not excluded from the sample. As a second language in order of dominance and acquisition, participants rated English, followed by Spanish, Arabic, Creole, and Portuguese. The most frequently indicated third language in both dominance and acquisition were Spanish, German, English, Italian, Arabic, and Portuguese. All the participants correctly translated English words “green”, “brown”, “pink” and “white”.

LEAP-Q		
Dominance French	1	0
Dominance English	2.26	0.056
Order French	1	0
Order English	2.19	0.052
French Use (%)	4.97 (1-5)	0.029
English Use (%)	1.73 (1-5)	0.160
French		
Acquisition	1.10 years	0.183
Fluent	3.03 years	0.228
Reading	5.54 years	0.147
Fluent Read	6.79 years	0.198
English		
Acquisition	9.85 years	0.351
Fluent	15.41 years	0.344
Reading	12.42 years	0.386
Fluent Read	14.75 years	0.404

432

433 ***Data Analysis***

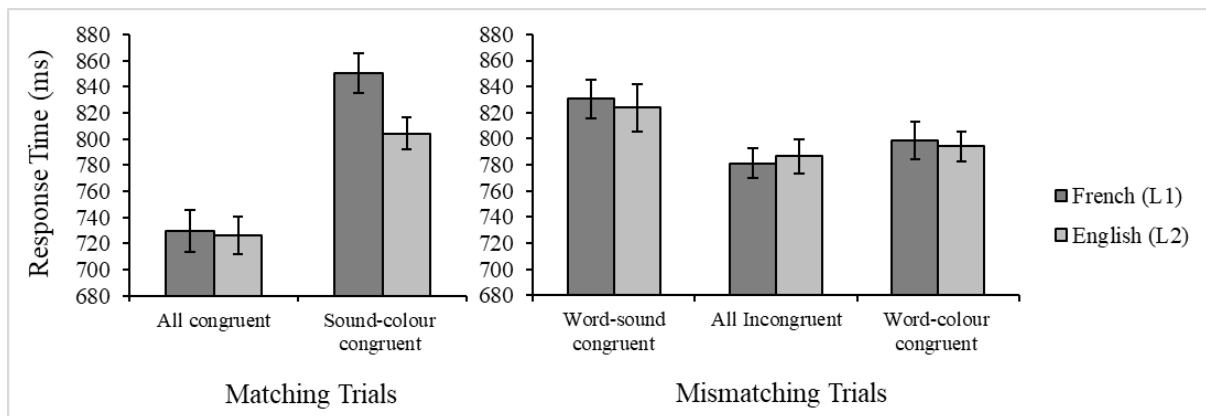
434 The mean correct response times (i.e., made during the 2000 ms response window)  
435 and mean percentage error were analysed. Response times were not trimmed (pre-planned  
436 analyses). However, we note that the direction and significance of all effects did not change in  
437 subsequent analyses with an Interquartile range (IQR) trim method, unless otherwise noted.  
438 No participants were excluded from the sample, as their individual accuracy rate was 86.35%  
439 or above. The congruency variable had different levels for “matching” and “mismatching”  
440 responses, and matching and mismatching trial types were analysed separately. One shared  
441 factor was a Distracter Language, with two levels: French (L1) and English (L2). Because the  
442 congruency variable had different levels for the “matching” and “mismatching” responses and  
443 because there are no relevant comparisons to make between the matching and mismatching  
444 trial types, two separate repeated measure analyses of variance with two within-subject factors  
445 were conducted. In the “matching” condition, 2 levels were analysed (*All congruent* and  
446 *Sound-colour congruent*), while in the “mismatching” condition, 3 levels were analysed  
447 (*Word-sound congruent*, *All incongruent* and *Word-colour congruent*).

448 **Response time (RT)**

449 Response times were recorded in milliseconds as the time elapsed from stimulus onset  
 450 to key press. A total of 5.98% trials were excluded from the analyses (5.77% incorrect and  
 451 .21% time out responses). Only RTs for correct responses in “matching” and “mismatching”  
 452 conditions were analysed and illustrated in Figure 5.

453 **Figure 5**

454 *Mean response times with standard errors for “matching” and “mismatching” trials*



455

456 **Matching trials**

457 There was a main effect of Trial Type;  $F(1,33) = 209.609$ ,  $MSE = 1606.534$ ,  $\eta_p^2 =$   
 458  $.864$ ,  $BF_{10} > 1000$ ,  $p < .001$ . Responses on *Sound-colour congruent* trials ( $M = 827$ ,  $SE =$   
 459  $13.30$ ) were slower than responses on *All congruent* trials ( $M = 728$ ,  $SE = 13.93$ ). The  
 460 significant main effect of Language was observed,  $F(1,33) = 11.638$ ,  $MSE = 1797.765$ ,  $\eta_p^2 =$   
 461  $.260$ ,  $BF_{10} = 1.124$ ,  $p = .001$ , with slower responses in French condition ( $M = 790$ ,  $SE =$   
 462  $14.71$ ) relative to English condition ( $M = 765$ ,  $SE = 12.53$ ). The interaction between Trial  
 463 Type and Language was also significant,  $F(1,33) = 9.272$ ,  $MSE = 1649.944$ ,  $\eta_p^2 = .219$ ,  $BF_{10}$   
 464  $= 11.021$ ,  $p < .01$ . There was no difference in response speed between French ( $M = 729$ ,  $SE =$   
 465  $16.06$ ) and English ( $M = 726$ ,  $SE = 14.45$ ) *All congruent* trials,  $t(33) = .286$ ,  $M_{diff} = 3$ ,  $BF_{10} =$   
 466  $.191$ ,  $BF_{01} = 5.236$ ,  $p = .776$ . However, responses were significantly slower on French ( $M =$

467 850,  $SE = 15.13$ ) *Sound-colour congruent* trials relative to English *Sound-colour congruent*  
 468 ( $M = 804$ ,  $SE = 12.14$ ) trials;  $t(33) = 6.847$ ,  $M_{diff} = 46$ ,  $BF_{10} > 1000$ ,  $p < .001$ .

#### 469 ***Mismatching trials***

470 The main effect of Trial Type was observed,  $F(2,66) = 36.205$ ,  $MSE = 926.505$ ,  $\eta_p^2 =$   
 471  $.523$ ,  $BF_{10} > 1000$ ,  $p < .001$ . Responses on *Word-sound congruent* ( $M = 827$ ,  $SE = 15.79$ )  
 472 trials were significantly slower than responses on *All incongruent* ( $M = 784$ ,  $SE = 12.01$ )  
 473 trials,  $t(33) = 7.156$ ,  $M_{diff} = 43$ ,  $BF_{10} > 1000$ ,  $p < .001$  and *Word-colour congruent* ( $M = 796$ ,  
 474  $SE = 12.44$ ) trials,  $t(33) = 5.085$ ,  $M_{diff} = 31$ ,  $BF_{10} > 1000$ ,  $p < .001$ . Responses on *Word-colour*  
 475 *congruent* trials were slower relative to responses on *All incongruent* trials,  $t(33) = 4.167$ ,  
 476  $M_{diff} = 12$ ,  $BF_{10} = 129.88$ ,  $p < .001$ . There was no main effect of Language<sup>3</sup>,  $F(1,33) = .278$ ,  
 477  $MSE = 727.161$ ,  $\eta_p^2 = .008$ ,  $BF_{10} = .161$ ,  $BF_{01} = 6.211$ ,  $p = .602$ , indicating that there is no  
 478 difference in response latencies between French and English trials. The interaction between  
 479 Trial Type and Language was also not significant,  $F(2,66) = .664$ ,  $MSE = 1031.101$ ,  $\eta_p^2 = .02$ ,  
 480  $BF_{10} = .179$ ,  $BF_{01} = 5.586$ ,  $p = .518$ .

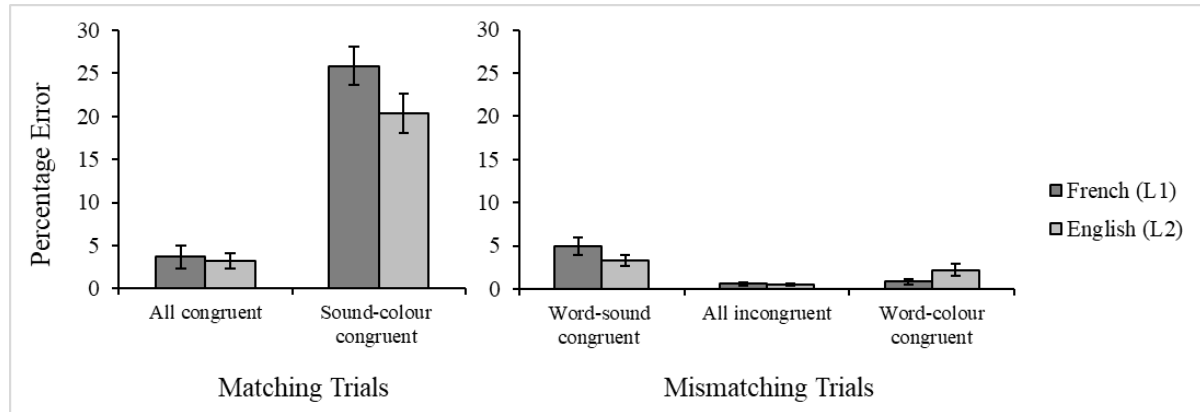
#### 481 ***Percentage error***

482 The mean percentage error data for all trial types and languages are presented in Figure 6.

---

<sup>3</sup> After trimming 512 outliers using the IQR method, the main effect of Language reached significance;  $F(1,33) = 6.243$ ,  $MSE = 581.77$ ,  $\eta_p^2 = .16$ ,  $p = .02$  for response times in Mismatching trials. Trials with French distracters ( $M = 781$ ;  $SE = 10.54$ ) were responded to slower than trials with English distracters ( $M = 773$ ,  $SE = 10.17$ ).



