

Applied Psycholinguistics, page 1 of 33, 2017
doi:[10.1017/S014271641700011X](https://doi.org/10.1017/S014271641700011X)

1 **Validity of new measures of implicit**
2 **knowledge: Distinguishing implicit**
3 **knowledge from automatized explicit**
4 **knowledge**

5 **YUICHI SUZUKI**
6 *Kanagawa University*

7 Received: July 28, 2016 Accepted for publication: March 27, 2017

ADDRESS FOR CORRESPONDENCE

Yuichi Suzuki, Department of Cross-Cultural Studies, Kanagawa University, 3-27-1, Rokkakubashi,
Kanagawa-ku, Yokohama-shi, Kanagawa, 221-8686, Japan. E-mail: szky819@kanagawa-u.ac.jp

8 **ABSTRACT**

9 Accumulating evidence suggests that time-pressured form-focused tasks like grammaticality judgment
10 tests (GJTs) can measure second language (L2) implicit knowledge. The current paper, however, pro-
11 poses that these tasks draw on automatized explicit knowledge. A battery of six grammar tests was
12 designed to distinguish automatized explicit knowledge and implicit knowledge. While three time-
13 pressured form-focused tasks (an auditory GJT, a visual GJT, and a fill in the blank test) were hypothe-
14 sized to measure automatized explicit knowledge, three real-time comprehension tasks (a visual-world
15 task, a word-monitoring task, and a self-paced reading task) were hypothesized to measure implicit
16 knowledge. One hundred advanced L2 Japanese learners with first language Chinese residing in Japan
17 took all six tests. Confirmatory factor analysis and multitrait-multimethod analysis provided an array
18 of evidence supporting that these tests assessed two types of linguistic knowledge separately with little
19 influence from the method effects. The results analyzed separately by length of residence in Japan (a
20 proxy for the amount of naturalistic L2 exposure) showed that learners with longer residence in Japan
21 can draw on implicit knowledge in the real-time comprehension tasks with more stability than those
22 with shorter residence. These findings indicate the potential of finely tuned real-time comprehension
23 tasks as measures of implicit knowledge.

24 The issue of implicit and explicit knowledge and learning mechanisms has attracted
25 attention from many second language acquisition (SLA) researchers because of its
26 theoretical and educational implications (e.g., Hulstijn, 2005). To tackle the issues
27 surrounding explicit and implicit knowledge and learning (e.g., interface issues),
28 the methodological problem of measuring implicit knowledge is crucial (Suzuki
29 & DeKeyser, 2017). Explicit and implicit knowledge are distinguished based on
30 awareness; implicit knowledge is deployed without awareness, whereas explicit
31 knowledge requires some level of awareness (DeKeyser, 2003; Williams, 2009).
32 Previous SLA studies have shown empirically that explicit and implicit knowledge
33 are distinct constructs that can be measured separately (Bowles, 2011; R. Ellis,
34 2005; Gutiérrez, 2013; Zhang, 2015). A recent study, however, employed a more

35 finely tuned psycholinguistic technique to examine real-time grammar processing
36 and cast doubt on the validity of existing implicit knowledge measures (Suzuki &
37 DeKeyser, 2015; Vafaei, Suzuki, & Kachinske, 2017).

38 The present paper reports a construct validation study of a new battery of finely
39 tuned tests for implicit knowledge: the eye-tracking-while-listening task, the word-
40 monitoring task, and the self-paced reading task. These real-time comprehension
41 tasks were compared with the existing tasks that have been claimed to measure im-
42 plicit knowledge, for example, timed grammaticality judgment tests (GJT), which
43 were hypothesized to draw more on automatized explicit knowledge in this study.

44 PROBLEMS IN PREVIOUS MEASURES OF “IMPLICIT” KNOWLEDGE IN 45 SLA

46 A seminal study by R. Ellis (2005, 2009) developed three tests that were hypoth-
47 esized to measure implicit knowledge: an oral narrative task, a timed GJT, and an
48 elicited imitation (EI) task. Since these tasks were performed under time pressure,
49 Ellis claimed that second language (L2) learners are more likely to draw on implicit
50 knowledge. He conducted a confirmatory factor analysis (CFA) and demonstrated
51 that these time-pressured tests were loaded onto a separate factor from an explicit
52 knowledge factor that untimed tests were loaded onto (i.e., an untimed GJT and
53 a metalinguistic knowledge test). This finding was essentially replicated in subse-
54 quent studies with different L2 populations (Ercetin & Alptekin, 2013; Gutiérrez,
55 Sarandi, 2015; Zhang, 2015) and a heritage learner population (Bowles,
56 2011).

57 The critical methodological factor that differentiated implicit knowledge from
58 “explicit” knowledge in those studies was imposing time pressure on the language
59 tests; however, time pressure cannot necessarily limit access to explicit knowledge
60 enough to ensure that implicit knowledge is drawn upon (DeKeyser, 2003; Suzuki
61 & DeKeyser, 2015; Vafaei et al., 2017). Proficient L2 learners may still access
62 explicit knowledge with awareness even if the execution is rapid (i.e., automa-
63 tized explicit knowledge), which is distinguished from the use of linguistic knowl-
64 edge *without* awareness (i.e., implicit knowledge). In other words, both implicit
65 knowledge and automatized explicit knowledge are accessed quickly, but they are
66 distinguished based on the awareness criterion. Highly automatized knowledge is
67 conscious knowledge that one can draw on quickly. It is functionally equivalent to
68 implicit knowledge in the sense that it is not easy to distinguish behaviorally (and
69 impossible to distinguish in mundane language use), but it remains knowledge
70 that one is aware of, and awareness is the defining criterion of explicit knowl-
71 edge. In cognitive psychology, automaticity (i.e., the end point of automatization)
72 is often characterized as lack of awareness (e.g., Jacoby, 1991; Posner & Snyder,
73 1975). However, automatization is a long process, and even highly automatized
74 skills do not always become 100% automatic, particularly with complex skills
75 like L2 learning. Automatization of explicit knowledge should be regarded as a
76 gradual development, not an all or nothing phenomenon (DeKeyser, 2015). In
77 the current paper, automatized explicit knowledge is thus defined as a body of
78 conscious linguistic knowledge including different levels of automatization. An
79 attempt is made to measure *partially* (not fully) automatized linguistic knowledge

80 with a conscious correlate, which can be theoretically distinguishable from implicit
81 knowledge.

82 A recent study provided evidence that it is possible to devise linguistic tasks that
83 can draw upon implicit knowledge separately from automatized explicit knowl-
84 edge (Suzuki & DeKeyser, 2015). Suzuki and DeKeyser (2015) demonstrated that
85 the EI task, which was the best measure of implicit knowledge in the test battery of
86 Ellis' (2009) study, did not measure implicit knowledge but drew on automatized
87 explicit knowledge. Even though time pressure was imposed and attention was di-
88 rected to meaning during the EI task, there appeared to be some room for accessing
89 automatized explicit knowledge for advanced L2 learners. Since EI tasks direct
90 learners' attention to meaning, it is certainly a better test of implicit knowledge
91 than form-focused tasks like the timed GJTs. These timed GJTs should be deemed
92 a measure of automatized explicit knowledge because they always require learn-
93 ers to pay attention to forms, which inevitably raises awareness of one's linguistic
94 knowledge (Vafaei et al., 2017). Of course, the level of awareness that each test
95 taker brings to the task may vary depending on his/her background. For instance,
96 native speakers and some L2 learners (e.g., heritage learners) with little experience
97 of learning of a L2 through formal instruction, who presumably possess little ex-
98 plicit knowledge, need to draw on implicit knowledge to perform a GJT, regardless
99 of its being timed or untimed. In contrast, learners with formal instruction may
100 tend to recourse to, or at least attempt to draw on, automatized explicit knowledge.
101 The present study focuses on L2 learners with some formal instruction and hy-
102 pothesizes that timed GJTs primarily draw on automatized explicit knowledge for
103 them.

104 THEORETICAL IMPORTANCE FOR DISTINGUISHING BETWEEN 105 IMPLICIT KNOWLEDGE AND AUTOMATIZED EXPLICIT KNOWLEDGE

106 The distinction between automatized explicit knowledge and implicit knowledge in
107 L2 learners has not been thoroughly researched. Differentiating linguistic knowl-
108 edge that has no conscious correlate (implicit) from that which involves con-
109 sciousness (explicit) but has been automatized bears important implications at
110 many levels. From an applied pedagogical perspective, the distinction may be triv-
111 ial practically (see further discussions in DeKeyser, 2015; Spada, 2015). From
112 a theoretical point of view, however, the distinction is indispensable for tackling
113 two related issues in explicit and implicit learning. The first problem concerns the
114 nature of linguistic knowledge types that L2 learners possess. By postulating au-
115 tomated explicit knowledge in addition to implicit and (nonautomated or less
116 automated) explicit knowledge, it allows for assessing L2 ability through more
117 scrutinized constructs. For instance, it is an empirical question of to what extent
118 L2 proficiency (e.g., measured by standardized tests) can be explained by implicit
119 knowledge, automatized explicit knowledge, and less automatized or nonautoma-
120 tized explicit knowledge (e.g., Elder & Ellis, 2009).

121 A more important point is that accurate identification/assessment of distinct
122 types of linguistic knowledge provides essential insight into the second problem,
123 that is, uncovering L2 learning processes. One of the central issues in the SLA field
124 is how explicit/implicit learning leads to the acquisition of implicit knowledge:

125 the interface issue. Many researchers express different positions as to whether
126 explicit knowledge facilitates the acquisition of implicit knowledge (DeKeyser,
127 2015; N. C. Ellis, 2005; Ellis, 2008; Hulstijn, 2002; Krashen, 1985; McLaughlin,
128 1987; Paradis, 2009). The theoretical distinction between automatized explicit
129 knowledge and implicit knowledge, along with valid measurements for them, can
130 advance our understanding of the interface issues in at least two related areas.

131 Explicit learning processes can be examined in more depth as L2 learning re-
132 sults in a large variability in the degree of automatization in L2 knowledge (e.g.,
133 DeKeyser, 1997, 2015). Implicit learning processes can be examined more closely
134 in relation to different types of explicit knowledge. A certain group of L2 learners
135 may first engage in explicit learning and succeed in attaining automaticity; they
136 may utilize automatized explicit knowledge to facilitate the acquisition of implicit
137 knowledge (Suzuki & DeKeyser, 2017). In contrast, a different type of L2 learner
138 possesses explicit knowledge, a large part of which is not automatized at all; these
139 learners may have to deploy implicit learning mechanisms in different ways from
140 the first group. It is also possible that the usefulness of explicit knowledge for im-
141 plicit learning varies depending on whether explicit knowledge is automatized or
142 not. The clearer operationalization and identification of the constructs are crucial
143 in revealing learning processes of different L2 learner populations through the lens
144 of explicit and implicit learning.

145 NEW MEASURES OF IMPLICIT KNOWLEDGE: REAL-TIME 146 COMPREHENSION TESTS

147 Following Suzuki and DeKeyser (2015), the current study proposes that implicit
148 knowledge is drawn upon when test takers *register*¹ specific grammatical structures
149 for real-time comprehension. Examining real-time grammar processing allows us
150 to capture whether learners can deploy their linguistic knowledge with very little
151 lag from the input; they are very unlikely to apply linguistic knowledge consciously
152 (Andringa & Curcic, 2015; Leung & Williams, 2012; Paradis, 2009; Suzuki &
153 DeKeyser, 2015). Only implicit knowledge makes it possible to operate at almost
154 the exact time of occurrence of targeted grammatical structures. More important,
155 measures of implicit knowledge should direct test takers' attention primarily to
156 meaning so that they do not raise awareness about grammatical structures to be
157 targeted. While form-focused tasks are direct measures of grammatical knowl-
158 edge, implicit knowledge tests are characterized as indirect measures. In what
159 follows, I will introduce three psycholinguistic measures that capture real-time
160 comprehension of grammatical structures, requiring no grammatical judgments
161 on the stimulus sentences. I will first discuss reaction-time measures with focus
162 on a word-monitoring task and a self-paced reading task. After that, I will intro-
163 duce a still newer method in the L2 field, an eye-tracking while-listening task (i.e.,
164 visual-world task).

