

VALIDATING GRAMMATICALITY JUDGMENT TESTS

Evidence from Two New Psycholinguistic Measures

QA 16 Payman Vafae, Yuichi Suzuki, and Ilna Kachisnke
17 *University of Maryland*

21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

Several previous factor-analytic studies on the construct validity of grammaticality judgment tests (GJTs) concluded that untimed GJTs measure explicit knowledge (EK) and timed GJTs measure implicit knowledge (IK) (Bowles, 2011; R. Ellis, 2005; R. Ellis & Loewen, 2007). It has also been shown that, irrespective of the time condition chosen, GJTs' grammatical sentences tap into IK, whereas their ungrammatical ones invoke EK (Gutiérrez, 2013). The current study examined these conclusions by employing two more fine-grained measures of IK: that is, a self-paced reading task and a word-monitoring task. The results of a confirmatory factor analysis revealed that manipulating GJTs' time conditions and/or the grammaticality of the sentences does not render them distinct measures of EK and IK. The current work shows that GJTs are too coarse to be measures of IK, and that the different types of GJTs measure different levels of EK.

45 We would like to thank professors Steve Ross, Mike Long, and Robert DeKeyser and the
46 *SSLA* editors and anonymous reviewers for their thoughtful suggestions and constructive
47 feedback. We are also grateful to the participants of our study. All errors and omissions
48 are, of course, our own.

49 Correspondence concerning this article should be addressed to Payman Vafae, 9348
Cherry Hill Road, Apt #708, College Park, MD 20740. E-mail: payman.vafae@gmail.com

1 INTRODUCTION

2
3
4 The current methodological study reexamined the second language
5 (L2) knowledge type that nonnative English speakers draw on to per-
6 form grammaticality judgment tests (GJTs). Previous factor-analytic
7 validity studies on GJTs employed elicited imitation (EI) and/or oral
8 narrative (ON) tasks as measures of implicit knowledge (IK). Their
9 authors concluded that manipulating GJTs' time conditions or the gram-
10 maticality of the sentences renders them distinct measures of IK and
11 explicit knowledge (EK). Several studies yielded findings indicating that
12 untimed GJTs measure EK, whereas timed GJTs measure IK (e.g., Bowles,
13 2011; R. Ellis, 2005; R. Ellis & Loewen, 2007). It has also been shown that,
14 regardless of time condition, participant responses to different GJT
15 stimulus types (i.e., grammatical and ungrammatical sentences) tap
16 into IK and EK, respectively (Gutiérrez, 2013).
17

18 Unlike previous studies, the current study employed two new psy-
19 cholinguistic measures of IK. More specifically, EI and ON were replaced
20 by a word-monitoring task (WMT) and a self-paced reading task (SPRT),
21 as these have been shown to be more valid measures of IK (Jiang, 2004,
22 2007; Jiang, Novokshanova, Masuda, & Wang, 2011; Suzuki, 2015; Suzuki
23 & DeKeyser, in press). For this reason, it was hypothesized that, by
24 including WMT and SPRT measures in a test battery, the results and
25 conclusions pertaining to the construct validity of GJTs would be dif-
26 ferent from previous studies. In addition, a metalinguistic knowledge
27 test (MKT), a well-established measure of EK, was included in the cur-
28 rent study.
29

30 Through the comparison of the performance of learners on the
31 WMT, SPRT, and MKT measures, as well as on several types of GJTs, it
32 was possible to ascertain whether manipulating time conditions and/
33 or grammaticality in GJT sentences can transform GJTs into distinct
34 measures of IK and EK.
35

36 To foreshadow the results and conclusions yielded by this research,
37 the aforementioned comparisons, carried out through confirmatory
38 factor analysis (CFA), revealed that GJTs are too coarse to be measures
39 of IK, and that manipulating their time conditions and sentence gram-
40 maticality does not render them distinct measures of IK and EK. Rather,
41 we concluded that, as behavioral measures such as GJTs are not pure
42 measures of IK or EK, on a continuum from being more explicit to more
43 implicit, GJTs fall closer to the explicit end.
44

45 The following section situates the significance of the current study
46 within a broader context of SLA research. It explains why more rigorous
47 validation studies on GJTs are needed. The next section provides a crit-
48 ical review of previous validation studies on GJTs, as they motivate the
49 research questions and design of the current study.

1 **Explicit and Implicit Knowledge and the Interface Issue**

2
3
4 The constructs of IK and EK are central to SLA theory construction and
5 evaluation. There is a consensus that SLA draws on both implicit and
6 explicit learning mechanisms (Bley-Vroman, 1991; DeKeyser, 2003; N. C.
7 Ellis, 2005; R. Ellis, 2004), which in turn may result in IK and/or EK
8 (Williams, 2009). Explicit knowledge is knowledge we are consciously
9 aware of, whereas IK is the knowledge that we have but are not aware of
10 (DeKeyser, 2009; Hulstijn, 2005; Williams, 2009). These two kinds of
11 knowledge also differ in the extent to which one can or cannot verbalize
12 them (R. Ellis, 2004, 2005). Our conscious access to EK allows us to ver-
13 balize it;¹ however, because IK is beyond awareness, it cannot be ver-
14 balized (DeKeyser, 2009).
15

16 It is currently believed that IK underlies the ability to use a L2 fluently
17 and confidently; therefore, development of IK should be the ultimate
18 goal of SLA (Doughty, 2003; N. C. Ellis, 1993; R. Ellis, 2005; Hulstijn,
19 2001).² Although there is solid evidence showing that implicit and
20 explicit representations are neurologically distinct (Paradis, 2009;
21 Ullman, 2011), the interaction between the two and how they influence
22 each other are still subjects of controversy. In SLA research, this con-
23 troversy is referred to as the *interface issue*. Central to the interface
24 issue is to what extent explicit learning and instruction impact implicit
25 learning and the development of IK (N. C. Ellis, 2011). There are three
26 positions with regard to the interface issue—the noninterface, strong-
27 interface, and weak-interface positions.
28

29 The noninterface position is often associated with Krashen (e.g.,
30 1994), who contended that conscious learning about L2s and subcon-
31 scious acquisition of L2s are two completely different phenomena that
32 result in distinct sources of knowledge (EK and IK, respectively) with no
33 interface between them. According to this position, subconscious
34 acquisition dominates L2 performance, learning can never convert into
35 acquisition, and conscious learning can only be used as a monitor
36 (editor) for performance. Proponents of the noninterface position
37 believe that EK and IK are located in different areas of the brain and are
38 thus accessed by different processes (Paradis, 1994).³
39