483 **Figure 6**484 *Mean percentage error with standard error for “matching” and “mismatching” trials*

485

486 **Matching trials**

487 There was a main effect of Trial Type,  $F(1,33) = 113.835$ ,  $MSE = 115.229$ ,  $\eta_p^2 = .775$ ,  
 488  $BF_{10} > 1000$ ,  $p < .001$ , indicating that participants made significantly more errors on *Sound-*  
 489 *colour congruent* ( $M = 23.07$ ,  $SE = 2.08$ ) than on *All congruent* trials ( $M = 3.43$ ,  $SE = .89$ ).

490 The main effect of Language was observed,  $F(1,33) = 8.034$ ,  $MSE = 37.752$ ,  $\eta_p^2 = .196$ ,  $BF_{10}$   
 491  $= .391$ ,  $BF_{01} = 2.557$ ,  $p = .01$ , with higher percentage errors on French ( $M = 14.75$ ,  $SE = 1.43$ )  
 492 than on English trials ( $M = 11.76$ ,  $SE = 1.39$ ). The interaction between Trial Type and

493 Language was marginally significant,  $F(1,33) = 4.272$ ,  $MSE = 49.6$ ,  $\eta_p^2 = .115$ ,  $BF_{10} = .987$ ,

494  $BF_{01} = 1.013$ ,  $p = .05$ . There was no significant difference in percentage error between French

495 ( $M = 3.68$ ,  $SE = 1.37$ ) and English ( $M = 3.19$ ,  $SE = .86$ ) *All congruent* trials,  $t(33) = .338$ ,  $M_{diff}$

496  $= .49$ ,  $BF_{10} = .194$ ,  $BF_{01} = 5.155$ ,  $p = .737$ . However, participants made significantly more

497 errors on French ( $M = 25.81$ ,  $SE = 2.23$ ) than on English ( $M = 20.33$ ,  $SE = 2.29$ ) *Sound-*

498 *colour congruent* trials,  $t(33) = 3.144$ ,  $M_{diff} = 5.483$ ,  $BF_{10} = 10.617$ ,  $p < .01$ , similar to the

499 response time data.

**500 Mismatching trials**

501           There was a main effect of Trial Type,  $F(2,66) = 19.381$ ,  $MSE = 11.884$ ,  $BF_{10} > 1000$ ,  
502  $\eta_p^2 = .37$ ,  $p < .001$ . That is, participants made significantly more mistakes in *Word-sound*  
503 *congruent* ( $M = 4.095$ ,  $SE = .69$ ) relative to *All incongruent* ( $M = .532$ ,  $SE = .118$ ) trials,  $t(33)$   
504  $= 5.524$ ,  $M_{diff} = 3.563$ ,  $BF_{10} > 1000$ ,  $p < .001$ ), and *Word-colour congruent* ( $M = 1.513$ ,  $SE =$   
505  $.456$ ) trials,  $t(33) = 3.826$ ,  $M_{diff} = 2.583$ ,  $BF_{10} = 54.49$ ,  $p = .001$ . The percentage error was  
506 larger in the *Word-colour congruent* than in the *All incongruent* condition,  $t(33) = 2.329$ ,  $M_{diff}$   
507  $= .98$ ,  $BF_{10} = 1.93$ ,  $p < .05$ . No significant main effect of Language was observed,  $F(1,33) =$   
508  $.102$ ,  $MSE = 6.423$ ,  $\eta_p^2 = .003$ ,  $BF_{10} = .154$ ,  $BF_{01} = 6.493$ ,  $p = .752$ . The interaction between  
509 Trial Type and Language was significant,  $F(2,66) = 5.112$ ,  $MSE = 7.647$ ,  $\eta_p^2 = .134$ ,  $BF_{10} =$   
510  $3.078$ ,  $p = .01$ . There were no significant differences in percentage errors between French and  
511 English *Word-sound congruent* trials,  $t(33) = 1.788$ ,  $M_{diff} = 1.645$ ,  $BF_{10} = .766$ ,  $BF_{01} = 1.305$ ,  
512  $p = .083$  and *All incongruent* trials,  $t(33) = .397$ ,  $M_{diff} = .08$ ,  $BF_{10} = .198$ ,  $BF_{01} = 5.05$ ,  $p =$   
513  $.694$ . However, participants made significantly more errors on English than French *Word-*  
514 *colour congruent* trials,  $t(33) = 2.223$ ,  $M_{diff} = 1.386$ ,  $BF_{10} = 1.587$ ,  $p < .05$ .

**515 Correlations**

516           As a supplementary analysis, we assessed the level to which language metric variables  
517 correlate with different types of trials with both French (L1) and English (L2) colour words  
518 used in the Stroop matching task. These analyses were purely exploratory and did not reveal  
519 any clear or significant results. However, we present these data in the Appendix for the  
520 interested reader.

**521 Discussion**

522           Experiment 1 had two aims: 1) compare the magnitude of between-language and  
523 within-language interference, and 2) investigate the source of interference in a bilingual

524 Stroop matching task with intermixed French (L1) and English (L2) colour word distracters.  
525 Within-language interference was larger than between-language interference, but only for  
526 *Sound-colour congruent* trials, with no significant difference between French and English  
527 word pairs across other trial types. That is, when a spoken word (e.g., “vert”, French for  
528 green) matched the ink colour of the written distracter, the French incongruent distracters  
529 (e.g., “marron”, French for brown printed in green) were responded to slower and less  
530 accurately than English incongruent distracters (e.g., “brown” in green). It is plausible that  
531 French written distracters lead to a strong task-irrelevant comparison (i.e., written word-  
532 spoken word) that impairs performance on a task-relevant comparison (i.e., ink colour-spoken  
533 word). *Sound-colour congruent* trials also had significantly higher percentage errors relative  
534 to all other trial types. This is probably due to the fact that both task-irrelevant comparisons  
535 activate the “mismatching” response in contrast to task-relevant comparison which activates  
536 the “matching” response. However, the observed pattern of results for both French and  
537 English “matching” trials clearly correspond to the assumptions of both stimulus and response  
538 conflict, with faster responses on *All congruent* relative to *Sound-colour congruent* trials.

539       Theoretically more interesting are the results for the mismatching trial types.  
540 Responses on *Word-sound congruent* trials were significantly slower and more error prone  
541 relative to *All incongruent* and *Word-colour congruent* trials (Bornstein, 2015). That is, both  
542 incongruent French (e.g., “vert” in brown) and English (e.g., “green” in brown) distracters  
543 slowed down responding when the word distracter corresponded to the auditory stimulus  
544 (e.g., hear “vert”). This contrasts with the results of Goldfarb and Henik (2006), who found  
545 the slowest “mismatching” responses for congruent distracters (i.e., *Word-colour congruent*  
546 trials). Interestingly, response latencies were almost identical in French and English condition,  
547 suggesting that responding to the spoken L1 word is equally affected by a written L1 word

548 (i.e., both spoken and written words are identical) and an L2 word (i.e., spoken and written  
549 words are not identical, but represent the same colour concept, e.g., “vert” and “green”).

550 The responses were the fastest in *All incongruent* condition, which confirms the  
551 assumptions of the response conflict account. This also aligns with the findings on  
552 behavioural data of Caldas and colleagues (2012) and Goldfarb and Henik (2006), thus  
553 confirming a role of response conflict in the Stroop matching task. In contrast, the semantic  
554 conflict account should have predicted that these trials would be the *slowest*, because the  
555 word, colour, and auditory stimulus are all incongruent with each other.