165 An increasing interest in a psycholinguistic approach to SLA has developed over
166 the decades, leading to an increasing use of reaction time (RT) to examine online
167 sentence processing in L2 (for a review, see Jiang, 2011). Representative tasks
168 include the word-monitoring task (Granena, 2013; Suzuki & DeKeyser, 2015)
169 and the self-paced reading task (Foote, 2011; Jiang, Novokshanova, Masuda, &

170 Wang, 2011; Roberts & Liszka, 2013). The advantage of using these RT methods,
171 over form-focused tasks like GJTs, is that we can indirectly measure grammatical
172 sensitivity without asking for grammaticality judgments. In the word-monitoring
173 task, participants listen for a monitoring word and respond to it as soon as they
174 hear it in an auditory sentence by pressing the key on the computer. They pay
175 attention to the meaning of the sentences, rather than to the grammatical forms,
176 because comprehension questions are presented after they hear the sentence. The
177 monitoring word is embedded in an auditory sentence and occurs right after a
178 target grammatical structure. When participants listen for a monitoring word (e.g.,
179 to) in an ungrammatical sentence (e.g., The man likes to play basketball), they
180 are expected to slow down to respond to the target word, compared to the one
181 in the grammatical counterpart (e.g., The man likes to play basketball). The RT
182 difference between grammatical and ungrammatical indicates the extent to which
183 participants detect the errors without awareness. The same logic of assessment
184 of online grammatical sensitivity applies to the self-paced reading task where
185 participants read a sentence word by word on the computer (see Instruments).

186 RT-based research has stood as a gold standard for psycholinguistic studies; how-
187 ever, a more fine-grained measurement technique, a visual-world task, has started
188 to be utilized to capture real-time L2 grammar processing (Grüter, Lew-Williams,
189 & Fernald, 2012; Hopp, 2013; Lew-Williams & Fernald, 2010; Trenkic, Mirkovic,
190 & Altmann, 2014). In the visual-world task, participants see a visual scene con-
191 sisting of pictures while listening to stimulus sentences with target grammatical
192 structures. By analyzing eye movements during the listening process, it can reveal
193 real-time comprehension of grammatical structures (Sedivy, 2010; Tanenhaus &
194 Trueswell, 2006). The advantages of applying the visual-world paradigm to L2
195 research are summarized as follows: (a) it requires no ungrammatical sentences
196 (little risk of raising awareness), (b) it is a direct measure of fast and ballistic (un-
197 stoppable) linguistic processing in real time, (c) it is simple and can be applied to
198 wider populations, and (d) it enjoys higher ecological validity than RT tasks. All
199 in all, the three tasks, by virtue of measuring real-time comprehension process,
200 should each measure implicit knowledge.

201 AMOUNT OF L2 EXPOSURE INFLUENCES RELIANCE OF IMPLICIT 202 KNOWLEDGE

203 In addition to examining the relationships among the language test scores, the
204 current study aims to obtain further evidence for validity of the new implicit
205 knowledge measures. It takes a very long time to acquire implicit knowledge
206 because a large amount of L2 input for specific grammatical forms is required to
207 develop this (Paradis, 2009). Individual differences in the amount of L2 exposure
208 have been found to be related to the acquisition of implicit knowledge (Suzuki
209 & DeKeyser, 2015). The work by Suzuki and DeKeyser (2015) revealed that per-
210 formance for a measure of implicit knowledge (i.e., the word-monitoring task)
211 was correlated with scores of the implicit learning aptitude (i.e., measured by the
212 serial-reaction time task) *only* among the L2 learners with longer length of resi-
213 dence (LOR) in the immersion context. It is conceivable that L2 learners with
214 more L2 experience are more likely to rely on implicit knowledge stably on the

Table 1. *Task features of the linguistic knowledge measurements*

	Indirect-Implicit Measures			Direct-Explicit Measures		
	Visual-World	Word Monit.	Self-Paced	Timed AGJT	Timed VGJT	Timed SPOT
Data	Fixation Proportion	RT	RT	Accuracy	Accuracy	Accuracy
Real-time processing	Yes	Yes	Yes	Yes/No	Yes/No	No
Focus	Meaning	Meaning ^a	Meaning	Form	Form	Form
Time pressure	No	Yes	Yes	Yes	Yes	Yes
Modality	Aural	Aural	Written	Aural	Written	Written

Note: AGJT, auditory judgment test; VGJT, visual judgment test; SPOT, Simple Performance-Oriented Test; RT, reaction time.

^aThe focus of attention is also directed to the monitoring word.

215 language tests. Following Suzuki and DeKeyser (2015), the current study recruited
 216 Japanese L2 learners who live in Japan and then divided them into two groups based
 217 on their LOR. By examining the LOR in the immersion context (a proxy for the
 218 amount of L2 exposure), the current study contributes to a better understanding
 219 of the measurement and development of implicit knowledge. It can also poten-
 220 tially inform a more stringent participants selection procedure for testing implicit
 221 knowledge.

222 THE CURRENT STUDY

223 The aim of the current study was to examine the validity of the behavioral mea-
 224 sures that could measure implicit knowledge and automatized explicit knowledge
 225 separately. Six language tests were developed to assess linguistic knowledge of
 226 three Japanese grammatical structures and were administered to 100 Japanese L2
 227 learners. The three indirect, real-time comprehension measures of grammatical
 228 knowledge (the visual-world task, the word-monitoring task, and the self-paced
 229 reading task) were hypothesized to assess implicit knowledge, whereas the other
 230 direct, form-focused measures (the timed auditory GJT, the timed visual GJT, and
 231 the timed fill in the blank test called Simple Performance-Oriented Test [SPOT]²)
 232 were hypothesized to draw on automatized explicit knowledge.

233 As shown in Table 1, the crucial differences between the two types of measures
 234 lie in (a) real-time sentence processing and (b) focus of attention. All three online
 235 measurements assess whether test takers can incrementally process the sentence
 236 while their attention is focused on the meaning of the sentences. They are less likely
 237 to use linguistic knowledge consciously because their real-time grammatical pro-
 238 cessing is examined within a time window of few hundred milliseconds (Andringa
 239 & Curcic, 2015; Leung & Williams, 2012; Paradis, 2009; Suzuki & DeKeyser,
 240 2015). In contrast, the two types of GJTs and the SPOT require them to focus on
 241 form or grammatical target points under time pressure. Even if the time pressure

242 is imposed on them, they are more likely to use explicit knowledge because the
243 tasks inherently predispose them to focus on form (Vafaei et al., 2017).

244 There are some differences between the GJTs and the SPOT. First, the amount of
245 attention to form may be greater in the SPOT than in the GJTs. In the SPOT (i.e., fill
246 in the blank test), learners have to focus on specific grammatical structures to fill in
247 the blanks, whereas in the GJTs they do not know whether and where grammatical
248 errors are embedded in each test item and need to search for grammatical errors.
249 Second, the SPOT might not have imposed as strong incentives to respond quickly
250 as the GJTs to complete the task, because a longer time limit was set on each
251 test item in the SPOT than in the GJTs (see Methods). It is still possible that
252 some learners make grammatical judgments in real time on the timed GJTs; the
253 requirement for “real-time processing” is less certain for the GJTs (see Table 1).

254 In order to validate the measurements for implicit knowledge and automatized
255 explicit knowledge, CFA was conducted to assess construct validity. In contrast
256 to exploratory factor analysis, CFA is a better approach to estimate relationships
257 among measured variables because it allows for identifying latent constructs by tak-
258 ing into account the measurement errors. The CFA procedures consist of (a) initial
259 model specification, (b) model evaluation, and (c) rival model comparison. In the
260 initial model specification, a CFA model is specified in advance based on prior the-
261 ories. Here, the two-factor model was hypothesized (see the left panel in Figure 1).

262 In model evaluation, the model is then tested with the gathered data and evalu-
263 ated by a goodness of fit. The identified model is then assessed for parameter
264 estimates such as factor loadings, error variances, and correlations between fac-
265 tors. Each parameter provides important information to examine validity. Factor
266 loadings represent the amount of variance in a measured variable (e.g., timed au-
267 ditory GJT) explained by the factor. For instance, high factor loadings of timed
268 auditory GJT, timed visual GJT, and timed SPOT suggest that these three mea-
269 sures measure a common theoretical construct (e.g., automatized explicit knowl-
270 edge). In other words, they serve as supporting evidence for *convergent validity*
271 or the extent to which measured variables are related (Campbell & Fiske, 1959).
272 In contrast, *discriminant validity* refers to the extent to which a latent factor (e.g.,
273 implicit knowledge) discriminates from other latent factors (e.g., automatized ex-
274 plicit knowledge). Discriminant validity can be evaluated by examining the relation
275 between the factors. A weak relationship between the two factors that were hy-
276 pothesized in the current study indicates the dissociation between implicit and
277 automatized explicit knowledge.

278 Further evidence for the discriminant validity can be evaluated by a rival model
279 comparison. As shown in the right panel in Figure 1, the one-factor model can
280 be specified as a rival model against the two-factor model.³ The one-factor model
281 can be plausible because all the language tests assess a single type of linguistic
282 knowledge. If the two-factor model is found to be better than the one-factor model,
283 it suggests that all six measures are not tapping into a single construct but two
284 distinct constructs, hence supporting the discriminant validity.

285 The current study takes a further step to evaluate the construct validity by con-
286 ducting multitrait-multimethod (MTMM) analysis (Widaman, 1985). The key ad-
287 vantage of MTMM analysis is that it assesses the extent to which the traits (i.e.,
288 latent constructs) were measured validly by taking into account the method effects

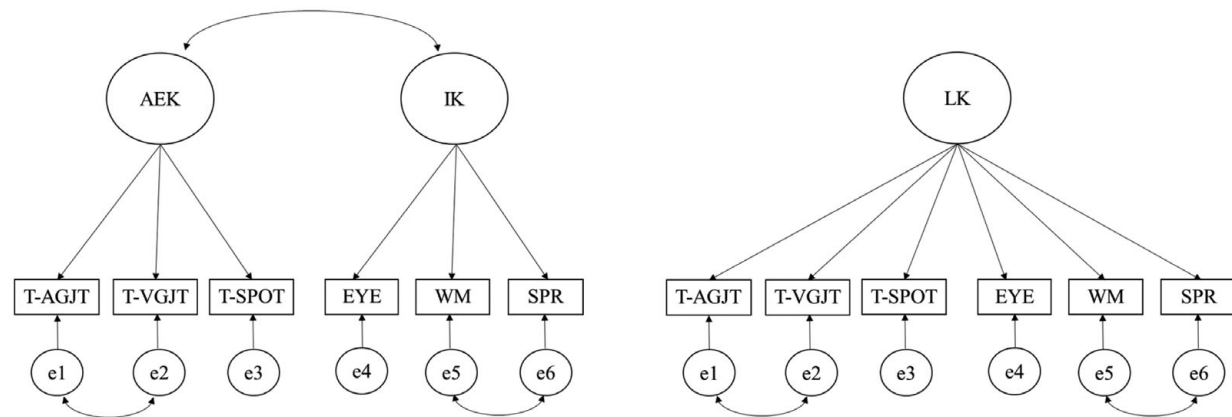


Figure 1. Confirmatory factor analysis models: (left) two-factor model and (right) one-factor model. AEK, automatized explicit knowledge; IK, implicit knowledge; LK, linguistic knowledge; T-AGJT, timed auditory grammaticality judgment test; T-VGJT, timed visual grammaticality judgment text; T-SPOT, timed SPOT; EYE, visual-world task; SPR, self-paced reading task; WM, word-monitoring task.