40 The strong-interface position, which is usually associated with
41 DeKeyser (e.g., 2007), offers an opposite view. His strong-interface
42 position should be understood within the models of skill acquisition, such
43 as ACT-R (Anderson & Lebiere, 1998). Within these models, a distinc-
44 tion is made between declarative and procedural knowledge. According
45 to DeKeyser (2009), “Declarative knowledge is knowledge THAT some-
46 thing is,” and procedural knowledge is “knowledge HOW to do something”
47 (p. 121). According to skill-acquisition models, learners first develop a
48 declarative encoding, whereby extensive practice is required to ensure
49

AQ2

AQ3

AQ4

AQ5

AQ6

1 that declarative knowledge leads to procedural knowledge. Further
2 practice leads to automatized knowledge, which may not require any
3 conscious processing (DeKeyser & Criado, 2013). This strong-interface
4 position implies a causal relationship between declarative knowledge
5 and proceduralization and automatization of the knowledge. In other
6 words, EK forms a prerequisite for the generation of IK (Segalowitz &
7 Hulstijn, 2005).
8

AQ7

9 The main claim of the weak-interface position is that EK does not
10 have a causal relationship with IK and only triggers or speeds up the
11 implicit learning process, which subsequently leads to the generation
12 of IK. For example, N. C. Ellis (e.g., 2005, 2008) contended that EK con-
13 tributes indirectly to the acquisition of IK by promoting some implicit
14 learning processes. According to the author, EK can make relevant lin-
15 guistic features salient, thus enabling learners to notice them and rec-
16 ognize the gap between the input and the linguistic knowledge they
17 possess (N. C. Ellis, 1994). N. C. Ellis (2008) further suggested that explicit
18 and implicit processes work in tandem and that there is a dynamic
19 interaction between them for consolidating IK.
20

21 The interface issue has been debated for decades and remains impor-
22 tant to SLA research for theoretical and pedagogical reasons (Hulstijn,
23 2005). However, dealing with EK and IK constructs and testing rival
24 interface hypotheses is challenging for several reasons. One of these
25 challenges is the lack of reliable and valid measurement tools, which is
26 particularly significant when attempting to measure IK. As both IK and
27 EK sources are involved in L2 performance, it is almost impossible to
28 devise pure behavioral measures of these two constructs (DeKeyser,
29 2009; R. Ellis, 2004, 2005). In addition, in constructing measures of IK,
30 the operationalization of the concept of awareness and other challenges
31 may be encountered. Consequently, researchers have been using a
32 variety of imperfect measures (e.g., metalinguistic tests for EK and
33 timed GJTs for IK). For this reason, rigorous validation studies on these
34 measures are required.
35

36 Thus far, GJTs have been among the most extensively used measures
37 in research on L2 acquisition. They have been subjected to several val-
38 idation studies. However, due to several methodological limitations
39 affecting previous studies, the construct validity of GJTs remains open
40 to scrutiny. The following section highlights the limitations of previous
41 studies and explains the motivation for the current project.
42
43
44

45 **Extant Validity Research on GJTs**

46
47

48 Previously, it was thought that GJTs provide a direct window into the
49 learners' linguistic competence. However, it is currently acknowledged

1 that GJTs only provide a measure of linguistic performance (R. Ellis,
2 1991). Many researchers have attempted to establish the type of linguistic
3 knowledge—whether implicit, explicit, or a combination of both—
4 learners draw on in their performance on GJTs. Earlier studies of this
5 type yielded the conclusion that the nature of learners' knowledge,
6 whether explicit or implicit, affects their judgment on GJTs (Chaudron,
7 1983). In other words, GJTs potentially lead participants to draw on
8 both IK and EK, depending on what source of knowledge they *mostly*
9 have at their disposal (R. Ellis, 2005).

10
11 Further research on GJTs aimed to investigate whether manipulating
12 GJT designs leads to differential participant performance. For instance,
13 it was hypothesized that if GJTs only ask participants to discriminate
14 between well-formed and deviant sentences, it is possible that they
15 evoke the use of pure intuition. On the other hand, the use of EK is very
16 likely if GJTs require locating the error and even editing/correcting
17 or describing the rule for judgment (R. Ellis, 1991). In addition, it was
18 hypothesized that the kind of knowledge GJTs prompt learners to rely
19 on to make their judgment depends on time conditions (i.e., whether
20 the test is timed or untimed). Time pressure may encourage test takers
21 to respond on the basis of their IK, whereas unlimited time may allow
22 them to rely on their EK (Bialystok, 1979). Moreover, it has been hypothesized
23 that the grammaticality of GJT sentences may induce differential
24 performance because learners rely on IK and EK for judging
25 grammatical and ungrammatical sentences, respectively (Gutiérrez,
26 2013).

27
28 Several factor-analytic validity investigations tested these hypotheses.
29 This series of factor-analytic studies commenced with the psychometric
30 work of Rod Ellis. R. Ellis (2005) conducted a psychometric study to
31 design several measures of EK and IK and evaluated their construct
32 validity with respect to EK and IK constructs. For operationalizing the
33 two constructs and distinguishing between them, Ellis proposed seven
34 criteria: degree of awareness, time available, focus of attention, systematicity
35 and certainty, metalanguage, and learnability. He subsequently
36 created five measures: namely, a timed GJT, an untimed GJT, an ON task,
37 an EI task, and a MKT. Based on the seven criteria, Ellis predicted that
38 the ON task, the EI task, and the timed GJT would tap into IK, whereas
39 the untimed GJT and the MKT would evoke the use of EK.