## 556 **Experiment 2**

557 Experiment 2 conceptually replicates Experiment 1 with one important modification.  
558 In particular, instead of the colour words used in Experiment 1, participants were presented  
559 with French and English colour associates. A complication with the matching task is that the  
560 predictions for the stimulus and response conflict accounts for mismatching trials are exactly  
561 in opposition. The response conflict account predicts that *All incongruent* trials should be the  
562 fastest of the three “mismatching” trial types (as observed), whereas the semantic conflict  
563 account predicts that they should be the slowest. Note that the predictions of both semantic  
564 and response conflict accounts for colour associates are identical to the predictions for colour  
565 words, already visualised in Figure 1. If both types of conflict exist, then it might be that the  
566 (larger) response conflict effect is concealing a (relatively smaller) semantic conflict effect.  
567 Therefore, one way to “reveal” the true effect of semantic conflict (assuming there is one, of  
568 course), would be to eliminate the response conflict. According to some, colour associates  
569 produce semantic conflict (e.g., (Glaser & Glaser, 1989; Schmidt & Cheesman, 2005), but not  
570 response conflict. If this logic is correct, it remains plausible that semantic conflict will be  
571 observed for colour associates. Although probably smaller, semantic conflict might emerge

572 due to strong conceptual links between colour associates and their corresponding colour  
573 words. For example, on a French *Sound-colour congruent* trial (e.g., see “ciel”, French for  
574 sky, printed in green, hear “vert”, French for green), a distracter “ciel”, associated with blue,  
575 should no longer interfere (or very little) with a relevant task comparison (i.e., “green”-  
576 “green”), simply because it does not belong to the same semantic category as a spoken word.  
577 Experiment 2 was therefore designed to further explore the role of semantic conflict that was  
578 possibly masked by response conflict in Experiment 1. Another question of interest concerns  
579 the distracter language. According to some models of bilingual memory, L2 words do not  
580 have strong direct access to semantics (Kroll & Stewart, 1994). Thus, while semantic conflict  
581 might be observed for L1 words, these models would predict the absence of a semantic  
582 conflict effect for L2 words.

### 583 **Method**

#### 584 *Participants*

585 A total of 33 (25 women) [removed for review] undergraduates ( $M_{age} = 20$ ;  $SD = 3.43$ )  
586 voluntarily participated in the experiment. The sample size was determined in the same way  
587 as in Experiment 1. All the selection criteria were identical to Experiment 1. Students who  
588 already participated in Experiment 1 were not allowed to participate in Experiment 2. Their  
589 average language background scores (mean age and standard errors) are presented in Table 2  
590 (see Results section).

#### 591 *Apparatus and materials, design, and procedure*

592 Experiment 2 was identical in all aspects to Experiment 1, with the following  
593 exceptions. First, colour words were replaced by colour associates in French (L1) and English  
594 (L2), which correspond to “blue”, “green”, “red”, and “yellow”, respectively (i.e.,  
595 “ciel”/“sky”, “herbe”/“grass”, “sang”/“blood”, and “citron”/“lemon”). These words were non-

596 cognates with a mean word length of 4.75 for French associates and 4.5 for English  
 597 associates. The colour associates from both languages were chosen based on: 1) their strong  
 598 association with a corresponding colour word (Nelson et al., 1998; Wilson et al., 1988) and 2)  
 599 their similarity in word length. Second, in line with used colour associates, spoken words  
 600 were “bleu” (blue), “vert” (green), “rouge” (red), and “jaune” (yellow). All trial timings were  
 601 identical to Experiment 1.

## 602 **Results**

603 Average language metric scores<sup>4</sup> are presented in Table 2. As in Experiment 1,  
 604 participants started acquiring French at early age (as it is a native language), while English  
 605 was learned as a first foreign language in schools (starting at around 10 years old), but again,  
 606 not to a very high level of mastery. Similar to Experiment 1, participants had rather low  
 607 objective English proficiency, as shown by the LexTale score, as well as low self-estimated  
 608 English level. All participants seemed to sufficiently fit our language dominance criteria.

## 609 **Table 2**

610 *Mean French and English language scores for and standard errors (in brackets)*

	<i>M</i>	<i>SE</i>
LexTale		
Years English	10 years	0.484
English level	2.79 (1-5)	0.155
Score	67.91 (0-100)	1.531
LEAP-Q		
Dominance French	1	0
Dominance English	2.1	0.049
Order French	1	0
Order English	2.06	0.045
French Use (%)	4.85 (1-5)	0.063
English Use (%)	1.82 (1-5)	0.147

<sup>4</sup> All the participants (33/33) indicated French as their first language in order of dominance and in order of acquisition. As a second language in order of dominance and in order of acquisition, participants mostly indicated English, followed by German, Spanish, and Vietnamese. The most frequent third language in dominance and order of acquisition was Spanish, followed by German, English, Italian, and Polish. The majority of participants correctly translated “sky” (31/33), “blood” (32/33), and “lemon” (32/33). However, only half of them correctly translated “grass” (17/33).

French		
Acquisition	0.59 years	0.179
Fluent	3.17 years	0.287
Reading	5.44 years	0.162
Fluent Read	6.61 years	0.252
English		
Acquisition	9.62 years	0.510
Fluent	15.2 years	0.458
Reading	11.39 years	0.486
Fluent Read	15.48 years	0.543

611

612 ***Data Analysis***

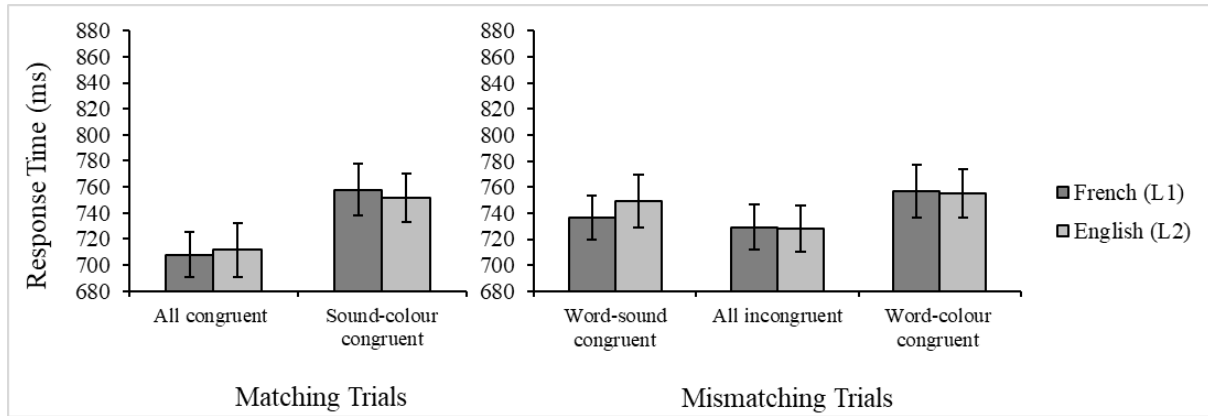
613 As in Experiment 1, the mean correct response times and mean percentage error <sup>5</sup>were  
 614 analysed. No participants were excluded from the sample, their individual accuracy rate  
 615 across the experiment was 89.84% or above. Two separate ANOVAs (one for Matching trials  
 616 and one for Mismatching trials) were conducted for both response times and percentage  
 617 errors.

618 ***Response time (RT)***

619 A total of 5.03% trials were excluded from the analyses (4.65% incorrect and .38%  
 620 time out responses). Only RTs for correct responses in Matching and Mismatching conditions  
 621 were analysed and illustrated in Figure 7.

---

<sup>5</sup> We note that subsequent analyses revealed that response time and error results were largely similar for all four words. It seems plausible that while recall (i.e., translation) was rather low for “grass”, participants were probably able to recognize the English word during the task.

622 **Figure 7**623 *Mean response times with standard errors for “matching” and “mismatching” trials*

624

625 ***Matching trials***

626 There was a main effect of Trial Type,  $F(1,32) = 32.467$ ,  $MSE = 2043.097$ ,  $\eta_p^2 = .504$ ,  
 627  $BF_{10} > 1000$ ,  $p < .001$ , suggesting that responses on *Sound-colour congruent* trials ( $M = 754$ ,  
 628  $SE = 18.71$ ) were significantly slower than responses on *All Congruent* trials ( $M = 710$ ,  $SE =$   
 629  $18.24$ ). However, there was no main effect of Language,  $F(1,32) = .041$ ,  $MSE = 1280.291$ ,  $\eta_p^2$   
 630  $= .001$ ,  $BF_{10} = .182$ ,  $BF_{01} = 5.494$ ,  $p = .840$ , indicating no overall difference in response speed  
 631 to French and English word trials. The interaction between Trial Type and Language was also  
 632 not significant,  $F(1,32) = .364$ ,  $MSE = 2425.755$ ,  $\eta_p^2 = .011$ ,  $BF_{10} = .348$ ,  $BF_{01} = 2.873$ ,  $p =$   
 633  $.550$ .

