

PAPER IN PRESS AT  
COGNITIVE RESEARCH: PRINCIPLES AND IMPLICATIONS

1 Can Corrections Spread Misinformation to New Audiences? Testing for the Elusive  
2 Familiarity Backfire Effect

3 Ullrich K. H. Ecker<sup>1</sup>, Stephan Lewandowsky<sup>2,1</sup>, & Matthew Chadwick<sup>1</sup>

4 <sup>1</sup> School of Psychological Science, University of Western Australia, 35 Stirling Hwy, Perth  
5 6009, Australia; ullrich.ecker@uwa.edu.au (UE), 21152437@student.uwa.edu.au (MC)

6 <sup>2</sup> School of Psychological Science, University of Bristol, 12a Priory Road, Bristol BS8 1TU,  
7 United Kingdom; stephan.lewandowsky@bristol.ac.uk

8 Word count: 11,814 (main text including footnotes, tables, and figure captions)

9 Address correspondence to: Ullrich Ecker, School of Psychological Science (M304),  
10 University of Western Australia, 35 Stirling Hwy, Perth 6009, Australia. Telephone: +618  
11 6488 3257; e-mail: ullrich.ecker@uwa.edu.au.

12 Abstract: Misinformation often continues to influence inferential reasoning after clear and  
13 credible corrections are provided; this effect is known as the continued influence effect. It has  
14 been theorized that this effect is partly driven by misinformation familiarity. Some  
15 researchers have even argued that a correction should avoid repeating the misinformation, as  
16 the correction itself could serve to inadvertently enhance misinformation familiarity and may  
17 thus backfire, ironically strengthening the very misconception it aims to correct. While  
18 previous research has found little evidence of such familiarity backfire effects, there remains  
19 one situation where they may yet arise: when correcting entirely novel misinformation, where  
20 corrections could serve to spread misinformation to new audiences who had never heard of it  
21 before. This article presents three experiments (total  $N = 1,718$ ) investigating the possibility  
22 of familiarity backfire within the context of correcting novel misinformation claims and after  
23 a one-week study-test delay. While there was variation across experiments, overall there was  
24 substantial evidence against familiarity backfire. Corrections that exposed participants to  
25 novel misinformation did not lead to stronger misconceptions compared to a control group  
26 never exposed to the false claims or corrections. This suggests that it is safe to repeat  
27 misinformation when correcting it, even when the audience might be unfamiliar with the  
28 misinformation.

29 **Keywords:** Continued influence effect; fact-checking; myth debunking; familiarity backfire  
30 effect; illusory truth effect; mere exposure effect

31 Significance statement: Misinformation often continues to influence people’s thinking and  
32 decision making even after they have received clear, credible corrections; this is known as the  
33 continued influence effect. It has been suggested that this effect is partly driven by the  
34 familiarity of false claims, such that people are particularly influenced by false claims that  
35 seem especially familiar (“I have heard that before, so there must be something to it!”). Some  
36 researchers have even recommended that a correction should avoid repeating the  
37 misinformation, out of concerns that the correction itself could inadvertently enhance the  
38 familiarity of the false claim. This could lead to corrections backfiring, ironically  
39 strengthening the very misconceptions they aim to correct. While previous research has found  
40 little evidence of such familiarity backfire effects, there remains one situation where they  
41 may yet arise: when correcting entirely novel misinformation. Such corrections might  
42 familiarize people with false claims they had never encountered before, and therefore, such  
43 corrections could serve to spread misinformation to new audiences. This article presents three  
44 online experiments (total  $N = 1,718$  participants) investigating the possibility of familiarity  
45 backfire within the context of correcting novel misinformation claims. While there was some  
46 variation across experiments, overall there was substantial evidence *against* familiarity  
47 backfire: Corrections that exposed participants to novel misinformation did not lead to  
48 stronger misconceptions compared to a control group never exposed to the false claims or  
49 corrections. This suggests that it is safe to repeat misinformation when correcting it, even  
50 when the audience might be unfamiliar with the misinformation.

51 Can Corrections Spread Misinformation to New Audiences? Testing for the Elusive  
52 Familiarity Backfire Effect

53 The advent of the internet and the subsequent rise of social media as a primary form  
54 of communication has facilitated the distribution of misinformation at unprecedented levels  
55 (Southwell & Thorson, 2015; Vargo, Guo, & Amazeen, 2018). Misinformation can have  
56 detrimental effects at a societal and individual level, as ill-informed decisions can have  
57 negative economic, social, and health-related consequences (Bode & Vraga, 2018; Lazer et  
58 al., 2018; Lewandowsky, Ecker, & Cook, 2017; MacFarlane, Hurlstone, & Ecker, 2020;  
59 Southwell & Thorson, 2015). This is concerning because there is a significant disparity  
60 between the ease of disseminating misinformation and the difficulty of correcting it.

61 Corrections can be ineffective, and individuals often continue to use corrected  
62 misinformation in their inferential reasoning, a phenomenon termed the continued influence  
63 effect (Chan, Jones, Hall Jamieson, & Albarracín, 2017; Johnson & Seifert, 1994;  
64 Lewandowsky, Ecker, Seifert, Schwarz, & Cook, 2012; Paynter et al., 2019; Rich &  
65 Zaragoza, 2016; Walter & Tukachinsky, 2020; Wilkes & Leatherbarrow, 1988).

66 One theoretical account of the continued influence effect assumes that it results from  
67 selective retrieval (Ecker, Lewandowsky, Swire, & Chang, 2011; Ecker, Lewandowsky, &  
68 Tang, 2010; Gordon, Quadflieg, Brooks, Ecker, & Lewandowsky, 2019; Swire, Ecker, &  
69 Lewandowsky, 2017). More specifically, in line with dual-processing models of memory  
70 (e.g., Diana, Reder, Arndt, & Park, 2006; Yonelinas & Jacoby, 2012; Zimmer & Ecker,  
71 2010), continued influence effects might arise when a reasoning task features a retrieval cue  
72 that automatically activates the misinformation, while recollection of the correction fails (see  
73 Ayers & Reder, 1998; Marsh & Fazio, 2006). According to this account, automatic  
74 misinformation activation is driven by familiarity, an automatic process that facilitates the  
75 rapid, context-free retrieval of previously encountered stimuli, whereby the degree of

76 activation of a memory representation depends upon the frequency with which the associated  
77 stimulus has been encountered in the past (Hintzman & Curran, 1994).

78         It follows that one driver of the continued influence effect may lie in the fact that  
79 misinformation is typically repeated within a correction, boosting its familiarity—for  
80 example, clarifying that vaccines do *not* cause autism all but requires repetition of the false  
81 vaccine-autism association (e.g., see Nyhan, Reifler, Richey, & Freed, 2014; Paynter et al.,  
82 2019). Apart from the fact that enhanced familiarity will facilitate automatic misinformation  
83 retrieval, familiarity has also been found to foster perceived truthfulness via metacognitive  
84 processes (Begg, Anas, & Farinacci, 1992; Dechêne, Stahl, Hansen, & Wänke, 2010; Parks &  
85 Toth, 2006)—either because enhanced familiarity indicates greater social consensus (Weaver,  
86 Garcia, Schwarz, & Miller, 2007; also see Arkes, Boehm, & Xu, 1991) or because familiar  
87 information is processed more fluently and the perceived fluency is misattributed to the  
88 information’s validity (Pennycook, Cannon, & Rand, 2018; Schwarz, Sanna, Skurnik, &  
89 Yoon, 2007; Unkelbach, 2007). Thus, corrections that repeat the misinformation might  
90 inadvertently increase the likelihood of it being retrieved and perceived as valid in  
91 subsequent reasoning tasks (Schwarz et al., 2007; Swire et al., 2017).

92         It has even been suggested that the boost in familiarity associated with the repetition  
93 of misinformation within a correction could be so detrimental that it could ironically *increase*  
94 belief in the corrected misinformation (Schwarz et al., 2007). This increase in post-correction  
95 belief in misinformation, relative to either a pre-correction baseline in the same sample of  
96 participants, or a no-misinformation-exposure baseline in a separate sample, has been termed  
97 the familiarity backfire effect (Cook & Lewandowsky, 2011; Lewandowsky et al., 2012). In  
98 order to avoid this effect, it is commonly suggested to educators, journalists, and science  
99 communicators that corrections should avoid repeating the targeted misinformation as much

100 as possible (Cook & Lewandowsky, 2011; Lewandowsky et al., 2012; Peter & Koch, 2016;  
101 Schwarz, Newman, & Leach, 2016; Schwarz et al., 2007).

102         However, despite familiarity backfire effects being prominently discussed in the  
103 literature, empirical evidence of such effects is scarce. In fact, the only clear demonstration of  
104 a familiarity backfire effect was reported in an unpublished manuscript by Skurnik, Yoon,  
105 and Schwarz (2007; discussed by Schwarz et al., 2007), who presented participants with a  
106 flyer juxtaposing “myths vs. facts” associated with the flu vaccine. It was found that after a  
107 30-min delay, a substantial proportion of myths were misremembered as facts, and that  
108 attitudes towards the flu vaccine became more negative compared to participants who had not  
109 been presented with the flyer. In a similar study, Skurnik, Yoon, Park, and Schwarz (2005)  
110 found that participants were more likely to misremember myths as facts after repeated versus  
111 singular retractions. However, these effects were only found with a three-day test delay and  
112 only in older adults (not after shorter delays and in younger adults, as in Skurnik et al., 2007),  
113 and the study also did not feature a baseline condition against which to assess actual  
114 “backfire”.

115         By contrast, a number of contemporary studies have failed to find evidence of  
116 familiarity backfire effects. For example, unlike Skurnik et al. (2005), Ecker et al. (2011)  
117 found that multiple retractions were more effective than singular retractions at reducing  
118 continued influence. Cameron et al. (2013) compared the effectiveness of flu-vaccine myth  
119 corrections that either avoided misinformation repetition (presenting facts only) or repeated  
120 misinformation (including one condition featuring Skurnik et al.’s [2007] “myths vs. facts”  
121 flyer). Flu-vaccine knowledge was measured prior to the manipulation and again after a  
122 week, together with post-intervention belief in the true and false claims. Cameron et al. found  
123 that all conditions were successful at reducing misconceptions, with the best outcomes in the  
124 “myths vs. facts” condition, and the worst outcomes in the facts-only condition that avoided

125 myth repetition. Likewise, Ecker, Hogan, and Lewandowsky (2017) found that repeating a  
126 piece of misinformation when correcting it actually led to stronger reduction of the continued  
127 influence effect than a correction that avoided misinformation repetition. They argued that  
128 misinformation repetition fosters co-activation of the misinformation and its correction,  
129 which in turn facilitates conflict detection and information integration when the correction is  
130 encoded, leading to stronger knowledge revision (see Kendeou, Walsh, Smith, & O'Brien,  
131 2014). Finally, Swire et al. (2017) presented participants with a series of true and false claims  
132 that were subsequently affirmed or corrected and measured the corresponding change in  
133 belief. They, too, failed to observe any familiarity backfire effects: post-correction belief in  
134 misinformation was always lower than pre-correction belief. This reduction in false-claim  
135 belief was observed even under conditions where the impact of familiarity (relative to  
136 recollection) should be maximal, viz. in elderly participants and after a long retention interval  
137 of up to three weeks. Swire et al. concluded that familiarity may contribute to continued  
138 influence effects (i.e., ongoing reliance on corrected misinformation, especially after a delay,  
139 when recollection of the correction fades but familiarity of the misinformation remains  
140 relatively intact; see Knowlton & Squire, 1995), but that misinformation familiarity is not  
141 typically associated with backfire effects (i.e., ironic boosts to false-claim beliefs relative to a  
142 pre-correction or no-exposure baseline).

143 In a recent study, Ecker, O'Reilly, Reid, and Chang (2020) found that presenting  
144 participants with only a correction (a brief retraction or a more detailed refutation) of a real-  
145 world false claim, without prior exposure to the false claim itself, decreased both false-claim-  
146 congruent reasoning and belief in the false claim relative to a control group who received no  
147 exposure to the claim. This demonstrated that mere exposure to a false claim within a  
148 correction did not cause a familiarity backfire effect. However, Ecker et al. highlighted one  
149 remaining situation where a familiarity backfire effect may yet occur: when *novel*

150 misinformation is introduced to a recipient through a correction. If a person's first encounter  
151 with a false claim is provided by a correction, the correction could inadvertently familiarize  
152 the person with the previously unfamiliar misinformation; corrections may thus potentially  
153 spread the misinformation to new audiences (as suggested by Schwarz et al., 2016). Indeed,  
154 the greatest boost to a claim's familiarity will be associated with the initial encounter, while  
155 additional encounters will bring about exponentially decreasing familiarity boosts (consistent  
156 with theoretical frameworks that propose novelty-dependent encoding; e.g., Oberauer,  
157 Lewandowsky, Farrell, Jarrold, & Greaves, 2012).

158         It is easy to see how social media could facilitate situations where an individual is  
159 exposed to a correction without previously having encountered the corresponding  
160 misinformation. Such exposure may not only familiarize the consumer with the novel  
161 misinformation, but may also lend some credibility to the false claim, in the sense that a  
162 correction may signal that someone actually believes the false claim to be true, thus  
163 warranting a correction. This makes the possibility of a familiarity backfire effect with novel  
164 misinformation a concerning notion. Thus, the main purpose of the present study was to  
165 investigate the possibility of a familiarity backfire effect within the context of correcting  
166 novel misinformation. To this end, the study aimed to replicate Ecker et al. (2020), using  
167 claims that were maximally novel to participants.

168         Except for the use of novel false claims, Experiment 1 was a straight replication of the  
169 brief-retraction conditions of Ecker et al. (2020; Experiment 2). Experiments 2 and 3 aimed  
170 to replicate Experiment 1, while manipulating factors that should influence the relative  
171 impact of familiarity, viz. retention interval (Experiment 2) and cognitive load during  
172 encoding (Experiment 3).

## Experiment 1

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

Experiment 1 presented participants with true and false claims and/or associated affirmative or corrective fact-checks. An example claim was “The national animal of Scotland is the unicorn” (see Figure 1 and Method section for further details). The experiment used a  $2 \times 2$  between-subjects design, fully crossing factors claim exposure (yes/no) and fact-check exposure (yes/no).<sup>1</sup> Conditions were no-exposure control (NE), claim-only (CO), fact-check-only (FCO), and claim-plus-fact-check (CFC; in this condition, participants first received all claims without any indication of validity, and then received the fact-checks separately). The experiment was designed to encourage participants to rely on familiarity during retrieval in order to maximize the possibility of observing familiarity-related backfire effects. Fact-checks in fact-check-only and claim-plus-fact-check conditions therefore simply stated the claim with a brief affirmation or correction (e.g., “The national animal of Scotland is the unicorn” followed by the word “TRUE” and a green tick mark; see Figure 2) but did not provide supporting, detailed information, since additional refutational information has been shown to increase the likelihood that the corrective message is later recollected (Chan et al., 2017; Ecker et al., 2020; Paynter et al., 2019; Swire et al., 2017). Additionally, a one-week retention interval between exposure and test was used, as the ability to engage in recollection diminishes over time, whilst familiarity remains relatively constant (Knowlton & Squire, 1995).

Belief in the claims at test was determined by direct claim-belief ratings, as well as a series of inference questions that indirectly measured claim belief by assessing claim-congruent reasoning. The inference questions were presented first because the inference score was determined a priori as the main dependent variable of interest, following ample precedent

---

<sup>1</sup> Technically, the design was a  $2 \times 2 \times 2$  mixed design with the within-subjects factor claim veracity (true/false); however, as the prime interest was on false claims, analyses were conducted separately for true and false claims.

196 (e.g., Ecker et al., 2017). The inference score provides a belief measure that is not  
197 “contaminated” by concurrent exposure to the core claim, whereas it is impossible to measure  
198 direct belief in a claim without at the same time exposing participants to it. Thus, only the  
199 inference score provides a “clean” baseline in the no-exposure condition. Moreover,  
200 presenting the claims for a direct belief rating first would have artificially increased claim  
201 familiarity across all conditions, and acted as a potent retrieval cue for recollection of the  
202 fact-checks. The core hypothesis ( $H1_{FIS}$ )<sup>2</sup> was that we would observe a familiarity backfire  
203 effect, that is, that mere exposure to corrective fact-checks would lead to increased inference  
204 scores relative to the no-exposure baseline (i.e.,  $NE < FCO$ ).

205 A series of secondary hypotheses was specified as follows (these are also  
206 summarized, together with the primary hypothesis, in Table 1 in the Results section):

207 Hypothesis  $H1_{FBR}$  was that mere corrections would also increase false-claim belief  
208 ratings relative to baseline (i.e.,  $NE < FCO$ ). Hypothesis  $H1_{TIS/TBR}$  was that mere affirmations  
209 would be effective and would thus increase inference scores and true-claim belief ratings  
210 relative to baseline (i.e.,  $NE < FCO$ ).

211 Hypothesis 2 investigated the illusory truth effect, whereby mere exposure to  
212 information renders it more likely to be evaluated as truthful (Dechêne et al., 2010). It was  
213 specified that mere exposure to claims would increase claim-congruent reasoning for both  
214 false claims ( $H2_{FIS}$ ) and true claims ( $H2_{TIS}$ ), and boost belief in both false ( $H2_{FBR}$ ) and true  
215 claims ( $H2_{TBR}$ ), relative to baseline (i.e.,  $NE < CO$ ).

216 Hypothesis 3 tested the effectiveness of fact-checking a claim that had already been  
217 encountered. It was specified that, relative to the claims-only condition, fact-checks of  
218 previously presented claims would decrease false-claim-congruent reasoning and false-claim

---

<sup>2</sup> Subscripts FIS, FBR, TIS, TBR will be used to refer to false-claim inference scores, false-claim belief ratings, true-claim inference scores, and true-claim belief ratings, respectively.

219 belief (i.e.,  $CFC < CO$ ;  $H3_{FIS}$  and  $H3_{FBR}$ ), while increasing true-claim-congruent reasoning  
220 and true-claim belief (i.e.,  $CFC > CO$ ;  $H3_{TIS}$  and  $H3_{TBR}$ ).

221 Finally, Hypothesis 4 tested if correcting previously presented false claims would  
222 reduce inference and belief scores back to or even below baseline. This is technically a test  
223 for continued influence, as previous research has found that corrections are often not able to  
224 eliminate misinformation influence down to baseline levels. However, in most continued-  
225 influence studies, the misinformation is initially presented as true and valid, whereas the  
226 initial presentation of false claims in the claim-plus-fact-check condition occurred without  
227 validation (i.e., the false claim was presented initially without being labelled a fact, which  
228 would have presumably increased initial belief, making it harder to subsequently bring belief  
229 back down to baseline). It was therefore not expected that inference scores would be greater  
230 in the claim-plus-fact-check condition than the no-exposure control. In fact, guided by the  
231 results of Ecker et al. (2020), we expected that corrections of previously presented false  
232 claims would decrease false-claim-congruent reasoning ( $H4_{FIS}$ ) and false-claim belief  
233 ( $H4_{FBR}$ ) back to or even below the level of the no-exposure control, and specified  
234 Hypothesis 4 as  $NE > CFC$ .