40
41 R. Ellis (2005) submitted his test battery data to exploratory factor
42 analysis (EFA). The results yielded a two-factor structure that confirmed
43 his predictions. More specifically, the ON task, the EI task, and the timed
44 GJT loaded on the first factor, whereas the untimed GJT and the MKT
45 loaded on the second. Ellis labeled the two factors IK and EK, respectively.
46 However, as Isemonger (2007) explained, Ellis's (2005) approach
47 suffered from a few major flaws from a methodological and analytical
48 perspective. As a result, the inferences and interpretations drawn were
49

AQ8

AQ9

1 not well supported. For example, from a methodological and analytical
2 point of view, because Ellis approached the factor analysis with an a
3 priori hypothesis, the use of EFA was not acceptable. Instead, CFA
4 should have been used because the prior hypothesis implied that the
5 measures would measure the distinct constructs of EK and IK.
6

7 To overcome the analytical flaws of the approach employed by R. Ellis
8 (2005), R. Ellis and Loewen (2007) reanalyzed the data used in that
9 study through a CFA. In the CFA, they tested the two-factor model from
10 the original EFA against a decision and production model, as a rival
11 model. This is also a two-factor model, with EI and ON tasks loading on
12 the production, and the two GJTs and the MKT loading on the decision
13 factor. The results yielded by this approach showed that only the first
14 hypothesized model produced adequate fit. However, whereas the CFA
15 approach is indeed a better option than EFA, the authors' execution of
16 the CFA could have been more thorough.
17

18 Confirmatory factor analysis is regarded as a process involving several
19 stages—namely, (a) initial model specification, (b) parameter identification
20 and estimation, (c) data-model fit assessment, (d) possible model
21 modification, and (e) rival model identification (which may jeopardize
22 causal inferences made from the original hypothesized model).
23 However, the research conducted by R. Ellis and Loewen (2007) lacked
24 a final CFA process—an adequate rival model identification. In addition,
25 either alternative models were not specified in the GJT validation
26 studies conducted to date or the tested alternative models were irrelevant
27 to the main purpose of the study. When rival CFA models are not
28 specified, conclusions about specific sets of measures can be highly
29 compromised. As Isemonger (2007) explained, “It is important that rival
30 models are tested because the fit of a particular model does not preclude
31 the possibility that other untested models fit better” (p. 109).
32

33 The only rival model R. Ellis and Loewen (2007) tested was the
34 decision and production model. The constructs in this rival model were,
35 however, irrelevant to the main constructs of the study. One important
36 rival model that could have provided a rebuttal to R. Ellis and Loewen's
37 claims is a one-factor model, which enables directly testing whether the
38 measures can actually distinguish the two constructs. Had this model
39 been tested and had it fit the data as well as their two-factor model, it
40 would not be possible for the researchers to conclude that their
41 measures tapped into two distinct constructs of IK and EK.
42

43 More recently, Kachinske and Vafaei (2014) set out to examine a one-
44 factor model as the alternative model to R. Ellis and Loewen's (2007)
45 original model. They reanalyzed R. Ellis and Loewen's data and found
46 that the one-factor model, which also accounted for the method effect
47 (by correlating the error terms between similar tasks), fit the data as
48 well as the authors' original model. As this finding suggests that the
49 one-factor model is statistically as acceptable as R. Ellis and Loewen's

1 two-factor model, the construct validity of the latter could not be sup-
2 ported. Kachinske and Vafae further pointed out similar flaws in the
3 factor-analytic approach adopted in several subsequent GJT validation
4 studies (i.e., Bowles, 2011; Gutiérrez, 2013).

5 Bowles (2011) created a battery of five tests of Spanish as a second
6 language by closely following R. Ellis's (2005) guidelines. Among other
7 analyses, Bowles conducted a CFA to examine the factorial structure
8 of her test battery and reported results that concurred with those
9 obtained by R. Ellis and Loewen (2007). However, Bowles did not exam-
10 ine any rival CFA models against her two-factor model. Furthermore,
11 the fitted two-factor model indicated high correlation between the two
12 hypothesized factors ($r = .87$). With such a high correlation between the
13 factors, it is hard to claim that the two factors are distinct. This reem-
14 phasizes the importance of testing rival models. If a single-factor model
15 had been tested, results could have shown a good model fit, and the
16 interpretations of the study could have changed.

17 Gutiérrez (2013) employed a similar test battery for a different L2
18 population—that is, Spanish learners in the United States. Gutiérrez fol-
19 lowed the guidelines proposed by R. Ellis (2005) and created a battery
20 of tests consisting of a timed GJT, an untimed GJT, and an MKT. The only
21 test in the battery hypothesized to be a measure of IK was the timed
22 GJT, whereas the untimed GJT and the MKT were considered measures
23 of EK. The author simultaneously examined the role of time pressure on
24 GJTs (timed and untimed) and the types of test items (grammatical and
25 ungrammatical) in order to scrutinize GJTs as measures of IK and EK.
26 It was hypothesized that, irrespective of time conditions, judging the
27 grammatical sentences in GJTs taps into IK, and judging ungrammatical
28 sentences engages EK. Gutiérrez conducted both an EFA and a CFA to
29 test his hypotheses.
30

31 In the CFA, Gutiérrez (2013) tested two rival models. In the first, both
32 the grammatical and ungrammatical sentences of the timed GJT loaded
33 on the construct of IK, and the grammatical and ungrammatical sen-
34 tences of the untimed GJT and the MKT loaded on the construct of EK.
35 In the second model, the grammatical sentences of both the timed and
36 untimed GJTs loaded on the construct of IK, and the ungrammatical
37 sentences of both types of GJTs and the MKT loaded on the EK con-
38 struct. The analyses yielded a better fit for the second model, implying
39 that, regardless of time pressure, grammaticality of the stimulus is what
40 distinguishes between the use of IK and EK in performing GJTs. How-
41 ever, Gutiérrez's study was not free from limitations.
42

43 First, having different types of GJTs as the only measures of IK, and a
44 MKT as the sole measure of EK, limited Gutiérrez's ability to test impor-
45 tant rival models. Second, bivariate correlations revealed statistically
46 significant coefficients between the MKT and all types of GJTs, irrespec-
47 tive of stimulus type and time condition. Although the correlations
48
49

1 between the MKT and the ungrammatical sentences of both types of
2 GJTs were the greatest, the correlations between the MKT and the other
3 measures cannot be ignored. Therefore, a one-factor model that accounted
4 for method effects could have been a plausible rival model.

