1	Detecting concealed language knowledge via response times
2	
3	Abstract
4	In the present study, we introduce a response-time-based test that can be used to detect
5	concealed language knowledge, for various potential applications (e.g., espionage, border
6	control, counter-terrorism). In this test, the examinees are asked to respond to repeatedly
7	presented items, including a real word in the language tested (suspected to be known by the
8	examinee) and several pseudowords. A person who understands the tested language
9	recognizes the real word and tends to have slower responses to it as compared to the
10	pseudowords, and, thereby, can be distinguished from those who do not understand the
11	language. This was demonstrated in a series of experiments including diverse participants
12	tested for their native language (German, Hungarian, Polish, Russian; $n = 312$), for second
13	language (English, German; $n = 66$), and several control groups ($n = 192$).
14	Keywords: deception, language, linguistic profiling, concealed information test,
15	response time

17 **1 Introduction**

Methods for discerning the truthfulness of a person who purports to be a native 18 speaker of a language have been recorded throughout history, from at least as early as the 19 11th century BCE to present day, in various and often crucial scenarios, such as ferreting out 20 infiltrators at battlefronts or verifying asylum claims (Reath 2004; Speiser 1942). Conversely, 21 there is no known test for discerning the truthfulness of a person who denies the knowledge 22 23 of a given language – short of the often repeated anecdote related to the Stroop task, typically recounting how Russian agents during the Cold War were detected based on slower decision 24 25 on the color of written Russian color words (e.g., Marx and Hillix 1987, p. 410; Peirce and MacAskill 2018, p. 111). 26

The Stroop task does not actually seem optimal for reliable deception detection,¹ but 27 the espionage scenario does occur in reality and is likely to continue posing a serious threat 28 (Riehle 2020) – undetected cases can be of historical importance (e.g., Black 1987; Kern et al. 29 2010). Hence, a reliable test for revealing concealed language knowledge could be of great 30 value for intelligence agencies, top-secret research facilities, and other highly confidential 31 organizations. Language tests could also be used for border control, and, in particular, to 32 verify asylum claims: When applicants lack documentation, determining their language 33 knowledge may be used to infer geographical origin (McNamara et al. 2016; Reath 2004). 34

¹ First, such a simple test is highly susceptible to faking (Boskovic et al. 2018; Hu et al. 2015; Verschuere et al. 2009). Second, the prevalence of color blindness throughout the world is estimated to be around 4-5%, which means that the Stroop effect is diminished in at least 4-5% of the cases to begin with. Third, basic colors are denoted by only a few simple words, thereby, restricting test material – moreover, color words are often similar between languages, and also often known to people who are even just vaguely familiar with a given language. Fourth, it would not be easy to standardize response times of verbal responses, let alone to implement an automatic analysis. Finally, no empirical validation for language detection purposes exists, but the generally moderate effect sizes of the Stroop task (e.g., Homack 2004) foreshadow poor diagnostic efficiency: Very large effect sizes are needed for reliable individual-level lie detection (cf. Suchotzki et al. 2017).

35	Since the actual origin is often suspected (e.g., Reath 2004), testing for the knowledge of the
36	corresponding language could be used for either preliminary screening or additional support
37	for previous conjectures. ² Other potential applications include scenarios where a person
38	deceitfully denies the knowledge of a language to (a) avoid cooperation with police or
39	military investigations (i.e., a suspect may deny understanding the language of the authorities,
40	or, alternatively, knowing the local language in a foreign operation); (b) justify, in a legal
41	case, not having understood rules or warnings; (c) claim insurance for language deficiency.
42	Less dramatically, it could also be useful for screening in psycholinguistic experiments in
43	which participants are not supposed to understand a certain language.
44	Finally, a test for concealed language knowledge could be used to detect not only
45	natural language, but also cants ("cryptolects") and similar coded language: Not only
46	organized crime members, but also terrorists are known to use secret jargon (Koskensalo
47	2015). Revealing that someone understands such a jargon would clearly warrant serious
48	further scrutiny. Hence, such a test could be valuable for screening national security
49	personnel, passengers at sensitive transport areas, detained terrorist suspects, etc.
50	All in all, a reliable method for detecting whether or not a person understands a given
51	language could be useful in a variety of high-stakes scenarios. In the present study, we
52	introduce the first method validated for this purpose.