## 235 **Method**

236 **Participants.** An a-priori power analysis using G\*Power3 (Faul, Erdfelder, Lang, &  
237 Buchner, 2007) indicated that a minimum sample size of 352 was needed to detect a small  
238 effect of  $f = .15$  between two groups with  $\alpha = .05$  and  $1 - \beta = .80$ . In order to account for  
239 attrition rates and ensure sufficient power, it was decided to recruit 440 participants—  
240 however, due to miscommunication, this sample size was used for the entire experiment even  
241 though its calculation was based on only two groups, and thus the experiment was somewhat  
242 underpowered. Participants were U.S.-based adult Amazon Mechanical Turk (MTurk)  
243 workers, who had completed at least 5,000 so-called human-intelligence tasks (HITs) with

244 97%+ approval. MTurk data are largely regarded as being of comparative quality to data  
245 from convenience samples (Berinsky, Huber, & Lenz, 2012; Hauser & Schwarz, 2015;  
246 Necka, Cacioppo, Norman, & Cacioppo, 2016).

247 A subset of 331 participants were randomly assigned to one of the three exposure  
248 conditions (CO, FCO, or CFC) of an experimental survey, with the constraint of  
249 approximately equal cell sizes. The retention rate between study and test was approximately  
250 80%, with 264 participants returning for the test phase. An additional 109 participants  
251 completed the NE control condition, which involved only a test phase and was therefore run  
252 separately (and concurrently with the test phase of the other conditions). Two participants  
253 were identified as erratic responders based on an a-priori exclusion criterion (see Results for  
254 details). The final sample size for analysis was thus  $N = 371$  (condition NE:  $n = 108$ ; CO:  
255  $n = 81$ ; FCO:  $n = 92$ ; CFC:  $n = 90$ ; age range: 20-71 years;  $M_{age} = 39.91$ ;  $SD_{age} = 11.99$ ;  
256 208 males, 160 females, and 3 participants of undisclosed gender). A post-hoc power analysis  
257 confirmed an achieved power in regards to the observed main effect of condition in the  
258 analysis of inference scores ( $\eta_p^2 = .022$ ; see Results below) of  $1 - \beta = .67$ . Participants were  
259 paid US\$0.40 for the study phase and US\$0.60 for the test phase.

## 260 **Materials.**

261 **Claims.** A total of 12 claims (six true, six false) were selected from an initial pool of  
262 48, with the intention of minimizing claim familiarity. To this end, prior to conducting the  
263 present study, the 48 claims were evaluated by a separate sample of  $N = 91$  participants via an  
264 MTurk survey (see Appendix for details). The familiarity and believability of each claim  
265 were rated on Likert scales ranging from 1 (low familiarity/believability) to 5 (high  
266 familiarity/believability). Claims with familiarity ratings  $> 2$  were excluded, as were  
267 excessively believable or unbelievable claims (believability ratings  $< 2$  or  $> 4$ ), resulting in a  
268 pool of 22 candidate claims. From this pool, the least familiar claims were then selected

269 while taking into account additional factors such as comprehensibility and the quality of  
270 corresponding inferential-reasoning questions that could be generated. All claims are  
271 provided in the Appendix. The average familiarity of selected false claims was  $M = 1.67$ ,  
272 with mean believability of  $M = 2.89$ ; average familiarity of selected true claims was  
273  $M = 1.63$ , with mean believability of  $M = 2.54$ .

274 Claims were presented in a format that mimicked a social-media post (see Figure 1).  
275 Each claim was associated with a different fictional account, and was displayed underneath  
276 the account name. A circular image with the first letter of the account handle was displayed  
277 instead of a traditional profile picture, similar to the default icon for a Google account.



278

279 *Figure 1.* Example of a true claim (left) and false claim (right).

280 **Fact-checks.** There were 12 fact-checks matched to the 12 claims; these were  
281 displayed in the same social-media format as the original claim (see Figure 2). Each fact-  
282 check repeated the corresponding claim along with an affirmation (a “TRUE” tag and a green  
283 tick) if the claim was true, or a correction (a “FALSE” tag and a red cross) if it was false. All  
284 fact-checks were associated with the fictional account “Facts First”, which was introduced as  
285 an independent and objective fact-checking group that verifies claims on social media.



286

287 *Figure 2.* Example of an affirmation (left) and correction (right).

288 **Measures.** Claim-related inferential reasoning was measured through a series of 24  
 289 inference questions designed to indirectly assess claim beliefs. There were two such  
 290 questions per claim, one of which was reverse-coded. Each item presented the participants  
 291 with a statement that was related to a claim, but did not repeat the claim itself. Statements  
 292 were designed such that agreeing or disagreeing with them would require reasoning that is  
 293 congruent or incongruent with belief in the original claim. An example item was “Facebook  
 294 is investing money into promoting inoffensive language on its platform”. Participants were  
 295 asked to rate their level of agreement with each statement on a Likert scale ranging from 0  
 296 (complete disagreement) to 10 (complete agreement). Inference questions are provided in the  
 297 Appendix. Claim belief was additionally measured through 12 direct belief ratings.  
 298 Participants were asked to indicate how much they believed each claim to be true or false on  
 299 a Likert scale ranging from 0 (certainly false) to 10 (certainly true).

300 **Procedure.** The experiment was administered using Qualtrics survey software  
 301 (Qualtrics, Provo, UT) via the CloudResearch platform (formerly TurkPrime; Litman,  
 302 Robinson, & Abberbock, 2017). After being presented with an ethically-approved  
 303 information sheet, participants answered demographic questions regarding their English  
 304 language proficiency, gender, age, and country of residence. In the study phase, depending on  
 305 experimental condition, participants read either a series of claims (claim-only condition CO),

306 a series of fact-checks (fact-check only, FCO), or a series of claims followed by a series of  
307 associated fact-checks (claim-plus-fact-check, CFC). All claims and/or fact-checks were  
308 presented individually for at least 3 s. After a one-week retention interval, participants who  
309 completed the study phase were invited by email to participate in the test phase. Participants  
310 in the no-exposure condition (NE) only completed the test phase. In the test phase,  
311 participants were first presented with the 24 inference questions. Inference questions were  
312 grouped by claim (i.e., paired questions were always presented together), but otherwise the  
313 sequence was randomized. Participants then answered the 12 direct belief questions in a  
314 random order. Finally, participants were asked if they had put in a reasonable effort and  
315 whether their data should be used for analysis (with response options “Yes, I put in  
316 reasonable effort”; “Maybe, I was a little distracted”; or “No, I really wasn’t paying any  
317 attention”), before being debriefed.

## 318 **Results**

319 Data from all experiments are available at <https://osf.io/69bq3/>. Before analysis, we  
320 applied a set of a-priori exclusion criteria. Three criteria were not met by any participants,  
321 namely English proficiency self-rated as “poor”, uniform responding (identified by a mean  
322  $SD < 0.5$  across all responses), and self-reported lack of effort (“no” response to the effort  
323 question). To identify erratic responding, we applied the following procedure: After inverting  
324 all reverse-keyed items such that greater inference scores reflected stronger claim-congruent  
325 reasoning, for each claim we calculated the mean absolute difference between the two  
326 inference-question responses (IQ1 and IQ2) and the belief rating (BR) as  $(|IQ1 - IQ2| + |IQ1$   
327  $- BR| + |IQ2 - BR|) / 3$ . The mean absolute differences across all 12 claims were then  
328 averaged to produce a final score, where entirely consistent responding would result in values  
329 approaching zero. This score was then used to identify and reject erratic responders, using the

330 inter-quartile outlier rule with a 2.2 multiplier (Hoaglin & Iglewicz, 1987). As mentioned  
331 earlier, we excluded  $n = 2$  erratic responders based on this procedure.

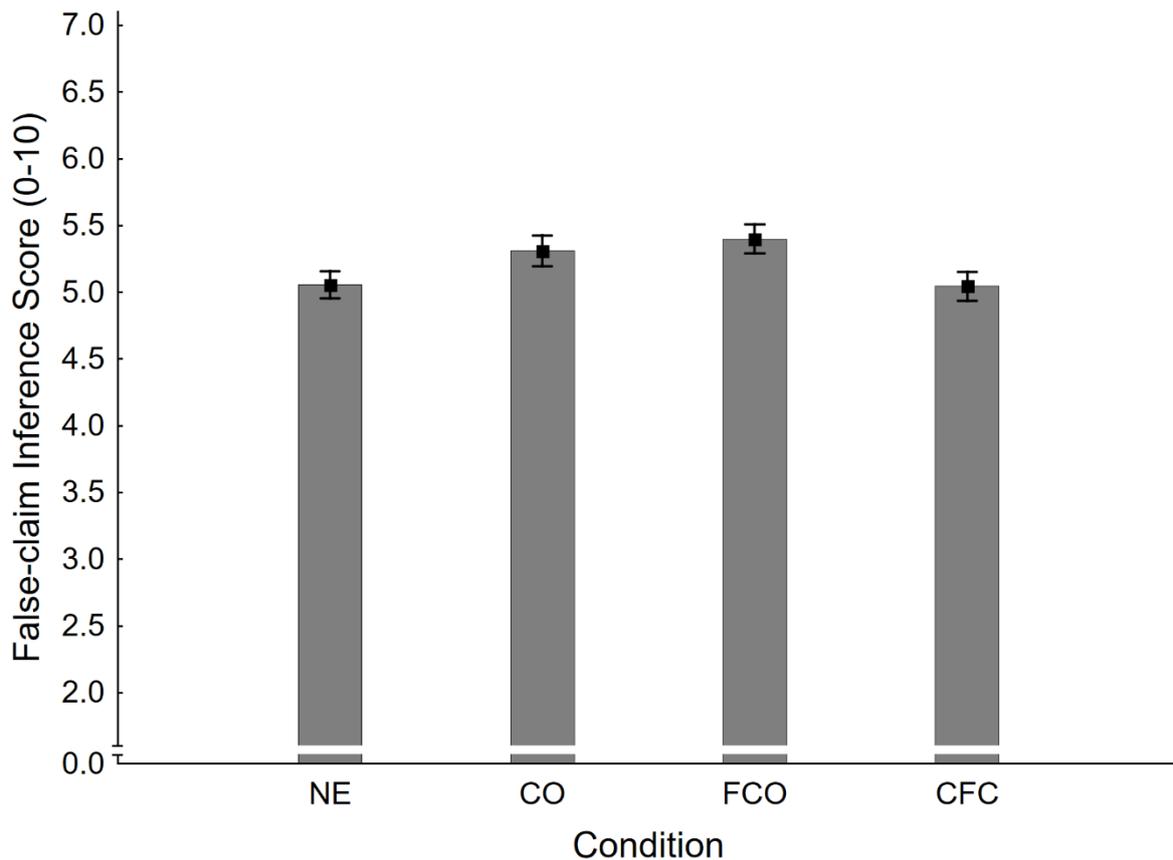
332 Mean false-claim and true-claim inference scores were calculated by averaging the  
333 scores associated with the 12 false-claim and 12 true-claim inference questions, respectively.  
334 Inference scores ranged from 0-10. The belief ratings associated with the six false claims  
335 were averaged to create a false-claim belief rating, and the ratings associated with the six true  
336 claims were averaged to create a true-claim belief rating. The scale was 0-10.

### 337 **False claims.**

338 *False-claim inference scores.* Mean false-claim inference scores across conditions  
339 are shown in Figure 3. A one-way ANOVA revealed a small but significant main effect of  
340 condition,  $F(3,367) = 2.73$ ,  $\eta_p^2 = .022$ ,  $p = .044$ . To test the primary hypothesis that  
341 corrections of novel myths would produce a familiarity backfire effect, a planned contrast  
342 compared no-exposure (NE:  $M = 5.06$ ,  $SE = 0.10$ ) and fact-check-only conditions (FCO:  
343  $M = 5.40$ ,  $SE = 0.11$ ). This contrast was significant,  $F(1,367) = 5.31$ ,  $\eta_p^2 = .014$ ,  $p = .022$ .  
344 Thus, a small familiarity backfire effect was observed, and  $H1_{FIS}$  was supported.

345 Next, three secondary planned contrasts were conducted on false-claim inference  
346 scores, applying the Holm-Bonferroni correction (Holm, 1979). The results of these contrasts  
347 are reported in the first panel of Table 1 (together with the primary contrast). In order to test  
348 for an illusory truth effect, we compared the claim-only (CO:  $M = 5.31$ ,  $SE = 0.12$ ) and no-  
349 exposure conditions. The difference was non-significant, and  $H2_{FIS}$  was rejected accordingly.

350 The effectiveness of correcting a previously encountered false claim was investigated  
351 by contrasting the claim-plus-fact-check (CFC:  $M = 5.04$ ,  $SE = 0.11$ ) and claim-only  
352 conditions. The difference was non-significant, and so  $H3_{FIS}$  was rejected.



353

354 *Figure 3.* Mean false-claim inference scores across conditions NE (no-exposure), CO (claim-  
 355 only), FCO (fact-check-only), and CFC (claim-plus-fact-check) in Experiment 1. Error bars  
 356 show standard errors of the mean.

357

In order to test if correcting previously presented false claims would reduce inference  
 358 scores below baseline, the no-exposure condition was contrasted with the claim-plus-fact-  
 359 check condition. The difference was clearly non-significant, so H<sub>4FIS</sub> was also rejected.

360

***False-claim belief ratings.*** Mean false-claim belief ratings across conditions are

361

shown in Figure 4. A one-way ANOVA found a significant main effect of condition,

362

$F(3,367) = 4.65, \eta_p^2 = .037, p = .003$ . A series of four planned contrasts was then conducted,

363

the results of which are reported in the second panel of Table 1.

364 Table 1

365 *Contrasts Run in Experiment 1*

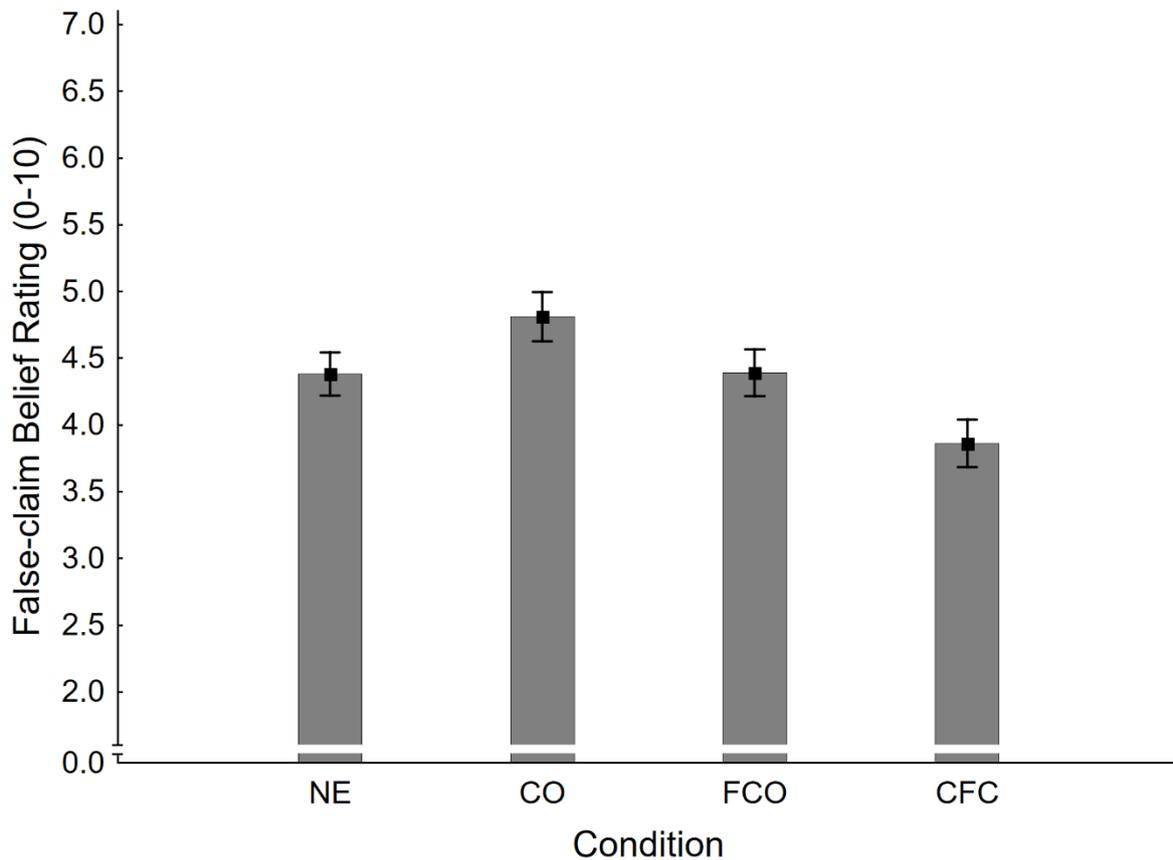
dV / Hypothesis	Effect tested	<i>F</i> (1,367)	<i>p</i>
False-claim inference scores			
<b>H1<sub>FIS</sub>: NE &lt; FCO</b>	<b>Familiarity backfire effect</b>	<b>5.31</b>	<b>.022*</b>
H2 <sub>FIS</sub> : NE < CO	Illusory truth effect	2.72	.100
H3 <sub>FIS</sub> : CFC < CO	Effect of claim+correction vs. claim-only	2.73	.099
H4 <sub>FIS</sub> : NE > CFC	Effect of claim+correction vs. baseline	0.01	.941
False-claim belief ratings			
H1 <sub>FBR</sub> : NE < FCO	Familiarity backfire effect	< 0.01	.971
H2 <sub>FBR</sub> : NE < CO	Illusory truth effect	3.03	.082
H3 <sub>FBR</sub> : CFC < CO	Effect of claim+correction vs. claim-only	13.75	<.001*
H4 <sub>FBR</sub> : NE > CFC	Effect of claim+correction vs. baseline	4.78	.029
True-claim inference scores			
H1 <sub>TIS</sub> : NE < FCO	Effect of affirmation vs. baseline	36.09	<.001*
H2 <sub>TIS</sub> : NE < CO	Illusory truth effect	4.23	.041*
H3 <sub>TIS</sub> : CFC > CO	Effect of claim+affirmation vs. claim-only	9.40	.002*
True-claim belief ratings			
H1 <sub>TBR</sub> : NE < FCO	Effect of affirmation vs. baseline	82.84	<.001*
H2 <sub>TBR</sub> : NE < CO	Illusory truth effect	5.32	.022*
H3 <sub>TBR</sub> : CFC > CO	Effect of claim+affirmation vs. claim-only	30.95	<.001*

366 *Note.* Hypotheses are numbered H1-4 (primary hypothesis in bold; see text for details);  
367 subscripts FIS, TIS, FBR, and TBR refer to false-claim and true-claim inference scores and  
368 belief ratings, respectively. Conditions are NE = no-exposure; CO = claim-only; FCO = fact-  
369 check-only; CFC = claim-plus-fact-check. \* indicates statistical significance (for secondary  
370 contrasts: after Holm-Bonferroni correction).