5 In summary, factor-analytic studies, including those conducted by
6 Bowles (2011), Ellis (2005), Ellis and Loewen (2007), and Gutiérrez (2013),
7 suffer from methodological issues, especially the failure to test rival
8 models. This omission in the approach compromises any conclusions
9 about the structural relations among measures and the factors in a
10 model. Therefore, in future studies, several important rival models (e.g.,
11 one-factor models) should be tested.
12
13
14
15

16 **Measuring IK through Online Processing**

17
18 Psycholinguistic research has revealed how L2 learners process gram-
19 matical structures in real time (e.g., Clahsen & Felser, 2006; Kaan, 2014).
20 Online grammatical processing has often been examined through reaction
21 time (RT) measures, such as self-paced reading tasks (e.g., Jiang, 2004;
22 Roberts & Liszka, 2013) and word-monitoring tasks (Granena, 2013;
23 Jiang, Hu, Lukyanchenko, & Cao, 2010; Suzuki & DeKeyser, in press).
24 These tasks can examine whether L2 learners are sensitive to grammat-
25 ical violations while they are reading or listening for comprehension.
26 In the self-paced reading task, for instance, participants read a sentence
27 containing a target grammatical structure word by word, and the RT to
28 each word read is recorded. Researchers examine whether participants
29 slow down to read the word(s) once they encounter a grammatical
30 error. For instance, when participants read a sentence with a subject-
31 verb agreement violation like “The boy in the room enjoy reading many
32 books,” they will slow down when reading the word after the violation
33 (i.e., *reading*), compared to their performance on the grammatical ver-
34 sion of the same sentence (The boy in the room enjoys reading many
35 books). By computing the RT difference between the grammatical and
36 ungrammatical sentences, the sensitivity to grammatical violation can
37 be estimated.
38
39

40 These online sentence-processing tasks are promising measures for
41 tapping into IK (Suzuki & DeKeyser, in press). First, the tasks can cap-
42 ture sensitivity to grammatical violation as the sentence unfolds in real
43 time. This minimizes the possibility that participants access their lin-
44 guistic knowledge consciously, that is, that they rely on EK (Paradis,
45 2009). Second, they can direct attention to meaning, as each sentence is
46 followed by a comprehension question asking about the content of the
47 sentence. This second feature contrasts with the design of GJTs because
48 GJTs of any kind draw participants’ attention to form. When learners are
49

1 instructed to decide whether a sentence is grammatical or ungrammat-
2 ical, they inevitably pay attention to the form, which potentially invokes
3 the use of EK. Although time pressure makes using EK harder, it does
4 not rule out the possibility that EK is accessed. Therefore, the validity
5 of GJTs can be further scrutinized if online processing measures that
6 draw learners' attention to meaning are included in the current study.
7

8 The present investigations hold promise for scrutinizing the validity
9 of GJTs because recent evidence suggests that online processing
10 measures tap into IK. Suzuki and DeKeyser (in press) compared the
11 word-monitoring task and the EI (a previously attested measure of IK)
12 as measures of IK. The tasks employed in their study—along with a MKT
13 as a measure of EK and a probabilistic serial reaction time (SRT) task,
14 which served as a measure of aptitude for implicit learning—were ad-
15 ministered to Japanese L2 learners who live in Japan. The study results
16 showed that the word-monitoring task was related to the SRT task only,
17 whereas the EI was associated with the MKT only. This pattern was
18 found only among the learners who had lived in Japan longer than a
19 certain number of years. This finding suggests that the word-monitoring
20 task can serve as an implicit processing measure among learners with
21 sufficient naturalistic L2 exposure.
22

23 The findings reported by Suzuki and DeKeyser also provide some
24 implications for participant selection. Participants in the validation
25 studies may need to have had sufficient naturalistic learning experi-
26 ence. For example, participants in the studies conducted by Gutiérrez
27 (2013) and Zhang (2015) were classroom learners with limited exposure
28 in L2 environments. As noted earlier, behavioral measures potentially
29 prompt learners to draw on both IK and EK, depending on what source
30 of knowledge they (mostly) have at their disposal (R. Ellis, 2005).
31 Therefore, more rigorous validation studies on GJTs should at least
32 recruit participants that are typically exposed to naturalistic, as well
33 as classroom-based, learning opportunities.
34
35

36 **THE STUDY**

37

38
39 The current study investigated the construct validity of GJTs. By keeping
40 modality constant, in the written mode, the study employed GJTs with
41 combinations of different stimulus types and time conditions. Gram-
42 maticality judgment tests with two stimulus types (grammatical vs.
43 ungrammatical sentences) and two time conditions (timed vs. untimed)
44 were developed.
45

46 In order to advance the current understanding of the methodological
47 problems in measuring EK and IK, the study incorporated two new psy-
48 cholinguistic measures, along with the GJTs. They were a SPRT and a
49 WMT, which should draw on IK to a greater extent. Explicit knowledge

1 was operationalized as the use of linguistic knowledge with attention to
2 form, requiring the use of metalinguistic knowledge. On the continuum
3 of implicit to explicit linguistic processing, a MKT was employed as the
4 most explicit test form.
5

7 RESEARCH QUESTIONS AND HYPOTHESES

10 The current study sought to address the following research questions:

12 *Question 1.* What is the relationship among performance on different types of
13 GJTs, a SPRT, a WMT, and a MKT?

14 *Question 2.* Does manipulating the time condition and stimulus type in GJTs
15 result in two distinct measures of EK and IK?
16

17 It was hypothesized that because GJTs draw attention to form, manip-
18 ulating their time condition or stimulus type does not transform them into
19 measures of implicit knowledge. Rather, online sentence-comprehension
20 tasks—such as WMTs and SPRTs, which draw attention to meaning—
21 are more valid measures of IK. It was further hypothesized that the
22 ungrammatical sentences of both GJTs are more valid measures of
23 knowledge of the target structures under examination. Therefore, it
24 was posited that a CFA model—which includes (a) only ungrammat-
25 ical sentences from both GJTs, as well as the MKT, as measures of
26 explicit knowledge and (b) the WMT and SPRT as measures of implicit
27 knowledge—would provide the best fit to the study data. Given the
28 need to answer the research questions and test the hypotheses, CFA
29 was chosen as the most suitable data analysis method.
30
31

33 METHOD

35 Participants

38 The main study participants were 79 learners of English as a second
39 language, who started learning English after about the age of 10⁴ in a
40 formal setting and subsequently moved to the United States. These
41 learners were Chinese international students who had lived in the
42 United States for a minimum of 1 year. Chinese international students
43 were chosen as study participants as this ensured that they all shared
44 a common first language (i.e., this element was constant). These partic-
45 ipants had scored a minimum of 90 on the TOEFL iBT test (or 6.5 on
46 IELTS). According to the ETS website, 90 is the minimum score of the
47 TOEFL iBT for the advanced level. In addition, according to ETS, an
48 IELTS score of 6.5 is equivalent to the TOEFL iBT score of 90. In terms of
49

1 the gender distribution of the sample, 52 of the participants were female,
 2 whereas the remaining 27 were male. Their educational attainment
 3 was mostly similar: 16 participants were undergraduates, whereas the
 4 remaining 63 were graduate students enrolled in various degree pro-
 5 grams at a U.S. mid-Atlantic university. Table 1 presents descriptive sta-
 6 tistics pertaining to students' demographic background.