1.1 Task design 53

54

In our language detection test, the main items are two real words in the given tested language, and four pseudowords that are graphemically similar to the real words. These items 55

² One method regularly applied in at least a dozen first-world countries is "language analysis for determination of origin" (McNamara et al. 2016). This method is controversial partly due to its questionable validity. For example, the use of a single Pakistani word could lead examiners to believe that the applicant is Pakistani (and, therefore, ineligible for asylum), although this could very well be accidental (Reath 2004, pp. 217-218). Testing for Pakistani language knowledge could provide valuable additional evidence.

are sequentially presented in a random order. The examinee is asked to press a key for each 56 item: One of the two real words is designated (selected randomly in the beginning of the task) 57 as *target* and requires pressing one response key (Key I on a standard keyboard), while all 58 other items (the other real word and all four pseudowords) require pressing another response 59 key (Key E). The other real word serves as probe. We assumed that only those who 60 understand the language would see the probe as saliently different from pseudowords, and 61 that they would respond slower to the probe as compared to pseudowords (which thus serve 62 as *control* items) – and, thereby, based on probe-control (real word versus pseudowords) 63 64 response time (RT) differences, they can be distinguished from those who do not understand the language. 65

This expected effect in such a design is supported by a series of related deception 66 detection studies for concealed information (Suchotzki et al. 2017; Verschuere and De 67 Houwer 2011: Verschuere and Meijer 2014), although there is no entirely certain or widely 68 accepted explanation for the underlying mechanism. In our view, it is decisive that the target 69 and probe share at least two interrelated key features (from the perspective of a person who 70 recognizes the probe among the controls): (a) Both target and probe meaning stand out as 71 task-relevant (the target because it requires a different key, and the probe because it pertains 72 to the deception scenario and is thereby semantically salient), and (b) both are, thus, 73 infrequent items compared to the controls – and yet the probe needs to be categorized 74 together with the controls - leading to the response conflict for probes (Lukács and Ansorge 75 2019; Seymour and Schumacher 2009; Verschuere et al. 2015). It follows that the greater the 76 similarity between target and probe items (relative to the controls), the larger the response 77 78 conflict (Suchotzki et al. 2018) - hence our choice of real words for both target and probe. Finally, apart from the main items (probe, target, controls), we included two kinds of 79 fillers that were (same as the general task instructions) always in the language acknowledged 80

81	to be understood by the examinee: ³ (a) expressions referring to meaningfulness and
82	genuineness (e.g., "MEANINGFUL," "TRUE," etc.) that had to be categorized with the same
83	key as the target (and, thus, opposite to the probe and the controls), and (b) expressions
84	referring to meaninglessness and fakeness (e.g., "UNTRUE," "FAKE," etc.) that had to be
85	categorized with the same key as the probe and controls. It is assumed (Lukács et al. 2017)
86	that fillers further slow down responses to the probes (when recognized by a person who
87	speaks the language) because the probes have to be categorized together with the
88	semantically incompatible expressions referring to meaninglessness (Nosek et al. 2007;
89	Rosch et al. 1976). In addition, by increasing the complexity of the otherwise excessively
90	simple task, fillers prevent strategically focusing on the target and thereby ignoring, to some
91	extent, the probe and its meaningfulness and relevance (Anderson 1991; Hu et al. 2013;
92	Reber 1989; Verschuere et al. 2015; Visu-Petra et al. 2013).
93	To establish not only conceptual (task-relevance, frequency) but also semantic

94 correspondence between the probe and the meaningfulness-referring fillers – and thereby
95 further enhance the probe response conflict – the probes (and therefore the targets too) were
96 meaningfulness-referring words as well.

97 **1.2 Study structure**

In the first two experiments, the test was performed only by speakers of the tested language (English and German; conducted in behavioral laboratory with university students), demonstrating real word versus pseudoword RT differences. In the subsequent three experiments, nonspeakers were tested too, so that classification accuracy could be assessed (with German, Hungarian, Polish, and Russian, as tested languages; in online experiments sampled from very diverse general populations of the respective countries).

³ A suspect might claim to only speak English, but is suspected to also speak German. In this case, task instructions and fillers in the test are in English. Only the probe and target items are German words (and the controls are German-like pseudowords).

104 **2 Experiments 1 and 2**

In Experiments 1 and 2, we tested native German speakers for English knowledge, and for German knowledge, respectively. At the same time, we also examined (a) in Experiment 1, whether meaninglessness-referring words (in the tested language) could serve as better controls than pseudowords, and (b), in Experiment 2, whether pseudowords could serve as better fillers than meaninglessness-referring words (in the instructions' language).

110 **2.1 Method**

111 2.1.1 Participants

Native German speaking students fluent in English participated for course credit at the
behavioral laboratory of [university name removed for masked review]. Sample size was
decided on by optional stopping using Bayes factor (BF) criterion (BF exceeding 5 for the
main within-subject comparison in each given experiment).