371 In order to test for a familiarity backfire effect in belief ratings, the no-exposure  
372 condition (NE:  $M = 4.38$ ,  $SE = 0.16$ ) was contrasted with the fact-check-only condition  
373 (FCO:  $M = 4.39$ ,  $SE = 0.17$ ). The difference was non-significant, and thus no additional  
374 evidence for familiarity backfire was obtained; H1<sub>FBR</sub> was rejected.

375 To test for an illusory truth effect, we compared no-exposure to claim-only (CO:  
376  $M = 4.81$ ,  $SE = 0.19$ ) conditions. Belief ratings were numerically higher in the claim-only  
377 condition, but the difference was non-significant; H2<sub>FBR</sub> was therefore rejected.

378 The effectiveness of corrections targeting a previously encountered false claim was  
 379 tested by contrasting the claim-plus-fact-check (CFC:  $M = 3.86$ ,  $SE = 0.18$ ) and claim-only  
 380 conditions. Belief ratings were significantly lower in the claim-plus-fact-check condition,  
 381 supporting  $H3_{FBR}$ .

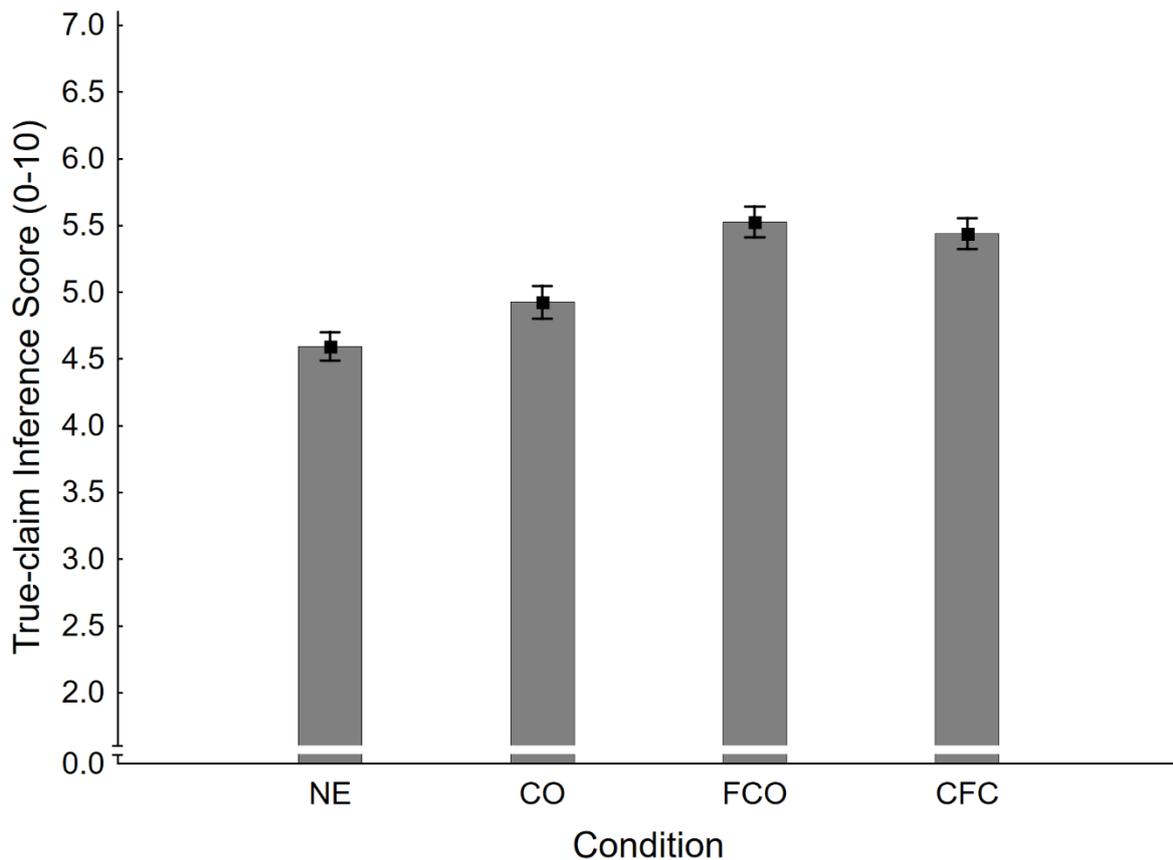


382

383 *Figure 4.* Mean false-claim belief ratings across conditions NE (no-exposure), CO (claim-  
 384 only), FCO (fact-check-only), and CFC (claim-plus-fact-check) in Experiment 1. Error bars  
 385 show standard errors of the mean.

### 386 **True claims.**

387 *True-claim inference scores.* Mean true-claim inference scores across conditions are  
 388 shown in Figure 5. A one-way ANOVA indicated a significant main effect of condition,  
 389  $F(3,367) = 15.92$ ,  $\eta_p^2 = .115$ ,  $p < .001$ . Three planned contrasts tested for specific condition  
 390 differences. Results are reported in the third panel of Table 1.



391

392 *Figure 5.* Mean true-claim inference scores across conditions NE (no-exposure), CO (claim-  
393 only), FCO (fact-check-only), and CFC (claim-plus-fact-check) in Experiment 1. Error bars  
394 show standard errors of the mean.

395

To test if mere affirmations increased inference scores relative to baseline, we  
396 compared no-exposure (NE:  $M = 4.59$ ,  $SE = 0.11$ ) and fact-check-only (FCO:  $M = 5.53$ ,  
397  $SE = 0.11$ ) conditions. This was a highly significant difference, so H1<sub>TIS</sub> was supported.

398

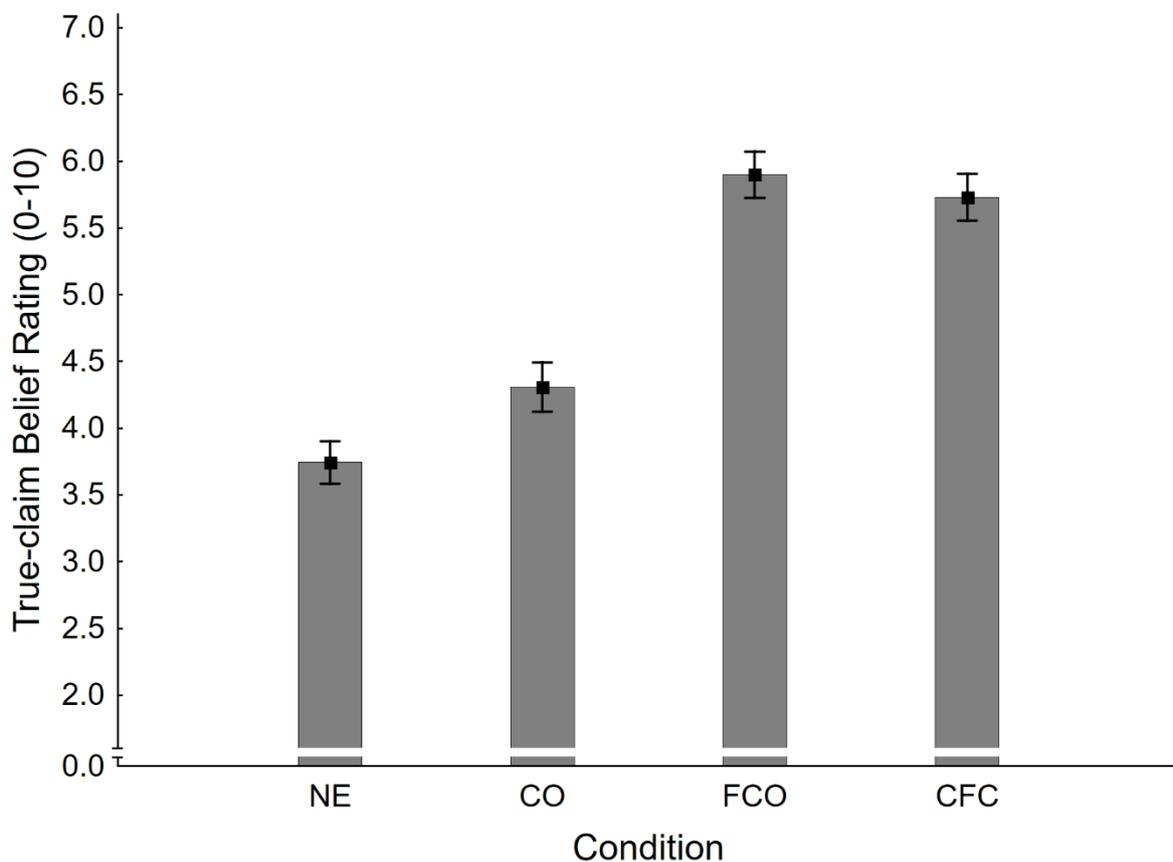
The illusory truth effect was tested for by contrasting no-exposure and claim-only  
399 (CO:  $M = 4.92$ ,  $SE = 0.12$ ) conditions. Inference scores were significantly greater in the  
400 claim-only condition; H2<sub>TIS</sub> was therefore supported.

401

The effectiveness of fact-checks affirming previously encountered claims was  
402 examined by contrasting claim-plus-fact-check (CFC:  $M = 5.44$ ,  $SE = 0.12$ ) and claim-only  
403 conditions. Inference scores were found to be significantly greater in the claim-plus-fact-  
404 check condition, so H3<sub>TIS</sub> was also supported.

405 To test the effectiveness of correcting a previously presented false claim relative to  
406 baseline, the claim-plus-fact-check condition was compared to no-exposure control. Belief  
407 ratings were numerically lower in the claim-plus-fact-check condition, but the contrast was  
408 non-significant after correcting for multiple tests;  $H_{4FBR}$  was thus rejected.

409 ***True-claim belief ratings.*** Mean true-claim belief ratings across conditions are shown  
410 in Figure 6. A one-way ANOVA returned a significant main effect of condition,  
411  $F(3,367) = 39.08$ ,  $\eta_p^2 = .242$ ,  $p < .001$ . Three planned contrasts were performed; results are  
412 presented in the fourth panel of Table 1.



413

414 *Figure 6.* Mean true-claim belief ratings across conditions NE (no-exposure), CO (claim-  
415 only), FCO (fact-check-only), and CFC (claim-plus-fact-check) in Experiment 1. Error bars  
416 show standard errors of the mean.

417 To test the effectiveness of a mere affirmation relative to baseline, we compared the  
418 no-exposure condition (NE:  $M = 3.74$ ,  $SE = 0.16$ ) with the fact-check-only condition (FCO:  
419  $M = 5.90$ ,  $SE = 0.17$ ). The difference was found to be highly significant, supporting  $H1_{TBR}$ .

420 To test for an illusory truth effect, we contrasted no-exposure and claim-only (CO:  
421  $M = 4.31$ ,  $SE = 0.19$ ) conditions. Belief ratings were significantly higher in the claim-only  
422 condition, providing evidence for an illusory truth effect and supporting  $H2_{TBR}$ .

423 Finally, we contrasted the claim-plus-fact-check (CFC:  $M = 5.73$ ,  $SE = 0.18$ ) and  
424 claim-only conditions to test whether an affirmation of a previously presented claim  
425 enhanced belief. Belief was higher in the claim-plus-fact-check condition, and so  $H3_{TBR}$  was  
426 supported.

## 427 **Discussion**

428 Experiment 1 found evidence for a small familiarity backfire effect on inference  
429 scores, supporting Skurnik et al. (2005). After a one-week study-test delay, participants who  
430 were exposed only to the corrective fact-check showed reasoning more in line with the false  
431 claim than participants never exposed to either the claim or the fact-check. This provides  
432 tentative evidence that corrections can backfire and ironically increase misinformed  
433 reasoning when they familiarize people with novel misinformation. However, no familiarity  
434 backfire effect was observed on direct belief ratings, suggesting that exposure to the  
435 previously corrected claim at test may have facilitated recollection of the correction. Given  
436 the small magnitude of the effect on inference scores, we aimed to replicate the result in  
437 Experiment 2 before drawing stronger conclusions; however, to foreshadow, the effect did  
438 not replicate.

439 Furthermore, Experiment 1 provided some additional evidence for illusory truth  
440 effects after just a single exposure (Begg et al., 1992; Dechêne et al., 2010; Pennycook et al.,

441 2018): Participants' claim-congruent reasoning and beliefs were stronger for claims they  
442 were previously exposed to, at least when the claims were actually true.

443 In general, it was found that fact-checks were effective when they targeted a claim  
444 that participants had already encountered before. Relative to the claim-only condition, the  
445 claim-plus-fact-check condition reduced false-claim beliefs and increased true-claim beliefs  
446 as well as true-claim-congruent reasoning (the reduction in false-claim inference scores was  
447 non-significant). These results replicate Ecker et al.'s (2020) finding that fact-checks tended  
448 to be more impactful if participants had previously been exposed to the relevant claim. The  
449 overall pattern also replicates Swire et al. (2017) in that affirmations tended to be more  
450 impactful than corrections, presumably because familiarity and recollection operate in unison  
451 for true claims (both driving acceptance) but stand in opposition with false claims (where  
452 claim familiarity will foster acceptance but correction recollection will drive rejection).  
453 However, correcting previously presented false claims did not reduce inference scores below  
454 the no-exposure baseline (the effect for belief ratings was marginal but non-significant). This  
455 contrasts to some extent with the findings of Ecker et al. (2020), although that study did not  
456 contrast no-exposure and claim-plus-fact-check conditions after a one-week delay. The  
457 absence of a stronger reduction is therefore again best explained by the tension between  
458 familiarity and recollection processes, with the latter more strongly compromised by the  
459 substantial retention interval.

## 460 **Experiment 2**

461 The aim of Experiment 2 was to replicate the familiarity backfire effect found in  
462 Experiment 1. Additionally, Experiment 2 manipulated retention interval, so the test was  
463 either immediate (henceforth indicated by lower-case i) or delayed by one week as in  
464 Experiment 1 (indicated by lower-case d). The rationale for this was that a familiarity  
465 backfire effect should arise only with a delayed test, not an immediate test, when recollection

466 of the correction will still be strong enough to avoid ironic correction effects. Experiment 2  
467 therefore replicated exactly the four experimental conditions of Experiment 1, but added  
468 claim-only, fact-check-only, and claim-plus-fact-check conditions with immediate test; it thus  
469 had a between-subjects design with the sole factor of condition (NE; COi; FCOi; CFCi; COd;  
470 FCOd; CFCd).

471 The design and analysis plan for Experiment 2 were pre-registered  
472 (<https://osf.io/69bq3/registrations>). As in Experiment 1, the core hypothesis regarded the  
473 familiarity backfire effect; it was hypothesized that false-claim inference scores would be  
474 higher in the delayed fact-check-only condition relative to no-exposure control ( $H1_{FISd}$ ;  $NE <$   
475  $FCOd$ ). A related secondary hypothesis was that in the immediate test, there should be no  
476 backfire and indeed a corrective effect ( $H1_{FISi}$ ;  $NE > FCOi$ ).

477 Supplementary hypotheses included the supplementary hypotheses of Experiment 1  
478 (we refrain from repeating these here, but they are specified again in Table 2); additional  
479 supplementary hypotheses were formulated regarding the effects of the delay manipulation  
480 on scores in the fact-check-only (H5) and claim-plus-fact-check (H6) conditions. It was  
481 assumed that significant forgetting would occur over time, implying that false-claim  
482 inference scores and belief ratings would be lower in the immediate fact-check-only (FCOi)  
483 and claim-plus-fact-check (CFCi) conditions than the respective delayed conditions ( $H5_{FIS}$   
484 and  $H5_{FBR}$ ; see Table 2), and that true-claim inference scores and belief ratings would be  
485 higher in the immediate fact-check-only (FCOi) and claim-plus-fact-check (CFCi) conditions  
486 than the respective delayed conditions ( $H5_{TIS}$  and  $H5_{TBR}$ ; see Table 2).

## 487 **Method**

488 **Participants.** A power analysis indicated that to detect an effect of the size observed  
489 in Experiment 1 (main effect of condition on false-claim inference scores,  $\eta_p^2 = .022$ ) with  
490  $\alpha = 0.05$  and  $1 - \beta = 0.80$  across the four replicated conditions would require a minimum

491 sample size of  $n = 123$  per condition. In Experiment 1, the lowest retention of any of the  
492 conditions was  $81/110 = 73.63\%$  (condition CO). It was thus decided to recruit  $n = 170$   
493 participants per condition in the delayed-test conditions and  $n = 130$  participants in the  
494 immediate-test conditions and the no-exposure condition, in the hope of achieving a test-  
495 phase sample size of  $n \approx 130$  per condition (i.e., total  $N = 3 \times 170 + 4 \times 130 = 1,030$ ).  
496 Participants were U.S.-based adult MTurk workers who had completed at least 5,000 HITs  
497 with 97%+ approval. Participants who had completed Experiment 1 were excluded from  
498 participation. The delayed-test conditions, the immediate-test conditions, and the no-exposure  
499 condition were again run separately due to differences in instructions and reimbursements,  
500 with random condition assignment in the delayed and immediate surveys. The immediate-test  
501 and no-exposure conditions were run concurrently with the delayed test; participants were not  
502 able to complete more than one condition.

503 A subset of 509 participants were randomly assigned to one of the three delayed-test  
504 conditions, with the constraint of approximately equal cell sizes. The retention rate between  
505 study and test was approximately 84%, with 427 participants returning for the test phase. An  
506 additional 521 participants completed the immediate-test and NE conditions. Nine  
507 participants were excluded based on a-priori criteria (see Results for details). The final  
508 sample size for analysis was thus  $N = 939$  (condition NE:  $n = 128$ ; COi:  $n = 129$ ; FCOi:  
509  $n = 129$ ; CFCi:  $n = 129$ ; COd:  $n = 140$ ; FCOd:  $n = 144$ ; CFCd:  $n = 140$ ; age range: 20-81  
510 years;  $M_{\text{age}} = 41.35$ ;  $SD_{\text{age}} = 11.97$ ; 469 males, 467 females, and 3 participants of undisclosed  
511 gender). Participants were paid US\$0.40 for the study phase and US\$0.60 for the test phase.

512 **Materials.** Claims, measures, and procedure were identical to Experiment 1, except  
513 that Experiment 2 also contained an immediate test, where participants just completed a  
514 1-min word puzzle between study and test.

**Results**

Before analysis, we applied a set of a-priori (pre-registered) exclusion criteria. Two criteria were not met by any participants, including English proficiency self-rated as “poor”, and self-reported lack of effort. Uniform responding and erratic responding were identified as in Experiment 1, which led to the exclusion of  $n = 5$  and  $n = 4$  participants, respectively.