289 (see, e.g., Bachman & Palmer, 1982, for applications of MTMM in L2 learner
290 populations). The current study utilized a pair of measurements that shared very
291 similar methods (e.g., the visual GJT and the auditory GJT). This leads to a po-
292 tential threat to the validity because some portion of correlations between similar
293 measures can be simply due to the similarity in methods (method effects) not to the
294 underlying common construct. MTMM analysis allows for assessing the extent to
295 which variance in the measurements could be attributed to traits versus to methods
296 (Podsakoff, MacKenzie, & Podsakoff, 2012). Specifically for the purpose of this
297 study, the present study addressed the construct validity as to whether the traits
298 (automatized explicit and implicit knowledge) could be measured rather than the
299 method effects (RT and GJT measures). Error covariance was imposed on two pairs
300 of measurements that shared similar methods (see Figure 1): the word-monitoring
301 task and the self-paced reading task, which utilized similar RT measures while
302 listening or reading for comprehension, and the visual GJT and the auditory GJT,
303 which shared the same procedure except for the modality difference.

304 RESEARCH QUESTION AND HYPOTHESES

305 The current study addressed the question whether the test battery measures the
306 constructs of implicit knowledge and automatized explicit knowledge separately.
307 The two-factor model was evaluated by five criteria for construct validity. First, the
308 two-factor model was examined as to whether it fits the current data set. Second,
309 it was investigated to what extent the measurements assessed either automatized
310 explicit knowledge or implicit knowledge construct (convergent validity). Third,
311 it was investigated the extent to which the set of measurements for implicit knowl-
312 edge and those for automatized explicit knowledge were dissociated (discriminant
313 validity). This discriminant validity for the two factors was tested by (a) computing
314 correlations between the two factors and (b) comparing the two-factor model and
315 the one-factor model. Fourth, a MTMM analysis was performed to assess traits
316 and method effects. Fifth, the study examined if the amount of L2 exposure in the
317 immersion setting, estimated by the LOR, moderated the results of the four criteria
318 above. For these criteria, five hypotheses were put forth:

- 319 1. *Hypothesis 1*: The data structure of the six measurements demonstrates a good fit
320 to the two-factor model.
- 321 2. *Hypothesis 2*: The factor loadings are strong and significant (systematic) for au-
322 tomatized explicit knowledge and implicit knowledge (convergent validity).
- 323 3. *Hypothesis 3a*: The relationship between the two latent factors is insubstantial
324 (discriminant validity).
325 *Hypothesis 3b*: The data structure of the six measurements demonstrates a poor
326 fit to the one-factor model (discriminant validity).
- 327 4. *Hypothesis 4*: The error covariance between the similar measurement methods is
328 nonsignificant or smaller than the covariance between the measurements for the
329 traits (method effects).
- 330 5. *Hypothesis 5*: The results from L2 learners who received long-term exposure in
331 the immersion setting confirm Hypotheses 1–4 more convincingly than the results
332 from L2 learners with less exposure.

Table 2. *Background Information of the second language learners*

	Age at Testing	Age of Arrival	Onset of Instruction	LOR (months)	Length of Instruction (months)
Whole group ($n = 100$)					
Mean	25.97	21.36	19.01	47.29	41.11
SD	4.47	2.66	2.25	27.71	17.44
Range	19–47	17–30	13–27	24–197	3–84
Short-LOR group ($n = 48$)					
Mean	23.88	21.21	18.69	30.13	41.54
SD	2.72	2.63	1.82	4.33	17.16
Range	19–32	17–29	13–24	24–38	6–72
Long-LOR group ($n = 52$)					
Mean	27.90	21.50	19.31	63.13	40.71
SD	4.91	2.72	2.57	30.66	17.84
Range	22–47	17–30	13–27	39–197	3–84

Note: LOR, length of residence.

333 METHODS

334 *Participants*

335 One hundred Japanese L2 learners (29 male, 71 female), whose first language was
 336 Chinese, were recruited in Tokyo and the surrounding Kanto area. Four require-
 337 ments had to be met by L2 learners in order to participate in the study: proficiency,
 338 age of arrival in Japan, LOR, and educational background. First, only advanced-
 339 level Japanese L2 learners were recruited. They were screened for Japanese pro-
 340 ficiency, which had to be equivalent to or higher than N1 in the standardized
 341 Japanese Language Proficiency Test.⁴ Second, I only focused on late L2 learners,
 342 who arrived in Japan after the age of 17. Third, their LOR in Japan was 2 years
 343 or longer. This cutoff point for LOR was roughly based on the previous findings
 344 that implicit knowledge seems to be exhibited most efficiently in online measure-
 345 ments (i.e., the word-monitoring task) when L2 learners have been immersed in
 346 the target country for 2.5 years of residence or longer (Suzuki & DeKeyser, 2015).
 347 Fourth, participants possessed at least a bachelor's degree or were enrolled in a
 348 4-year college at the time of testing. The sampled population consisted of under-
 349 graduate students ($n = 34$), MA students ($n = 40$), PhD students ($n = 14$), and
 350 office workers ($n = 12$) at the time of testing. Forty-three out of 100 participants
 351 majored in Japanese language studies (i.e., Japanese or Japanese education as a
 352 foreign/second language) in undergraduate studies; 27 out of 61 participants with an
 353 MA degree or currently seeking one in Japanese language studies; and 5 out of 14
 354 participants with a doctorate degree or who are pursuing one in Japanese language
 355 studies. The rest of the participants' major varied (e.g., economic, architecture, en-
 356 gineering, management, law, psychology, physics, and liberal arts). Background
 357 information about the L2 learners is presented in Table 2.

358 The whole group was split in half by using the median LOR of 39 months
 359 (see Table 2). According to independent t tests, the two groups were significantly

360 different in terms of LOR and age at testing ($p < .001$). The other factors (age of
361 arrival, onset of instruction, and length of instruction) were not different ($p > .05$).
362 Fifty-one native speakers (NSs) were also recruited to serve as a baseline for the
363 linguistic knowledge tasks (see the Analysis section).

364 *Target structures*

365 Three Japanese linguistic structures were tested across the six language tests:
366 transitive/intransitive verb pairs, classifiers, and locative particles (*ni/de*). These
367 structures were chosen because they generate some prediction of upcoming infor-
368 mation, which can be demonstrated by the visual-world task. All target structures
369 are usually explicitly taught in beginner-level Japanese classes.

370 *Transitive–intransitive verb pairs.* Sixteen transitive–intransitive verb pairs were
371 chosen (Jacobsen, 1992). The pairs share the stem, but morphological markings
372 distinguish transitive from intransitive. For instance, the transitive verb *war-u* (to
373 break) has the intransitive counterpart *war-eru* (to be broken). A theme is dis-
374 cernible by the object-marking particle *o* for the transitive verb (e.g., *sara-o waru*;
375 someone breaks the dish). For the intransitive verb, the theme should be marked
376 with the subject-marking particle *ga* rather than *o* (e.g., *sara-ga wareru*; the dish
377 got broken). Note that action doer is implied in the transitive verb.

378 *Classifiers.* Eight classifiers were chosen and matched with 4 nouns; there were
379 32 classifier–noun pairs. For instance, *chaku* is a counter for clothes as in *go-chaku*
380 *no doresu* (*five-CHAKU-Genitive dress*; *five dresses*). Although some classifiers
381 are shared between Japanese and Chinese, we chose the classifier–noun pairs that
382 were not shared in order to avoid mere transfer from Chinese to Japanese (see
383 online-only supplementary material Appendix A).

384 *Locative particles: Ni/De.* The particles *ni* and *de* are multifunctional case markers,
385 and the usage for locations was focused on in the current study. In particular, *de*
386 indicates the place where an action takes place (e.g., *toshokan-de benkyousuru*;
387 study in the library), whereas *ni* is used to indicate the place where a thing or
388 a person exists (e.g., *toshokan-ni iru*; I will be in the library). It has been found
389 that Chinese speakers tend to overuse *ni* for *de* (Hasuike, 2004). Not all of the
390 usage for *ni* is difficult, and a relatively easier usage is expressing destination with
391 motion verbs (e.g., *cafe ni hairu*; enter the cafe). In sum, action verbs agree with
392 the location particle *de*, static verbs with the location particle *ni*, and motion verbs
393 with the destination particle *ni*.

394 *Instruments*

395 *Visual-world paradigm.* In the visual-world task, participants were first presented
396 with a scene consisting of four pictures on the computer screen for 5.5 s. They
397 then listened to sentences while their eye movements were being tracked, using
398 an EyeLink-II system (SR Research, Osgoode, Ontario, Canada) with a sampling
399 rate of 500 Hz. Participants were presented with a total of 64 trials: 48 critical

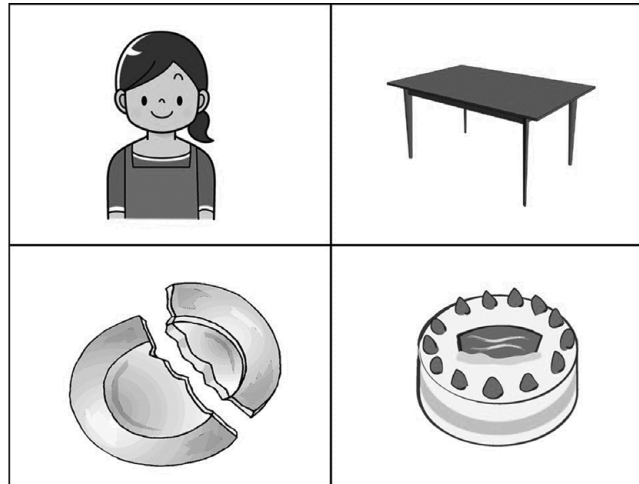
400 trials and 16 filler trials. Sixteen trials were prepared for each of the three lin-
 401 guistic structures tested, and participants heard two sentences for each trial. The
 402 critical sentence was always presented as the first sentence so that participants
 403 were not influenced by any information from the previous sentence (16 trials \times
 404 3 structures = 48 sentences). The second sentence acted as a filler to divert the
 405 participants' attention from the critical sentence and to avoid revealing the pur-
 406 pose of the study (48 filler grammatical sentences). There were also 16 filler trials,
 407 resulting in 32 filler grammatical sentences. All trials were presented in semiran-
 408 domized order such that the same trial type never occurred more than twice in a
 409 row. The location of the four objects on display was rotated across trials. After
 410 each trial, a yes/no comprehension question was asked to ensure that participants'
 411 attention was focused on the meaning of the sentence (cf. Dussias et al., 2013).
 412 Half of the questions asked about the critical sentence, and the other half about the
 413 filler sentence. Two practice trials were given to familiarize the participant with
 414 the procedure of the task.

415 The display always involved a target object and a competitor object. There
 416 were two types of trials for each target structure: target trials (where the target
 417 object was mentioned in the critical sentence) and the competitor trials (where
 418 the competitor object was mentioned). As shown in Figure 2, each display for
 419 transitive/intransitive structures consisted of a person (e.g., the mother), a contrast
 420 object (e.g., the table), a theme (e.g., the broken dish), and a distractor. The person
 421 was defined as a target, whereas the contrast object was defined as a competitor.
 422 Two types of critical trials were created: transitive and intransitive trials.

423 The first part of both sentences always followed the same form: *NP1-ACC-*
 424 *transitive verb-iru-no-wa-adverb-NP2* (It is NP2 that TRANSITIVE-VERB NP1)
 425 and *NP1-SUB-intransitive verb-iru-no-wa-adverb-NP2* (The reason is NP2 why
 426 NP1 INTRANSITIVE-VERB), where NP1, ACC, and NP2 are noun phrase 1,
 427 accusative, and noun phrase 2, respectively. NP2 was always a person (e.g., the
 428 mother) in the transitive trials (defined as target trials), whereas it was always
 429 a contrast object (e.g., the table) in the intransitive trials (defined as competitor
 430 trials). The eye movements were analyzed from the onset of the case marker (*ga*
 431 or *o*). If participants were sensitive to the transitivity of the verb, then looks to
 432 the target (e.g., mother) would be greater in the target trials than in the competitor
 433 trials. This is because a segment of NP-ACC and *te*-form of a transitive verb (i.e.,
 434 *osara-wo watte*) implied an action doer. The task design for the other two structures
 435 is described in online-only supplementary material Appendix B.