In Experiment 1, the initial sample of 60 participants already fulfilled our criteria. The 116 data of 20 participants had to be excluded (6 due to technical issues, 2 due to too low 117 accuracy, 12 for not selecting all four English words correctly during the verification task at 118 the end of the test; as per preregistration), leaving 40 participants (age = 21.2 ± 3.3 ; 8 male). 119 In Experiment 2, the initial sample of 55 participants did not fulfill our criteria, hence 120 20 more participants were invited three times, at which point the criterion was fulfilled with 121 100 completed tests (as not all invitations were answered). Twenty participants' data had to 122 be excluded (1 due to too low accuracy, 6 for low LexTale score – see below), leaving 93 123

124 (age = 21.4 ± 3.1 ; 38 male).

125 **2.1.2** Procedure

Participants were told about the purpose of the experiment, and they were asked to imagine themselves, during the testing, in a scenario where it would be crucial for them to conceal the knowledge of the tested language (English or German).

129	The main task, in each test, contained four blocks, each with its own unique set of
130	probe, target, and four controls. Fillers were placed among these items in a random order, but
131	with the restrictions that each of the 9 fillers (3 meaningfulness-referring, 6 meaninglessness-
132	referring) preceded each of the 4 probes, 4 targets, and 16 controls exactly one time.
133	Participants had to press Key I when the target appeared, and Key E when the probe or a
134	control appeared. Whenever a meaningfulness-referring filler appeared, participants had to
135	press the Key I (same as for targets), while whenever a meaninglessness-referring filler
136	appeared, they had to press Key E (same as for probe and controls).
137	The probes and targets were always real and meaningfulness-referring words in the
138	tested language (English in Experiment 1, with German instructions; and German in
139	Experiment 2, with English instructions). In each block, each probe, target, and control was
140	repeated 18 times. In Experiment 1, for each participant, two blocks had pseudowords as
141	controls, and two other blocks had meaninglessness-referring English words as controls; see
142	Table 1.
143	
144	
145	
146	
1.47	
14/	
148	
149	
150	
151	

152 **Table 1**

Item Type	Example 1	Example 2	Correct Key
Target	meaningful	proper	#I
Probe	genuine	true	#E
Control	onscaft, wrute, sieringlest, deborent	unknown, wrong, fake, untrue	#E
Filler-T	bedeutsam, ve	rtraut, wahr ⁴	#I
Filler-NT	unbedeutend, unve unbekannt, and	ertraut, gefälscht, ere, sonstiges ⁵	#E

153 Item Types Examples for Experiment 1

154

Note. Each example depicts a possible set of all items in a single block. Example 1 shows 155 possible items in a block with pseudoword controls, while Example 2 shows possible items in 156 a block with meaninglessness-referring controls: The only difference between the two 157 conditions concerned these *controls* – the probe and target items are interchangeable (i.e., 158 they are randomly assigned in each condition from the same pool of words), and the fillers 159 are always identical. Filler-T: "target-side" meaningfulness-referring fillers; Filler-NT: 160 161 "nontarget-side" meaninglessness-referring fillers. 162 In Experiment 2, two blocks had, analogously to Experiment 1, meaninglessness-163

- referring English words (instruction language) as fillers to be categorized together with the probe and controls (Key *E*), and two other blocks had pseudoword fillers instead (Table 2).

⁴ Meaningful, familiar, true.

⁵ Meaningless, unfamiliar, fake, unknown, other, miscellaneous.

166

167 **Table 2**

168 Item Types Examples for Experiment 2

Item Type	Example 1	Example 2	Correct Key			
Target	bekannt ⁶	sinnvoll ⁷	#I			
Probe	vertraut	bedeutsam	#E			
Control	glätisch, redengig, pauflich, schlinst	plaucklos, hokisch, tintzlich, klotselig	#E			
Filler-T	true, meaningj	#I				
Filler-NT	untrue, fake, foreign, random, unfamiliar, invalid	ontreg, dake, saneign, mindaw, unamidiar, imbodal	#E			
Note. Example 1 sh	ows possible items in a blo	ock with meaninglessness	-referring Filler-NT			
items, while Example 2 shows possible items in a block with pseudoword Filler-NT items:						
The only difference is in Filler-NT; the probe, target (real German words), and controls						
(pseudowords) are interchangeable, and Filler-T items are always identical. Filler-T: "target-						
side" meaningfulness-referring fillers; Filler-NT: "nontarget-side" meaninglessness-referring						

- 174
- 175

169

170

171

172

173

The inter-trial interval randomly varied between 0.5 and 0.8 s. In case of an incorrect response or no response within 1 s, the caption "Inkorrekt!" ("incorrect!") or "Zu langsam!" ("too slow!") appeared in red color, respectively, for 0.5 s, followed by the next trial. The

⁶ Known.

fillers.