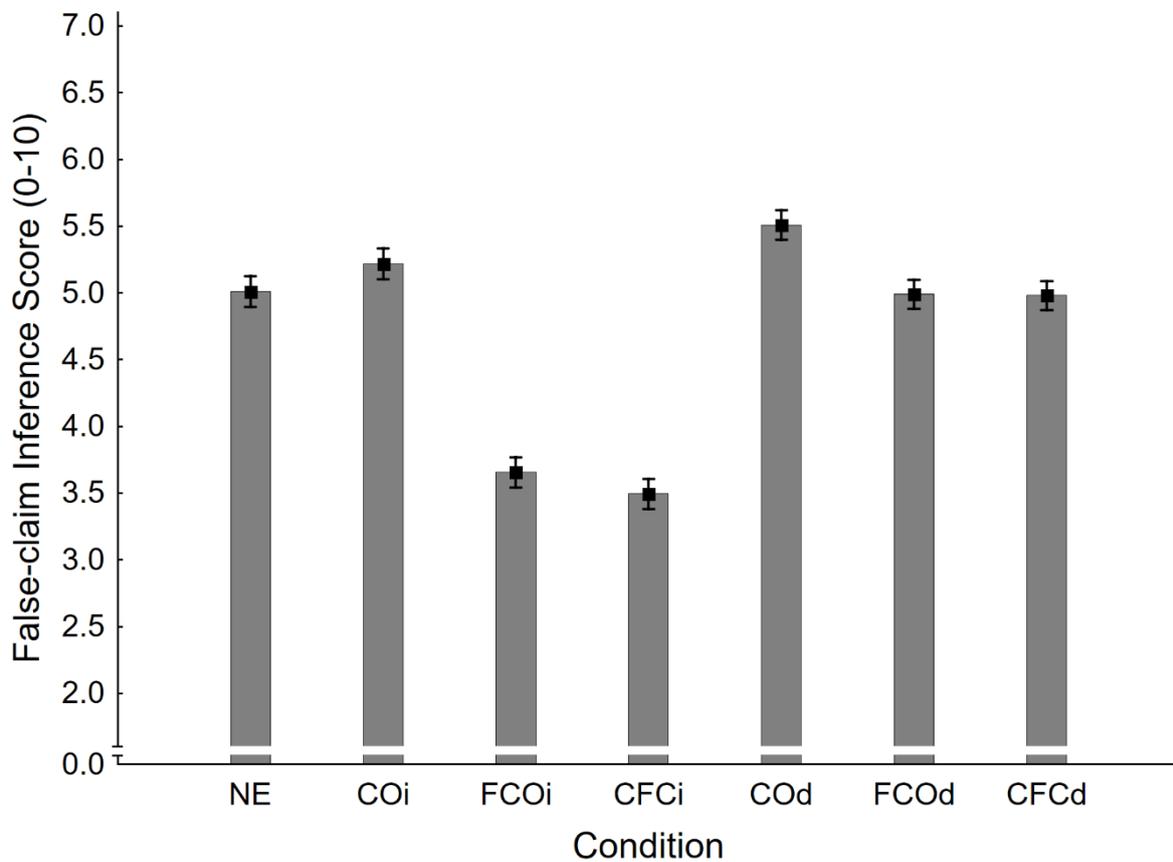
Inference and belief scores were calculated as in Experiment 1.

**False claims.**

*False-claim inference scores.* Mean false-claim inference scores across conditions are shown in Figure 7. A one-way ANOVA revealed a significant main effect of condition,  $F(6,932) = 48.66$ ,  $\eta_p^2 = .239$ ,  $p < .001$ . To test the primary hypothesis that corrections of novel myths would produce a familiarity backfire effect, a planned contrast compared no-exposure (NE:  $M = 5.01$ ,  $SE = 0.11$ ) and delayed fact-check-only (FCOd:  $M = 4.99$ ,  $SE = 0.11$ ) conditions. This was clearly non-significant,  $F(1,932) = 0.02$ ,  $\eta_p^2 < .001$ ,  $p = .895$ . Thus, no familiarity backfire effect was observed, and  $H1_{FISd}$  was not supported. However, the inference score in the immediate fact-check-only condition (FCOi:  $M = 3.65$ ,  $SE = 0.11$ ) was significantly lower than no-exposure control, supporting secondary hypothesis  $H1_{FISi}$ .

Next, the supplementary planned contrasts were conducted on false-claim inference scores. Results are reported in the first panel of Table 2 (together with the primary contrast). To summarize, we found evidence of an illusory truth effect in the delayed (COd:  $M = 5.51$ ,  $SE = 0.11$ ) but not the immediate test (COi:  $M = 5.22$ ,  $SE = 0.11$ ), rejecting  $H2_{FISi}$  and supporting  $H2_{FISd}$ . Corrections of previously presented false claims (CFCi:  $M = 3.49$ ,  $SE = 0.11$ ; CFCd:  $M = 4.98$ ,  $SE = 0.11$ ) were found effective relative to the claim-only condition at both delays (supporting  $H3_{FISi}$  and  $H3_{FISd}$ ). However, compared against the no-exposure baseline, corrections of previously presented false claims were effective immediately but not after a delay (supporting  $H4_{FISi}$  and rejecting  $H4_{FISd}$ ). As expected, the delay had a significant

540 impact on correction effectiveness in both fact-check-only and claim-plus-fact-check  
 541 conditions (supporting H5<sub>FIS</sub> and H6<sub>FIS</sub>).



542

543 *Figure 7.* Mean false-claim inference scores across conditions NE (no-exposure), COi/d  
 544 (claim-only, immediate/delayed test), FCOi/d (fact-check-only, immediate/delayed test), and  
 545 CFCi/d (claim-plus-fact-check, immediate/delayed test) in Experiment 2. Error bars show  
 546 standard errors of the mean.

547 ***False-claim belief ratings.*** Mean false-claim belief ratings across conditions are  
 548 shown in Figure 8. A one-way ANOVA revealed a significant main effect of condition,  
 549  $F(6,932) = 56.14, \eta_p^2 = .265, p < .001$ . Planned contrasts were run to test specific hypotheses;  
 550 results are provided in the second panel of Table 2.

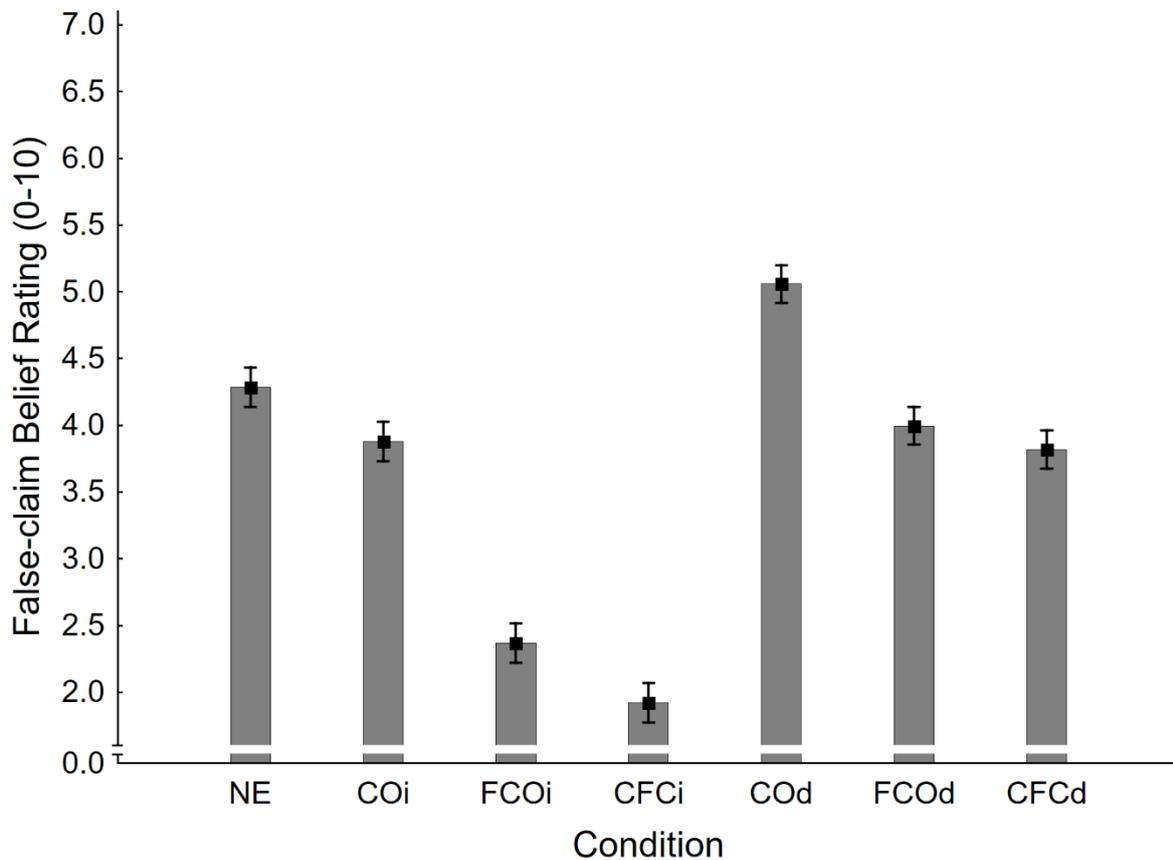
551 Table 2

552 *Contrasts Run in Experiment 2*

dV / Hypothesis	Effect tested	<i>F</i> (1,932)	<i>p</i>
False-claim inference scores			
H1 <sub>FISi</sub> : NE > FCOi	Effect of correction vs. baseline	70.33	<.001*
H2 <sub>FISi</sub> : NE < COi	Illusory truth effect	1.63	.202
H3 <sub>FISi</sub> : CFCi < COi	Effect of claim+correction vs. claim-only	114.14	<.001*
H4 <sub>FISi</sub> : NE > CFCi	Effect of claim+correction vs. baseline	88.08	<.001*
<b>H1<sub>FISd</sub>: NE &lt; FCOd</b>	<b>Familiarity backfire effect</b>	<b>0.02</b>	<b>.895</b>
H2 <sub>FISd</sub> : NE < COd	Illusory truth effect	9.89	.002*
H3 <sub>FISd</sub> : CFCd < COd	Effect of claim+correction vs. claim-only	11.64	<.001*
H4 <sub>FISd</sub> : NE > CFCd	Effect of claim+correction vs. baseline	0.04	.850
H5 <sub>FIS</sub> : FCOi < FCOd	Delay effect on correction	72.21	<.001*
H6 <sub>FIS</sub> : CFCi < CFCd	Delay effect on claim+correction	88.43	<.001*
False-claim belief ratings			
H1 <sub>FBRi</sub> : NE > FCOi	Effect of correction vs. baseline	83.13	<.001*
H2 <sub>FBRi</sub> : NE < COi	Illusory truth effect	3.75	.053 <sup>^</sup>
H3 <sub>FBRi</sub> : CFCi < COi	Effect of claim+correction vs. claim-only	87.36	<.001*
H4 <sub>FBRi</sub> : NE > CFCi	Effect of claim+correction vs. baseline	126.87	<.001*
H1 <sub>FBRd</sub> : NE < FCOd	Familiarity backfire effect	2.02	.155
H2 <sub>FBRd</sub> : NE < COd	Illusory truth effect	14.09	<.001*
H3 <sub>FBRd</sub> : CFCd < COd	Effect of claim+correction vs. claim-only	37.91	<.001*
H4 <sub>FBRd</sub> : NE > CFCd	Effect of claim+correction vs. baseline	5.12	.024
H5 <sub>FBR</sub> : FCOi < FCOd	Delay effect on correction	63.33	<.001*
H6 <sub>FBR</sub> : CFCi < CFCd	Delay effect on claim+correction	85.50	<.001*
True-claim inference scores			
H1 <sub>TISi</sub> : NE < FCOi	Effect of affirmation vs. baseline	148.92	<.001*
H2 <sub>TISi</sub> : NE < COi	Illusory truth effect	2.21	.137
H3 <sub>TISi</sub> : CFCi > COi	Effect of claim+affirmation vs. claim-only	121.24	<.001*
H1 <sub>TISd</sub> : NE < FCOd	Effect of affirmation vs. baseline	18.24	<.001*
H2 <sub>TISd</sub> : NE < COd	Illusory truth effect	0.19	.666
H3 <sub>TISd</sub> : CFCd > COd	Effect of claim+affirmation vs. claim-only	15.07	<.001*
H5 <sub>TIS</sub> : FCOi > FCOd	Delay effect on affirmation	23.51	<.001*
H6 <sub>TIS</sub> : CFCi > CFCd	Delay effect on claim+affirmation	72.57	<.001*
True-claim belief ratings			
H1 <sub>TBRi</sub> : NE < FCOi	Effect of affirmation vs. baseline	45.71	<.001*
H2 <sub>TBRi</sub> : NE < COi	Illusory truth effect	3.79	.052 <sup>^</sup>
H3 <sub>TBRi</sub> : CFCi > COi	Effect of claim+affirmation vs. claim-only	108.75	<.001*
H1 <sub>TBRd</sub> : NE < FCOd	Effect of affirmation vs. baseline	57.16	<.001*
H2 <sub>TBRd</sub> : NE < COd	Illusory truth effect	6.49	.011*
H3 <sub>TBRd</sub> : CFCd > COd	Effect of claim+affirmation vs. claim-only	25.31	<.001*
H5 <sub>TBR</sub> : FCOi > FCOd	Delay effect on affirmation	0.38	.536
H6 <sub>TBR</sub> : CFCi > CFCd	Delay effect on claim+affirmation	1.37	.243

553 *Note.* Hypotheses are numbered H1-6 (primary hypothesis in bold; see text for details);  
554 subscripts FISi/d, TISi/d, FBRi/d, and TISi/d refer to false-claim and true-claim inference  
555 scores and belief ratings in immediate and delayed tests, respectively. Conditions are NE =  
556 no-exposure; COi/d = claim-only with immediate/delayed test; FCOi/d = fact-check-only

557 with immediate/delayed test; CFCi/d = claim-plus-fact-check with immediate/delayed test.  
 558 \* indicates statistical significance after Holm-Bonferroni correction. ^ indicates an effect in  
 559 the opposite of hypothesized direction.



560

561 *Figure 8.* Mean false-claim belief ratings across conditions NE (no-exposure), COi/d (claim-  
 562 only, immediate/delayed test), FCOi/d (fact-check-only, immediate/delayed test), and CFCi/d  
 563 (claim-plus-fact-check, immediate/delayed test) in Experiment 2. Error bars show standard  
 564 errors of the mean.

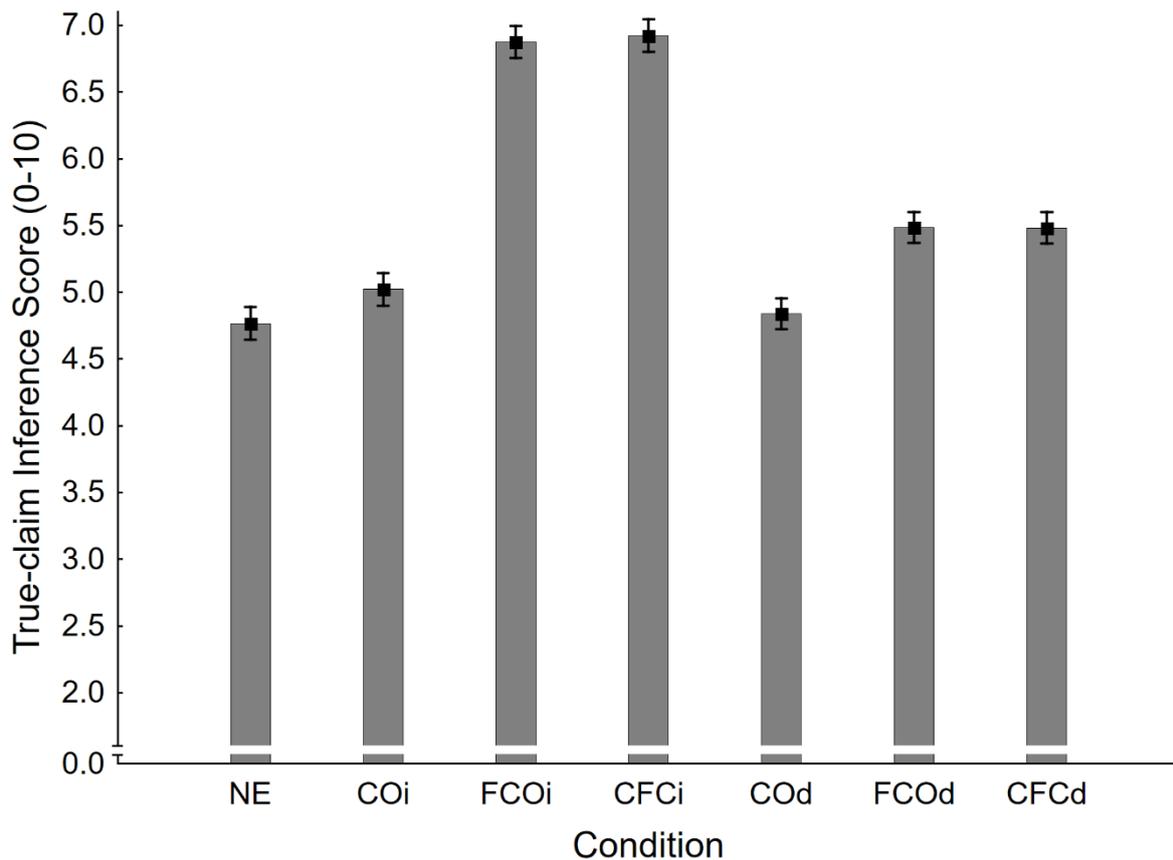
565 There was no evidence of familiarity backfire in belief ratings, as the delayed fact-  
 566 check-only condition (FCOd:  $M = 3.99$ ,  $SE = 0.14$ ) did not differ significantly from no-  
 567 exposure control (NE:  $M = 4.29$ ,  $SE = 0.15$ );  $H1_{FBRd}$  was thus rejected. However, a mere  
 568 correction was effective in the immediate test (FCOi:  $M = 2.37$ ,  $SE = 0.15$ ), supporting  
 569  $H1_{FBRi}$ . There was mixed evidence regarding illusory truth effects, with no-exposure differing  
 570 significantly from the claim-only condition in the delayed test (COd:  $M = 5.06$ ,  $SE = 0.14$ )  
 571 but not the immediate test (COi:  $M = 3.88$ ,  $SE = 0.15$ ), supporting  $H2_{FBRd}$  but rejecting

572 H2<sub>FBRi</sub>. Corrections of previously presented false claims (CFCi:  $M = 1.92$ ,  $SE = 0.15$ ; CFCd:  
573  $M = 3.82$ ,  $SE = 0.14$ ) were found effective relative to the claim-only condition at both delays  
574 (supporting H3<sub>FBRi</sub> and H3<sub>FBRd</sub>). However, mirroring the inference-score results, compared  
575 against the no-exposure baseline, corrections of previously presented false claims were  
576 effective immediately but not after a delay (supporting H4<sub>FBRi</sub> and rejecting H4<sub>FBRd</sub>). Delay  
577 again had a significant impact on correction effectiveness in both fact-check-only and claim-  
578 plus-fact-check conditions (supporting H5<sub>FBR</sub> and H6<sub>FBR</sub>).