634 ***Mismatching trials***

635 The main effect of Trial Type was observed,  $F(2,64) = 21.143$ ,  $MSE = 589.472$ ,  $\eta_p^2 =$   
 636  $.398$ ,  $BF_{10} > 1000$ ,  $p < .001$ . *Word-colour congruent* trials ( $M = 756$ ,  $SE = 18.87$ ) were  
 637 responded to slower than *All incongruent* ( $M = 729$ ,  $SE = 17.46$ ) trials,  $t(32) = 6.293$ ,  $M_{diff} =$   
 638  $27$ ,  $BF_{10} > 1000$ ,  $p < .001$ , and *Word-sound congruent* ( $M = 743$ ,  $SE = 17.99$ ) trials,  $t(32) =$   
 639  $3.004$ ,  $M_{diff} = 13$ ,  $BF_{10} = 7.70$ ,  $p = .01$ . Responses were slower on *Word-sound congruent*  
 640 relative to *All incongruent* trials,  $t(32) = 3.663$ ,  $M_{diff} = 14$ ,  $BF_{10} = 35.69$ ,  $p < .01$ . There was



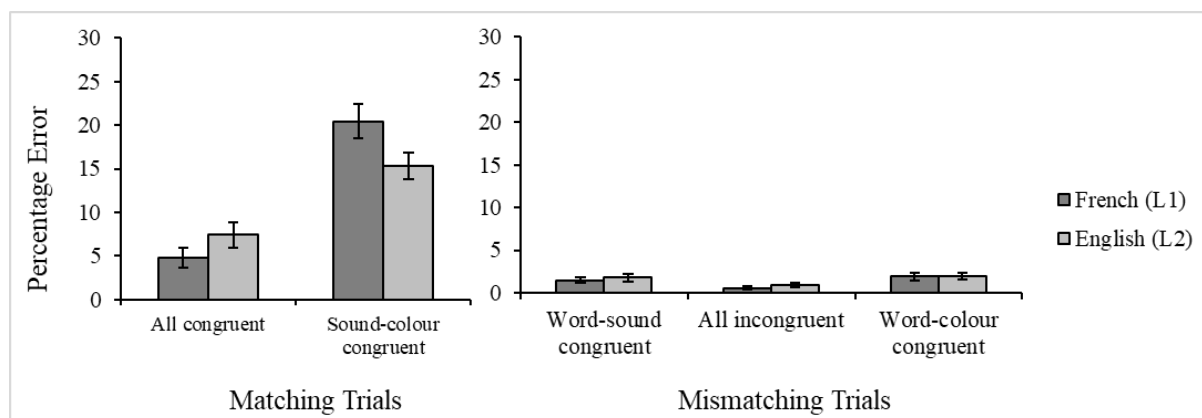
641 no main effect of Language,  $F(1,32) = .581$ ,  $MSE = 882.089$ ,  $\eta_p^2 = .018$ ,  $BF_{10} = .193$ ,  $BF_{01} =$   
 642  $5.181$ ,  $p = .451$ , suggesting no overall difference in response speed between French and  
 643 English word trials. The interaction between Trial Type and Language was also not  
 644 significant,  $F(2,64) = 1.073$ ,  $MSE = 1043.801$ ,  $\eta_p^2 = .032$ ,  $BF_{10} = .25$ ,  $BF_{01} = 4$ ,  $p = .348$ .

#### 645 *Percentage error*

646 The mean percentage error data for all trial types and languages are presented in  
 647 Figure 8.

#### 648 **Figure 8**

649 *Mean percentage errors with standard errors for “matching” and “mismatching” trials*



650

#### 651 *Matching trials*

652 There was a main effect of Trial Type,  $F(1,32) = 77.71$ ,  $MSE = 58.774$ ,  $\eta_p^2 = .708$ ,  
 653  $BF_{10} > 1000$ ,  $p < .001$ , suggesting that *Sound-colour congruent* trials ( $M = 17.859$ ,  $SE =$   
 654  $1.498$ ) were significantly more error-prone relative to *All congruent* trials ( $M = 6.095$ ,  $SE =$   
 655  $.969$ ). No main effect of Language was observed,  $F(1,32) = 1.32$ ,  $MSE = 38.6$ ,  $\eta_p^2 = .04$ ,  $BF_{10}$   
 656  $= .233$ ,  $BF_{01} = 4.292$ ,  $p = .259$ , suggesting no overall difference in percentage error between  
 657 French and English word trials. An interaction between Trial Type and Language was  
 658 significant,  $F(1,32) = 7.839$ ,  $MSE = 61.967$ ,  $\eta_p^2 = .197$ ,  $BF_{10} = 12.331$ ,  $p = .01$ . That is, there

659 was no difference in percentage error between French ( $M = 4.798$ ,  $SE = 1.149$ ) and English  
660 ( $M = 7.392$ ,  $SE = 1.422$ ) *All congruent* trials,  $t(32) = 1.516$ ,  $M_{diff} = 2.594$ ,  $BF_{10} = .525$ ,  $BF_{01} =$   
661  $1.905$ ,  $p = .139$ . However, participants made significantly more errors on French ( $M = 20.399$ ,  
662  $SE = 1.966$ ) than on English ( $M = 15.32$ ,  $SE = 1.486$ ) *Sound-colour congruent* trials,  $t(32) =$   
663  $2.854$ ,  $M_{diff} = 5.079$ ,  $BF_{10} = 5.56$ ,  $p = .01$ .

#### 664 ***Mismatching trials***

665 A main effect of Trial Type was significant,  $F(2,64) = 7.53$ ,  $MSE = 3.182$ ,  $\eta_p^2 = .19$ ,  
666  $BF_{10} = 34.428$ ,  $p = .001$ . Participants made significantly more errors on *Word-colour*  
667 *congruent* trials ( $M = 1.91$ ,  $SE = .32$ ) relative to *All incongruent* ( $M = .75$ ,  $SE = .18$ ) trials,  
668  $t(32) = 4.06$ ,  $M_{diff} = 1.16$ ,  $BF_{10} = 96.42$ ,  $p < .001$ . There was no difference in percentage error  
669 between *Word-colour congruent* and *Word-sound congruent* ( $M = 1.61$ ,  $SE = .33$ ) trials,  $t(32)$   
670  $= .873$ ,  $MEAN_{diff} = .30$ ,  $BF_{10} = .26$ ,  $BF_{01} = 3.85$ ,  $p = .389$ . Participants made more errors on  
671 *Word-sound congruent* relative to *All incongruent* trials,  $t(32) = 2.86$ ,  $M_{diff} = .862$ ,  $BF_{10} =$   
672  $5.63$ ,  $p < .05$ . There was no significant main effect of Language,  $F(1,32) = 1.179$ ,  $MSE =$   
673  $2.931$ ,  $\eta_p^2 = .035$ ,  $BF_{10} = .243$ ,  $BF_{01} = 4.115$ ,  $p = .286$ . An interaction between Trial Type and  
674 Language was also not significant,  $F(2,64) = .154$ ,  $MSE = 3.435$ ,  $\eta_p^2 = .005$ ,  $BF_{10} = .105$ ,  $BF_{01}$   
675  $= 9.524$ ,  $p = .858$ .

#### 676 ***Correlations***

677 As in Experiment 1, we assessed the level to which language metric variables correlate  
678 with different trial types with both French (L1) and English (L2) colour associates used in the  
679 Stroop matching task. Similar to Experiment 1, there were no significant correlations.  
680 However, we present these data in the Appendix.

681 **Discussion**

682 Experiment 2 aimed to 1) compare the magnitude of between-language and within-  
683 language interference produced by French (L1) and English (L2) colour associates, and 2)  
684 investigate the source of this interference. In line with the predictions of both semantic and  
685 response conflict accounts, *Sound-colour congruent* trials are responded to slower and with  
686 more errors relative to *All congruent* trials. Interestingly, a lack of interaction suggests that  
687 participants were equally fast in responding to French and English distracters. This contrasts  
688 the assumption of larger within-language (i.e., produced by French distracters) relative to  
689 between-language (i.e., produced by English distracters) interference.

690 Concerning the “mismatching” trials, *Word-colour congruent* trials were responded to  
691 slower than *Word-sound congruent* and *All incongruent* trials, suggesting that congruent  
692 colour associates (e.g., “ciel” in blue or “sky” in blue) interfere with “mismatching”  
693 responses, as observed by Goldfarb and Henik (2016) and Caldas et al. (2012) with colour  
694 words. It is plausible that participants take additional time to process the congruency of the to-  
695 be-ignored written colour associates, which slows down responding. Interestingly, almost  
696 equal response times were observed with both French and English distracters, suggesting that  
697 first and second language distracters might be processed in a similar way.

698 Finally, responses were again the fastest on *All incongruent* trials, which aligns with  
699 the assumption of the response conflict account. That is, even for colour associate distracters,  
700 participants perform all three task comparisons, which suggest the same, “mismatching”  
701 response alternative. Thus, contrary to expectations, the use of colour associates did not  
702 eliminate response conflict, allowing us to observe a potential true (but small) semantic  
703 conflict effect. Instead, colour associates (like colour words) seemingly produced response  
704 conflict.