436 We were primarily interested in looks to the two possible locations, coded as tar-
 437 gets and competitors. Before the primary analyses, each time window was shifted
 438 200 ms after the linguistic cues in the speech stream to account for the time it takes
 439 to generate a saccadic eye movement (Matin, Shao, & Boff, 1993). In order to com-
 440 pute a "sensitivity index" for individuals, "target advantage (TA) scores" were first
 441 computed separately for target trials and competitor trials as follows: target looks
 442 divided by the sum of target looks and competitor looks. TA scores were then stan-
 443 dardized (z transformed) across the three structures, and the sensitivity index was
 444 computed by the TA difference scores as follows: TA in the target trials – TA in the
 445 competitor trials. The higher sensitivity score indicated more developed linguistic
 446 knowledge. These sensitivity scores were computed after time locking to 200 ms

Q1

**(1) Transitive trials (Target trials):**

Osara wo watte iru no wa soko ni iru okaasan desu.
Dish-ACC breaking be NOMINALIZER TOP there-LOC exist mother be.
(It is the mother that is breaking the dish.)

(2) Intransitive trials (Competitor trials):

Osara ga warete iru no wa soko ni aru teeburu kara ochite shimatta kara desu.
Dish-SUB breaking be NOM. TOP there-LOC exist table from fall off because.
(The dish is broken because it fell off the table.)

Figure 2. Visual scene and critical sentences for transitive/intransitive structure.

447 from the data-drive onset (see online-only supplementary material Appendix C for
448 details).

449 *Word-monitoring task.* In the word-monitoring task, participants were instructed
450 to listen to a sentence for a target word and to press a button as soon as they
451 identified it in the spoken sentence. The target word remained on the screen until
452 the response was made. A yes/no comprehension question appeared on the screen,
453 so that participants' attention was directed to the sentence meaning as well as
454 the target word. For instance, a sample sentence targeting transitive structure is
455 presented below.

[Target Word: *mazeru*]

456 *Ao to kiiri no enogu o/*ga mazeru to, kirei na mimdori ni naru.*
Blue and yellow paint-ACC/SUB mix if, beautiful green become
When you mix blue and yellow paints, it becomes beautiful green.

457 The target sentence always included a segment of the case-marking particle
458 (*ga* or *o*) and a verb (transitive or intransitive). The target word was always the

459 verb following the particles (*ga* or *o*). The “sensitivity index” was computed by
460 the RT difference scores (ungrammatical RT – grammatical RT) across the three
461 structures, after the average RTs of grammatical and ungrammatical items were
462 standardized (*z* scores) in order to treat the sensitivity across the target structures
463 equally. The magnitude of this sensitivity index is used to index how developed
464 one’s implicit knowledge is (Granena, 2013; Suzuki & DeKeyser, 2015).

465 The list of stimulus sentences included 48 target sentences (16 for each structure,
466 half grammatical and half ungrammatical) and 48 grammatical filler sentences.
467 Half of the items for each condition were followed by a yes/no comprehension
468 question. The ratio was kept equal between a positive response and a negative re-
469 sponse. Sample stimulus sentences for the other structures are presented in online-
470 only supplementary material Appendix D.

471 *Self-paced reading task.* In the self-paced reading task, participants were asked
472 to read a sentence word by word as quickly as possible while paying attention to
473 its meaning to answer a comprehension question accurately. The first word of a
474 sentence appeared on the left side of the screen, and when the button was pressed,
475 the next word appeared to the right of the preceding word, which disappeared
476 upon the presentation of the following word (moving-window presentation). When
477 participants read the final word followed by the period, they pressed a second key
478 to continue to either the next test item or a comprehension question. Words were
479 presented in Japanese characters in chunks consisting of a clause or *bunsetsu* (i.e.,
480 content word + function word). For example, a sample sentence with the transitive
481 structure is presented below (a slash indicates a unit of presentation).

483 [Region 1 = *mazeru to*, Region 2 = *ii*]
482 *Uta no gurupu o/ tsukuru tokini/ danshi to/ joshi o(*ga)/*
Singing group-OBJ/ make when boy and girls-OBJ
mazeru to/ ii/ baransu ni/ naru to omou/
484 mix if good balance becomes think
When you form a singing group, I think it makes a good balance if you mix boys and
girls.

485 The region of interest where RTs were compared between grammatical and
486 ungrammatical sentences was at the critical word where the error occurred in the
487 ungrammatical sentences (Region 1). This word was located in the same position
488 as that in the word-monitoring task so that the effects could be compared fairly
489 between the word-monitoring task and the self-paced reading task. RTs of the word
490 immediately following the critical word (Region 2) were also included to capture
491 spillover effects (Mitchell, 1984). In a similar way to the word-monitoring task,
492 the sensitivity index was computed for individuals as *z*-standardized RT scores
493 (ungrammatical RT – grammatical RT) at Regions 1 and 2 combined.

494 As in the word-monitoring task, a list of stimulus sentences included 48 tar-
495 get sentences (16 for each structure, half grammatical and half ungrammatical)
496 and 48 grammatical filler sentences. Half of the items for each condition were
497 followed by a yes/no comprehension question. The ratio was kept equal be-
498 tween a positive response and a negative response. Sample stimulus sentences

499 with the other structures are presented in the online-only supplementary material
500 Appendix E.

501 *Timed auditory GJT.* In the computer-delivered timed auditory GJT, participants
502 listened to an aural stimulus sentence and indicated whether each sentence was
503 grammatical or ungrammatical by pressing a response button. They were asked to
504 press a key as soon as an error was detected in the sentence. The time limit imposed
505 on the task was 10 s for each item. Responses that were longer than certain time
506 limits were then dealt with after administering the test (see Data Analysis section
507 for details). The stimulus sentences consisted of 48 target sentences (16 for each
508 structure, half grammatical and half ungrammatical) as well as 16 grammatical
509 filler sentences. Before the actual test, participants took a practice session. The
510 percentage accuracy score was calculated for all the items. One item in the auditory
511 GJT was excluded from the analyses because the sentence was not unambiguously
512 grammatical or ungrammatical (58% in NS accuracy rate).

513 *Timed written GJT.* As in the timed auditory GJT, the timed visual GJT was also
514 administered on a computer. The procedure was identical to the one in the timed
515 auditory GJT except for the modality of the stimulus sentences. Participants were
516 presented with a written sentence on a screen and asked to indicate whether each
517 sentence was grammatical or ungrammatical by pressing a response button as
518 quickly as possible. They were allowed to press the key while the sentence was
519 played when the error was detected within the sentence. The time limit imposed
520 on the task was 10 s for each item. The stimulus sentences consisted of 48 target
521 sentences (16 for each structure, half grammatical and half ungrammatical) as well
522 as 16 grammatical filler sentences. The percentage accuracy score was calculated
523 for all the items. One item in the visual GJT was excluded from the analyses because
524 the sentence was not unambiguously grammatical or ungrammatical (68% in NS
525 accuracy rate).

526 *Timed SPOT (fill in the blank test).* In the timed SPOT, the participants were
527 presented with a single sentence with some blanks on the computer screen. Then,
528 they had to fill in the blank with Japanese characters on the answer sheet as quickly
529 as possible. A blank was left in each sentence to specifically target one of the
530 linguistic structures. Once they filled in the answer on the sheet, they pressed a
531 computer button to move on to the next item. Participants were told to respond
532 as quickly as possible. The time limit for each test item was accidentally set to
533 100 s, instead of 10 s (see Data Analysis section). The number of characters to be
534 filled in the sentence was indicated by the number of blank circles in the sentence
535 (see sample items in online-only supplementary material Appendix F). A syllabic
536 *hiragana* character was used to fill in the blanks. The stimulus set consisted of
537 48 target sentences (16 for each structure) and 16 filler sentences. The percentage
538 accuracy score was calculated over all items for the target sentences.

539 *Procedure*

540 Participants were tested individually in a soundproof booth. After the consent
541 form and the background questionnaire, the linguistic tasks were administered

542 in fixed order from the most implicit linguistic tasks to the more explicit: the
543 visual-world task, the word-monitoring task, the self-paced reading task, the timed
544 auditory GJT, the timed visual GJT, and the timed SPOT. Before taking each task,
545 participants were presented with several practice items to familiarize them with
546 the procedure. All the stimulus sentences were different across the tasks in order
547 to reduce practice effects. They were presented in a fixed semirandom order in
548 each task, interspersing different types of stimulus sentences, in order to conceal
549 the purpose of the study. It took approximately 2 hr to complete the tasks, and
550 participants were given two 3-min breaks, one after the visual-world task and
551 another after the self-paced reading task.

552 *Data analysis*

553 *Real-time comprehension tasks.* For all three implicit knowledge tests (real-time
554 comprehension tasks), data cleaning procedures were conducted. Specifically, ac-
555 curacy of the comprehension questions was computed. A participant whose error
556 rate was higher than 25% would be excluded from further analysis to ensure that
557 each individual was paying attention to meaning (Jiang et al., 2011). None of the
558 participants scored below 75%; all participants' eye-movement and RT data were
559 analyzed. More detailed results from data cleaning procedures are presented in
560 online-only supplementary material Appendix G.

561 *Timed form-focused tasks.* Previous studies like R. Ellis (2005) and Bowles
562 (2011) set the time limit for *presenting* each sentence based on the NSs' aver-
563 age response time plus an additional 20% of the time for each sentence. A more
564 lenient time pressure was imposed on the current tasks: 10 s across all the test
565 items. Instead of imposing a strict time-out for duration of sentence presentation,
566 L2 learners' responses were screened after the data was collected. If the response
567 time was not within a certain time limit based on the NSs' RTs, those responses
568 were scored as incorrect. Initial review of data revealed that around 15%–30%
569 of the responses would be discarded *even for NSs' responses* in the three form-
570 focused tasks when we imposed the 20% + NSs' RT for each item. It seemed
571 more reasonable to impose time pressure in which most NSs can perform the task
572 accurately. We decided to identify a different percentage value so that 90% of the
573 NSs' responses were scored correct. In other words, percentages to be added to
574 the NSs' mean RT were determined such that the NSs' mean error rate of the total
575 score was kept less than 10%. The cutoff percentages that retained 90% of NSs
576 data were mean RTs + 50% for the auditory GJT, mean RTs + 120% for the visual
577 GJT, and mean RTs + 50% for the SPOT. These cutoff points were used to score
578 the responses of L2 learners in the three tests.

579 *Data summary: Missing data and data transformation*

580 Before presenting the results for L2 learners, native speakers' performance on
581 the six language tests was checked (see online-only supplementary material Ap-
582 pendix H). They showed sensitivity to the manipulation of stimulus sentences in
583 the meaning-focused tests (visual-world, word-monitoring, and self-paced reading

Table 3. *Descriptive statistics for the language tests by second language learners*

	<i>N</i>	Possible			Min	Max	95% CI	Cronbach α
		Max	<i>M</i>	<i>SD</i>				
Eye ^a	100	—	0.01	0.09	-0.26	0.24	[-0.01, 0.03]	—
WM ^a	100	—	22	54	-111	162	[11, 33]	0.91
SPR ^a	100	—	36	90	-198	351	[18, 54]	0.96
T-AGJT ^b	100	100	43.43	12.12	14.58	76.19	[0.41, 0.46]	0.67
T-VJGT ^b	100	100	30.64	16.28	2.08	82.74	[0.27, 0.34]	0.85
T-SPOT ^b	99	100	27.13	23.37	0	91.67	[0.20, 0.29]	0.95

Note: Eye, visual-world task; WM, word-monitoring task; SPR, self-paced reading task; T-AGJT, timed auditory judgment test; T-VJGT, timed visual judgment test; T-SPOT, timed Simple Performance-Oriented Test.

^aThe values for the online comprehension tasks indicate sensitivity index.

^bThe values for the form-focused tasks indicate percentage accuracy score.

584 tasks) and high accuracy (all above 90% in accuracy) in the form-focused tasks
585 (auditory and written GJTs and SPOT).