⁷ Sensible.

main task was preceded by three short practice rounds that included all items from the
upcoming first block, and participants had to repeat any round on which they had too few
correct responses in time (for further details see the analogous task in, e.g., Lukács and
Ansorge 2019). For analysis, only trials with a correct response between 0.15 s and 1 s were
used.

At the end of the language test, as a verification task, participants were shown all probes and controls, and were asked to select the probes (ensuring that they understood the language). The data of those who did not select all four probes correctly were excluded in Experiment 1 (but not in Experiment 2, since all participants were already verified native German speakers). Finally, all participants completed a LexTALE test for English language comprehension (Lemhöfer and Broersma 2012): The data of those with a score below 60% (minimum score for B2 level) were excluded.

To calculate illustrative areas under the curves (AUCs) for probe-control RT mean differences as predictors, we simulated nonspeaker groups for the RT data using 1,000 normally distributed values with a mean of zero and an *SD* derived from the corresponding empirical data as $SD_{real} \times 0.5 + 7$ ms (which has been shown to very closely approximate actual data; Lukács and Specker 2020).

For all five experiments, preregistrations, all testing material (working PsychoPy or JavaScript/HTML codes for each task), the lists of all tested real words and pseudowords in each given language (and detailed description of their origin, creation, and the corresponding selection mechanisms during testing), analysis scripts, collected data, and an online appendix with supplementary analyses are available via

201 <u>https://osf.io/p78u3/?view_only=b581a7a9af7c4a9f91377c920c5731a3</u>. [Temporarily

202 masked preregistration links: Exp. 1:

203 https://osf.io/fq42x/?view_only=6eb975b4221c403187f26ea989e94417, Exp. 2:

10

204	https://osf.io/tqr6j/?view_only=93e059a914434c44b180ed79d5602f19, Exp 3:
205	https://osf.io/sg32f/?view_only=6d3f1f66072345eca61372d797b75762, Exp. 4:
206	https://osf.io/2g76c/?view_only=c02550fde63841d6b04bcc6bab12aec9, Exp. 5:
207	https://osf.io/gdk92/?view_only=c64e367ee8e6434ea2080c7d31f6d2e0]
208	2.2 Results
209	Large differences (ranging from 43.3 to 156.7 ms) were found between probe and
210	control RTs in both experiments, indicating potential for high classification accuracy; see
211	Table 3. In the within-subject comparison of Experiment 1, blocks with pseudoword controls
212	proved to have 40.11 ms larger probe-control differences than those with meaninglessness-
213	referring controls, 95% CI [19.11, 61.12], <i>d</i> = 0.61, 95% CI [0.27, 0.95], <i>t</i> (39) = 3.86, <i>p</i>
214	$< .001$, $BF_{10} = 67.54$, indicating higher potential for pseudoword controls. In the within-
215	subject comparison of Experiment 2, blocks with pseudoword fillers and meaninglessness-
216	referring fillers had probe-control differences of only 5.20 ms difference in their magnitudes,
217	95% CI [-7.33, 17.74], $d = 0.09$, 95% CI [-0.12, 0.29], $t(92) = 0.82$, $p = .412$, $BF_{01} = 6.28$.
218	
219	
220	
221	
222	
223	
224	
225	
226	
227	
228	

229 **Table 3**

	Probe	Control	Target	Filler-NT	Filler-T	P - C	AUC _{sim}
Experiment 1							
Pseudoword	528±70	445±33	569±48	518±51	589±47	83.4±52.0	.928 [.886, .969]
Real word	520±66	477±47	556±46	498±56	589±51	43.3±38.9	.828 [.765, .890]
Experiment 2							
Pseudoword	603±70	446±43	569±45	441±44	592±52	156.7±47.5	.998 [.995, 1]
Real word	606±78	454±42	579±48	498±49	605±54	151.5±57.4	.991 [.982, 1]

230 Response Times and Simulated AUCs in Experiments 1 and 2

231

232 Note. Means and SDs for individual RT means (ms) for different item types, and for probe-

233 control differences (P - C), and corresponding simulated AUCs (as AUC_{sim}, with 95% CIs in

234 brackets). *Pseudoword* denotes pseudoword controls in Experiment 1, and pseudoword Filler-

235 NT items in Experiment 2; *Real word* denotes meaningfulness-referring controls in

236 Experiment 1, and meaningfulness-referring Filler-NT items in Experiment 2. Filler-T:

237 "target-side" meaningfulness-referring fillers; Filler-NT: "nontarget-side" meaninglessness-

referring fillers; AUC: area under the curve.