579 **True claims.**

580 *True-claim inference scores.* Mean true-claim inference scores across conditions are  
581 shown in Figure 9. A one-way ANOVA indicated a significant main effect of condition,  
582  $F(6,932) = 56.62$ ,  $\eta_p^2 = .267$ ,  $p < .001$ . Planned contrasts tested for specific condition  
583 differences; results are reported in the third panel of Table 2.

584 It was found that a mere affirmation increased inference scores relative to the no-  
585 exposure baseline (NE:  $M = 4.77$ ,  $SE = 0.12$ ) in both immediate (FCOi:  $M = 6.87$ ,  $SE = 0.12$ )  
586 and delayed (FCOd:  $M = 5.48$ ,  $SE = 0.12$ ) tests, supporting H1<sub>TISi</sub> and H1<sub>TISd</sub>. There was no  
587 evidence for illusory truth effects, with no significant difference between claim-only and no-  
588 exposure conditions in either the immediate (COi:  $M = 5.02$ ,  $SE = 0.12$ ) or delayed (COd:  
589  $M = 4.84$ ,  $SE = 0.12$ ) test; H2<sub>TISi</sub> and H2<sub>TISd</sub> were thus rejected. Affirmations of previously  
590 presented true claims (CFCi:  $M = 6.92$ ,  $SE = 0.12$ ; CFCd:  $M = 5.48$ ,  $SE = 0.12$ ) were found  
591 effective relative to the claim-only condition at both delays (supporting H3<sub>TISi</sub> and H3<sub>TISd</sub>).  
592 Again, delay had a significant impact on affirmation effectiveness in both fact-check-only  
593 and claim-plus-fact-check conditions (supporting H5<sub>TIS</sub> and H6<sub>TIS</sub>).



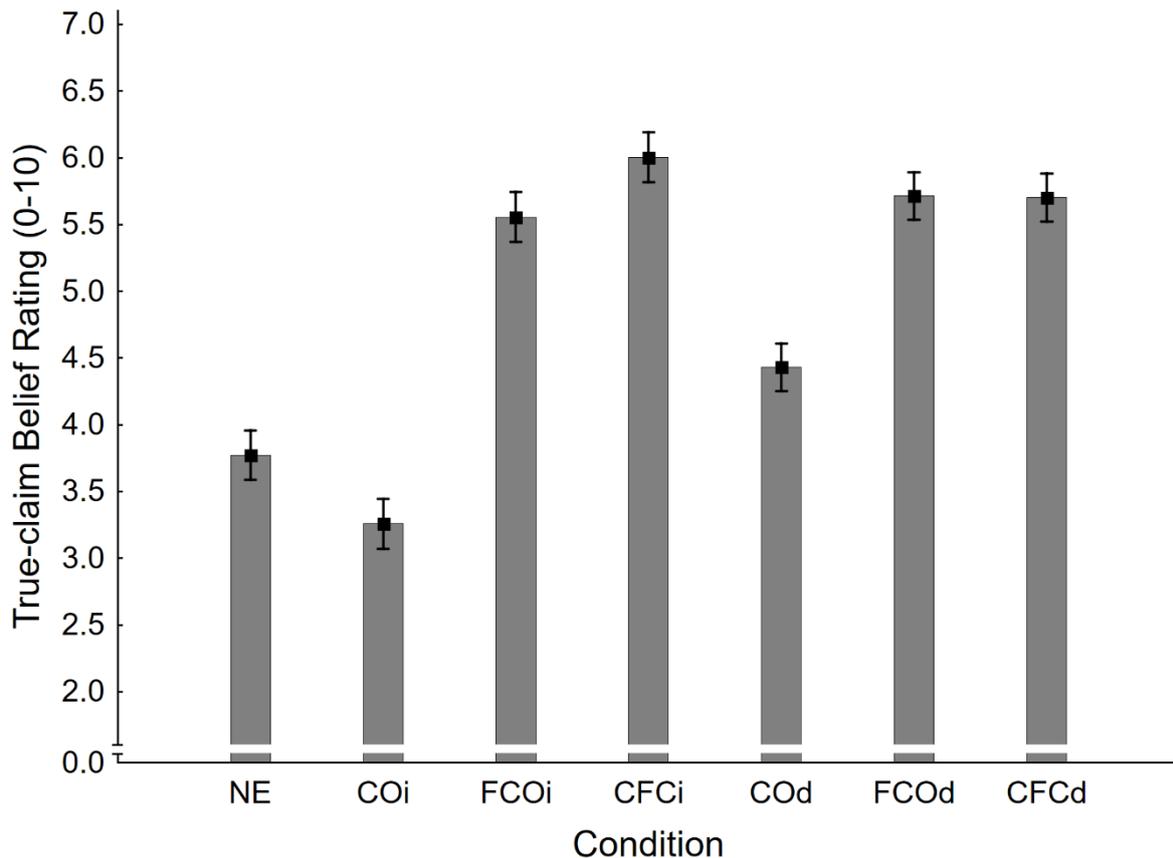
594

595 *Figure 9.* Mean true-claim inference scores across conditions NE (no-exposure), COi/d  
 596 (claim-only, immediate/delayed test), FCOi/d (fact-check-only, immediate/delayed test), and  
 597 CFCi/d (claim-plus-fact-check, immediate/delayed test) in Experiment 2. Error bars show  
 598 standard errors of the mean.

599 ***True-claim belief ratings.*** Mean true-claim belief ratings across conditions are shown  
 600 in Figure 10. A one-way ANOVA returned a significant main effect of condition,  
 601  $F(6,932) = 34.98, \eta_p^2 = .184, p < .001$ . Planned contrasts tested for specific condition  
 602 differences; results are reported in the fourth panel of Table 2.

603 It was found that a mere affirmation increased belief ratings relative to the no-  
 604 exposure baseline (NE:  $M = 3.77, SE = 0.19$ ) in both immediate (FCOi:  $M = 5.56, SE = 0.19$ )  
 605 and delayed (FCOd:  $M = 5.71, SE = 0.18$ ) tests, supporting  $H1_{TBRi}$  and  $H1_{TBRd}$ . There was  
 606 mixed evidence for illusory truth effects, with a significant difference between claim-only  
 607 and no-exposure conditions in the delayed (COD:  $M = 4.43, SE = 0.18$ ) but not the immediate  
 608 (COi:  $M = 3.26, SE = 0.19$ ) test, supporting  $H2_{TBRd}$  and rejecting  $H2_{TBRi}$ . Affirmations of

609 previously presented true claims (CFCi:  $M = 6.00$ ,  $SE = 0.19$ ; CFCd:  $M = 5.70$ ,  $SE = 0.18$ )  
 610 were found effective relative to the claim-only condition at both delays (supporting H3<sub>TBRi</sub>  
 611 and H3<sub>TBRd</sub>). In contrast to the inference scores, delay had no significant impact on  
 612 affirmation effectiveness in fact-check-only and claim-plus-fact-check conditions (rejecting  
 613 H5<sub>TBR</sub> and H6<sub>TBR</sub>).



614

615 *Figure 10.* Mean true-claim belief ratings across conditions NE (no-exposure), COi/d (claim-  
 616 only, immediate/delayed test), FCOi/d (fact-check-only, immediate/delayed test), and CFCi/d  
 617 (claim-plus-fact-check, immediate/delayed test) in Experiment 2. Error bars show standard  
 618 errors of the mean.

## 619 Discussion

620 The primary aim of Experiment 2 was to replicate the familiarity backfire effect  
 621 observed in Experiment 1. The effect did not replicate; there was no evidence for familiarity  
 622 backfire in either the false-claim inference scores or the false-claim belief scores. This is  
 623 consonant with the results Ecker et al. (2020) obtained with non-novel claims, and suggests

624 that the familiarity boost effected by exposure to a false claim within a correction may be  
625 sufficient to offset the corrective effect of a mere fact-check after a one-week delay (thus  
626 resulting in the observed null effect), but not sufficient to cause ironic misconception-  
627 strengthening effects.

628 Evidence for illusory truth effects was again mixed: False-claim inference scores and  
629 belief ratings, as well as true-claim belief ratings, were greater in the claim-only condition  
630 compared to the no-exposure baseline in a delayed test. This stands in contrast to Experiment  
631 1, where illusory truth effects were found only for true claims. Given that participants were  
632 unable to reliably differentiate between true and false claim prior to fact-checks being  
633 provided, we suspect that the best explanation for the overall pattern is that illusory truth  
634 effects after a single exposure are small, and whether or not a statistically significant effect is  
635 obtained is partially down to random variation. There were no significant illusory truth  
636 effects in the immediate test, suggesting that illusory truth effects may be delay-dependent  
637 and thus occur only if memory is relatively more reliant on familiarity.

638 As in Experiment 1, fact-checks were generally effective when they targeted a claim  
639 that participants had already encountered before. Relative to the claim-only condition, the  
640 claim-plus-fact-check condition reduced false-claim beliefs and false-claim-congruent  
641 reasoning and increased true-claim beliefs and true-claim-congruent reasoning across both  
642 retention intervals. This again replicates Ecker et al.'s (2020) finding that fact-checks are  
643 more impactful if participants had previously been exposed to the relevant claim. However,  
644 replicating Experiment 1, correcting previously presented false claims did not reduce  
645 inference scores or belief ratings below the no-exposure baseline after a delay. This is again  
646 best explained by the fact that familiarity and recollection processes stand in opposition when  
647 it comes to delayed appraisals of corrected false claims. Additional support for this  
648 theoretical notion comes from the pattern of delay effects observed: While both fact-check-



674 The design and analysis plan for Experiment 3 were pre-registered  
675 (<https://osf.io/69bq3/registrations>). As in Experiments 1 and 2, the core hypothesis pertained  
676 to the familiarity backfire effect; it was hypothesized that false-claim inference scores would  
677 be higher in the delayed fact-check-only condition under high load than no-exposure control  
678 ( $H1_{FISI+}$ ;  $NE < FCOI+$ ). We also hypothesized that familiarity backfire would occur without  
679 load ( $H1_{FISI-}$ ;  $NE < FCOI-$ ), as in Experiment 1, even though based on Experiment 2 we did  
680 not expect to support this hypothesis.

681 Supplementary hypotheses included some of the supplementary hypotheses of  
682 Experiments 1 and 2; these are not repeated here but specified again in Table 3. Additional  
683 supplementary hypotheses were formulated regarding the effects of the cognitive-load  
684 manipulation on scores in the fact-check-only conditions. It was assumed that load would  
685 reduce correction effects. We therefore expected that false-claim inference scores and belief  
686 ratings would be greater in the load condition than the no-load condition (i.e.,  
687  $FCOI+ > FCOI-$ ;  $H7_{FIS}$  and  $H7_{FBR}$ , respectively; see Table 3), while true-claim inference  
688 scores and belief ratings would be greater in the no-load condition than the load condition  
689 ( $FCOI+ < FCOI-$ ;  $H7_{TIS}$  and  $H7_{TBR}$ , respectively; see Table 3).

## 690 **Method**

691 **Participants.** Participants were U.S.-based adult MTurk workers who had completed  
692 at least 5,000 HITs with 97%+ approval. Participants who had completed Experiment 1 or 2  
693 were excluded from participation. The (two-phase) fact-check-only conditions were again run  
694 separately from the no-exposure condition, with random load-condition assignment within  
695 the fact-check-only conditions. The no-exposure condition was run concurrently with the  
696 delayed fact-check-only test; participants were not able to complete more than one condition.

697 Sampling decisions were guided by the power analysis presented in Experiment 2. A  
698 total of 400 participants were randomly assigned to one of the two fact-check-only

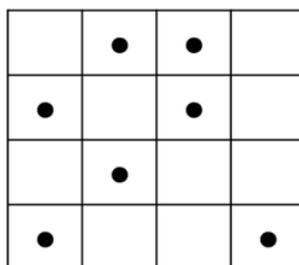
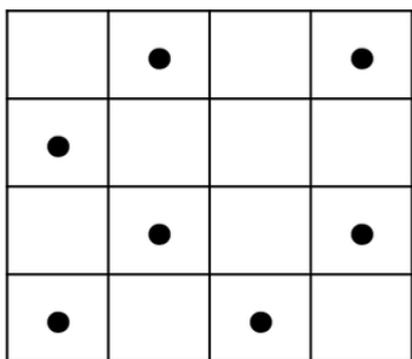
699 conditions, with the constraint of approximately equal cell sizes. Failure to complete the  
700 secondary task above chance level led to the exclusion of  $n = 17$  participants from the test  
701 phase. The retention rate between study and test was approximately 68%, with 260  
702 participants returning for the test phase. An additional 151 participants completed the no-  
703 exposure condition. Three participants were excluded based on a-priori criteria (see Results  
704 section for details). The final sample size for analysis was thus  $N = 408$  (condition NE:  
705  $n = 150$ ; FCOI-:  $n = 128$ ; FCOI+:  $n = 130$ ; age range: 20-74 years;  $M_{\text{age}} = 40.86$ ;  
706  $SD_{\text{age}} = 12.10$ ; 180 males, 227 females, and 1 participant of undisclosed gender). Participants  
707 were paid US\$0.40 for the study phase and US\$0.60 for the test phase.

708 **Materials.** Claims, measures, and procedure were identical to Experiment 1, with the  
709 exception of the secondary task—a dot-pattern-recognition task—used to manipulate  
710 cognitive load (following de Neys & Schaeken, 2007). Participants were presented with a dot  
711 matrix preceding each fact-check (2 s presentation time) and had to perform a 2AFC  
712 recognition test immediately after reading the fact-check. The to-be-remembered pattern was  
713 complex in the FCOI+ condition (seven dots in random locations, with no more than two  
714 (three) dots in any vertical/horizontal (diagonal) line; 2-4 dots overlap between the two test  
715 alternatives) but trivial in the FCOI- condition (four dots in a vertical/horizontal line; four  
716 random positions in test lure; see Figure 11). Above-chance performance was defined as at  
717 least 8 out of 12 correct (cumulative probability when guessing  $p = .194$ ).

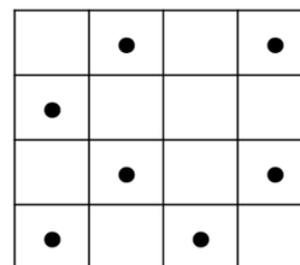
## 718 **Results**

719 Before analysis, we applied the same pre-registered exclusion criteria as in  
720 Experiments 1 and 2. The criterion of “poor” English proficiency was not met by any  
721 participant, but  $n = 1$  participant was excluded due to self-reported lack of effort. Uniform  
722 and erratic responding each led to the exclusion of  $n = 1$  participant. Inference and belief  
723 scores were calculated as in Experiments 1 and 2.

Please memorize this pattern!



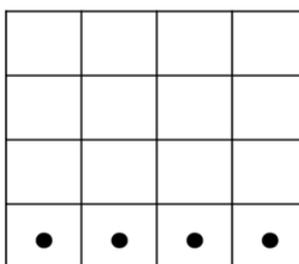
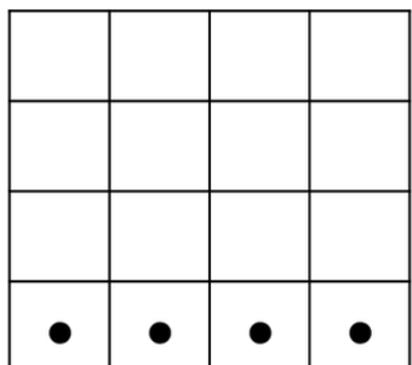
A



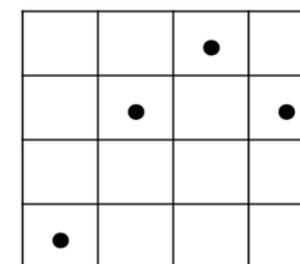
B

Which pattern did you just memorize?

Please memorize this pattern!



A



B

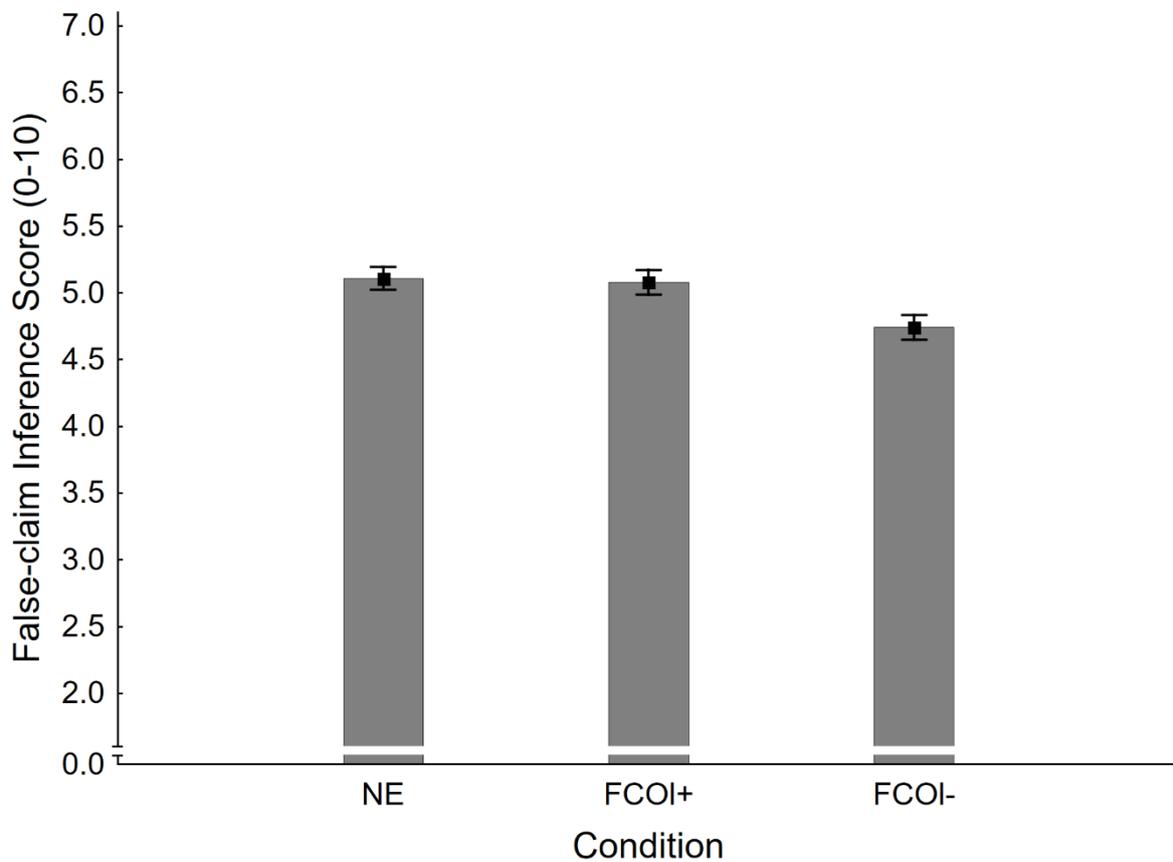
Which pattern did you just memorize?