586 Descriptive statistics for all the measures performed by L2 learners are presented
587 in Table 3. L2 learners showed no sensitivity to the target structures in the visual-
588 world task, whereas they demonstrated some sensitivity in the word-monitoring
589 and the self-paced reading tasks (see online-only supplementary material Appendix
590 I for details). Their performance on the form-focused tasks was low; they scored
591 highest on the timed auditory GJT, followed by the timed visual GJT, and then the
592 timed SPOT.

593 Reliability indices were all above .65 and deemed acceptable (Loewenthal,
594 2004). The timed auditory GJT showed lower reliability (.67) than the other
595 form-focused tasks in the test battery perhaps because the test takers had only
596 one chance to listen to a spoken stimulus sentence. The internal consistency
597 (e.g., Cronbach α) was not computed for the visual-world task in the current
598 study because no standard procedure exists for estimating internal consistency
599 of the visual-world task; one promising approach is to examine test-retest reli-
600 ability (Farris-Trimble & McMurray, 2013). Since the test-retest reliability was
601 not available in the current study, the current findings should be interpreted with
602 caution.

603 Prior to the CFA and MTMM analyses, tests of univariate normality were exam-
604 ined for the six test scores. The total scores of the T-SPOT were positively skewed;
605 square root transformation was applied to reduce skewness. Based on the stan-
606 dardized coefficients of skewness and kurtosis (*z* scores), all the variables met the
607 assumption of univariate normality ($p > .05$). Multivariate normality of the score
608 distribution was examined by Mardia's coefficient. The coefficients (chi-square)
609 were 1.648 ($p = .439$) for all the six tests and 0.007 ($p = .996$) for the five tests,
610 both of which met the assumption of multivariate normality. Out of 100 partici-
611 pants, only 1 participant had missing cases in T-SPOT. Since this person was the
612 only one who had a missing case in the language tests, this person was excluded
613 from the analyses.

614 *CFA and MTMM analysis*

615 The two hypothesized CFA models (one-factor and two-factor models) were en-
616 tered into the CFA analyses (Figure 1). All the analyses were implemented in the
617 software package LISREL 9.1 (Jöreskog & Sörbom, 2013). Five hypotheses were
618 evaluated. The models were evaluated statistically with a maximum likelihood
619 method to estimate the model parameters (Hypothesis 1). Multiple fit indices were
620 jointly used to assess the model fit in addition to the chi-square statistics (Brown,
621 2006; Hoyle & Panter, 1995). The following three categories of fit indices were
622 utilized to assess the overall goodness of fit of the CFA models: absolute fit in-
623 dices (standardized root mean square), incremental fit indices (the comparative
624 fit index and the Bentler–Bonnet nonnormed fit index), and fit indices adjusting
625 for model parsimony (root mean square error of approximation). According to
626 the findings of simulation studies conducted by Hu and Bentler (1999), a good
627 fit between the target model and the observed data (maximum likelihood esti-
628 mation) was obtained in instances where standardized root mean square residual
629 values were below 0.09, root mean square error of approximation values were
630 below 0.06, and comparative fit index and Bentler–Bonnet nonnormed fit index
631 were above 0.96. Based on these empirically derived criteria, each of the mod-
632 els was assessed to exhibit one of three levels of fit: good fit, marginal fit, and
633 poor fit. When the indices in two or three out of three categories met the criteria
634 above, the model was considered to be a good fit (Hu & Bentler, 1999). When
635 none of the fit indices reach the criteria, the model was considered to be a poor
636 fit.

637 In order to seek evidence for convergent validity (Hypothesis 2), the magni-
638 tudes and significance of the factor loadings were examined. The discriminant
639 validity was assessed by the correlation between the two latent factors (Hypoth-
640 esis 3a). In addition, the discriminant validity was also evaluated by compar-
641 ing the one-factor and two-factor models by the goodness of fit testing indexed
642 by the chi-square statistics as the two models were nested (Hypothesis 3b). A
643 correlated uniqueness model, which is an alternative MTMM approach (Brown,
644 2006), was constructed to determine the extent to which variance in the mea-
645 surements could be attributed to latent constructs of linguistic knowledge (traits)
646 and to specific methods (Hypothesis 4). This model correlated the error be-
647 tween the timed visual GJT and the timed auditory GJT and the one between
648 the word-monitoring task and the self-paced reading task.⁵ Since the model in-
649 volved the two traits and two methods, the factor loadings on the same trait factor
650 were constrained to equality (Brown, 2006, p. 220). Finally, the same analyses
651 above were conducted separately for the short-LOR and the long-LOR groups
652 (Hypothesis 5).

653 RESULTS

654 Pearson's correlation coefficients will be presented first among the six language
655 test scores, followed by the results from the two competing CFA models with the
656 whole group, short-LOR group, and the long-LOR group. After that, results from
657 MTMM analyses will be presented.

Table 4. *Intercorrelations of the language tests (whole group, n = 99)*

	Eye	WM	SPR	T-AGJT	T-VGJT	T-SPOT
Eye	—	.093	.129	.153	.185	.212*
WM		—	.261**	.060	-.074	.057
SPR			—	.164	.073	.102
T-AGJT				—	.681**	.508**
T-VGJT					—	.553**
T-SPOT						—

Note: Eye, visual-world task; WM, word-monitoring task; SPR, self-paced reading task; T-AGJT, timed auditory judgment test; T-VGJT, timed visual judgment test; T-SPOT, timed Simple Performance-Oriented Test.

* $p < .05$. ** $p < .01$.

658 CFA

659 Table 4 shows the correlation matrix for the six linguistic test scores for the whole
660 group of L2 learners. Significantly moderate relationships were found among the
661 timed form-focused tasks ($.508 < r < .681$), whereas the correlations among the
662 three online tests were weak, and the only significant relationship among the on-
663 line measures, between the word-monitoring and the self-paced reading tasks, was
664 weak ($r = .261, p = .009$). Unexpectedly, the visual world task was significantly
665 correlated only with T-SPOT, possibly because both tests did not use any ungram-
666 matical sentences.

667 Both two-factor and one-factor models produced a good fit (see Table 5). A chi-
668 square difference test was conducted to compare the two-factor and the one-factor
669 models. The two-factor model fit better than the one-factor model at the descriptive
670 level, but the difference was not significant, $\chi^2_{\text{difference}} = 0.897, df = 1, p = .344$.

671 Figure 3 presents both models with factor loadings and significant correlated
672 errors. In the two-factor model, the two latent factors were moderately correlated
673 ($r = .47, p = .069$). Factor loadings for automatized explicit knowledge were high
674 and significant, whereas those for implicit knowledge were much lower and the path
675 to the visual-world task (EYE) was only marginally significant. For the one-factor
676 model, the factor loadings for automatized explicit knowledge were identical to
677 the two-factor model, but all the factor loadings for implicit knowledge were lower
678 than the two-factor model. This partially supported that the two-factor model was
679 more plausible than the one-factor model, and the latent factor largely contributes
680 to the form-focused tasks.

681 In order to investigate how the amount of L2 experience changes the validity
682 of the tests, CFAs were conducted separately for the two subsets. The correlation
683 matrix is presented for the short-LOR and long-LOR groups in Table 6. The form-
684 focused tasks converged to a similar extent for the whole group both in the short-
685 LOR group ($.515 < r < .691$) and in the long-LOR groups ($.534 < r < .626$). While
686 there were no meaningful relationships among the three online tasks in the short-
687 LOR group ($-.129 < r < .100$), the online measures were correlated more highly
688 with each other in the long-LOR group than in the whole group ($.237 < r < .343$).

Table 5. Fit indices for confirmatory factor analysis models (two-factor and one-factor models) and MTMM models

	Model	df	χ^2	p	NNFI	CFI	SRMR	RMSEA [90% CI]	Fit
Whole (n = 99)	Two factor	7	6.043	.535	1.019	1	0.036	0 [0–0.113]	Good
	One factor	8	6.940	.543	1.018	1	0.044	0 [0–0.107]	Good
	MTMM	9	9.060	.432	0.999	0.999	0.070	0.008 [0–0.114]	Good
Short LOR(n = 47)	Two factor				Improper solution				Poor
	One factor	9	4.894	.844	1.159	1	0.055	0 [0–0.094]	Good
	MTMM				Improper solution				Poor
Long LOR(n = 52)	Two factor	8	7.527	.481	1.015	1	0.071	0 [0–0.156]	Good
	One factor	9	17.328	.044	0.758	0.855	0.116	0.133 [0.022–0.227]	Poor
	MTMM	9	10.622	.303	0.953	0.972	0.097	0.059 [0–0.173]	Good

Note: MTMM, multitrait-multimethod; NNFI, Bentler–Bonnet nonnormed fit index; CFI, comparative fit index; SRMR, standardized root mean square residual; RMSEA, root mean square error of approximation; LOR, length of residence. The cutoff values for good fit: SRMR < 0.09, RMSEA < 0.06, and CFI and NNFI > 0.96.

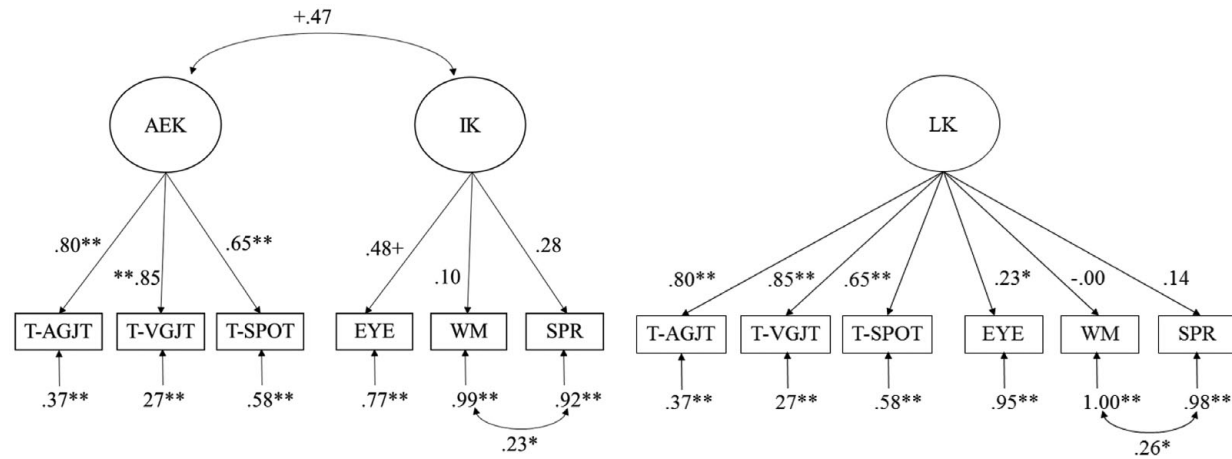


Figure 3. (Left) Two-factor model and (right) one-factor model in the whole second language group ($n = 99$). Standardized coefficients: $+p < .10$, $*p < .05$, $**p < .01$.

Table 6. *Intercorrelations of the language tests for the short-LOR group (n = 47) and long-LOR group (n = 52)*

	Eye	WM	SPR	T-AGJT	T-VGJT	T-SPOT
Short-LOR group						
Eye	—	-.129	-.057	.146	.217	.170
WM		—	.100	.130	-.010	.165
SPR			—	.142	.137	.128
T-AGJT				—	.691**	.515**
T-VGJT					—	.539**
T-SPOT						—
Long-LOR group						
Eye	—	.237	.343*	.158	.131	.266
WM		—	.270	.010	-.157	.018
SPR			—	.173	.012	.077
T-AGJT				—	.626**	.534**
T-VGJT					—	.567**
T-SPOT						—

Note: LOR, length of residence; Eye, visual-world task; WM, word-monitoring task; SPR, self-paced reading task; T-AGJT, timed auditory judgment test; T-VGJT, timed visual judgment test; T-SPOT, timed Simple Performance-Oriented Test.

* $p < .05$. ** $p < .01$.