239

3 Experiments 3-5

In Experiments 3, 4, and 5, we tested Hungarian natives for German (as a second language) and for Hungarian (native language), and Polish and Russian native speakers for their native languages – and, for all these cases, we also tested respective nonspeaker control groups.

245 **3.1 Method**

246 **3.1.1 Participants**

Participants for Experiments 3-5 were recruited via the online crowdsourcing platform Prolific (https://www.prolific.co/). The information regarding native and second languages were self-reported by participants on Prolific, and we invited only those who fulfilled our required criteria (e.g., "Hungarian native" and "fluent in German," for testing Hungarian natives for German as a second language).

For Experiments 3 and 4, the preregistered sample sizes were based on the estimated 252 available participants on Prolific. In Experiment 3, the sample was also limited by the 253 254 actually participating Hungarian participants (hence, collection was stopped, despite not having reached the goal of 50 participants, after 15 days, as preregistered): 41 German 255 speakers and 33 nonspeakers participated out of which 19 had to be excluded (14 for too low 256 German LexTALE score, 5 for too low accuracy), leaving 26 speakers (age = 28.4 ± 6.9 ; 17 257 male), and 29 nonspeakers (age = 29.0 ± 8.0 ; 9 male). Participants were paid 3.10 GBP for the 258 20-25 min experiment, and a potential 1.55 GBP bonus if they were not detected as 259 understanding German.⁸ 260

In Experiment 4, 50 Hungarian natives and 50 Polish natives participated, both tested for Hungarian as well as Polish (simulating a scenario where two different concealed languages are suspected), hence serving as each other's control groups – out of which 5 had to be excluded (3 for not selecting probes correctly, 2 for too low accuracy), leaving 49 Hungarian (age = 26.2±6.8; 38 male) and 46 Polish speakers (age = 25.1±6.9; 37 male). Participants were paid 4.88 GBP for the 20-25 min experiment, and a potential 0.50 GBP bonus for not having been detected in either language (hence altogether max. 1.00 GBP).

⁸ Successful detection for this purpose and for automatic feedback, was based on a d = 0.3 (standardized mean difference) between probe and control RTs, a higher level than in previous studies to favor participants (Noordraven and Verschuere 2013).

In Experiment 5, we again used optional stopping (see Footnote 10), which was fulfilled after 130 Russian native speakers (two additions of 30 following an initial 70) and 70 English monolinguals participated, out of which 8 had to be excluded (1 for not selecting probes correctly, 7 for too low accuracy), leaving 124 Russian natives (age = 31.8 ± 10.4 [1 unknown]; 46 male) and 68 English monolinguals (age = 31.1 ± 9.4 ; 36 male). Participants were paid 3.28 GBP for the 25-30 min experiment, and a potential 0.50 GBP bonus for not having been detected as understanding Russian.

275 **3.1.2** Procedure

The procedure and tasks were the same as in Experiment 1 and 2, unless otherwise noted.

Participants were told about the purpose of the experiment and were asked to imagine 278 themselves, during the testing, in a scenario where it would be crucial for them to conceal any 279 knowledge of the tested language (or both tested languages, in case of Experiment 4). 280 In Experiment 3, Hungarian native speakers were tested for German and had 281 Hungarian task instructions and fillers. In Experiments 4 and 5 (with participants tested for 282 their native languages), instructions and fillers were in English. Probes and targets were 283 always meaningfulness-referring words in the respective tested language, while the controls 284 were corresponding pseudowords. The target-side fillers were always meaningfulness-285 referring expressions. 286

The nontarget-side fillers were meaninglessness-referring expressions in Experiment 3. For a random half of participants in Experiment 4 and for all in Experiment 5, the nontargetside fillers were meaninglessness-referring expressions in two of the four blocks, but shuffled-letter items in the other two blocks. Preceding filler type change in either case, one short practice round had to be passed before commencing the given block with the new fillers. The six unique probe, target, and controls in each given block served as the basis of the six

293	nontarget shuffled-letter fillers. The given item's letters were reshuffled for each
294	presentation.9
295	At the end of the test, participants were shown all probes and controls, and were asked
296	to select the probes. As a precaution, in Experiments 4 and 5, the data of those who did not
297	select at least three probes correctly in their native language were excluded. In Experiment 3,
298	participants completed a LexTALE for German, and the data of those with a score below
299	60% were excluded. (In Experiments 4 and 5, participants completed LexTALE for English
300	for potential exploratory analysis.)
301	3.2 Results
302	For all three experiments, AUCs and related data are shown in Table 4.
303	
304	
305	
306	
307	
308	

⁹ We hypothesized that shuffled-letter items may be a better mental representation of meaninglessness (nonwords, nonsense words, pseudowords) than words that refer to meaninglessness (and yet are actually meaningful, i.e., existing words), and thereby lead to larger probe-control differences. The BF for comparing the two versions was also used for optional stopping of participant collection. The detailed report on this manipulation, and our related test length analyses, are available via https://osf.io/p78u3/?view_only=b581a7a9af7c4a9f91377c920c5731a3, and are planned to be published in a separate paper. In short, the results seem in line with our preregistered expectations, although their benefit for classification accuracy is limited. They have no important role in the results reported below: The changes in AUCs in particular, depending on different conditions, are relatively small (within 5%).