705

**General Discussion**

706           The present study aimed to explore the effects of bilingual colour word and colour  
707 associate distracters on matching stimuli presented in auditory (i.e., spoken word) and visual  
708 (i.e., ink colour) formats. In Experiment 1, participants were presented with either congruent  
709 or incongruent colour words in French (L1) and English (L2), accompanied with a spoken  
710 French colour word. Experiment 2 followed the same logic, but French and English colour  
711 associates appeared as distracters. In both experiments, participants were explicitly instructed  
712 to ignore the colour word and to respond based on whether ink colour and spoken word match  
713 or mismatch. This manipulation allowed comparisons between two matching trial types (*All*  
714 *congruent* and *Sound-colour congruent*) and three mismatching trial types (*Word-sound*  
715 *congruent*, *All incongruent*, and *Word-colour congruent*).

716           The first question of interest concerns the language of distracters. Since only French  
717 colour words were used as spoken stimuli, French distracters should produce within-language  
718 interference, whereas English distracters should produce between-language interference. As  
719 already discussed in the Introduction, previous findings suggest that within-language  
720 interference is usually larger than between-language interference (Fang et al., 1981; Kiyak,  
721 1982; MacLeod, 1991). We observed this pattern with the matching trial types, where there  
722 was evidence for a larger congruency effect for L1 than L2. No language differences were  
723 found for mismatching trial types, however. This makes the findings similar to those expected  
724 for more balanced bilinguals. It is important to note that participants were tested only on a  
725 small set of words (i.e., colour words), which are often learned in early phases of second  
726 language learning. It would be interesting to test these finding with less balanced bilinguals or  
727 by using a larger set of distracting words, which might reveal clearer differences between L1  
728 and L2 items. Future work may also make use of mixed modelling of individual-trial response  
729 times, as traditional methods of data analysis do not always account for individual differences

730 across bilingual participants (Privitera et al., 2023). Alternatively, L2 words might possess a  
731 strong link with their corresponding conceptual representations, similarly to L1 words (Šaban  
732 & Schmidt, 2021; Schmidt et al., 2018). As discussed in the Introduction, L2 words could  
733 possess strong semantic connections, lexical connections, or a combination of both.  
734 Therefore, the nature of L2 connections and their strength towards lexical and semantic  
735 representations should help elucidate the similarities/differences observed in patterns for both  
736 L1 and L2 words.

737         However, it seems that the difference in magnitudes of within- and between-language  
738 interference is even smaller for colour associates (Experiment 2) relative to colour words  
739 (Experiment 1). That is, overall response times were faster for colour associates than for  
740 colour words (Schmidt & Cheesman, 2005). Moreover, no difference was observed between  
741 French and English trials, thus suggesting that the first and second language colour associates  
742 seem to interfere less with L1 spoken colour words relative to colour word distracters. This  
743 can be due to the fact that colour associates, although semantically related to colour words, do  
744 not correspond to the spoken colour words. This finding thus revealed that the meaning of the  
745 written distracter, either from L1 or L2, cannot be completely ignored, resulting in a decrease  
746 of the response speed within which ink colour and spoken words were judged as “matching”  
747 or “mismatching”. This interference produced by written distracters seems to increase  
748 proportionally with its similarity to the spoken word. That is, in both experiments, spoken  
749 words were French colour words. Responses were generally slower in Experiment 1 when the  
750 same set of French colour words was used as distracters. That is, written, to-be-ignored colour  
751 word distracters also served as potential targets. In contrast, responses were faster in  
752 Experiment 2 when colour associates were used as distracters. Although these colour  
753 associates were semantically related to spoken colour words, they were not targets. This  
754 aligns with the assumptions of the *response set membership* account (Klein, 1964; Risko et

755 al., 2006), which refers to a larger interference for words (e.g., distracters) that are potential  
756 targets (e.g., or a to-be-attended stimulus dimension, such as a spoken word in the Stroop  
757 matching task). This has been confirmed with both colour words and colour associates (Klein,  
758 1964; Risko et al., 2006; Schmidt & Cheesman, 2005; Sharma & McKenna, 1998) in the  
759 literature and in the present series of experiments.

760 A second question of interest is the source of interference produced in the Stroop  
761 matching task. The semantic conflict account suggests that responses will be the slowest on  
762 trials in which task dimensions activate multiple colour concepts. For instance, larger  
763 interference is expected on trials in which two contrasting colour representations are  
764 simultaneously activated (e.g., *Sound-colour congruent* trials) relative to trials in which only  
765 one colour representation is activated (e.g., *All congruent* trials). In contrast, the response  
766 conflict account focuses on task comparisons and assumes that responses will be slowest on  
767 trials in which task-relevant and task-irrelevant comparisons suggest different responses. That  
768 is, responses should be faster on trials in which all three task comparisons suggest the same  
769 response option (e.g., “match” or “mismatch”, for *All congruent* and *All incongruent* trials,  
770 respectively), relative to those trials in which one comparison activates one response option,  
771 whereas two other comparisons activate contrasting response option (e.g., on *Word-sound*  
772 *congruent* or *Word-colour congruent* trials). The interplay between semantic and response  
773 conflict is also possible. For instance, these two conflict effects might be in opposition in the  
774 matching task. That is, the larger response conflict is “masking” the smaller semantic conflict.  
775 One way to measure the true effect of semantic conflict would be to eliminate the response  
776 conflict. To do so, colour associates (which are assumed to produce semantic conflict  
777 exclusively) were used as alternative to colour words in Experiment 2.

778 As expected, the response times were slower for incongruent trials (i.e., *Sound-colour*  
779 *congruent*) relative to congruent trials (i.e., *All congruent*) with “matching” response.

780 However, previous findings suggest that the response times are slower for congruent relative  
781 to incongruent trials with “mismatching” responses (Bornstein, 2015; Caldas et al., 2012;  
782 Goldfarb & Henik, 2006). That is, *Word-colour congruent* trials (e.g., “green” in green, hear  
783 “pink”) are assumed to be responded to slower than *All incongruent* (e.g., “green” in brown,  
784 hear “pink”) and *Word-sound congruent* (e.g., “green” in brown, hear “green”). This has been  
785 replicated in Experiment 2 using colour associates, when *Word-colour congruent* trials (e.g.,  
786 “sky” in blue, hear “green”) produced the slowest response latencies as compared to other two  
787 types of trial. However, this pattern was not observed in Experiment 1 which made use of  
788 colour words. In Experiment 1, the responses were slowest on *Word-sound congruent* trials  
789 (e.g., “green” in brown, hear “green”). That is, instead of focusing on congruency of the  
790 written stimuli exclusively, as suggested by previous studies, participants tend to compare a  
791 written, to-be-ignored stimulus, with a spoken word, thus engaging a task-irrelevant  
792 comparison.

793 Navon (1985) introduced the notion of outcome-conflict to reflect a state where the  
794 output of one task modifies (and potentially interferes) a variable that is relevant to the  
795 performance of a concurrent task (Navon, 1985; Navon & Miller, 1987). In this  
796 conceptualization, performance in the Stroop matching task is determined by a conflict of  
797 outcomes between three separate dimensions, each one resulting in either a “matching” or  
798 “mismatching” outcome. It is possible that outcome-conflicts occurred whenever the relevant  
799 matching task and the two mistakenly performed matching tasks produced conflicting  
800 outcomes (i.e., “matching” vs. “mismatching”). Interference effects were large and significant  
801 only in conditions that featured such a conflict. For instance, outcome-conflict does not  
802 predict any interference in *All congruent* and *All incongruent* conditions because all three  
803 comparisons between colour representations indicate the same response, “matching” and  
804 “mismatching”, respectively. According to this account, when one irrelevant matching

805 outcome conflicted with the response (e.g., on *Word-sound congruent* and *Word-colour*  
806 *congruent* trials, when a correct response was “mismatch”, and two irrelevant comparisons  
807 suggest “match” and “mismatch”), the interference should be smaller than on trials in which  
808 both irrelevant outcomes conflicted with the response (e.g., on *Sound-colour congruent* trials  
809 when a correct response was “match”, but both irrelevant comparisons suggest “mismatch”).  
810 In sum, as the number of outcome-conflicts becomes larger, performance is more prone to  
811 errors. Our results align with this: the percentage error was extremely high in the *Sound-*  
812 *colour congruent* condition relative to remaining four trial types (in both Experiment 1 and  
813 Experiment 2). Consequently, to achieve higher accuracy, participants probably focus on  
814 serial processing of separate comparisons, which in turn might have produced additional  
815 response delays. This is also observable in the present results, with *Sound-colour congruent*  
816 trials being slower relative to all other trial types.