724

725 *Figure 11.* Example study and test patterns in fact-check-only conditions with high load  
 726 (FCOI+; top) and low load (FCOI-; bottom).

727 ***False-claim inference scores.*** Mean false-claim inference scores across conditions  
 728 are shown in Figure 12. A one-way ANOVA revealed a significant main effect of condition,  
 729  $F(2,405) = 5.00$ ,  $\eta_p^2 = .024$ ,  $p = .007$ . To test the primary hypothesis that corrections of novel  
 730 myths would produce a familiarity backfire effect, a planned contrast compared the no-  
 731 exposure condition (NE:  $M = 5.11$ ,  $SE = 0.09$ ) with the fact-check-only condition with load  
 732 (FCOI+:  $M = 5.08$ ,  $SE = 0.09$ ). This was clearly non-significant,  $F(1,405) = 0.06$ ,  $\eta_p^2 < .001$ ,  
 733  $p = .810$ . We also contrasted the no-exposure condition with the fact-check-only condition  
 734 with no load (FCOI-:  $M = 4.74$ ,  $SE = 0.09$ ), which mirrors the test for familiarity backfire in  
 735 Experiments 1 and 2. This was significant,  $F(1,405) = 8.45$ ,  $\eta_p^2 = .020$ ,  $p = .004$ , but  
 736 constituted a *corrective* effect (i.e.,  $NE > FCOI-$ ). Thus, no familiarity backfire effect was

737 observed, and  $H1_{FIS+}$  and  $H1_{FIS-}$  were rejected. A supplementary planned contrast found a  
 738 significant effect of cognitive load, supporting  $H7_{FIS}$  (see top panel of Table 3).



739

740 *Figure 12.* Mean false-claim inference scores across conditions NE (no-exposure) and  
 741 FCOI+/- (fact-check-only, with/without cognitive load) in Experiment 3. Error bars show  
 742 standard errors of the mean.

743 ***False-claim belief ratings.*** Mean false-claim belief ratings across conditions are  
 744 shown in Figure 13. A one-way ANOVA revealed a significant main effect of condition,  
 745  $F(2,405) = 3.49$ ,  $\eta_p^2 = .017$ ,  $p = .032$ . Planned contrasts were run to test specific hypotheses;  
 746 results are provided in the second panel of Table 3.

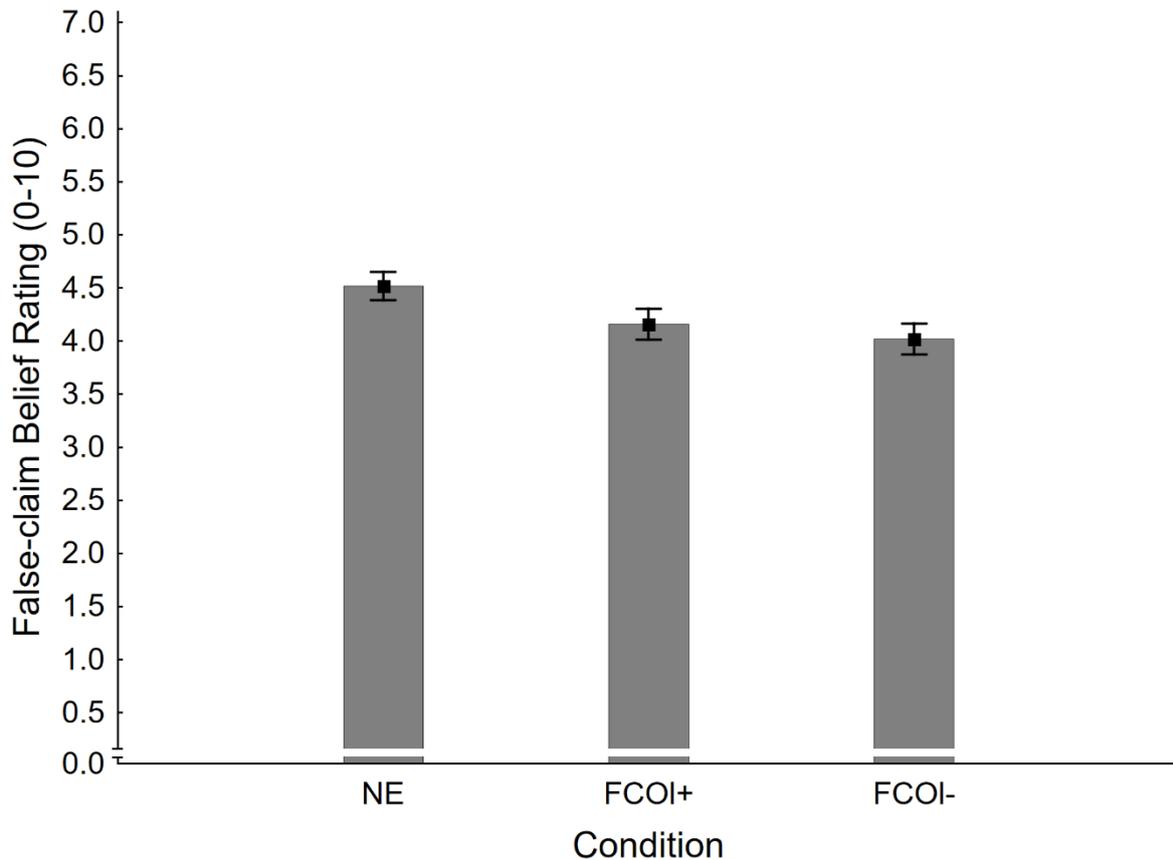
747 Table 3

748 *Contrasts Run in Experiment 3*

dV / Hypothesis	Effect tested	<i>F</i> (1,405)	<i>p</i>
False-claim inference scores			
<b>H1<sub>FISl+</sub>: NE &lt; FCOl+</b>	<b>Familiarity backfire effect</b>	<b>0.06</b>	<b>.810</b>
<b>H1<sub>FISl-</sub>: NE &lt; FCOl-</b>	<b>Familiarity backfire effect</b>	<b>8.45</b>	<b>.004*<sup>^</sup></b>
H7 <sub>FIS</sub> : FCOl- < FCOl+	Load effect on correction	6.65	.010*
False-claim belief ratings			
H1 <sub>FBRl-</sub> : NE < FCOl-	Familiarity backfire effect	6.40	.012* <sup>^</sup>
H1 <sub>FBRl+</sub> : NE < FCOl+	Familiarity backfire effect	3.39	.066 <sup>^</sup>
H7 <sub>FBR</sub> : FCOl- < FCOl+	Load effect on correction	0.45	.501
True-claim inference scores			
H1 <sub>TISl-</sub> : NE < FCOl-	Effect of affirmation vs. baseline	19.21	<.001*
H1 <sub>TISl+</sub> : NE < FCOl+	Effect of affirmation vs. baseline	15.69	<.001*
H7 <sub>TIS</sub> : FCOl- > FCOl+	Load effect on affirmation	0.18	.671
True-claim belief ratings			
H1 <sub>TBRl-</sub> : NE < FCOl-	Effect of affirmation vs. baseline	40.61	<.001*
H1 <sub>TBRl+</sub> : NE < FCOl+	Effect of affirmation vs. baseline	22.32	<.001*
H7 <sub>TBR</sub> : FCOl- > FCOl+	Load effect on affirmation	2.60	.108

749 *Note.* Hypotheses are numbered H1 and H7 (primary hypotheses in bold; see text for details);  
750 subscripts FIS, TIS, FBR, and TBR refer to false-claim and true-claim inference scores and  
751 belief ratings, respectively; no-load and load conditions are indicated by l- and l+. Conditions  
752 are NE = no-exposure; FCOl-/+ = fact-check-only with no load or with load. \* indicates  
753 statistical significance (for secondary contrasts: after Holm-Bonferroni correction).  
754 <sup>^</sup> indicates effect in the opposite of hypothesized direction.

755 It was found that a mere correction with no load at encoding (FCOl-:  $M = 4.02$ ,  
756  $SE = 0.14$ ) reduced false-claim belief relative to no-exposure control (NE:  $M = 4.52$ ,  
757  $SE = 0.13$ ); this rejects familiarity backfire hypothesis H1<sub>FBRl-</sub>. The fact-check-only condition  
758 with load (FCOl+:  $M = 4.16$ ,  $SE = 0.14$ ) did not differ significantly from either of the two  
759 other conditions; this rejects H1<sub>FBRl+</sub> and H7<sub>FBR</sub>.

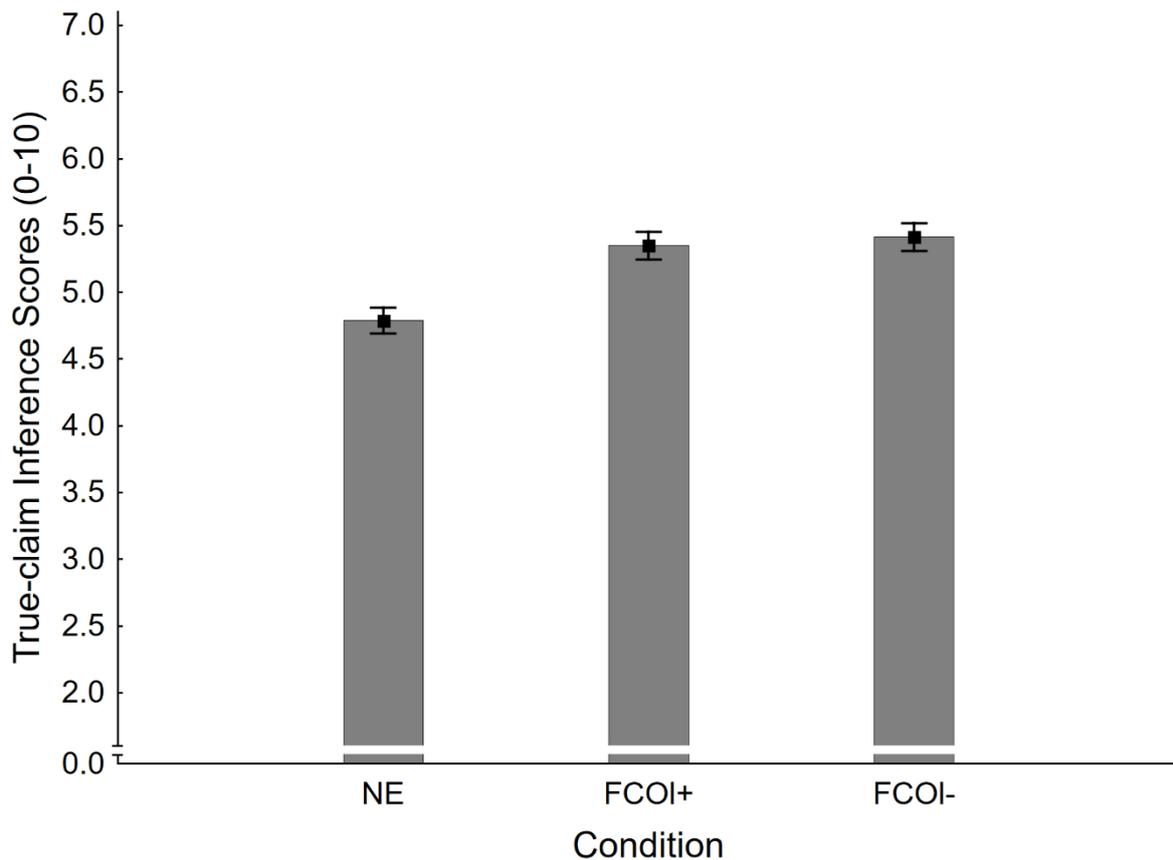


760

761 *Figure 13.* Mean false-claim belief ratings across conditions NE (no-exposure) and FCOI+/-  
 762 (fact-check-only, with/without cognitive load) in Experiment 3. Error bars show standard  
 763 errors of the mean.

764

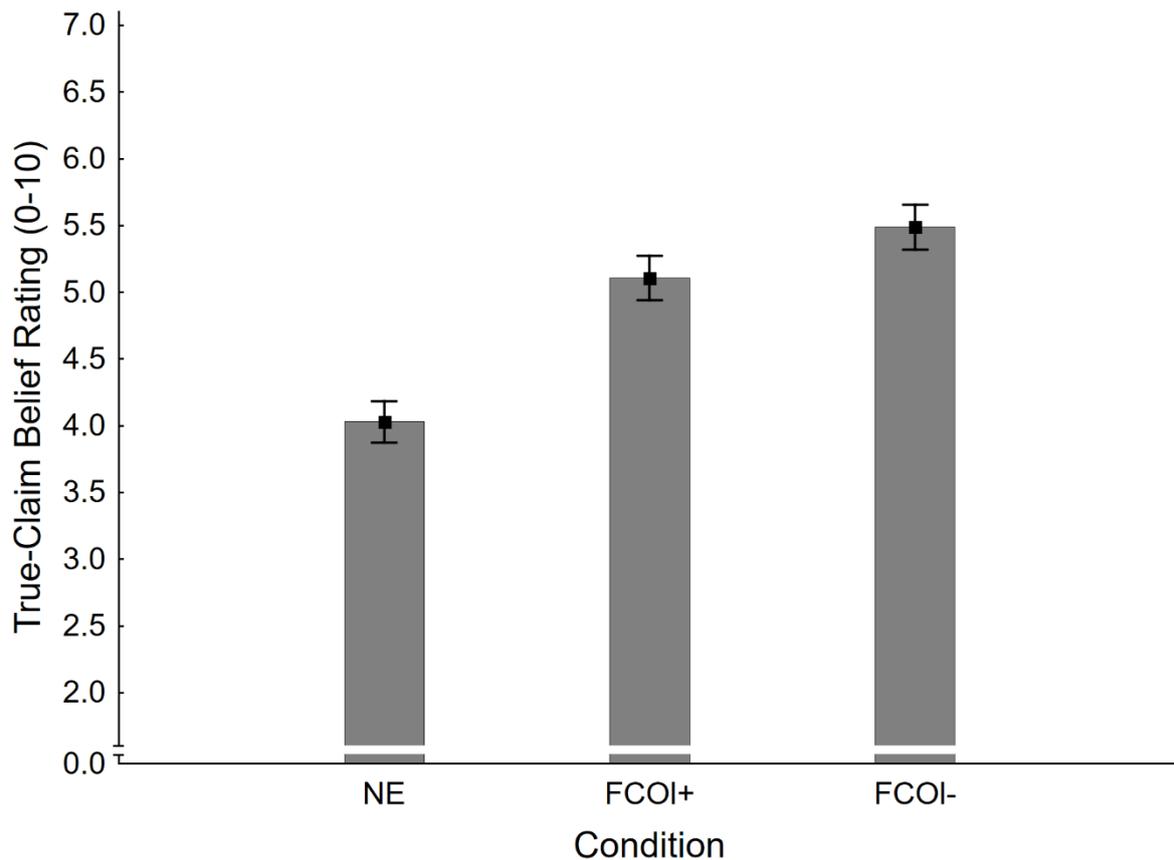
***True-claim inference scores.*** Mean true-claim inference scores across conditions are  
 765 shown in Figure 14. A one-way ANOVA indicated a significant main effect of condition,  
 766  $F(2,405) = 11.99, \eta_p^2 = .056, p < .001$ . Planned contrasts tested for specific condition  
 767 differences; results are reported in the third panel of Table 3. It was found that affirmations  
 768 were equally effective across load conditions (NE:  $M = 4.79, SE = 0.10$ ; FCOI-:  $M = 5.41,$   
 769  $SE = 0.10$ ; FCOI+:  $M = 5.35, SE = 0.10$ ); this supports H1<sub>TISL-</sub> and H1<sub>TISL+</sub>, and rejects H7<sub>TIS</sub>.



770

771 *Figure 14.* Mean true-claim inference scores across conditions NE (no-exposure) and  
 772 FCOI+/- (fact-check-only, with/without cognitive load) in Experiment 3. Error bars show  
 773 standard errors of the mean.

774 ***True-claim belief ratings.*** Mean true-claim belief ratings across conditions are shown  
 775 in Figure 15. A one-way ANOVA yielded a significant main effect of condition,  
 776  $F(2,405) = 22.32, \eta_p^2 = .099, p < .001$ . Planned contrasts tested for specific condition  
 777 differences; results are reported in the fourth panel of Table 3. A mere affirmation increased  
 778 true-claim belief equally in both load conditions (NE:  $M = 4.03, SE = 0.16$ ; FCOI-:  $M = 5.49,$   
 779  $SE = 0.17$ ; FCOI+:  $M = 5.11, SE = 0.17$ ); this supports  $H1_{TBR-}$  and  $H1_{TBR+}$ ; it rejects  $H7_{TBR}$ .



780

781 *Figure 15.* Mean true-claim belief ratings across conditions NE (no-exposure) and FCOI+/-  
782 (fact-check-only, with/without cognitive load) in Experiment 3. Error bars show standard  
783 errors of the mean.

## 784 Discussion

785 Experiment 3 again found no evidence for familiarity backfire effects in either  
786 inference scores or belief ratings. In fact, the no-load condition of Experiment 3 found  
787 evidence that a mere correction of a novel claim significantly *reduced* false-claim-congruent  
788 reasoning and false-claim belief in a delayed test. Under cognitive load at encoding, a mere  
789 correction was unable to reduce misinformed reasoning and beliefs relative to no-exposure  
790 control, but also did no harm. In sum, Experiment 3 found no evidence of familiarity backfire  
791 and is thus more in line with the findings of Ecker et al. (2020) than the findings of Skurnik et  
792 al. (2007). The fact that cognitive load at fact-check encoding reduced the impact of a  
793 correction on false-claim inference scores but did not influence the effects of affirmations can

794 be seen as additional evidence that avoiding false-claim-congruent reasoning relies on  
795 recollection of the correction, which would have been impaired by the cognitive load  
796 (however, no such effect was observed for false-claim belief ratings). Moreover, mere  
797 affirmations were generally found to increase true-claim-congruent reasoning and true-claim  
798 belief after a delay irrespective of load at encoding, in line with Ecker et al. (2020) and  
799 Experiments 1 and 2.

### 800 **Bayesian Analyses**

801 To further corroborate the evidence for or against familiarity backfire effects, we  
802 employed supplementary Bayesian analyses; these have the advantage that evidence in  
803 support of a null hypothesis can be quantified (e.g., see Wagenmakers, Love et al., 2018).  
804 Specifically, Bayesian ANOVAs were run on inference scores and belief ratings from the no-  
805 exposure and fact-check-only conditions of Experiments 1-3 (separately and conjointly; the  
806 analysis on Experiment 3 data and the conjoint analysis were pre-registered before running  
807 Experiment 3). These tested whether there was evidence for a model including a condition  
808 factor over a null model. Mean inference scores across experiments were  $M = 5.06$  ( $SE = .05$ )  
809 for the no-exposure condition and  $M = 5.01$  ( $SE = .06$ ) for the fact-check only condition (or  
810  $M = 5.12$  [ $SE = .06$ ] when using the load condition of Experiment 3). Mean belief ratings  
811 across experiments were  $M = 4.40$  ( $SE = .08$ ) for the no-exposure condition and  $M = 4.10$   
812 ( $SE = .09$ ) for the fact-check only condition (or  $M = 4.15$  [ $SE = .08$ ] when using the load  
813 condition of Experiment 3).