309 **Table 4**

	Probe	Control	Target	Filler-NT	Filler-T	P – C	AUC	TPR	TNR
Exp. 3 (GE)									
Speaker	519±48	492±39	598±50	573±50	605±53	27.0±25.3	.708	50	76
Nonspeaker	· 503±38	493±36	592±47	589±43	609±52	9.4±17.6	[.571, .845]	.38	./0
Exp. 4 (HU)									
Speaker	565±70	469±44	577±57	500±48	607±55	95.6±43.8	.992	02	1
Nonspeaker	455±37	461±38	580±39	526±56	587±39	-6.1±13.1	[.982, 1]	.92	1
Exp. 4 (PL)									
Speaker	549±64	489±47	584±53	494±52	601±47	60.0±32.8	.980	01	.98
Nonspeaker	487±54	488±55	595±61	521±60	590±49	-1.3±10.4	[.959, 1]	.91	
Exp. 5 (RU)									
Speaker	586±79	500±54	616±57	517±54	621±53	85.8±52.1	.939	.86	.96
Nonspeaker	497±62	496±61	616±62	539±63	596±57	1.0±17.1	[.906, .972]		
Note. Means	and SDs f	for indivi	dual RT	means (ms	s) for diff	erent item	types, and fo	r prob	e-
control differ	ences (P -	– C), and	, most in	nportantly,	correspo	nding AU	Cs (95% CIs	in	
brackets). TP	R: true po	ositive rat	tes (ratio	of correct	ly detecte	ed speakers), TNR: true	negat	ive
rates (ratio of correctly detected speakers), using arbitrary optimal cutoffs (maximal									
Youden's index) for classification. Filler-T: "target-side" meaningfulness-referring fillers:									
Filler-NT [•] "nontarget-side" meaninglessness-referring fillers AUC [•] area under the curve [•] GF [•]									
Cormon: HII: Hungarian: DI : Daligh: DII: Duggian									
Exclu	ding part	icipants	who are	suspected	of not c	omplying	with the requ	uiremo	ent is

310 *Response Times and AUCs in Experiments 3-5*

reasonable, but it may be argued that it limits the generalizability of our results and restricts conclusions. Therefore, we exploratorily recalculated AUCs using all participants in all three experiments without any exclusions. The changes in the AUCs (cf. Table 4) are negligible: .693, 95% CI [.572, .813] (TPR = .59, TNR = .79; 41 speakers, 33 nonspeakers) in Experiment 3; .992, 95% CI [.981, 1] (TPR = .92, TNR = 1; 50 speakers, 50 nonspeakers) in Experiment 4 for the Hungarian language test; .979, 95% CI [.958, 1] (TPR = .90, TNR = .98; 50 speakers, 50 nonspeakers) in Experiment 4 for the Polish language test; .931, 95% CI [.896, .966] (TPR = .85, TNR = .96; 130 speakers, 70 nonspeakers) in Experiment 5.

329 4 General discussion

First and foremost, we have demonstrated, based on testing three different languages, that our test can detect concealed *native* language knowledge with very high classification accuracy. We have also found strong evidence that the test provides classification accuracy well above chance level for second languages too – although it remains to be shown to what extent the knowledge of specific words and general fluency in the tested language might affect the outcomes.

The complex design of the test offers a number of opportunities for improvement and 336 fine-tuning: As one of many possibilities, in Experiment 1, we have shown that different 337 control items may affect classification accuracy. The test could also be specifically improved 338 for any given language by finding the optimal set of probe, target, and control items. For 339 example, while in our study the probes and targets were assigned randomly (for general proof 340 of concept), it seems likely that pairing close synonyms (as probe and target) would work 341 even better. Testing for concealed language knowledge, as compared to other kinds of 342 deception, is particular in that it does not require experimental setup, such as a mock-crime, 343 to simulate an appropriate scenario. Ground truth is relatively easy to establish (e.g., via a 344 preliminary interview in the examinee's native language), and it is likely that real suspects 345 are no more difficult to detect than experimental participants (Kleinberg, and Verschuere 346

- 348 would be crucial before real-life application.
- 349 We invite independent replications and further related research using freely available
- easy-to-use software for testing and evaluation (Lukács 2019). As explained in the
- introduction (Section 1), this novel method has wide-ranging potential for screening or for
- 352 providing additional evidence in various situations such as spotting spies, criminals,
- 353 terrorists, or detecting suspicious language-related inconsistencies in legal cases.