817         The present findings also align with the confluence model proposed by Eviatar and  
818 colleagues (1994) based on their findings from a visual matching task. According to this  
819 model, in matching tasks, all stimulus dimensions are processed automatically and  
820 simultaneously regardless of task relevance. This processing and an interference produced by  
821 the outputs between all task dimensions precede response selection. In the present study,  
822 visual and auditory stimuli were processed until their representations could be compared. The  
823 “matching” or “mismatching” among the outputs of these comparisons determined the  
824 response speed and the likelihood of selecting the correct response alternative. This  
825 interpretation is similar to the one proposed by Navon’s (1985) outcome-conflict account.  
826 However, this confluence model is more specifically oriented toward matching tasks and  
827 more explicit regarding the processing stage to which interference is attributed (Eviatar et al.,  
828 1994).



829           The present findings with colour word distracters (Experiment 1) align with  
830 behavioural data of Caldas and colleagues (2012) and those of Goldfarb and Henik (2006),  
831 providing an additional support for the response conflict account. Interestingly, we observed  
832 response conflict effect even with colour associates, which we assumed (incorrectly) would  
833 eliminate the response conflict component. However, the electrophysiological data of Caldas  
834 and colleagues (2012) supported a semantic conflict account. This data showed that conflict  
835 related brain activity, as indicated by a greater frontal negativity (N450), was not observed for  
836 a “mismatching” condition that featured conflicting irrelevant “matching” output. Rather,  
837 N450 amplitude was greater in *Word-colour congruent* and *All incongruent* conditions than in  
838 the *Word-sound congruent* condition. This discrepancy between behavioural and  
839 electrophysiological data suggests that interference produced in the Stroop matching task  
840 could be due to contributions of both semantic and response conflict. It is plausible that the  
841 role of semantic conflict in explaining the Stroop matching interference could be evidenced  
842 exclusively by using more subtle measures, such as electrophysiology. Another possibility is  
843 that there still might be a semantic conflict effect observable in behavioural studies, however,  
844 it is still being masked by response conflict.

845           The present results clearly indicate a role for response conflict in the Stroop matching  
846 task, for colour words and colour associates and in the first and less-fluent second language.  
847 However, the role of semantic conflict is less clear. As highlighted in this manuscript, one  
848 peculiarity of the matching task is that it can only provide evidence for either response  
849 conflict or semantic conflict, but not both, as the two are pitted against each other. As such, it  
850 is not currently clear whether semantic conflict was absent in our studies, or rather merely  
851 smaller than (and therefore concealed by) response conflict. Future research could help  
852 answering these inquiries. Indeed, as indicated in the Introduction, one of the goals of the  
853 present manuscript was to assess some competing models of bilingual memory. According to

854 certain models, stimulus conflict should only occur for L1 words in early language learners,  
855 but not for L2 words, whereas other models suggest that stimulus conflict should occur for  
856 both. Given the absence of stimulus conflict in the present task, even for L1 words, we were  
857 unable to assess such competing models with the current data. In sum, despite the fact that  
858 response conflict plays an important role in the interference produced in the Stroop matching  
859 task, this does not discard the possibility that some other, non-response (i.e., semantic)  
860 conflict also contributes to this effect, which remains a focus of debate (Caldas et al., 2020;  
861 Dittrich & Stahl, 2017; Green et al., 2016; Luo, 1999).

862

### **Conclusion**

863 The present experiments explored how different types of first and second language words  
864 influence audio-visual matching performance. The findings suggest that, regardless of the  
865 distracting language (L1 vs. L2), responses were the fastest on trials in which task  
866 comparisons activate fewer response alternatives, supporting the assumption of the response  
867 conflict account. That is, performance is faster when no competition between response  
868 alternatives occurs. The present work serves as a good starting point in understanding how  
869 simultaneous audio-visual processing affects response selection across languages and word  
870 types.

871

872 **References**

- 873 Altarriba, J., & Mathis, K. M. (1997). Conceptual and lexical development in second  
874 language acquisition. *Journal of Memory and Language*, 36(4), 550–568.  
875 <https://doi.org/10.1006/jmla.1997.2493>
- 876 Augustinova, M., & Ferrand, L. (2014). Automaticity of word reading: Evidence from the  
877 semantic Stroop paradigm. *Current Directions in Psychological Science*, 23(5), 343–  
878 348. <https://doi.org/10.1177/0963721414540169>
- 879 Besner, D., Stolz, J. A., & Boutilier, C. (1997). The stroop effect and the myth of  
880 automaticity. *Psychonomic Bulletin & Review*, 4(2), 221–225.  
881 <https://doi.org/10.3758/BF03209396>
- 882 Blumenfeld, H. K., & Marian, V. (2014). Cognitive control in bilinguals: Advantages in  
883 Stimulus–Stimulus inhibition. *Bilingualism: Language and Cognition*, 17(3), 610–  
884 629. <https://doi.org/10.1017/S1366728913000564>
- 885 Bornstein, I. S. (2015). *Behavioural measures of interference and facilitation in an*  
886 *audiovisual colour-word Stroop matching task*. The University of British Columbia.
- 887 Caldas, A. L., Machado-Pinheiro, W., Daneyko, O., & Riggio, L. (2020). The Stroop-  
888 matching task as a tool to study the correspondence effect using images of graspable  
889 and non-graspable objects. *Psychological Research*, 84(7), 1815–1828.  
890 <https://doi.org/10.1007/s00426-019-01191-5>
- 891 Caldas, A. L., Machado-Pinheiro, W., Souza, L. B., Motta-Ribeiro, G. C., & David, I. A.  
892 (2012). The Stroop matching task presents conflict at both the response and  
893 nonresponse levels: An event-related potential and electromyography study: Stroop  
894 matching task and conflict. *Psychophysiology*, 49(9), 1215–1224.  
895 <https://doi.org/10.1111/j.1469-8986.2012.01407.x>

- 896 Dittrich, K., & Stahl, C. (2017). Two distinct patterns of interference in between-attribute  
897 Stroop matching tasks. *Attention, Perception, & Psychophysics*, *79*(2), 563–581.  
898 <https://doi.org/10.3758/s13414-016-1253-x>
- 899 Dyer, F. N. (1971). Color-naming interference in monolinguals and bilinguals. *Journal of*  
900 *Verbal Learning and Verbal Behavior*, *10*(3), 297–302. [https://doi.org/10.1016/S0022-](https://doi.org/10.1016/S0022-5371(71)80057-9)  
901 [5371\(71\)80057-9](https://doi.org/10.1016/S0022-5371(71)80057-9)
- 902 Dyer, F. N. (1973). Interference and facilitation for color naming with separate bilateral  
903 presentations of the word and color. *Journal of Experimental Psychology*, *99*(3), 314–  
904 317. <https://doi.org/10.1037/h0035245>
- 905 Egeth, H. E., Blecker, D. L., & Kamlet, A. S. (1969). Verbal interference in a perceptual  
906 comparison task. *Perception & Psychophysics*, *6*(6), 355–356.  
907 <https://doi.org/10.3758/BF03212790>
- 908 Eviatar, Z., Zaidel, E., & Wickens, T. (1994). Nominal and physical decision criteria insame-  
909 different judgments. *Perception & Psychophysics*, *56*(1), 62–72.  
910 <https://doi.org/10.3758/BF03211691>
- 911 Fang, S.-P., Tzeng, O. J. L., & Alva, L. (1981). Intralanguage vs. Interlanguage Stroop effects  
912 in two types of writing systems. *Memory & Cognition*, *9*(6), 609–617.  
913 <https://doi.org/10.3758/BF03202355>
- 914 Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G\*Power 3: A flexible statistical  
915 power analysis program for the social, behavioral, and biomedical sciences. *Behavior*  
916 *Research Methods*, *39*(2), 175–191. <https://doi.org/10.3758/BF03193146>
- 917 Ferrand, L., & Augustinova, M. (2014). Differential effects of viewing positions on standard  
918 versus semantic Stroop interference. *Psychonomic Bulletin & Review*, *21*(2), 425–431.  
919 <https://doi.org/10.3758/s13423-013-0507-z>

- 920 Flowers, J. H. (1975). "Sensory" interference in a word-color matching task. *Perception &*  
921 *Psychophysics*, 18(1), 37–43. <https://doi.org/10.3758/BF03199364>
- 922 Glaser, M. O., & Glaser, W. R. (1989). Context effects in Stroop-like word and picture  
923 processing. *Journal of Experimental Psychology: General*, 118, 13–42.
- 924 Goldfarb, L., & Henik, A. (2006). New Data Analysis of the Stroop Matching Task Calls for a  
925 Reevaluation of Theory. *Psychological Science*, 17(2), 96–100.  
926 <https://doi.org/10.1111/j.1467-9280.2006.01670.x>
- 927 Goldfarb, L., & Henik, A. (2007). Evidence for task conflict in the Stroop effect. *Journal of*  
928 *Experimental Psychology: Human Perception and Performance*, 33(5), 1170–1176.  
929 <https://doi.org/10.1037/0096-1523.33.5.1170>
- 930 Green, M. L., Locker, L., Boyer, T. W., & Sturz, B. R. (2016). Stroop-like interference in a  
931 match-to-sample task: Further evidence for semantic competition? *Learning and*  
932 *Motivation*, 56, 53–64. <https://doi.org/10.1016/j.lmot.2016.09.003>
- 933 Hamers, J. F., & Lambert, W. E. (1972). Bilingual interdependencies in auditory perception.  
934 *Journal of Verbal Learning and Verbal Behavior*, 11, 303–310.
- 935 Kiyak, H. A. (1982). Interlingual interference in naming color words. *Journal of Cross-*  
936 *Cultural Psychology*, 13(1), 125–135.
- 937 Klein, G. S. (1964). Semantic power measured through the interference of words with color-  
938 naming. *The American Journal of Psychology*, 77(4), 576.  
939 <https://doi.org/10.2307/1420768>
- 940 Kroll, J. F., & Stewart, E. (1994). Category interference in translation and picture naming:  
941 Evidence for asymmetric connections between bilingual memory representations.  
942 *Journal of Memory and Language*, 33, 149–174.