689 The two CFA models were statistically evaluated. For the short-LOR group,
690 the two-factor model failed to converge, and the one-factor model fit the data set
691 well with acceptable fit indices (see Table 5). For the long-LOR group, in contrast,
692 the two-factor model fit the data significantly better than the one-factor model,
693 $\chi^2_{\text{difference}} = 9.801$, $df = 1$, $p = .002$. While the one-factor model yielded a poor
694 fit in all the indices, the two-factor model indicated a good fit (see Table 5). The
695 factor loadings for the good-fit models are presented for the short-LOR (one-factor
696 model) and long-LOR group (two-factor model) in Figure 4.

697 For the one-factor model of short-LOR group, factor loadings from the measure-
698 ments hypothesized to assess automatized explicit knowledge were consistently
699 high, but the loadings from the measurements hypothesized to assess implicit
700 knowledge were as low as the whole-group results, suggesting that the L2 learners
701 relied on automatized explicit knowledge more. For the two-factor model of the
702 long-LOR group, factor loadings for the implicit knowledge factor were higher
703 than in the model for the whole group, in addition to the high factor loadings for
704 the automatized explicit knowledge. The covariance between automatized explicit
705 knowledge and implicit knowledge was lower in the long-LOR group ($r = .22$,
706 $p = .258$), suggesting that the two latent factors were more distinct in the long-LOR
707 group than in the whole group.

708 *MTMM analysis*

709 The fit indices of the correlated uniqueness model indicated a good fit for the whole
710 group of L2 learners (see Table 5). As shown in Figure 5, the model results showed
711 that all the trait (factor) loadings were statistically significant ($p < .05$). As in the

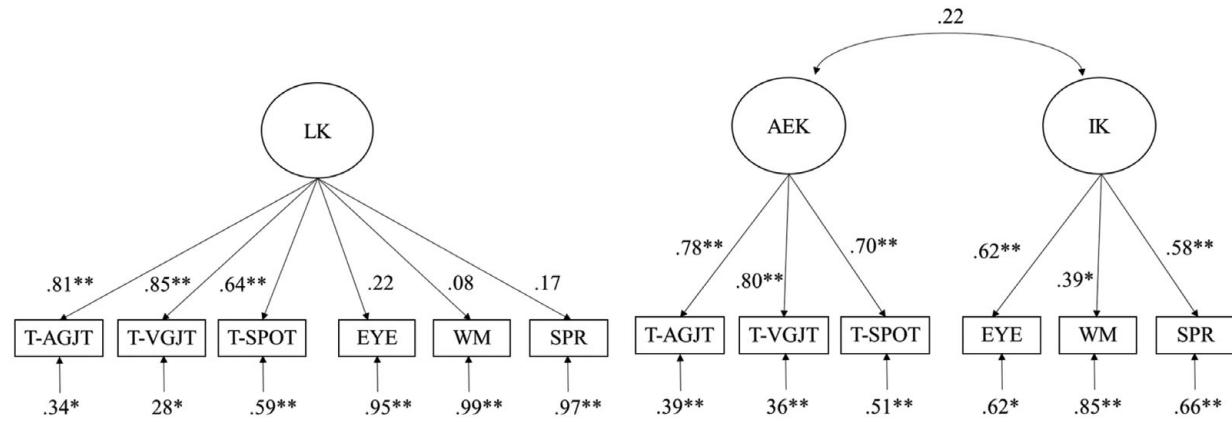


Figure 4. (Left) One-factor model for short-length of residence group ($n = 47$) and (right) two-factor model for long-length of residence group ($n = 52$). Standardized coefficients: * $p < .05$, ** $p < .01$.

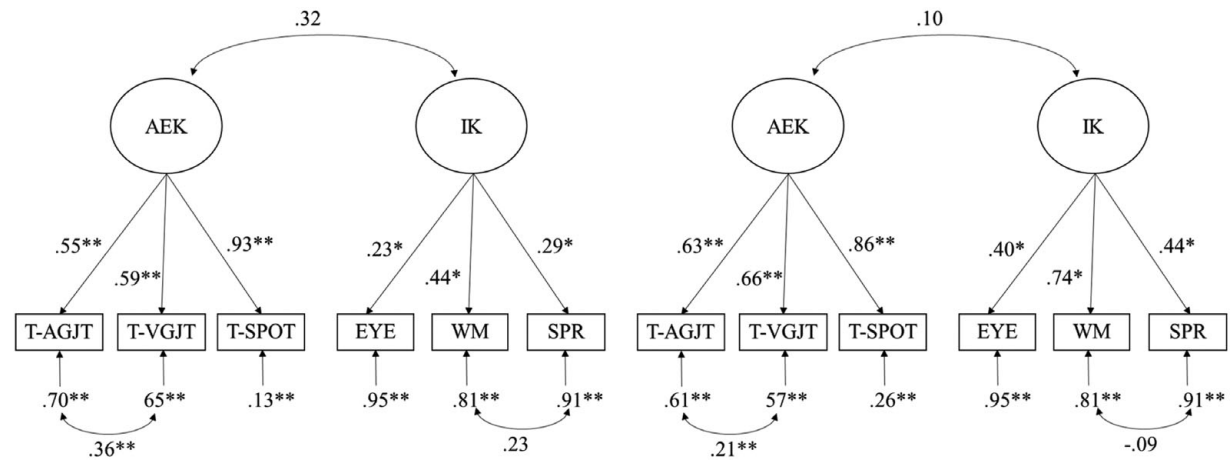


Figure 5. (Left) Multivariate-multimethod models for whole group ($n = 99$) and (right) two-factor model for long-length of residence group ($n = 52$). Standardized coefficients: * $p < .05$, ** $p < .01$.

712 CFA models, the factor loadings were moderate to large in the automatized explicit
713 knowledge measures (range = .55–.93), whereas the trait loadings for implicit
714 knowledge were small to moderate (range = .23–.44). A small and nonsignificant
715 correlation between the two traits was found ($r = .32, p = .175$). The presence
716 of method effects was examined by the correlated uniqueness (errors) among the
717 similar methods. Although the correlated uniqueness was significant between the
718 visual GJT and the auditory GJT ($r = .36, p < .001$), its magnitude was smaller than
719 any of the trait (factor) loadings from the two GJTs (.55 and .59). The correlated
720 uniqueness between the word-monitoring task and the self-paced reading task was
721 not significant ($r = .23, p = .364$), and its magnitude was also smaller than either
722 of the trait loadings (.44 and .29). Method effects evaluated in the MTMM model
723 are marginal, indicating that the set of measurements estimated traits reliably with
724 little influence from the method effects.

725 The same analysis was conducted on the short-LOR group and the long-LOR
726 group, respectively. The model resulted in an improper solution for the short-LOR
727 group; the model for the long-LOR group indicated a good fit of the model with two
728 of the three types of acceptable fit indices (see Table 5).⁶ As shown in Figure 5, the
729 model results showed that all the trait factor loadings were statistically significant
730 ($p < .01$). The magnitude of the trait loadings was medium to large, both for the
731 automatized explicit knowledge measures (range = .63–.86) and for the implicit
732 knowledge (range = .40–.74). A nonsignificant negligible correlation between the
733 two traits also constitutes evidence for discriminant validity ($r = .10, p = .175$).

734 The presence of method effects was investigated through the correlated unique-
735 ness among the similar methods. Although the correlated uniqueness was sig-
736 nificant between the visual GJT and the auditory GJT ($r = .21, p < .001$), the
737 magnitude was smaller than the trait factor loadings from the two GJTs (.63 and
738 .66, $p < .001$). The correlated uniqueness between the word-monitoring task and
739 the self-paced reading task was not significant ($r = -.09, p = .364$), and the mag-
740 nitude of the trait loadings was larger than the correlated uniqueness (.44 and .74,
741 $p < .001$). Method effects estimated in the long-LOR group were smaller than the
742 whole group, providing stronger support for the stability of traits.

743 DISCUSSION

744 The current study addressed whether the three online psycholinguistic measures
745 tap the distinct construct from other time-pressured form-focused tests. Over-
746 all, the results of CFA confirmed that the two-factor model fit the data well
747 (Hypothesis 1). Results of subset analysis demonstrated a different pattern for the
748 two L2 groups varying in the amount of L2 experience (LOR). For the short-LOR
749 group, the two-factor model did not converge, but the one-factor model produced
750 a good fit. In contrast, the two-factor model, but not the one-factor model, fit the
751 data well for the long-LOR group.

752 *Construct validity of measures for automatized explicit and implicit* 753 *knowledge*

754 With regard to Hypothesis 2, although the factor loadings for automatized explicit
755 knowledge were high and statistically significant (range = .65–.85), the loadings

756 for implicit knowledge were much lower (range = .10–.48) in CFA. These low
757 loadings underscore the challenges to devise measurements for implicit knowledge,
758 indicating weak convergent validity for the measurements for implicit knowledge.
759 Nevertheless, supporting evidence was provided for the discriminant validity, given
760 a factor correlation below .80 (Brown, 2006; Hypotheses 3a and 3b). Although
761 the one-factor model fit the data as well as the two-factor model, the substantial
762 loadings in the one-factor model were all from the form-focused tasks. Moreover,
763 the factor loadings from the three online measurements were all lower in the one-
764 factor model than in the two-factor model. The MTMM analysis further showed
765 stronger trait factors than method effects for both GJTs and reaction-time measures
766 (Hypothesis 4).

767 Although a good deal of evidence has been provided for the construct validity of
768 the hypothesized two-factor model, uncertainty is inevitably involved as to whether
769 these two factors are automatized explicit and implicit knowledge. As proposed at
770 the outset of this study (see Table 1 above), however, the tests for automatized ex-
771 plicit knowledge and implicit knowledge were designed to maximally differentiate
772 the two types of tests in terms of the level of awareness involved during the test.
773 Form-focused tasks such as timed GJTs and SPOT *directly* ask participants to pay
774 attention to grammatical structures in stimulus sentences, raising the awareness of
775 linguistic knowledge (i.e., explicit knowledge). Having said that, it is impossible to
776 completely rule out the possibility that some L2 learners draw on implicit knowl-
777 edge to perform timed GJTs. If learners were able to register the error online to
778 make judgments in the timed GJT and had little reflection on their judgments, they
779 might have relied primarily on implicit knowledge (Godfroid et al., 2015). With
780 this inevitable ambiguous nature of GJTs in mind, however, if behavioral language
781 tests are considered on a continuum spectrum from more explicit to more implicit,
782 timed form-focused tasks like GJTs are probably considered closer to the explicit
783 end of the continuum (DeKeyser, 2003; Vafae et al., 2017).

784 In contrast, indirect real-time comprehension measures hypothesized to assess
785 implicit knowledge never ask participants to detect errors. Instead, participants
786 are asked to pay attention to the meaning of a sentence so that they can answer
787 the comprehension question. This indirect nature of the grammar tests can prevent
788 learners from becoming aware of their linguistic knowledge use and thus minimize
789 the involvement of (automatized) explicit knowledge (Andringa & Curcic, 2015;
790 Leung & Williams, 2012; Paradis, 2009; Suzuki & DeKeyser, 2015). Given these
791 rationales of the test design, the current evidence suggests that the two factors
792 should be labeled as automatized explicit knowledge and implicit knowledge. It
793 casts doubt on the construct validity of the previous test battery of explicit and
794 implicit knowledge developed by R. Ellis (2005) and further utilized by others
795 (Bowles, 2011; Ercetin & Alptekin, 2013; Gutiérrez, 2013; Sarandi, 2015; Zhang,
796 2015). Although previous research is cautious in the interpretation that timed GJTs
797 are a less pure measure for implicit knowledge (e.g., Loewen, 2009), time-pressure
798 cannot guarantee the inaccessibility of automatized explicit knowledge (DeKeyser,
799 2003; Suzuki & DeKeyser, 2015; Vafae et al., 2017).