814 The Bayes factors ( $BF_{10}$ ) in Table 4 quantify the evidence for or against inclusion of  
815 the condition factor. A  $BF_{10} > 1$  suggests evidence in favor of including a condition factor  
816 (which can be interpreted as a main effect of condition); a  $BF_{10} < 1$  suggests evidence in  
817 favor of the null model. For example,  $BF_{10} = 10$  would suggest that the data are 10 times  
818 more likely to have occurred under the alternative hypothesis than the null hypothesis;

819  $BF_{10} = 0.10$  would suggest that the data are 10 times more likely to occur under the null  
 820 hypothesis.  $BF$  values between 0.33 and 3 are taken to only provide anecdotal evidence;  $BF$   
 821 values between 0.1 and 0.33, or 3 and 10 constitute moderate/substantial evidence;  $BF$  values  
 822  $< 0.1$  or  $> 10$  provide strong to very strong evidence (Jeffreys, 1961; Wagenmakers, Love et  
 823 al., 2018).

824 Table 4

825 *Results from Bayesian Analyses across Experiments 1-3*

dV	Effect direction	$BF_{10}$
Experiment 1		
FIS	NE < FCO (familiarity backfire)	2.801
FBR	NE = FCO (no familiarity backfire)	0.154*
Experiment 2		
FISd	NE = FCO (no familiarity backfire)	0.135*
FBRd	NE = FCO (no familiarity backfire)	0.363
Experiment 3		
FISl-	NE > FCO (corrective effect)	11.757**
FBRI-	NE > FCO (corrective effect)	3.065*
FISl+	NE = FCO (no familiarity backfire)	0.135*
FBRI+	NE = FCO (no familiarity backfire)	0.774
Experiments 1-3		
FIS(l-)	NE = FCO (no familiarity backfire)	0.104*
FBR(l-)	NE > FCO (corrective effect)	1.799
FIS(l+)	NE = FCO (no familiarity backfire)	0.112*
FBR(l+)	NE = FCO (no familiarity backfire)	0.760

826 *Note.* FIS and FBR: False-claim inference scores and belief ratings from delayed test. As test  
 827 delay was manipulated in Experiment 2, only the delayed test variables (FISd and FBRd)  
 828 were entered into analysis. No-load (FISl-; FBRI-) and load (FISl+; FBRI+) conditions of  
 829 Experiment 3 were included in separate analysis of Experiment 3, and also in separate  
 830 conjoint analyses. The condition factor includes only conditions NE (no-exposure) and FCO  
 831 (fact-check-only). \* indicates substantial and \*\* indicates strong evidence for or against the  
 832 null.

833 As can be seen in Table 4, the evidence for a familiarity backfire effect from the  
 834 inference scores in Experiment 1 was only anecdotal, while Experiment 2 provided  
 835 substantial evidence against a familiarity backfire effect, and Experiment 3 yielded strong

836 evidence for a *corrective* effect in the no-load condition (which matched the conditions of  
837 Experiments 1 and 2), while providing substantial evidence against familiarity backfire in the  
838 load condition. Likewise, the secondary belief measures suggested substantial evidence  
839 against backfire in Experiment 1 and substantial evidence for a corrective effect in the no-  
840 load condition of Experiment 3. However, the main conclusion to be drawn, from the  
841 conjoint analyses, is that the experiments reported in this paper overall yielded substantial to  
842 strong evidence against familiarity backfire effects: Across experiments, while the secondary  
843 belief-rating data remained inconclusive, the primary inferential reasoning data were found to  
844 be approximately nine times more likely to have occurred under the null hypothesis.

### 845 **General Discussion**

846 The main focus of this paper was to investigate whether mere exposure to a correction  
847 could familiarize people with a novel piece of misinformation such that it would negatively  
848 affect their reasoning and beliefs. In other words, we tested whether corrections of novel  
849 misinformation could elicit a familiarity-driven backfire effect, which may ironically  
850 strengthen misconceptions and spread misinformation to new audiences (Schwarz et al.,  
851 2007, 2016).<sup>3</sup> Experiment 1 found some evidence for a familiarity backfire effect, but the  
852 evidence was statistically weak and the result failed to occur in an exact replication with  
853 greater experimental power (Experiment 2) as well as a close replication that added only a  
854 trivial secondary task (the no-load condition of Experiment 3). In fact, both Experiments 2  
855 and 3 yielded substantial evidence *against* the presence of a familiarity backfire effect, even

---

<sup>3</sup> We note that Kessler, Braasch, & Kardash (2019) recently reported a backfire effect with vaccination misinformation, which they observed only in people with high “flexible thinking” scores. Kessler et al. speculated that in flexible thinkers—those who open-mindedly consider new information—corrections might thus spread novel misinformation. However, they did not measure misinformation novelty, and only prior vaccination *beliefs* and not prior vaccination *knowledge* predicted the backfire effects they observed; it therefore seems more likely that these effects were driven by worldview rather than familiarity (see Lewandowsky et al., 2012; Ecker & Ang, 2019).

856 under conditions that should maximize reliance on familiarity and thus facilitate occurrence  
857 of familiarity backfire, viz. the combination of novel claims that maximized the familiarity  
858 boost conveyed by first exposure, a relatively long one-week retention interval, and  
859 correction encoding under cognitive load (the load condition of Experiment 3). Thus, while  
860 there was some variability across experiments, the overall evidence was in support of the null  
861 hypothesis. This meshes well with previous studies failing to find evidence for familiarity  
862 backfire with more familiar claims (Ecker et al., 2017, 2020; Swire et al., 2017).

863         However, this does not rule out misinformation familiarity as an important driver of  
864 continued influence effects. This is because we also found consistent evidence that after a  
865 delay of one week, affirmations of true claims were more effective than corrections of false  
866 claims. This closely mirrors the pattern observed by Swire et al. (2017)<sup>4</sup> and thus  
867 corroborates their conclusion that misinformation familiarity can be a counterproductive  
868 force when correcting false claims. That is, the overall evidence observed here suggests, in  
869 line with Swire et al., that acceptance of false claims can be driven by claim familiarity, in  
870 particular when the ability to recollect the correction is reduced (e.g., due to delay-related  
871 forgetting or cognitive load). This can offset the correction entirely, such that endorsement of  
872 a false claim and false-claim-congruent reasoning can return to baseline after a one-week  
873 delay, which essentially means that even a correction that is reasonably effective in the short  
874 term can lose its impact within a week, relative to a no-exposure control condition (as  
875 demonstrated in Experiment 2; see Figures 7 and 8; note that corrections were still somewhat  
876 effective relative to the claim-only condition). However, the boost to claim familiarity  
877 through claim repetition within the correction is typically not substantial enough to cause  
878 actual backfire. Broadly speaking, these results support the view that memory-based

---

<sup>4</sup> Peter and Koch (2016) also observed this asymmetry, although they referred to this asymmetry itself as a familiarity backfire effect, which in our view is a misnomer.

879 evaluation processes determine inferential reasoning and endorsement of claims much more  
880 than metacognitive judgments of fluency (cf. Schwarz et al., 2007). The conflicting results  
881 from Experiment 1 can only serve as a reminder that one should never place too much  
882 emphasis on the findings of a single experiment (e.g., see Murayama, Pekrun, & Fiedler,  
883 2014), and that significant *p*-values can translate to only “anecdotal” evidence under a  
884 Bayesian framework (see Wagenmakers, Marsman et al., 2018, for a detailed discussion). We  
885 speculate that some of the variability in findings arose due to the use of novel claims. While it  
886 was necessary for the present project to use novel claims for the theoretical and practical  
887 reasons outlined earlier, the claims we used are not generally representative of claims  
888 encountered in the real world, which are typically grounded in contextual world knowledge.  
889 Ratings of such novel claims may be inherently less reliable than ratings of familiar claims  
890 that can tap into pre-existing knowledge and beliefs (Swire-Thompson, DeGutis, & Lazer,  
891 2020).

892         Additional evidence obtained in the present set of experiments regards the illusory  
893 truth effects conveyed by mere exposure (Begg et al., 1992; Dechêne et al., 2010; Parks &  
894 Toth, 2006; Unkelbach, 2007; Weaver et al., 2007). While some research has found that even  
895 a single exposure to a false claim can have measurable impact on claim endorsement (e.g.,  
896 Pennycook et al., 2018), the evidence here was somewhat mixed. Experiment 1 found some  
897 evidence for illusory truth effects with true but not false claims, whereas Experiment 2 found  
898 evidence for illusory truth effects after a delay with false claims (and also on true-claim belief  
899 ratings but not inference scores). This pattern was observed despite the fact that participants  
900 could not reliably differentiate between true and false claims, and control-group (no-  
901 exposure) belief ratings were generally lower for true claims in both experiments. The fact  
902 that illusory truth effects were only observed in the delayed test of Experiment 2 but not in  
903 the immediate test suggests that these effects were indeed driven by familiarity rather than

904 perceived social consensus (see Pennycook et al., 2018; Unkelbach, 2007; Weaver et al.,  
905 2007). However, apart from that, we can only conclude from these results that a single  
906 exposure to a claim can lead to enhanced subsequent endorsement, but that this is not always  
907 the case. Thus, to some extent, this mirrors our conclusions regarding the role of familiarity  
908 for continued influence, in that the evidence regarding illusory truth effects we obtained is  
909 somewhat inconsistent, but generally suggests that familiarity likely impacts reasoning and  
910 endorsement of claims (we also note that evidence for illusory truth effects in general is much  
911 more solid than the evidence for familiarity backfire effects; e.g., see De keersmaecker et al.,  
912 2020).

913         The practical implications of this research are clear: Recommendations to front-line  
914 educators and communicators to entirely avoid repeating misinformation when debunking  
915 (Cook & Lewandowsky, 2011; Lewandowsky et al., 2012; Peter & Koch, 2016; Schwarz,  
916 Newman, & Leach, 2016; Schwarz et al., 2007) were unwarranted. Recent research indicates  
917 that repeating misinformation when correcting it can have a positive effect, enhancing a  
918 correction in the short term (presumably by increasing the salience of the correction and  
919 facilitating conflict resolution and knowledge revision processes; see Ecker et al., 2017;  
920 Kendeou et al., 2014). There is also evidence that exposure to a correction that repeats a piece  
921 of (non-novel) misinformation does not lead to backfire effects relative to either a pre-  
922 correction or no-exposure baseline (Ecker et al., 2020). Finally, the present study suggests  
923 that exposure to a correction does not cause familiarity backfire relative to a no-exposure  
924 control even with novel claims, and thus corrections do not seem to spread misinformation to  
925 new audiences easily.

926         That being said, recommendations to avoid *unnecessary* misinformation repetition  
927 should arguably remain in place—while one repetition in the context of a correction may  
928 have benefits for correction salience, additional repetition of the misinformation runs the risk

929 of enhancing familiarity without any added benefit. Moreover, while we have demonstrated  
930 that corrections do not backfire when it comes to specific beliefs about a proposition, one  
931 needs to differentiate this from the over-arching framing that is achieved by stating  
932 something that is false (see Lakoff, 2010). For example, a government official stating that  
933 there are “no plans for a carbon tax” may achieve a reduction in the specific belief that a  
934 carbon tax rollout is being prepared, but at the same time using the word “tax” may make  
935 people who oppose new taxes for ideological or pragmatic reasons think about climate  
936 change as a threat rather than an opportunity (also see Fletcher, 2009; Kahan, 2010;  
937 Lewandowsky et al., 2017). Therefore, communicators should perhaps focus their  
938 considerations more on the framing of their corrections, as repeating the misinformation  
939 *frame* might do more damage than repetition of the misinformation itself. Investigating the  
940 effects of frame repetition within corrections is therefore an important target for future  
941 research.

## References

942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964

Arkes, H. R., Boehm, L. E., & Xu, G. (1991). Determinants of judged validity. *Journal of Experimental Social Psychology, 27*, 576-605. doi:10.1016/0022-1031(91)90026-3

Ayers, M. S., & Reder, L. M. (1998). A theoretical review of the misinformation effect: Predictions from an activation-based memory model. *Psychonomic Bulletin & Review, 5*, 1-21. doi:10.3758/BF03209454

Begg, I. M., Anas, A., & Farinacci, S. (1992). Dissociation of processes in belief: Source recollection, statement familiarity, and the illusion of truth. *Journal of Experimental Psychology: General, 121*, 446-458. doi:10.1037/0096-3445.121.4.446

Berinsky, A., Huber, G., & Lenz, G. (2012). Evaluating online labor markets for experimental research: Amazon.com's Mechanical Turk. *Political Analysis, 20*, 351-368. doi:10.1093/pan/mpr057

Bode, L., & Vraga, E. K. (2018). See something, say something: Correction of global health misinformation on social media. *Health Communication, 33*, 1131-1140.

Cameron, K. A., Roloff, M. E., Friesema, E. M., Brown, T., Jovanovic, B. D., Hauber, S., & Baker, D. W. (2013). Patient knowledge and recall of health information following exposure to "facts and myths" message format variations. *Patient Education and Counseling, 92*, 381-387. doi:10.1016/j.pec.2013.06.017

Chan, M. S., Jones, C. R., Hall Jamieson, K., & Albarracin, D. (2017). Debunking: A meta-analysis of the psychological efficacy of messages countering misinformation. *Psychological Science, 28*, 1531-1546. doi:10.1177/0956797617714579

Cook, J., & Lewandowsky, S. (2011). *The debunking handbook*. Retrieved from [http://www.skepticalscience.com/docs/Debunking\\_Handbook.pdf](http://www.skepticalscience.com/docs/Debunking_Handbook.pdf)

- 965 Craik, F. I. M., Govoni, R., Naveh-Benjamin, M., & Anderson, N. D. (1996). The effects of  
966 divided attention on encoding and retrieval processes in human memory. *Journal of*  
967 *Experimental Psychology: General*, *125*, 159-180. doi:10.1037/0096-3445.125.2.159
- 968 De keersmaecker, J., Dunning, D., Pennycook, G., Rand, D. G., Sanchez, C., Unkelbach, C.,  
969 & Roets, A. (2020). Investigating the robustness of the illusory truth effect across  
970 individual differences in cognitive ability, need for cognitive closure, and cognitive  
971 style. *Personality and Social Psychology Bulletin*, *46*, 204-215.  
972 doi:10.1177/0146167219853844
- 973 de Neys, W., & Schaeken, W. (2007). When people are more logical under cognitive load:  
974 Dual task impact on scalar implicature. *Experimental Psychology*, *54*, 128-133.  
975 doi:10.1027/1618-3169.54.2.128
- 976 Dechêne, A., Stahl, C., Hansen, J., & Wänke, M. (2010). The truth about the truth: A meta-  
977 analytic review of the truth effect. *Personality and Social Psychology Review*, *14*,  
978 238-257. doi:10.1177/1088868309352251
- 979 Diana, R. A., Reder, L. M., Arndt, J., & Park, H. (2006). Models of recognition: A review of  
980 arguments in favor of a dual-process account. *Psychonomic Bulletin & Review*, *13*, 1-  
981 21. doi:10.3758/BF03193807
- 982 Ecker, U. K. H., & Ang, L. C. (2019). Political attitudes and the processing of misinformation  
983 corrections. *Political Psychology*, *40*, 241-260. doi:10.1111/pops.12494
- 984 Ecker, U. K. H., Hogan, J. L., & Lewandowsky, S. (2017). Reminders and repetition of  
985 misinformation: Helping or hindering its retraction? *Journal of Applied Research in*  
986 *Memory and Cognition*, *6*, 185-192. doi:10.1016/j.jarmac.2017.01.014

- 987 Ecker, U. K. H., Lewandowsky, S., Swire, B., & Chang, D. (2011). Correcting false  
988 information in memory: Manipulating the strength of misinformation encoding and its  
989 retraction. *Psychonomic Bulletin & Review*, *18*, 570-578. doi:10.3758/s13423-011-  
990 0065-1
- 991 Ecker, U. K. H., Lewandowsky, S., & Tang, D. T. W. (2010). Explicit warnings reduce but  
992 do not eliminate the continued influence of misinformation. *Memory & Cognition*, *38*,  
993 1087-1100. doi:10.3758/MC.38.8.1087
- 994 Ecker, U. K. H., O'Reilly, Z., Reid, J. S., & Chang, E. P. (2020). The effectiveness of short-  
995 format refutational fact-checks. *British Journal of Psychology*, *111*, 36-54.  
996 doi:10.1111/bjop.12383
- 997 Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G\*Power 3: A flexible statistical  
998 power analysis program for the social, behavioral, and biomedical sciences. *Behavior*  
999 *Research Methods*, *39*, 175-191. doi:10.3758/BF03193146
- 1000 Fletcher, A. L. (2009). Clearing the air: The contribution of frame analysis to understanding  
1001 climate policy in the United States. *Environmental Politics*, *18*, 800-816.  
1002 doi:10.1080/09644010903157123
- 1003 Gordon, A., Quadflieg, S., Brooks, J. C. W., Ecker, U. K. H., & Lewandowsky, S. (2019).  
1004 Keeping track of 'alternative facts': The neural correlates of processing  
1005 misinformation corrections. *NeuroImage*, *193*, 46-56.  
1006 doi:10.1016/j.neuroimage.2019.03.014
- 1007 Hauser, D. J., & Schwarz, N. (2016). Attentive Turkers: MTurk participants perform better  
1008 on online attention checks than do subject pool participants. *Behavior Research*  
1009 *Methods*, *48*, 400-407. doi:10.3758/s13428-015-0578-z