- 943 Lemhöfer, K., & Broersma, M. (2012). Introducing LexTALE: A quick and valid lexical test  
944 for advanced learners of English. *Behavior Research Methods*, *44*(2), 325–343.  
945 <https://doi.org/10.3758/s13428-011-0146-0>
- 946 Lorentz, E., McKibben, T., Ekstrand, C., Gould, L., Anton, K., & Borowsky, R. (2016).  
947 Disentangling Genuine Semantic Stroop Effects in Reading from Contingency Effects:  
948 On the Need for Two Neutral Baselines. *Frontiers in Psychology*, *7*.  
949 <https://doi.org/10.3389/fpsyg.2016.00386>
- 950 Luo, C. R. (1999). Semantic competition as the basis of Stroop interference: Evidence from  
951 color-word matching tasks. *Psychological Science*, *10*(1), 35–40.  
952 <https://doi.org/10.1111/1467-9280.00103>
- 953 MacLeod, C. M. (1991). Half a century of research on the Stroop effect: An integrative  
954 review. *Psychological Bulletin*, *109*(2), 163–203. [https://doi.org/10.1037/0033-](https://doi.org/10.1037/0033-2909.109.2.163)  
955 [2909.109.2.163](https://doi.org/10.1037/0033-2909.109.2.163)
- 956 Mägiste, E. (1982). Automaticity and interference in bilinguals. *Psychological Research*,  
957 *44*(1), 29–43. <https://doi.org/10.1007/BF00308553>
- 958 Marian, V., Blumenfeld, H. K., & Kaushanskaya, M. (2007). The Language Experience and  
959 Proficiency Questionnaire (LEAP-Q): Assessing language profiles in bilinguals and  
960 multilinguals. *Journal of Speech, Language, and Hearing Research*, *50*(4), 940–967.  
961 [https://doi.org/10.1044/1092-4388\(2007\)067](https://doi.org/10.1044/1092-4388(2007)067)
- 962 Morton, J. (1969). Categories of interference: Verbal mediation and conflict in card sorting.  
963 *British Journal of Psychology*, *60*(3), 329–346. [https://doi.org/10.1111/j.2044-](https://doi.org/10.1111/j.2044-8295.1969.tb01204.x)  
964 [8295.1969.tb01204.x](https://doi.org/10.1111/j.2044-8295.1969.tb01204.x)
- 965 Navon, D. (1985). Attention division or attention sharing? In *Attention and performance XI*  
966 (pp. 133–146). Erlbaum.

- 967 Navon, D., & Miller, J. (1987). Role of outcome conflict in dual-task interference. *Journal of*  
968 *Experimental Psychology: Human Perception and Performance*, 13(3), 435–448.
- 969 Nelson, D. L., McEvoy, C. L., & Schreiber, T. A. (1998). *The University of South Florida*  
970 *word association, rhyme, and word fragment norms*.  
971 <http://www.usf.edu/FreeAssociation/>
- 972 Posner, M. I., & Snyder, C. R. R. (1975). Attention and cognitive control. In *Information*  
973 *Processing and cognition: The Loyola Symposium* (pp. 55–85). Hillsdale, NJ:  
974 Erlbaum.
- 975 Preston, M. S., & Lambert, W. E. (1969). Interlingual interference in a bilingual version of  
976 the Stroop color-word task. *Journal of Verbal Learning and Verbal Behavior*, 8(2),  
977 295–301. [https://doi.org/10.1016/S0022-5371\(69\)80079-4](https://doi.org/10.1016/S0022-5371(69)80079-4)
- 978 Privitera, A. J., Momenian, M., & Weekes, B. S. (2023). Modeling the bilingual advantage:  
979 Do results differ between analysis methods? *Ampersand*, 11, 100134.  
980 <https://doi.org/10.1016/j.amper.2023.100134>
- 981 Risko, E. F., Schmidt, J. R., & Besner, D. (2006). Filling a gap in the semantic gradient: Color  
982 associates and response set effects in the Stroop task. *Psychonomic Bulletin & Review*,  
983 13(2), 310–315. <https://doi.org/10.3758/BF03193849>
- 984 Šaban, I., & Schmidt, J. R. (2021). Stimulus and response conflict from a second language:  
985 Stroop interference in weakly-bilingual and recently-trained languages. *Acta*  
986 *Psychologica*, 218, 103360. <https://doi.org/10.1016/j.actpsy.2021.103360>
- 987 Schmidt, J. R., & Cheesman, J. (2005). Dissociating stimulus-stimulus and response-response  
988 effects in the Stroop task. *Canadian Journal of Experimental Psychology/Revue*  
989 *Canadienne de Psychologie Expérimentale*, 59(2), 132–138.  
990 <https://doi.org/10.1037/h0087468>

- 991 Schmidt, J. R., Crump, M. J. C., Cheesman, J., & Besner, D. (2007). Contingency learning  
992 without awareness: Evidence for implicit control. *Consciousness and Cognition*,  
993 *16*(2), 421–435. <https://doi.org/10.1016/j.concog.2006.06.010>
- 994 Schmidt, J. R., Hartsuiker, R. J., & De Houwer, J. (2018). Interference in Dutch–French  
995 bilinguals: Stimulus and response conflict in intra- and interlingual Stroop.  
996 *Experimental Psychology*, *65*(1), 13–22. <https://doi.org/10.1027/1618-3169/a000384>
- 997 Seymour, P. H. K. (1977). Conceptual encoding and locus of the Stroop effect. *Quarterly*  
998 *Journal of Experimental Psychology*, *29*(2), 245–265.
- 999 Sharma, D., & McKenna, F. P. (1998). Differential components of the manual and vocal  
1000 Stroop tasks. *Memory & Cognition*, *26*(5), 1033–1040.  
1001 <https://doi.org/10.3758/BF03201181>
- 1002 Simon, J. R., & Berbaum, K. (1988). Effect of irrelevant information on retrieval time for  
1003 relevant information. *Acta Psychologica*, *67*(1), 33–57. [https://doi.org/10.1016/0001-](https://doi.org/10.1016/0001-6918(88)90023-6)  
1004 [6918\(88\)90023-6](https://doi.org/10.1016/0001-6918(88)90023-6)
- 1005 Stoet, G. (2010). PsyToolkit: A software package for programming psychological  
1006 experiments using Linux. *Behavior Research Methods*, *42*(4), 1096–1104.  
1007 <https://doi.org/10.3758/BRM.42.4.1096>
- 1008 Stoet, G. (2017). PsyToolkit: A novel web-based method for running online questionnaires  
1009 and reaction-time experiments. *Teaching of Psychology*, *44*(1), 24–31.  
1010 <https://doi.org/10.1177/0098628316677643>
- 1011 Stroop, J. R. (1935). Studies on interference in serial verbal reactions. *Journal of*  
1012 *Experimental Psychology*, *18*, 643–661.
- 1013 Tanaka, J. W., & Presnell, L. M. (1999). Color diagnosticity in object recognition. *Perception*  
1014 *& Psychophysics*, *61*(6), 1140–1153. <https://doi.org/10.3758/BF03207619>



- 1015 Treisman, A., & Fearnley, S. (1969). The Stroop Test: Selective Attention to Colours and  
1016 Words. *Nature*, 222(5192), 437–439. <https://doi.org/10.1038/222437a0>
- 1017 Tzelgov, J., Henik, A., & Leiser, D. (1990). Controlling Stroop interference: Evidence from a  
1018 bilingual task. *Journal of Experimental Psychology: Learning, Memory, and*  
1019 *Cognition*, 16(5), 760–771. <https://doi.org/10.1037/0278-7393.16.5.760>
- 1020 Virzi, R. A., & Egeth, H. E. (1985). Toward a translational model of Stroop interference.  
1021 *Memory & Cognition*, 13(4), 304–319. <https://doi.org/10.3758/BF03202499>
- 1022 Wilson, M., Kiss, G., & Armstrong, C. (1988). *University of Oxford, EAT: the Edinburgh*  
1023 *associative corpus*. <http://hdl.handle.net/20.500.12024/1251>
- 1024
- 1025