800 The visual-world task is probably a superior measure to the RT measures be-
801 cause it requires no ungrammatical sentences, which makes the task most im-
802 plicit. Furthermore, it directly captures real-time grammar processing via eye

803 movements without any mediation such as through button presses. These advan-
804 tages were empirically supported by the results from the CFA models (Figures 3
805 and 4) indicating that the factor loadings from the visual-world task were the high-
806 est in the current test battery. In contrast, the RT measures (self-paced reading and
807 word-monitoring tasks) necessitate ungrammatical items, which may potentially
808 raise awareness of the linguistic form. However, they are still useful assessment
809 tools of implicit knowledge because L2 learners are unlikely to apply linguistic
810 knowledge consciously when their grammatical processing is time locked within a
811 few hundred milliseconds. When errors are registered without awareness, the level
812 of noticing may sometimes rise up to consciousness as maintenance rehearsal is
813 carried out in working memory. Regardless of the fact that registration may lead to
814 further awareness *after the point of ungrammaticality*, implicit knowledge should
815 be deployed for the registration of the errors *at the exact time of occurrence*, which
816 is captured by the RT measures (Suzuki & DeKeyser, 2015). Furthermore, although
817 the RT measures required button presses, the MTMM findings showed negligible
818 method effects (Figure 5).

819 *Further evidence: More L2 exposure leads to more stable use of implicit*
820 *knowledge*

821 There was a striking difference between the two LOR groups; the one-factor model
822 produced a good fit for the short-LOR group, as opposed to the two-factor model for
823 the long-LOR group (Hypothesis 1). Regarding Hypothesis 2, the factor loadings
824 for the latent factor in the one-factor model (i.e., linguistic knowledge) suggest
825 that L2 learners in the short-LOR group primarily rely on automatized explicit
826 knowledge, as all the loadings for online measures were very low. Inspecting the
827 results from the two-factor model in the long-LOR group, the factor loadings for
828 automatized explicit knowledge were as good as for the whole group (range = .70–
829 .80). It was critical that the factor loadings for implicit knowledge were higher and
830 statistically significant: the two moderate loadings (.58 for the self-paced reading
831 task and .62 for the visual-world task) and one weak loading (.39 for the word-
832 monitoring task) in CFA.

833 The discriminant validity was further supported only for the long-LOR group
834 (Hypotheses 3a and 3b), suggesting that both automatized explicit knowledge and
835 implicit knowledge had been assessed more distinctively for them. In regard to
836 the method effects (Hypothesis 4), the MTMM analysis for the long-LOR group
837 further indicated that the correlated error of the two GJTs was significant but less
838 than the trait factor loadings, and that of the word-monitoring task and the self-
839 paced reading task was of nonsignificant small negative value. The traits seemed
840 to be assessed more reliably with negligible method effects.

841 In sum, the overall findings from the long-LOR group supported all the hypothe-
842 ses (except for the fit indices, Hypothesis 1, probably due to the smaller number of
843 participants) more strongly, including the convergent validity of implicit knowl-
844 edge. Even though implicit knowledge is much harder to assess, compared to
845 automatized explicit knowledge, it is possible to tap into implicit knowledge with
846 more stability, particularly when more experienced L2 learners performed the test
847 battery. This corroborated the previous findings in Suzuki and DeKeyser (2015)

848 and is consistent with Paradis's (2009) claim that explicit and implicit knowledge
849 coexist in the L2 system, and the reliance of implicit knowledge increases over
850 time through more L2 experience. Furthermore, regardless of the three analyzed
851 groups, the factor loadings for automatized explicit knowledge were high (range =
852 .63–.93). This suggests that late L2 learners with some formal instruction, as was
853 the case for the present study, tend to rely on explicit knowledge very consistently
854 (DeKeyser, 2007; Paradis, 2009). Results might be different when different L2
855 populations such as heritage learners and learners who received early foreign lan-
856 guage instruction were administered with the current test set (e.g., Bowles, 2011;
857 Phillip, 2009).

858 The caveat for the current subset analysis is that since each model consists of six
859 indicators, even the rough estimation of the necessary sample size (10 participants
860 \times 6 indicators = 60) indicates that the sample size was less than ideal. The results
861 from the subset analyses should be interpreted cautiously; however, it highlights a
862 methodological gap in the previous studies. Most of the previous validation studies
863 recruited classroom learners with limited immersion experience (Bowles, 2011;
864 Gutiérrez, 2013; Zhang, 2015). The initial study by R. Ellis (2005, 2009) recruited
865 L2 learners in the immersion context, but their LOR was relatively short (1.9 years).
866 The present findings underscore that the amount of L2 exposure in the immersion
867 context should be taken seriously for future research in the validation studies of
868 measures of explicit and implicit knowledge.

869 *Suggestions for further research*

870 The current study opens several avenues for future research. More rigorous vali-
871 dation studies are needed for developing implicit knowledge measures. First, the
872 reliability of the visual-world task was not assessed in the current study. Farris-
873 Trimble and McMurray (2013) examined test–retest reliability of the visual-world
874 task by requiring participants to complete the visual-world task for spoken word
875 recognition twice (Day 1 and Day 2, which were separated by a week). The re-
876 sults showed that eye-movement patterns were closely related between Day 1 and
877 2, suggesting that the visual-world task is stable enough to index an individual's
878 language processing. The present study could not assess the reliability of the task;
879 it should be examined rigorously in future research.

880 Second, another point concerns the generalizability of the present findings. The
881 current study focused on Japanese L2 learners experienced with both formal in-
882 struction and naturalistic environment. Further research is clearly needed in other
883 L2 learner populations, different first language-second language combinations,
884 other linguistic structures tested, and so on.

885 Third, the timed GJT tasks used here do not exactly follow the methodology
886 of previous studies (e.g., R. Ellis, 2005; Zhang, 2015). In prior research, the time
887 limit for each test item in timed GJTs was fixed to average response times by
888 native speakers plus an extra 20% of the time, whereas the present study attempted
889 to analyze the data with post hoc procedures. In the previous study where the
890 time limit was imposed on the responses for each item on the test, test takers
891 must have experienced more pressure and exhibited stronger motivation to make
892 judgments quickly compared to the current format of GJTs. Further research should

893 follow the exact design of the previous tasks in order to collectively advance the
894 understandings in the previous measurements.

895 *Conclusions*

896 The present study set out to investigate the validity of more finely tuned implicit
897 knowledge measures that are distinguishable from automatized explicit knowl-
898 edge. An array of validity evidence supported that the six measurements assessed
899 two distinct linguistic knowledge types, suggesting that indirect real-time compre-
900 hension tasks measured implicit knowledge, which was distinguished from autom-
901 atized explicit knowledge. Although automatized explicit knowledge was assessed
902 relatively easily by the conventional time-pressured form-focused tasks, it seems
903 to be much harder to measure implicit knowledge through behavioral measures,
904 particularly for L2 learners with less L2 exposure. Evidently, further investigations
905 are needed as the theoretical distinction between automatized explicit knowledge
906 and implicit knowledge can bear important implications for understanding SLA.

907 ACKNOWLEDGMENTS

908 This study was supported by IGERT: Biological and Computational Foundations of Lan-
909 guage Diversity (NSF DGE-0801465), the Office of the Graduate Dean for a Summer
910 Research Fellowship, the Language Learning Dissertation Grant Program, and the PhD
911 program in Second Language Acquisition at the University of Maryland. This study is
912 based on part of the author's doctoral dissertation, which was submitted to the University
913 of Maryland College Park. I thank Robert DeKeyser, Yi Ting Huang, Steven Ross, and Nan
914 Jiang for their advice on this project; Yuki Hirose and Edson Tadashi Miyamoto for their gen-
915 erous support with data collection; Kaoru Koyanagi, Yukiko Okuno, Tomomi Nishikawa,
916 Hiromi Ozeki, and Kiyoko Tadokoro for their assistance in recruiting participants; and Kei
917 Harata and Jun Fujita for their assistance in developing materials.

918 NOTES

- 919 1. The word *register* is used in the technical term in the current paper, meaning that
920 cognitive registration of linguistic input that does not require awareness (see Suzuki &
921 DeKeyser, 2015; Tomlin & Villa, 1994, for details).
- 922 2. Since this procedure was similar to the format of existing tests in the Japanese education
923 system, where it is called the SPOT, this task is called the timed SPOT here (Kobayashi,
924 Ford-Niwa, & Yamoto, 1996).
- 925 3. A third model was also constructed in terms of modality of measurements: a written-
926 aural model. It never converged for the current data sets, however.
- 927 4. Japanese Language Proficiency Test (JLPT) N1 (which corresponds to the previous
928 JLPT Levels 1) is roughly equivalent to the ACTFL Superior on the OPI scale (Kanno,
929 Hasegawa, Ikeda, & Ito, 2005). JLPT Level 1 is the minimum requirement for accep-
930 tance into a regular college undergraduate/graduate program in Japan.
- 931 5. These two correlated errors were also imposed on the confirmatory factor analysis mod-
932 els; for parsimony, however, only the correlated errors that were statistically significant
933 were retained in the final model.

- 934 6. Although the standardized root mean square and Benter–Bonnet nonnormed fit indices
 935 did not pass the criteria for the long-LOR group, they were close to the criteria. The
 936 overall assessment of the fit was deemed acceptable.

937 SUPPLEMENTARY MATERIAL

938 To view supplementary material for this article, please visit <https://doi.org/10.1017/S014271641700011X>
 939

940 REFERENCES

- 941 Andringa, S., & Curcic, M. (2015). How explicit knowledge affects online L2 processing. *Studies in*
 942 *Second Language Acquisition*, 37, 237–268. doi:10.1017/S0272263115000017
- 943 Bachman, L. F., & Palmer, A. S. (1982). The construct validation of some components of communica-
 944 tive proficiency. *TESOL Quarterly*, 16, 449–465.
- 945 Bowles, M. A. (2011). Measuring implicit and explicit linguistic knowledge. *Studies in Second Lan-*
 946 *guage Acquisition*, 33, 247–271. doi:10.1017/S0272263110000756
- 947 Brown, T. A. (2006). *Confirmatory factor analysis for applied research*. New York: Guilford Press.
- 948 Campbell, D. T., & Fiske, D. W. (1959). Convergent and discriminant validation by the multitrait-
 949 multimethod matrix. *Psychological Bulletin*, 56, 81–105.
- 950 DeKeyser, R. M. (1997). Beyond explicit rule learning. *Studies in Second Language Acquisition*, 19,
 951 195–221.
- 952 DeKeyser, R. M. (2003). Implicit and explicit learning. In C. J. Doughty & H. M. Long (Eds.), *The*
 953 *handbook of second language acquisition* (pp. 312–348). Oxford: Blackwell.
- 954 DeKeyser, R. M. (2007). Skill acquisition theory. In B. VanPatten & J. Williams (Eds.), *Theories in*
 955 *second language acquisition* (pp. 97–114). Mahwah, NJ: Erlbaum.
- 956 DeKeyser, R. M. (2015). Skill acquisition theory. In B. VanPatten & J. Williams (Eds.), *Theo-*
 957 *ries in second language acquisition: An introduction* (2nd ed., pp. 94–112). New York:
 958 Routledge.
- 959 Elder, C., & Ellis, R. (2009). Implicit and explicit knowledge of an L2 and language proficiency. In
 960 R. Ellis, S. Loewen, C. Elder, R. Erlam, J. Philp, & H. Reinders (Eds.), *Implicit and explicit*
 961 *knowledge in second language learning, testing and teaching* (pp. 167–193). Tonawanda, NY:
 962 Multilingual Matters.
- 963 Ellis, N. C. (2005). At the interface: Dynamic interactions of explicit and implicit language knowledge.
 964 *Studies in Second Language Acquisition*, 27, 305–352. doi:10.1017/S027226310505014X
- 965 Ellis, R. (2005). Measuring implicit and explicit knowledge of a second language. *Studies in Second*
 966 *Language Acquisition*, 27, 141–172. doi:10.1017/S0272263105050096
- 967 Ellis, R. (2008). *The study of second language acquisition* (2nd ed.). Oxford: Oxford University Press.
- 968 Ellis, R. (2009). Measuring implicit and explicit knowledge of a second language. In R. Ellis,
 969 S. Loewen, C. Elder, R. Erlam, J. Philp, & H. Reinders (Eds.), *Implicit and explicit knowledge*
 970 *in second language learning, testing and teaching* (pp. 31–64). Tonawanda, NY: Multilingual
 971 Matters.
- 972 Ellis, R., Loewen, S., Elder, C., Erlam, R., Philp, J., & Reinders, H. (2009). *Implicit and explicit*
 973 *knowledge in second language learning, testing and teaching*. Tonawanda, NY: Multilingual
 974 Matters.
- 975 Ercetin, G., & Alptekin, C. (2013). The explicit/implicit knowledge distinction and working memory:
 976 Implications for second-language reading comprehension. *Applied Psycholinguistics*, 34, 727–
 977 753. doi:10.1017/S0142716411000932
- 978 Farris-Trimble, A., & McMurray, B. (2013). Test–retest reliability of eye tracking in the visual world
 979 paradigm for the study of real-time spoken word recognition. *Journal of Speech, Language,*
 980 *and Hearing Research*, 56, 1328–1345. doi:10.1044/1092-4388