- 1010 Hicks, J. L., & Marsh, R. L. (2000). Toward specifying the attentional demands of  
1011 recognition memory. *Journal of Experimental Psychology: Learning, Memory, and*  
1012 *Cognition*, 26, 1483-1498. doi:10.1037/0278-7393.26.6.1483
- 1013 Hintzman, D. L., & Curran, T. (1994). Retrieval dynamics of recognition and frequency  
1014 judgments: Evidence for separate processes of familiarity and recall. *Journal of*  
1015 *Memory and Language*, 33, 1-18. doi:10.1006/jmla.1994.1001
- 1016 Hoaglin, D. C., and Iglewicz, B. (1987). Fine tuning some resistant rules for outlier labeling.  
1017 *Journal of American Statistical Association*, 82, 1147-1149. doi:10.2307/2289392
- 1018 Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian*  
1019 *Journal of Statistics*, 6, 65-70. doi:10.2307/4615733
- 1020 Jeffreys, H. (1961). *Theory of probability*. Oxford: Oxford University Press.
- 1021 Johnson, H. M., & Seifert, C. M. (1994). Sources of the continued influence effect: When  
1022 misinformation in memory affects later inferences. *Journal of Experimental*  
1023 *Psychology: Learning, Memory, and Cognition*, 20, 1420-1436. doi:10.1037/0278-  
1024 7393.20.6.1420
- 1025 Kendeou, P., Walsh, E. K., Smith, E. R., & O'Brien, E. J. (2014). Knowledge revision  
1026 processes in refutation texts. *Discourse Processes*, 51, 374-397.  
1027 doi:10.1080/0163853X.2014.913961
- 1028 Kessler, E. D., Braasch, J. L. G., & Kardash, C. M. (2019) Individual differences in revising  
1029 (and maintaining) accurate and inaccurate beliefs about childhood vaccinations.  
1030 *Discourse Processes*, 56, 415-428. doi:10.1080/0163853X.2019.1596709
- 1031 Knowlton, B. J., & Squire, L. R. (1995). Remembering and knowing: Two different  
1032 expressions of declarative memory. *Journal of Experimental Psychology: Learning,*  
1033 *Memory, and Cognition*, 21, 699-710. doi:10.1037//0278-7393.21.3.699

- 1034 Lakoff, G. (2010). Why it matters how we frame the environment. *Environmental*  
1035 *Communication*, 4, 70-81. doi:10.1080/17524030903529749
- 1036 Lazer, D. M., Baum, M. A., Benkler, Y., Berinsky, A. J., Greenhill, K. M., Menczer, F., ...  
1037 Rothschild, D. (2018). The science of fake news. *Science*, 359, 1094-1096.  
1038 doi:10.1126/science.aao2998
- 1039 Lewandowsky, S., Ecker, U. K. H., & Cook, J. (2017). Beyond misinformation:  
1040 Understanding and coping with the “post-truth” era. *Journal of Applied Research in*  
1041 *Memory and Cognition*, 6, 353-369. doi:10.1016/j.jarmac.2017.07.008
- 1042 Lewandowsky, S., Ecker, U. K. H., Seifert, C. M., Schwarz, N., & Cook, J. (2012).  
1043 Misinformation and its correction: Continued influence and successful debiasing.  
1044 *Psychological Science in the Public Interest*, 13, 106-131.  
1045 doi:10.1177/1529100612451018
- 1046 Litman, L., Robinson, J., & Abberbock, T. (2017). TurkPrime.com: A versatile  
1047 crowdsourcing data acquisition platform for the behavioral sciences. *Behavior*  
1048 *Research Methods*, 49, 433-442. doi:10.3758/s13428-016-0727-z
- 1049 MacFarlane, D., Hurlstone, M. J., & Ecker, U. K. H. (2020). Protecting consumers from  
1050 fraudulent health claims: A taxonomy of psychological drivers, interventions, barriers,  
1051 and treatments. *Social Science & Medicine*. doi:10.1016/j.socscimed.2020.112790
- 1052 Marsh, E. J., & Fazio, L. K. (2006). Learning errors from fiction: Difficulties in reducing  
1053 reliance on fictional stories. *Memory & Cognition*, 34, 1140-1149.  
1054 doi:10.3758/BF03193260
- 1055 Murayama, K., Pekrun, R., & Fiedler, K. (2014). Research practices that can prevent an  
1056 inflation of false-positive rates. *Personality and Social Psychology Review*, 18, 107-  
1057 118. doi:10.1177/1088868313496330

- 1058 Necka, E. A., Cacioppo, S., Norman, G. J., & Cacioppo, J. T. (2016). Measuring the  
1059 prevalence of problematic respondent behaviors among MTurk, campus, and  
1060 community participants. *PLOS ONE*, *11*, e0157732.  
1061 doi:10.1371/journal.pone.0157732
- 1062 Nyhan, B., Reifler, J., Richey, S., & Freed, G. L. (2014). Effective messages in vaccine  
1063 promotion: A randomized trial. *Pediatrics*, *133*, e835-e842. doi:10.1542/peds.2013-  
1064 2365
- 1065 Oberauer, K., Lewandowsky, S., Farrell, S., Jarrold, C., & Greaves, M. (2012). Modeling  
1066 working memory: An interference model of complex span. *Psychonomic Bulletin &  
1067 Review*, *19*, 779-819. doi:10.3758/s13423-012-0272-4
- 1068 Parks, C. M., & Toth, J. P. (2006). Fluency, familiarity, aging, and the illusion of truth.  
1069 *Aging, Neuropsychology, and Cognition*, *13*, 225-253. doi:10.1080/138255890968691
- 1070 Paynter, J. M., Luskin-Saxby, S., Keen, D., Fordyce, K., Frost, G., Imms, C., ... Ecker, U. K.  
1071 H. (2019). Evaluation of a template for countering misinformation: Real-world autism  
1072 treatment myth debunking. *PLOS ONE*, *14*, e0210746.  
1073 <https://doi.org/10.1371/journal.pone.0210746>
- 1074 Pennycook, G., Cannon, T. D., & Rand, D. G. (2018). Prior exposure increases perceived  
1075 accuracy of fake news. *Journal of Experimental Psychology: General*, *147*, 1865-  
1076 1880. doi:10.1037/xge0000465
- 1077 Peter, C., & Koch, T. (2016). When debunking scientific myths fails (and when it does not):  
1078 The backfire effect in the context of journalistic coverage and immediate judgments  
1079 as prevention strategy. *Science Communication*, *38*, 3-25.  
1080 doi:10.1177/1075547015613523

- 1081 Rich, P. R., & Zaragoza, M. S. (2016). The continued influence of implied and explicitly  
1082 stated misinformation in news reports. *Journal of Experimental Psychology:*  
1083 *Learning, Memory, and Cognition*, 42, 62-74. doi:10.1037/xlm0000155
- 1084 Schwarz, N., Newman, E., & Leach, W. (2016). Making the truth stick & the myths fade:  
1085 Lessons from cognitive psychology. *Behavioral Science & Policy*, 2, 85-95.
- 1086 Schwarz, N., Sanna, L. J., Skurnik, I., & Yoon, C. (2007). Metacognitive experiences and the  
1087 intricacies of setting people straight: Implications for debiasing and public  
1088 information campaigns. *Advances in Experimental Social Psychology*, 39, 127-161.  
1089 doi:10.1016/S0065-2601(06)39003-X
- 1090 Skurnik, I., Yoon, C., Park, D. C., & Schwarz, N. (2005). How warnings about false claims  
1091 become recommendations. *Journal of Consumer Research*, 31, 713-724.  
1092 doi:10.1086/426605
- 1093 Skurnik, I., Yoon, C., & Schwarz, N. (2007). *Myths and facts about the flu: Health education*  
1094 *campaigns can reduce vaccination intentions*. Unpublished manuscript available from  
1095 [http://webuser.bus.umich.edu/yoonc/research/Papers/Skurnik\\_Yoon\\_Schwarz\\_2005\\_](http://webuser.bus.umich.edu/yoonc/research/Papers/Skurnik_Yoon_Schwarz_2005_Myths_Facts_Flu_Health_Education_Campaigns_JAMA.pdf)  
1096 [Myths\\_Facts\\_Flu\\_Health\\_Education\\_Campaigns\\_JAMA.pdf](http://webuser.bus.umich.edu/yoonc/research/Papers/Skurnik_Yoon_Schwarz_2005_Myths_Facts_Flu_Health_Education_Campaigns_JAMA.pdf)
- 1097 Southwell, B. G., & Thorson, E. A. (2015). The prevalence, consequence, and remedy of  
1098 misinformation in mass media systems. *Journal of Communication*, 65, 589-595.  
1099 doi:10.1111/jcom.12168
- 1100 Swire, B., Ecker, U. K. H., & Lewandowsky, S. (2017). The role of familiarity in correcting  
1101 inaccurate information. *Journal of Experimental Psychology: Learning, Memory, and*  
1102 *Cognition*, 43, 1948-1961. doi:10.1037/xlm0000422
- 1103 Swire-Thompson, B., DeGutis, J., & Lazer, D. (2020). Searching for the backfire effect:  
1104 Measurement and design considerations. doi:10.31234/osf.io/ba2kc

- 1105 Unkelbach, C. (2007). Reversing the truth effect: Learning the interpretation of processing  
1106 fluency in judgments of truth. *Journal of Experimental Psychology: Learning,*  
1107 *Memory, and Cognition*, *33*, 219-230. doi:10.1037/0278-7393.33.1.219
- 1108 Vargo, C. J., Guo, L., & Amazeen, M. A. (2018). The agenda-setting power of fake news: A  
1109 big data analysis of the online media landscape from 2014 to 2016. *New Media and*  
1110 *Society*, *20*, 2028-2049. doi:10.1177/1461444817712086
- 1111 Wagenmakers, E.-J., Love, J., Marsman, M., Jamil, T., Ly, A., Verhagen, J., ... Morey, R. D.  
1112 (2018). Bayesian inference for psychology. Part II: Example applications with JASP.  
1113 *Psychonomic Bulletin & Review*, *25*, 58-76. doi:10.3758/s13423-017-1323-7
- 1114 Wagenmakers, E.-J., Marsman, M., Jamil, T., Ly, A., Verhagen, J., Love, J., ... Morey, R. D.  
1115 (2018). Bayesian inference for psychology. Part I: Theoretical advantages and  
1116 practical ramifications. *Psychonomic Bulletin & Review*, *25*, 35-57.  
1117 doi:10.3758/s13423-017-1343-3
- 1118 Walter, N., & Tukachinsky, R. (2020). A meta-analytic examination of the continued  
1119 influence of misinformation in the face of correction: How powerful is it, why does it  
1120 happen, and how to stop it? *Communication Research*, *47*, 155-177.  
1121 doi:10.1177/0093650219854600
- 1122 Weaver, K., Garcia, S. M., Schwarz, N., & Miller, D. T. (2007). Inferring the popularity of an  
1123 opinion from its familiarity: A repetitive voice can sound like a chorus. *Journal of*  
1124 *Personality and Social Psychology*, *92*, 821-833. doi:10.1037/0022-3514.92.5.821

**1125 Declarations**

1126 Ethics approval and consent to participate: Approval to conduct this research was granted by  
1127 the Human Ethics Office of the University of Western Australia under RA/4/1/8104. All  
1128 participants consented to participation after receiving an approved information sheet.

1129 Consent for publication: Not applicable (no identifiable data published).

1130 Availability of data and materials: Materials are provided in the Appendix. Data are available  
1131 on the Open Science Framework website at <https://osf.io/69bq3/>.

1132 Competing interests: The authors declare no competing interests.

1133 Funding: The research was supported by the Australian Research Council under grant  
1134 DP160103596, awarded to the first and second authors.

1135 Authors' contributions: UE and SL conceptualized and designed the experiments; UE and  
1136 MC created the experimental surveys, analyzed the data, and wrote the initial manuscript  
1137 draft; SL contributed to the writing.

1138 Acknowledgements: We thank Hannah Beach, Rosie Carbutt, Charles Hanich, and Ee Pin  
1139 Chang for research assistance.

**Appendix**

1140

**Claim Pilot Rating**

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

A total of 100 MTurk participants partook in the online pilot survey; participants could not participate in the main study. Data from participants were excluded for the following two a-priori reasons: (1) uniform responding and (2) completing the survey in less than five minutes. Two participants were classified as uniform responders; across all responses, they showed  $SD < 0.467$ , the lower outlier criterion of the inter-quartile rule with a 2.2 multiplier. Eight participants completed the survey in less than five minutes. One participant met both criteria. Consequently,  $n = 9$  participants were excluded, resulting in a final sample size of  $N = 91$  (age range: 20-64 years;  $M_{\text{age}} = 36.48$ ;  $SD_{\text{age}} = 10.30$ ; 50 males, 40 females, and one participant of undisclosed gender). Claim ratings from the pilot study are presented in Tables A1 and A2.

Table A1. *Familiarity and Believability Ratings of False Claims in Pilot Study*

Claim	Familiarity		Believability	
	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
*Facebook is about to launch a “no swearing” campaign.	1.52	0.95	2.69	1.07
Frequently wearing silk garments in direct contact with the skin can cause spontaneous lactation.	1.53	0.96	3.65	0.98
*50,000 men were raped in South Africa last year.	1.56	1.04	3.66	0.92
NASA is predicting six consecutive days of darkness in the Northern hemisphere in 2022 due to a rare astronomical event.	1.57	1.07	3.25	1.06
*The ratio of male:female CEOs in Manchester, UK is 1:1.	1.60	0.94	3.59	1.03
*The outer skin of a pineapple emits a dangerous toxin into the environment when it breaks down.	1.60	1.02	2.19	1.19
The Cinderella Castle at Disneyland Florida can be disassembled during hurricanes.	1.64	1.04	3.24	1.11
Fibers found in cow skin are now being added to Botox injections.	1.76	1.09	2.85	1.00
“Camo” the German shepherd is the only dog in history to become an Officer of the British Empire.	1.76	1.07	2.52	1.06
*Hugh Hefner donated a fifth of his will to the Planned Parenthood charity.	1.84	1.08	2.47	0.99
*Placing a car battery on a cement floor can drain it and lead to its decay.	1.90	1.28	2.74	1.39
The first artificial intelligence robot has been appointed as a teaching assistant in Japan.	2.00	1.11	2.21	0.97
Nike footwear has to meet a quota of containing at least 20% recycled materials.	2.10	1.17	3.41	1.32
The motor-vehicle accident rate regularly surges after the Super Bowl in the home state of the losing team.	2.15	1.26	3.15	1.06
Bitcoin is used by the American government as a way to keep track of online criminal activity.	2.18	1.41	3.03	0.99
Wireless signals have a direct negative impact on plant growth.	2.22	1.28	2.15	1.20
The “redhead gene” is becoming extinct.	2.42	1.45	3.60	1.02
Drinking cold water can be bad for your health.	2.43	1.48	3.81	1.32
Antibacterial mouthwash helps cure colds and sore throats.	2.47	1.41	2.75	1.21
An at-home administration kit to screen for type-1 diabetes is currently being introduced.	2.48	1.33	2.52	1.09
Hospitals are busier on full moons.	3.10	1.62	2.93	1.36
St Bernard’s dogs once carried brandy barrels around their necks while rescuing people lost in the mountains.	3.43	1.63	2.29	1.01
If you pluck a grey hair, more grey hairs will arrive in its place.	3.53	1.54	2.27	0.91
Turkey meat makes you sleepy.	3.97	1.58	3.22	0.98

*Note.* \* indicates claims used in Experiments 1-3.

Table A2. *Familiarity and Believability Ratings of True Claims in Pilot Study*

Claim	Familiarity		Believability	
	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
*A man in Canada was still allowed to board his flight after a pipe bomb was found in his bag.	1.40	0.87	2.29	1.09
In Turkey, people do not chew gum at night due to a superstition that it represents chewing the flesh of the dead.	1.40	0.87	2.62	1.03
*The color of a chicken's egg is related to the chicken's earlobe color.	1.51	0.97	2.59	1.03
Hippopotamus milk is pink. [item excluded as actually found to be false.]	1.54	1.08	2.33	0.87
*Chicken carcasses can be used for renewable energy.	1.57	1.12	3.38	1.01
2014 was the deadliest year for flying on a plane, with 992 fatalities globally.	1.62	0.96	3.27	1.10
*The national animal of Scotland is the unicorn.	1.68	1.23	2.78	1.08
*Saudi Arabia has revealed plans for a \$500 billion "no fossil fuels" mega-city.	1.69	1.08	1.98	0.94
Exposure to microwaves can open the blood-brain barrier.	1.81	1.10	2.84	1.23
In 2015, Sweden imported nearly 1.3 million tons of waste from Norway, the UK, Ireland and others.	1.84	1.19	2.47	1.28
*Pessimism may be inherited due to a genetic mutation.	1.93	1.17	2.20	1.01
Dandelion root extract is being tested as a cancer treatment.	2.13	1.31	3.43	1.03
After the release of "The Hunger Games" in 2012, women's participation in archery rose by 105%.	2.20	1.37	2.74	1.03
Germany has officially removed any tuition fees for both local and international college students.	2.25	1.34	3.57	0.90
Honeybee stings are used in the treatment of arthritis.	2.29	1.49	2.75	1.22
The plague is still active in the US today.	2.30	1.49	3.00	0.95
26 civilians died in the conflict along the Ukraine-Russia border in the summer of 2017 alone.	2.35	1.31	2.78	1.05
The heart of a blue whale is so massive that a human being can swim through its arteries.	2.49	1.58	2.97	1.22
Coca Cola single-bottle production is over 110 billion per year.	2.52	1.34	3.63	1.37
A mattress doubles its weight after 10 years of usage.	2.57	1.56	3.52	1.21
Roughly 800 journalists have been killed globally over the last 10 years.	2.59	1.51	2.22	1.23
Nano-robots are being tested in the treatment of cancer.	2.78	1.45	3.20	1.32
Acne is a hereditary condition.	2.79	1.31	2.75	1.17
China is implementing a citizen ranking system to determine who is a good citizen.	3.47	1.71	2.21	1.29

*Note.* \* indicates claims used in Experiments 1-3.

**Inferential-reasoning Questions (R = reverse-coded)**

**False claims.**

1. Facebook does not care about the language used on its platform. (R)
2. Facebook is investing money into promoting inoffensive language on its platform.
3. Men in South Africa generally do not need to be concerned about sexual assault. (R)
4. The rising number of reported HIV cases in South Africa is partially due to a large number of male rape cases.
5. Industries concerned about gender equity can look to Manchester, UK, for solutions.
6. In Manchester's corporate environment, males are much more likely than females to be promoted to senior managerial positions. (R)
7. There should be an awareness campaign to educate consumers about the environmental risks associated with pineapple skin.
8. There is no need to worry about how to dispose of pineapple skin. (R)
9. Hugh Hefner was bankrupt when he died. (R)
10. The late Hugh Hefner was a philanthropist, supporting various charities.
11. Concrete has no impact on electronics. (R)
12. When taking out the battery of your car, it is important not to place it on a concrete floor.

**True claims.**

1. Canada is lax regarding security on commercial flights.
2. In Canada, suspected terrorists are given an immediate ban on flying. (R)
3. A farmer can look at a chicken and predict the colour of the egg (white or brown) it will lay.
4. Whether a chicken's egg is white or brown is completely random. (R)
5. The only profitable use of chickens lies in meat and egg production. (R)

6. In the future, it is likely that some of our energy will come from bio-matter such as animal remains.
7. Souvenir shops in Scotland are likely to stock unicorn figures.
8. Scotland's national animal can be found in most zoos. (R)
9. Saudi Arabia is investing billions of dollars into environmental sustainability.
10. Investment in renewable energy technology in Saudi Arabia is virtually non-existent.  
(R)
11. In the future, genetic testing will be able to tell you if your baby will grow into a pessimistic person.
12. Whether people become pessimistic depends entirely on their life experiences. (R)