1026 **Declaration**

1027 **Conflict of interest**

1028 The authors have no conflicts of interest to declare.

1029

- 1030 **Replication package**
- 1031 Research materials
- 1032 All research materials, including participant recruitment material, questionnaire, task  
1033 instructions and debriefing form are available at  
1034 [https://osf.io/48q2p/?view\\_only=3c4cdec3f832446984291fc5f22f6392](https://osf.io/48q2p/?view_only=3c4cdec3f832446984291fc5f22f6392) under the section  
1035 “Research materials”
- 1036 Data
- 1037 Data is available at [https://osf.io/48q2p/?view\\_only=3c4cdec3f832446984291fc5f22f6392](https://osf.io/48q2p/?view_only=3c4cdec3f832446984291fc5f22f6392)  
1038 under the section “Data”.
- 1039 Analysis code
- 1040 Instructions and code required to reproduce all analyses are available at  
1041 [https://osf.io/48q2p/?view\\_only=3c4cdec3f832446984291fc5f22f6392](https://osf.io/48q2p/?view_only=3c4cdec3f832446984291fc5f22f6392) under the section  
1042 “Analysis code”.
- 1043

## Appendix

1044

1045 Table A1 presents the non-parametric rank-based Spearman's correlation coefficients  
 1046 between the behavioural measures (i.e., response times and error rates) and language metric  
 1047 scores for Experiment 1. We observed that only percentage error, but not response speed,  
 1048 correlated with certain language metric variables (e.g., age of development of English reading  
 1049 skills or percentage of English exposure). Note however that after applying a Holm-  
 1050 Bonferroni correction for multiple comparisons, none of the correlations were significant at  $\alpha$   
 1051 = .05, so these correlations should be interpreted with caution.

1052 **Table A1**1053 *Correlations between behavioural and language metric scores in Experiment 1*

	French (L1)										English (L2)									
	Matching				Mismatching						Matching				Mismatching					
	All congruent		Sound-colour congruent		Word-sound congruent		All incongruent		Word-colour congruent		All congruent		Sound-colour congruent		Word-sound congruent		All incongruent		Word-colour congruent	
	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR
LexTALE	-.186	.022	-.211	.157	-.090	-.054	-.093	-.054	-.107	-.121	-.248	-.126	-.123	.157	-.086	-.054	-.132	-.015	-.132	.212
English Level	-.128	-.296	-.027	-.077	.113	.294	.098	-.385	.053	.047	.021	-.002	-.029	.034	.103	.077	.047	.161	.099	-.116
Years English	.032	.145	.043	-.035	.171	.033	.183	-.081	.139	.166	-.117	-.113	.197	-.128	.318	-.005	.179	-.098	.176	-.194
% French Exposure	.027	.096	.098	-.098	.062	-.193	.009	.110	.062	-.390	-.044	-.227	.027	.196	-.009	-.028	.115	.111	.062	-.311
% English Exposure	.032	-.214	-.086	-.430	.025	-.068	.123	-.318	.014	-.125	.046	-.032	-.074	-.185	.152	.066	.056	-.242	.086	-.193
FRENCH																				
Acquisition	.093	.136	.090	-.146	-.061	-.339	.020	-.218	-.082	-.297	.234	.058	.068	-.139	-.101	<b>-.456</b>	.005	-.123	.127	-.238
Fluent	-.202	.023	-.090	.304	-.068	.080	-.202	.069	-.240	-.083	-.150	.073	-.158	.160	-.136	-.004	-.119	.069	-.032	.070
Reading	-.019	-.009	.044	.138	.049	.241	.026	.131	-.057	-.134	.068	-.153	.068	.240	.089	.059	.060	.151	.184	.067
Fluent Reading	.062	.307	.029	.206	-.015	.034	-.071	.057	-.080	-.364	.148	-.027	-.025	.327	-.004	-.017	.026	-.053	.069	-.089
ENGLISH																				
Acquisition	-.113	-.063	-.088	.017	.052	.076	-.134	.059	-.070	.009	.009	-.073	-.161	-.173	-.134	.213	-.084	.032	-.134	.129
Fluent	.139	.213	.135	.180	-.082	.051	.067	-.028	-.025	-.105	.075	.014	.242	.188	.082	-.097	.019	-.189	.068	.012
Reading	-.021	<b>.470</b>	-.141	.081	.006	.061	-.109	.200	-.083	.021	-.240	-.114	-.193	.114	-.024	.393	-.069	-.090	-.146	.185
Fluent Reading	-.079	.327	-.187	.128	-.026	.156	-.057	.052	-.073	.052	-.203	-.087	-.198	-.024	.077	.388	-.023	-.256	-.109	.179

1054 Note. Italic =  $p < .05$ , Bold =  $p < .01$ ; no tests were significant after Holm-Bonferroni correction.

1055

1056 Table A2 presents the same correlation for the Experiment 2 data. As in Experiment 1,  
 1057 none of the correlations were significant at  $\alpha = .05$  after applying the Holm-Bonferroni  
 1058 correction for multiple comparisons. As such, the following should be interpreted with  
 1059 caution. We observed that the response speed for all trial types (both French and English)  
 1060 were negatively correlated with the age of reading in French. That is, the earlier participants  
 1061 started reading in French, the slower their responses were. This seems reasonable because  
 1062 reading is often considered as an automatic skill (Augustinova & Ferrand, 2014) acquired  
 1063 early in life. However, in this task, participants were explicitly instructed to avoid reading a  
 1064 distracter since it represents a task-irrelevant dimension and impairs matching/mismatching  
 1065 responses.

1066 **Table A2**

1067 *Correlations between behavioural and language metric scores in Experiment 2*

	French (L1)										English (L2)												
	Matching					Mismatching					Matching					Mismatching							
	All congruent		Sound-colour congruent		Word-sound congruent		All incongruent		Word-colour congruent			All congruent		Sound-colour congruent		Word-sound congruent		All incongruent		Word-colour congruent			
	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	RT	ERR	
LexTALE	.293	.016	.109	.077	.112	.101	.054	-.285	.111	-.109	.202	.172	.183	.067	.166	.048	.084	.024	-.001	.043			
English Level	.121	.323	.039	-.008	.026	.082	-.040	-.369	.033	-.237	.096	.141	.022	.057	.018	-.143	-.039	.098	.011	-.050			
Years English	.161	.181	.095	-.188	.198	-.116	.246	-.117	.135	.065	.140	.189	.200	-.325	.112	-.282	.173	.099	.172	-.357			
% French Exposure	.124	-.213	.275	.286	.302	.228	.186	.256	.266	.094	.435	.327	.284	.254	.266	.214	.266	.232	.337	.171			
% English Exposure	-.141	-.094	.052	-.129	-.043	.237	-.039	-.415	-.010	-.203	-.009	.048	-.011	-.235	-.078	-.153	-.078	.147	-.075	.116			
FRENCH																							
Acquisition	-.141	.202	-.072	-.061	-.117	-.136	-.123	.107	-.271	.110	-.095	.041	-.228	-.097	-.142	-.078	-.135	-.045	-.061	-.092			
Fluent	.055	-.017	.248	.084	.262	.066	.286	.152	.238	.262	.203	.209	.120	.199	.196	.059	.176	.089	.226	.280			
Reading	-.411	.056	-.497	.030	-.454	.100	-.505	.150	-.467	-.064	-.465	-.159	-.422	.165	-.464	-.063	-.482	.210	-.445	.136			
Fluent Reading	-.116	.255	-.189	-.078	-.158	-.196	-.087	.187	-.108	.040	-.196	-.106	-.150	-.003	-.111	-.239	-.127	.133	-.097	-.189			
ENGLISH																							
Acquisition	.149	.123	.148	-.001	.110	.043	.085	-.057	.101	-.101	.283	.118	.155	.092	.163	.085	.129	-.069	.067	.224			
Fluent	.329	-.120	.107	-.178	.147	.031	.159	-.545	.216	-.313	.080	-.234	.262	-.164	.221	-.130	.099	-.263	-.023	.112			
Reading	.040	.039	.085	-.281	.020	-.225	.065	-.312	.116	-.408	.036	-.284	.135	-.238	.092	-.328	.128	-.030	.041	-.017			
Fluent Reading	.004	-.176	.061	-.237	-.012	-.070	.077	-.520	.029	-.335	-.137	-.382	.064	-.293	.008	-.230	.013	-.238	-.092	.088			

1068 Note. Italic =  $p < .05$ , Bold =  $p < .01$ ; no tests were significant after Holm-Bonferroni correction.

1069