- 981 Foote, R. (2011). Integrated knowledge of agreement in early and late English–Spanish bilinguals.
982 *Applied Psycholinguistics*, 32, 187–220. doi:10.1017/S0142716410000342
- 983 Godfroid, A., Loewen, S., Jung, S., Park, J., Gass, S., & Ellis, R. (2015). Timed and untimed gram-
984 maticity judgements measure distinct types of knowledge. *Studies in Second Language Ac-*
985 *quisition*, 37, 269–297. doi:10.1017/S0272263114000850
- 986 Granena, G. (2013). Individual differences in sequence learning ability and second language acquisition
987 in early childhood and adulthood. *Language Learning*, 63, 665–703. doi:10.1111/lang.12018
- 988 Grüter, T., Lew-Williams, C., & Fernald, A. (2012). Grammatical gender in L2: A pro-
989 duction or a real-time processing problem? *Second Language Research*, 28, 191–215.
990 doi:10.1177/0267658312437990
- 991 Gutiérrez, X. (2013). The construct validity of grammaticality judgment tests as measures of
992 implicit and explicit knowledge. *Studies in Second Language Acquisition*, 35, 423–449.
993 doi:10.1017/S0272263113000041
- 994 Hasuike, I. (2004). Basho wo arawasu kakujoshi “ni” no kajou shiyō ni kansuru ichikousatsu:
995 Chuukyū reberu no chuugokuowasha no joshi sentaku sutōratejii bunseki [Investigation on the
996 overuse of the particle “Ni” for location: Analysis of particle choice strategies by intermediate
997 Chinese speakers]. *Nihongo Kyōiku*, 122, 52–61.
- 998 Hopp, H. (2013). Grammatical gender in adult L2 acquisition: Relations between lexical and syntactic
999 variability. *Second Language Research*, 29, 33–56. doi:10.1177/0267658312461803
- 1000 Hoyle, R. H., & Panter, A. T. (1995). Writing about structural equation models. In R. H. Hoyle (Ed.),
1001 *Structural equation modeling: Concepts, issues, and applications* (pp. 158–176). Thousand
1002 Oaks, CA: Sage.
- 1003 Hu, L. T., & Bentler, P. M. (1999). Cutoff criteria for fit indexes in covariance structure analy-
1004 sis: Conventional criteria versus new alternatives. *Structural Equation Modeling*, 6, 1–55.
1005 doi:10.1080/10705519909540118
- 1006 Hulstijn, J. H. (2002). Towards a unified account of the representation, processing and ac-
1007 quisition of second language knowledge. *Second Language Research*, 18, 193–223.
1008 doi:10.1191/0267658302sr207oa
- 1009 Hulstijn, J. H. (2005). Theoretical and empirical issues in the study of implicit and ex-
1010 plicit second-language learning. *Studies in Second Language Acquisition*, 27, 129–140.
1011 doi:10.1017/S0272263105050084
- 1012 Jacobsen, W. M. (1992). *The transitive structure of events in Japanese*. Tokyo: Kuroshio.
- 1013 Jacoby, L. L. (1991). A process dissociation framework: Separating automatic from intentional uses of
1014 memory. *Journal of Memory and Language*, 30, 513–541. doi:10.1016/0749-596X(91)90025-F
- 1015 Jiang, N. (2011). *Conducting reaction time research in second language studies*. New York: Routledge.
- 1016 Jiang, N., Novokshanova, E., Masuda, K., & Wang, X. (2011). Morphological congruency and
1017 the acquisition of L2 morphemes. *Language Learning*, 61, 940–967. doi:10.1111/j.1467-
1018 9922.2010.00627.x
- 1019 Jöreskog, K., & Sörbom, D. (2013). *LISREL 9.1 for Windows*. Skokie, IL: Scientific Software Interna-
1020 tional.
- 1021 Kanno, K., Hasegawa, T., Ikeda, K., & Ito, Y. (2005). Linguistic profiles of heritage bilingual learners of
1022 Japanese. In J. Cohen, K. T. McAliser, K. Rolstad, & J. MacSwan (Eds.), *ISB4: Proceedings of*
1023 *the 4th International Symposium on Bilingualism* (pp. 1139–1151). Somerville, MA: Cascadilla
1024 Press.
- 1025 Kobayashi, N., Ford-Niwa, J., & Yamoto, H. (1996). SPOT: A new testing method of Japanese language
1026 proficiency [Nihongo nouryoku no atarashii sokuteihou: SPOT]. *Japanese-Language Education*
1027 *Around the Globe*, 6, 201–218.
- 1028 Krashen, S. D. (1985). *The input hypothesis: Issues and implications*. New York: Longman.
- 1029 Leung, J. H. C., & Williams, J. N. (2012). Constraints on implicit learning of grammati-
1030 cal form-meaning connections. *Language Learning*, 62, 634–662. doi:10.1111/j.1467-
1031 9922.2011.00637.x

- 1032 Lew-Williams, C., & Fernald, A. (2010). Real-time processing of gender-marked articles by native
1033 and non-native Spanish speakers. *Journal of Memory and Language*, *63*, 447–464.
1034 doi:[10.1016/j.jml.2010.07.003](https://doi.org/10.1016/j.jml.2010.07.003)
- 1035 Loewen, S. (2009). Grammaticality judgment tests and the measurement of implicit and explicit L2
1036 knowledge. In R. Ellis, S. Loewen, C. Elder, R. Erlam, J. Philp, & H. Reinders (Eds.), *Implicit and explicit knowledge in second language learning, testing and teaching* (pp. 94–112).
1037 Tonawanda, NY: Multilingual Matters.
- 1038 Loewenthal, K. M. (2004). *An introduction to psychological tests and scales* (2nd ed.). Hove: Psychology
1039 Press.
- 1040 Matin, E., Shao, K., & Boff, K. R. (1993). Saccadic overhead: Information-processing time with and
1041 without saccades. *Perception & Psychophysics*, *53*, 372–380. doi:[10.3758/BF03206780](https://doi.org/10.3758/BF03206780)
- 1042 McLaughlin, B. (1987). *Theories of second-language learning*. London: Routledge.
- 1043 Mitchell, D. C. (1984). An evaluation of subject-paced reading tasks and other methods for investigating
1044 immediate processes in reading. In D. E. Kieras & M. A. Just (Eds.), *New methods in reading*
1045 *comprehension research* (pp. 69–89). Hillsdale, NJ: Erlbaum.
- 1046 Paradis, M. (2009). *Declarative and procedural determinants of second languages*. Philadelphia, PA:
1047 Benjamins.
- 1048 Philip, J. (2009). Pathways to proficiency: Learning experiences and attainment in implicit and explicit
1049 knowledge of English as a second language. In R. Ellis, S. Loewen, C. Elder, R. Erlam, J. Philp,
1050 & H. Reinders (Eds.), *Implicit and explicit knowledge in second language learning, testing and*
1051 *teaching* (pp. 194–215). Tonawanda, NY: Multilingual Matters.
- 1052 Podsakoff, P. M., MacKenzie, S. B., & Podsakoff, N. P. (2012). Sources of method bias in social science
1053 research and recommendations on how to control it. *Annual Review of Psychology*, *63*, 539–569.
1054 doi:[10.1146/annurev-psych-120710-100452](https://doi.org/10.1146/annurev-psych-120710-100452)
- 1055 Posner, M. I., & Snyder, C. R. (1975). Attention and cognitive control. In R. L. Solso (Ed.), *Information*
1056 *processing in cognition: The Loyola symposium* (pp. 55–85). Hillsdale, NJ: Erlbaum.
- 1057 Roberts, L., & Liszka, S. A. (2013). Processing tense/aspect-agreement violations on-line in the second
1058 language: A self-paced reading study with French and German L2 learners of English. *Second*
1059 *Language Research*, *29*, 413–439. doi:[10.1177/0267658313503171](https://doi.org/10.1177/0267658313503171)
- 1060 Sarandi, H. (2015). Reexamining elicited imitation as a measure of implicit grammatical knowledge and
1061 beyond . . . ? *Language Testing*. Advance online publication. doi:[10.1177/0265532214564504](https://doi.org/10.1177/0265532214564504)
- 1062 Sedivy, J. C. (2010). Using eyetracking in language acquisition research. In E. Blom & S. Unsworth
1063 (Eds.), *Experimental methods in language acquisition research* (pp. 115–138). Philadelphia:
1064 PA: Benjamins.
- 1065 Spada, N. (2015). SLA research and L2 pedagogy: Misapplications and questions of relevance. *Language*
1066 *Teaching*, *48*, 69–81. doi:[10.1017/S026144481200050X](https://doi.org/10.1017/S026144481200050X)
- 1067 Suzuki, Y., & DeKeyser, R. M. (2015). Comparing elicited imitation and word monitoring as measures
1068 of implicit knowledge. *Language Learning*, *65*, 860–895. doi:[10.1111/lang.12138](https://doi.org/10.1111/lang.12138)
- 1069 Suzuki, Y., & DeKeyser, R. M. (2017). The interface of explicit and implicit knowledge in a second
1070 language: Insights from individual differences in cognitive aptitudes. *Language Learning*.
1071 Advance online publication.
- 1072 Tanenhaus, M. K., & Trueswell, J. C. (2006). Eye movements and spoken language comprehension.
1073 In M. J. Traxler & M. A. Gernsbacher (Eds.), *Handbook of psycholinguistics* (2nd ed.). New
1074 York: Elsevier.
- 1075 Tomlin, R. S., & Villa, V. (1994). Attention in cognitive science and second language acquisition.
1076 *Studies in Second Language Acquisition*, *16*, 183–203. doi:[10.1017/S0272263100012870](https://doi.org/10.1017/S0272263100012870)
- 1077 Trenkic, D., Mirkovic, J., & Altmann, G. T. M. (2014). Real-time grammar processing by native and
1078 non-native speakers: Constructions unique to the second language. *Bilingualism: Language*
1079 *and Cognition*, *17*, 237–257. doi:[10.1017/S1366728913000321](https://doi.org/10.1017/S1366728913000321)
- 1080 Vafaei, P., Suzuki, Y., & Kachinske, I. (2017). Validating grammaticality judgment tests: Evidence
1081 from two new psycholinguistic measures. *Studies in Second Language Acquisition*, *39*, 59–95.
1082 doi:[10.1017/S0272263115000455](https://doi.org/10.1017/S0272263115000455)
- 1083

- 1084 Widaman, K. F. (1985). Hierarchically nested covariance structure models for multitrait-multimethod
1085 data. *Applied Psychological Measurement*, 9, 1–26. doi:[10.1177/014662168500900101](https://doi.org/10.1177/014662168500900101)
- 1086 Williams, J. N. (2009). Implicit learning in second language acquisition. In W. C. Ritchie & T. K. Bhatia
1087 (Eds.), *The new handbook of second language acquisition* (pp. 319–353). London: Emerald
1088 Group.
- 1089 Zhang, R. (2015). Measuring university-level L2 learners' implicit and explicit linguistic knowledge.
1090 *Studies in Second Language Acquisition*, 37, 457–486. doi:[10.1017/S0272263114000370](https://doi.org/10.1017/S0272263114000370)