7 In addition, the language tests employed in the study were piloted
 8 with a group of English native speakers (NSs), all of whom were under-
 9 graduate students, on two separate occasions. First, 20 NSs took part in
 10 the initial item analyses. In light of the results of the first phase, some of
 11 the materials and test items were modified, and the revised version was
 12 offered to 10 further NSs.
 13
 14
 15
 16

17 Target Structures

18
 19 Four English target structures—present hypothetical conditional, third-
 20 person *s*, simple past/present perfect, and mass/count nouns—were
 21 used to construct the GJTs and the SPRT, WMT, and MKT measures.
 22 The reasons for choosing these four structures were twofold. First, past
 23 research (e.g., R. Ellis, 2009) suggested that these four target structures
 24 are among the most difficult structures in the English language for EFL
 25 and ESL learners to master. In addition, these structures could easily be
 26 incorporated into SPRT and WMT items.
 27
 28
 29
 30

31 Instruments

32
 33 Before describing the details of each of the tasks, it should be men-
 34 tioned that four parallel sets of sentences (Set 1, 1', 2, and 2') were cre-
 35 ated for each of the following tasks: the untimed GJT, the timed GJT, the
 36
 37

38 **Table 1.** Descriptive statistics for background information
 39

| | Mean | Median | <i>SD</i> | Min | Max | Range |
|---------------------|-------|--------|-----------|-----|-----|-------|
| 40 Age | 24.45 | 24 | 3.5 | 18 | 36 | 18 |
| 41 LOR ^a | 31.85 | 26 | 25.4 | 12 | 146 | 134 |
| 42 AOA ^b | 21.74 | 22 | 2.86 | 17 | 31 | 14 |
| 43 AOS ^c | 9.45 | 10 | 2.96 | 1 | 18 | 17 |
| 44 TOEFL | 98.46 | 99 | 5.43 | 90 | 110 | 20 |

45 ^a Length of residence in the United States in months.

46 ^b Age of arrival in the United States in years.

47 ^c Age of starting to learn English in China in years.
 48
 49

1 SPRT, and the WMT. The sets were equal in terms of length (length of
2 sentences was kept between 9 and 13 words) and complexity (all sentences
3 were simple in structure, with no embedded clauses), as well as
4 in the frequency and density of their lexicons (the frequency of the
5 words was checked in a corpus).⁵ The following samples targeting
6 mass/count nouns illustrate how sentences in these four sets were constructed.
7 In Set 1 and 1', the sentences were very similar, with minimal
8 changes to some of the words:
9

10
11 1: Mary added a lot of sugar(s) to her coffee.⁶

12 1': Tom added a lot of sugar(s) to his tea.
13

14 The same level of similarity existed between Set 2 and 2':
15

16 2: Mary likes to put a lot of sugar(s) to her coffee.

17 2': Tom likes to put a lot of sugar(s) in his tea.
18
19

20 However, when the differences between the sets were compared,
21 the difference between Set 1 and Set 2 was found to be greater than
22 that between 1' and 1 and 2 and 2', respectively. Thus, Sets 1 and 2
23 were used for the SPRT and WMT, whereas Sets 1' and 2' were employed
24 in the untimed and timed GJTs. By using less similar sentences (e.g.,
25 Sets 1 and 2) for more similar tasks (e.g., the SPRT and WMT), spurious
26 correlations between scores on the measures were avoided. In
27 other words, relationships between the scores, if found, should indicate
28 the commonality of the task designs (and constructs measured)
29 rather than the similarity among the sentences used in each task.
30 The tests were programmed and delivered through DMDX (Forster &
31 Forster, 2003).
32
33
34

35 **Timed and Untimed GJTs** 36 37

38 Each of the timed and untimed GJTs were composed of 96 sentences.
39 Sixteen items were presented for each target structure, half of which
40 were grammatical and the other half ungrammatical. Similarly, among
41 the 32 filler sentences (testing other target structures) included, 16 were
42 grammatical and 16 ungrammatical. The results pertaining to the
43 filler items were not included in the analyses. For each of the GJTs,
44 two counterbalanced lists of sentences were created. In List 1, half of
45 the target sentences were grammatical, and half ungrammatical. The
46 grammaticality of the sentences was reversed in List 2, to ensure
47 that no target sentence was presented twice in the same condition in
48 one list.
49

1 For the untimed GJT, participants were instructed to decide whether
2 the sentences were correct or incorrect and were reminded that there
3 was no time constraint. Sentences in the timed GJT appeared on the
4 screen for 2–5.5 s.⁷ The time limit for each item for the main part of
5 study was established based on the reaction time (RT) of the NSs in the
6 pilot. More specifically, the time limit for each of the sentences of the
7 GJTs was equal to $1.2 \times$ NSs' RT.
8

11 **SPRT**

14 The SPRT assessed online grammatical sensitivity while participants
15 were reading sentences for comprehension. In this task, participants
16 were asked to read a sentence word by word as quickly and accurately
17 as possible. The first word in a sentence appeared on the left-hand side
18 of the screen, and when the keyboard button was pressed, the next
19 word appeared to the right of the preceding word, which disappeared
20 on the presentation of the following word (moving-window presentation).
21 When participants read the final word followed by the period,
22 they pressed a second key to continue to a comprehension question.
23

24 To develop the SPRT, two lists of stimulus sentences were con-
25 structed, each consisting of 64 target sentences (16 for each structure).
26 The two lists were counterbalanced, whereby one half of the target sen-
27 tences were grammatical and the other half ungrammatical in List 1. In
28 List 2, the grammaticality of the sentences was reversed, so that no
29 target sentence was presented twice in the same condition in one list.
30 As before, 32 grammatical sentences were also included in each list as
31 filler sentences. All the sentences were followed by a comprehension
32 question that required a simple yes/no response. The ratio between yes
33 and no responses was kept equal. Once again, RT on the filler sentences
34 was excluded in the analyses.
35