354

355	References
356	Anderson, J. R. 1991. The adaptive nature of human categorization. Psychological Review,
357	98(3), 409-429. https://doi.org/10.1037/0033-295X.98.3.409
358	Black, I. 1987. The origins of Israeli intelligence. Intelligence and National Security, 2(4),
359	151-156. https://doi.org/10.1080/02684528708431920
360	Boskovic, I., Biermans, A. J., Merten, T., Jelicic, M., Hope, L., and Merckelbach, H. 2018.
361	The modified Stroop Task is susceptible to feigning: Stroop performance and
362	symptom over-endorsement in feigned test anxiety. Frontiers in Psychology, 9, 1195.
363	https://doi.org/10.3389/fpsyg.2018.01195
364	Homack, S. 2004. A meta-analysis of the sensitivity and specificity of the Stroop Color and
365	Word Test with children. Archives of Clinical Neuropsychology, 19(6), 725-743.
366	https://doi.org/10.1016/j.acn.2003.09.003
367	Hu, X., Bergström, Z. M., Bodenhausen, G. V., and Rosenfeld, J. P. 2015. Suppressing
368	unwanted autobiographical memories reduces their automatic influences: Evidence
369	from electrophysiology and an Implicit Autobiographical Memory Test.
370	Psychological Science, 26(7), 1098-1106. https://doi.org/10.1177/0956797615575734
371	Hu, X., Evans, A., Wu, H., Lee, K., and Fu, G. 2013. An interfering dot-probe task facilitates
372	the detection of mock crime memory in a reaction time (RT)-based concealed
373	information test. Acta Psychologica, 142(2), 278–285.
374	https://doi.org/10.1016/j.actpsy.2012.12.006
375	Kern, G., Schecter, J. L., and Ransom Clark, J. 2010. The trouble with atomic spies.
376	Intelligence and National Security, 25(5), 705–724.
377	https://doi.org/10.1080/02684527.2010.537125
378	Kleinberg, B., and Verschuere, B. 2016. The role of motivation to avoid detection in reaction
379	time-based concealed information detection. Journal of Applied Research in Memory

- *and Cognition*, 5(1), 43–51. https://doi.org/10.1016/j.jarmac.2015.11.004
- 381 Koskensalo, A. 2015. Secret language use of criminals: Their implications to legislative
- institutions, police, and public social practices. *Sino-US English Teaching*, 12(7),
- 383 497–509. https://doi.org/10.17265/1539-8072/2015.07.005
- Lemhöfer, K., and Broersma, M. 2012. Introducing LexTALE: A quick and valid Lexical
- Test for Advanced Learners of English. *Behavior Research Methods*, 44(2), 325–343.
 https://doi.org/10.3758/s13428-011-0146-0
- 387 Lukács, G. 2019. CITapp A response time-based Concealed Information Test lie detector

388 web application. *Journal of Open Source Software*, *4*(34), 1179.

- 389 https://doi.org/10.21105/joss.01179
- 390 Lukács, G., and Ansorge, U. 2019. Information leakage in the response time-based Concealed

391 Information Test. *Applied Cognitive Psychology*, *33*(6), 1178–1196.

- 392 https://doi.org/10.1002/acp.3565
- Lukács, G., Kleinberg, B., and Verschuere, B. 2017. Familiarity-related fillers improve the
 validity of reaction time-based memory detection. *Journal of Applied Research in*
- 395 *Memory and Cognition*, 6(3), 295–305. https://doi.org/10.1016/j.jarmac.2017.01.013
- 396 Lukács, G., and Specker, E. 2020. Dispersion matters: Diagnostics and control data computer
- 397 simulation in Concealed Information Test studies. *PLOS ONE*, *15*(10), e0240259.
- 398 https://doi.org/10.1371/journal.pone.0240259
- Marx, M. H., and Hillix, W. A. 1987. *Systems and theories in psychology* (4th ed, Vol. 2.
 McGraw-Hill.
- McNamara, T., Van Den Hazelkamp, C., and Verrips, M. 2016. LADO as a language test:
 Issues of validity. *Applied Linguistics*, *37*(2), 262–283.
- 403 https://doi.org/10.1093/applin/amu023
- 404 Noordraven, E., and Verschuere, B. 2013. Predicting the sensitivity of the reaction time-