36 As in the study conducted by Jiang (2007), the region of interest,
37 where RTs were compared between grammatical and ungrammatical
38 sentences, was set at three different word positions (see the under-
39 lined words in Table 2): at the critical word (i.e., where the error occurred
40 in the ungrammatical sentences) and at the two words immediately
41 following the critical word (to capture spillover effects). The word
42 preceding the critical region was also used as a baseline in order to
43 ascertain that the reading time of the word before the critical region
44 did not differ between grammatical and ungrammatical sentences. If
45 participants were sensitive to the grammatical error that preceded
46 the critical region, their reading time would be delayed at (some of)
47 these three positions. Table 2 shows some examples of each of the
48 target structures.
49

1 **Table 2.** Sample sentences with critical regions
2

| 3 Target | Sample sentence | Critical word |
|-------------------------|--|---------------|
| 5 Count/mass | Mary added a lot of <u>sugar(s)</u> to her coffee. | Sugar(s) |
| 6 Third-person -s | The boy in the room <u>enjoy(s)</u> reading many books. | Enjoy(s) |
| 8 Present perfect | Last spring he <u>(has)</u> planted many roses in the garden. | Has |
| 10 Present hypothetical | If I lived in Miami, I <u>can/could</u> have a house near the beach. | Can/could |

13
14 **WMT**

17 Similarly to the SPRT, the WMT also assessed online grammatical sensitivity during reading comprehension. Instead of self-paced reading, in this test participants were instructed to read a sentence presented automatically word by word on the screen. Their task was to respond to a target word or a monitoring word that appeared at one of the locations within the sentence. First, they were presented with a monitoring word in the center of the screen for 2 s. Next, each word in the sentence appeared on the screen for 1 s. The respondents were instructed to press the keyboard button as soon as they saw the target word in the sentence.

28 The monitoring word was always placed after the relevant target structure in the critical stimulus sentences. The difference in the RT to the target word between grammatical and ungrammatical sentences provided the index for online grammatical sensitivity. The monitoring word was located in the same position as the critical word in the SPRT, allowing the effects to be compared fairly between the WMT and the SPRT. Table 3 provides some examples of each of the target structures.

37
38 **Table 3.** Sample sentences with critical words and monitoring words
39

| 40 Target | Sample sentence | Monitoring word |
|-------------------------|--|-----------------|
| 41 Count/mass | Mary added a lot of <u>sugar(s)</u> to her coffee. | To |
| 43 Third-person -s | The boy in the room <u>enjoy(s)</u> reading many books. | Reading |
| 45 Present perfect | Last spring he <u>(has)</u> planted many roses in the garden. | Planted |
| 48 Present hypothetical | If I lived in Miami, I <u>can/could</u> have a house near the beach. | Have |

1 As in the SPRT, two lists of stimulus sentences were constructed
2 for the WMT. Each list consisted of 64 target sentences (16 for each
3 structure). The two lists were counterbalanced, whereby List 1 was
4 composed of an equal number of grammatical and ungrammatical tar-
5 get sentences, and the grammaticality of the sentences was reversed in
6 List 2. Once again, no target sentence was presented twice in the same
7 condition in one list. Similarly, 32 grammatical sentences were also
8 included in each list as filler sentences. All sentences were followed by
9 a comprehension question requiring a yes/no response, with an equal
10 ratio between the two. Reaction times for the filler sentences were
11 excluded in the analyses.
12
13
14

15 **MKT**

16
17
18 The MKT was constructed with 20 items pertaining to the same target
19 structures as in the GJTs, the SPRT, and the WMT. Five sentences were
20 used for each target structure. All sentences in this task were ungram-
21 matical and were similar to the ungrammatical sentences used in the
22 GJTs. For each item, a sentence appeared on the screen. Participants
23 were informed that the sentences were all ungrammatical, and their
24 task was to state the reason for the ungrammaticality and then provide
25 the correct form. There was no time constraint in this task. A rubric
26 (Appendix A) was developed for rating the learners' performance on
27 the MKT. The rubric detailed all of the possible acceptable and unac-
28 ceptable responses. According to this rubric, partial credit could be
29 assigned to each response, with 1 point for correct explanation and
30 1 point for correction. A total score of 2 was assigned for a response
31 that included both a correct explanation and correction. Two researchers
32 used the rubric and rated the responses independently. Their ratings
33 were subjected to Rasch analysis, and the ability logit produced by
34 Rasch was used as MKT data for further analyses.
35
36
37
38
39

40 **Procedure**

41
42 The five linguistic measures (timed GJT, untimed GJT, SPRT, WMT,
43 and MKT) were administered, starting with more implicit tasks and
44 progressing to more explicit tasks. Although other cognitive measures
45 were administered in the study, the results are not reported here. The
46 participants took these cognitive tests between the linguistic measures
47 to delay exposure to the target structures and to minimize any prac-
48 tice effect. All the measures were administered in a 2-hr session.
49

1 However, the battery was divided into two 1-hr blocks with a 15-min
 2 break between the two. Learners were paid \$25 for their participation.
 3 Table 4 shows the order and time for each of the measures.
 4

6 ANALYSES

8 Pilot Study Involving NSs

10 All linguistic measures employed in the study were piloted with English
 11 NSs. As noted previously, the pilot consisted of two separate assess-
 12 ments. In the first phase, 20 NSs were asked to check the test items for
 13 ease of comprehension and correctness. Item and reliability analyses
 14 were subsequently conducted to diagnose the problematic items in the
 15 GJTs and MKT. The RTs of the NSs on the timed GJT items were also
 16 computed, to set the time limit for individual sentences in the task for
 17 the main study. The SPRT and WMT data were then analyzed to ensure
 18 that the tasks captured NSs' sensitivity to the incorrect target struc-
 19 tures. The first phase of the pilot study identified some problematic
 20 items in the timed GJT, the untimed GJT, and the SPRT, prompting
 21 appropriate revisions. The revised tasks were subsequently adminis-
 22 tered to another group of the NSs ($n = 10$), who took part in the second
 23 phase of the pilot study. Because the WMT and MKT data from the first
 24 phase showed that the task functioned as expected, they were not given
 25 to the second group of NSs.
 26