- 405 based Concealed Information Test: Detecting deception with the Concealed
- 406 Information Test. *Applied Cognitive Psychology*, 27(3), 328–335.
- 407 https://doi.org/10.1002/acp.2910
- 408 Nosek, B. A., Greenwald, A. G., and Banaji, M. R. 2007. The Implicit Association Test at
- 409 Age 7: A methodological and conceptual review. In *Social psychology and the*
- 410 *unconscious: The automaticity of higher mental processes* (pp. 265–292. Psychology
- 411 Press.
- 412 Peirce, J., and MacAskill, M. 2018. Building experiments in PsychoPy. Sage.
- 413 Reath, A. 2004. Language analysis in the context of the asylum process: Procedures, validity,
- 414 and consequences. *Language Assessment Quarterly*, *1*(4), 209–233.
- 415 https://doi.org/10.1207/s15434311laq0104_2
- 416 Reber, A. S. 1989. Implicit learning and tacit knowledge. Journal of Experimental
- 417 *Psychology: General*, *118*(3), 219–235. https://doi.org/10.1037/0096-3445.118.3.219
- Riehle, K. P. 2020. Russia's intelligence illegals program: An enduring asset. *Intelligence and National Security*, *35*(3), 385–402.
- 420 https://doi.org/10.1080/02684527.2020.1719460
- 421 Rosch, E., Mervis, C. B., Gray, W. D., Johnson, D. M., and Boyes-Braem, P. 1976. Basic
- 422 objects in natural categories. *Cognitive Psychology*, 8(3), 382–439.
- 423 https://doi.org/10.1016/0010-0285(76)90013-X
- 424 Seymour, T. L., and Schumacher, E. H. 2009. Electromyographic evidence for response
- 425 conflict in the exclude recognition task. *Cognitive, Affective, and Behavioral*
- 426 *Neuroscience*, *9*(1), 71–82. https://doi.org/10.3758/CABN.9.1.71
- 427 Speiser, E. A. 1942. The Shibboleth Incident (Judges 12:6. *Bulletin of the American Schools*428 *of Oriental Research*, 85, 10–13. https://doi.org/10.2307/1355052
- 429 Suchotzki, K., De Houwer, J., Kleinberg, B., and Verschuere, B. 2018. Using more different

- and more familiar targets improves the detection of concealed information. *Acta Psychologica*, 185, 65–71.
- Suchotzki, K., Kakavand, A., and Gamer, M. 2019. Validity of the reaction time Concealed
 Information Test in a prison sample. *Frontiers in Psychiatry*, *9*, Article 745.
- 434 https://doi.org/10.3389/fpsyt.2018.00745
- 435 Suchotzki, K., Verschuere, B., Van Bockstaele, B., Ben-Shakhar, G., and Crombez, G. 2017.
- 436 Lying takes time: A meta-analysis on reaction time measures of deception.

437 *Psychological Bulletin*, *143*(4), 428–453. https://doi.org/10.1037/bul0000087

- 438 Verschuere, B., and De Houwer, J. 2011. Detecting concealed information in less than a
- 439 second: Response latency-based measures. In B. Verschuere, G. Ben-Shakhar, and E.
- 440 Meijer (Eds.), Memory detection: Theory and application of the Concealed

441 *Information Test* (pp. 46–62). Cambridge University Press.

- Verschuere, B., Kleinberg, B., and Theocharidou, K. 2015. RT-based memory detection: Item
 saliency effects in the single-probe and the multiple-probe protocol. *Journal of Applied Research in Memory and Cognition*, 4(1), 59–65.
- $144 \qquad 11 ppilea Research in Memory and Cognition, <math>4(1), 55, 65.$
- 445 Verschuere, B., and Meijer, E. H. 2014. What's on your mind?: Recent advances in memory
- detection using the Concealed Information Test. European Psychologist, 19(3), 162–
- 447 171. https://doi.org/10.1027/1016-9040/a000194
- Verschuere, B., Prati, V., and Houwer, J. D. 2009. Cheating the lie detector: Faking in the
 autobiographical implicit association test. *Psychological Science*, *20*(4), 410–413.
- 450 https://doi.org/10.1111/j.1467-9280.2009.02308.x
- 451 Visu-Petra, G., Varga, M., Miclea, M., and Visu-Petra, L. 2013. When interference helps:
- 452 Increasing executive load to facilitate deception detection in the Concealed
- 453 Information Test. *Frontiers in Psychology*, *4*, Article 146.
- 454 https://doi.org/10.3389/fpsyg.2013.00146