27 **GJTs.** All the items from both GJTs were scored dichotomously, and
 28 the results were recorded as zero or one. Through item and reliability
 29 analyses, items with error rates higher than 25% were flagged for revision.
 30 These items were revised to ensure that grammatical and ungram-
 31 matical sentences clearly functioned in the expected manner. When
 32 the second phase of the pilot, with 10 NSs, was carried out, the results
 33 revealed that none of the items had an error rate higher than 25%.
 34

35 **Table 4.** Order and time of the measures in each of the two blocks

| Block 1 | | Block 2 | |
|---------------------|--------|--------------------------|--------|
| Task | Time | Task | Time |
| Consent form | 5 min | Timed GJT | 10 min |
| WMT | 20 min | Cognitive measure 2 | 10 min |
| Cognitive measure 1 | 15 min | Untimed GJT | 20 min |
| SPRT | 20 min | MKT | 15 min |
| | | Background questionnaire | 5 min |

1 The NSs' RTs in the second round were used to set the time limit for the
2 individual sentences for the learners. Following the work of R. Ellis
3 (2005), the NSs' RT to each individual sentence was increased by 20% to
4 set the time limit. Depending on the length and complexity of the sen-
5 tences, the time limit varied across items. The average RT for the entire
6 test was 3.39 s, and RTs ranged from 2 to 5.5 s.
7

8
9 **SPRT.** After the first phase of the pilot ($n = 20$), some items on the
10 SPRT were revised. Based on the statistical results and the NSs' feed-
11 back, several sentences were revised to make them sound more nat-
12 ural or unambiguously grammatical or ungrammatical. The revised
13 SPRT was then given to another group of NSs ($n = 10$), and their RTs
14 to the word prior to the target structure (Region 0) and the average
15 RT in the critical region (the target word and the following two words)
16 for both grammatical and ungrammatical sentences were measured.
17 Two assumptions were tested for the SPRT as a linguistic measure—
18 namely, for the NSs (a) the average RT to Region 0 should be statisti-
19 cally the same across grammatical and ungrammatical sentences,
20 and (b) the average RT to the critical region for the ungrammatical
21 sentences should be statistically greater than that measured for the
22 grammatical sentences. The results showed that (a) participants read
23 the words prior to the target structure similarly in both grammatical
24 and ungrammatical sentences, and that (b) NSs, who have IK of the
25 target structures, slowed down when they came across the incorrect
26 use of target structures. Through four separate sets of paired-samples
27 t -tests, RTs to Region 0 and the critical region for the four target struc-
28 tures in the grammatical and ungrammatical sentences were com-
29 pared. Table 5 summarizes the results from the second phase of the
30 pilot, with 10 NSs.
31

32
33 The results revealed that there was no statistically significant differ-
34 ence in RTs to Region 0 across grammatical and ungrammatical sen-
35 tences. On the other hand, for all four target structures, the RTs to the
36 critical region were statistically greater in the ungrammatical sentences
37 than in the grammatical sentences.
38

39
40 **WMT.** In order to establish whether NSs showed online sensitivity to
41 the target structures, the RTs obtained by the group of NSs in the first
42 phase of the pilot ($n = 20$) were analyzed. Through four separate sets of
43 paired-samples t -tests, RTs to the target word across grammatical and
44 ungrammatical items were compared, and the results are summarized
45 in Table 6.
46

47 The results revealed that RTs to the target word in the ungrammat-
48 ical sentences were statistically greater than RTs in the grammatical
49 sentences. These results confirm that the WMT captured the online
sensitivity.

Table 5. Results of the paired-samples *t*-tests and descriptive statistics for the SPRT (NS data, second phase of the pilot)

| Outcome | Grammatical | | Ungrammatical | | <i>n</i> | <i>t</i> | <i>df</i> |
|--------------------------------------|-------------|-----------|---------------|-----------|----------|----------|-----------|
| | <i>M</i> | <i>SD</i> | <i>M</i> | <i>SD</i> | | | |
| Third-person Region 0 | 255.5 | 21.13 | 258.5 | 24.2 | 10 | -.35 | 9 |
| Mass/count nouns Region 0 | 257.2 | 14.49 | 256.9 | 9.59 | 10 | .07 | 9 |
| Past/perfect Region 0 | 276.6 | 16.26 | 293.4 | 32.09 | 10 | -2.1 | 9 |
| Present hypothetical Region 0 | 326.8 | 69.89 | 311.2 | 58.85 | 10 | .47 | 9 |
| Third-person critical region | 668.9 | 54.1 | 999.6 | 312.29 | 10 | -3.02* | 9 |
| Mass/count nouns critical region | 614.2 | 63.2 | 1,337.3 | 285.82 | 10 | -8.37* | 9 |
| Past/perfect critical region | 600.9 | 73.52 | 990.3 | 227.6 | 10 | -6.31* | 9 |
| Present hypothetical critical region | 720.4 | 121.66 | 1,599.3 | 408.22 | 10 | -6.31* | 9 |

* $p < .05$.

Main Study on L2 Learners

After the two-stage pilot testing, the revised test battery was administered to the 80 L2 learners. Even though 80 learners took all the tasks, the data pertaining to one of the participants had to be excluded because the computer failed to record the required information for several tasks. Thus the following analyses were conducted on a sample size of 79.

Item and Reliability Analysis. To begin with, item analyses were conducted on all the items in both the timed and untimed GJTs, and analyses were repeated for the grammatical and ungrammatical sentences separately. Henceforth, for brevity, the following acronyms will be used for different combinations of GJTs: total timed GJT (T-GJT), timed GJT

Table 6. Results of the *t*-tests and descriptive statistics for the WMT

| Outcome | Grammatical | | Ungrammatical | | <i>n</i> | <i>t</i> | <i>df</i> |
|----------------------|-------------|-----------|---------------|-----------|----------|----------|-----------|
| | <i>M</i> | <i>SD</i> | <i>M</i> | <i>SD</i> | | | |
| Third person | 527.6 | 117.58 | 670.65 | 237.47 | 20 | -3.05* | 19 |
| Mass/count nouns | 508.75 | 97.26 | 625.1 | 255.44 | 20 | -2.63* | 19 |
| Past/perfect | 421.3 | 88.14 | 621.15 | 352.91 | 20 | -2.55* | 19 |
| Present hypothetical | 415.95 | 71.61 | 589.3 | 316.34 | 20 | -2.62* | 19 |

* $p < .05$.

1 grammatical sentences (T-GJT-G), timed GJT ungrammatical sentences
 2 (T-GJT-U), total untimed GJT (Un-GJT), untimed GJT grammatical sen-
 3 tences (Un-GJT-G), and untimed GJT ungrammatical sentences (Un-GJT-U).
 4 For each of the GJTs, items with a zero or negative item discrimination
 5 (ID) index (item-total correlation) were deleted to improve the reliability.
 6 Eight items were deleted from the timed GJT, and twelve from the untimed
 7 GJT. Next, reliability was estimated (Cronbach’s alpha) for the aforemen-
 8 tioned GTJ measures, the SPRT, and the WMT. The interrater reliability of
 9 the two raters for the MKT was also computed. Two independent ratings
 10 for the MKT were subjected to a Rasch analysis, and the ability logit for
 11 individual learners was computed to generate data for subsequent
 12 analyses. Table 7 summarizes the reliability estimates for all the tasks of
 13 the study. The reliability estimates for the GJTs and the MKT in the cur-
 14 rent study were within the acceptable range (above .70, as suggested by
 15 Nunnally, 1978), and they all were within the ranges reported in the pre-
 16 vious research (Bowles, 2011; R. Ellis, 2005; Gutiérrez, 2013; Zhang, 2015).
 17 Cronbach’s alpha for the SPRT was high, whereas the WMT’s internal
 18 consistency was lower but close to .70.
 19
 20

21
 22 **Descriptive Statistics.** Following item analysis and exclusion of items
 23 with an inappropriate ID index, descriptive statistics were computed
 24 for all measures. Using multiple sources, such as statistical tests
 25 (Kolmogorov–Smirnov and Shapiro–Wilk tests) and graphs, the univar-
 26 iate normality of the measures was assessed. Depending on the severity
 27 of the skewness, two types of transformation are usually conducted.
 28 A square root transformation is carried out for data that differ moder-
 29 ately from the normal distribution, whereas a log transformation is
 30 more appropriate for data exhibiting substantial deviation (Tabachnick &
 31 Fidell, 2001). The T-GJT, T-GJT-U, Un-GJT-U, SPRT, and WMT data did not
 32 require any transformation. A square root transformation was carried
 33 out for the remaining measures because they differed from the normal
 34
 35
 36

37 **Table 7.** Reliability estimates

| 38 Task | 39 Number of items | 40 Reliability estimate |
|-------------|--------------------|------------------------------|
| 41 T-GJT | 42 56 | .75 |
| 43 T-GJT-G | 44 26 | .74 |
| 45 T-GJT-U | 46 30 | .78 |
| 47 Un-GJT | 48 52 | .83 |
| 49 Un-GJT-G | 22 | .74 |
| Un-GJT-U | 30 | .88 |
| WMT | 32 | .65 |
| SPRT | 32 | .95 |
| MKT | 20 | Interrater reliability = .91 |

distribution moderately. No log transformation was conducted. Table 8 summarizes descriptive statistics for all the measures. As can be seen, the univariate skewness and kurtosis values for all measures were within the acceptable range of +/- 1, with the exception of kurtosis for SPR. However, when a multivariate normality test was conducted, the assumption of multivariate normality was met (skewness: z score = .76, p value = .447; kurtosis: z score = 1.31, p value = .19; skewness and kurtosis: chi-squared = 2.3, p value = .32).

Correlational Analysis. Before conducting CFA, in order to explore the relationships among the measures of interest for the present study, a Pearson product moment correlation analysis was conducted. Table 9 summarizes the results, presenting only the statistically significant correlations, for clarity. For the complete table of correlations, see Appendix B.

As can be seen in Table 9, the WMT and SPRT correlated with each other, but not with any other measures. The T-GJT correlated only with the grammatical and ungrammatical sentences contained within. The Un-GJT also correlated with the grammatical and ungrammatical sentences it was composed of, as well as with the MKT. The T-GJT-G correlated with the Un-GJT-G only, but the T-GJT-U correlated with the Un-GJT, Un-GJT-U, and MKT. Finally, the Un-GJT-U correlated not only with the Un-GJT and T-GJT-U but also with the MKT.

Confirmatory Factor Analysis. Confirmatory factor analysis was chosen as the main method of data analysis in order to test the prior hypotheses. Unlike in previous studies, 20 different CFA models were tested. These models included the model hypothesized in this study and models employed in previous studies (i.e. Bowles, 2011; R. Ellis & Loewen, 2007; Gutiérrez, 2013; Zhang, 2015), as well as several other rival models. Table 10 summarizes information pertaining to all tested models. It also

Table 8. Descriptive statistics

| Task | Mean | SD | Min | Max | Skew | Kurt |
|----------|--------|--------|--------|------|------|------|
| T-GJT | 27.25 | 6.52 | 12 | 43 | -.08 | -.2 |
| T-GJT-G | 23.27 | .76 | 21.42 | 25 | -.05 | -.16 |
| T-GJT-U | 9.3 | 4.85 | 0 | 21 | .4 | -.35 |
| Un-GJT | 47.3 | .94 | 45.34 | 50 | .25 | .54 |
| Un-GJT-G | 20.88 | .715 | 19.26 | 22 | -.42 | -.55 |
| Un-GJT-U | 18.42 | 6.15 | 0 | 29 | -.68 | .08 |
| MKT | 1.86 | .36 | .93 | 2.59 | .02 | .44 |
| SPRT | 10.72 | 298.91 | -1,274 | 749 | -.82 | 3.78 |
| WMT | -13.71 | 188.29 | -486 | 448 | .03 | .22 |

Table 9. Correlational matrix for all the measures

| Task | T-GJT | T-GJT-G | T-GJT-U | Un-GJT | Un-GJT-G | Un-GJT-U | MKT | SPRT | WMT |
|----------|--------------|---------|--------------|--------|----------|----------|--------------|------|-----|
| T-GJT | — | | | | | | | | |
| T-GJT-G | .65** | — | | | | | | | |
| T-GJT-U | | | — | | | | | | |
| Un-GJT | .33** | | .49** | — | | | .22* | | |
| Un-GJT-G | | | .48** | | — | | .45** | | |
| Un-GJT-U | | | | | | — | .47** | | |
| MKT | | | | | | | — | | |
| SPRT | | | | | | | | — | |
| WMT | | | | | | | | | — |

** Correlation is significant at the 0.01 level (two-tailed).

* Correlation is significant at the 0.05 level (two-tailed).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Table 10. Summary of the tested CFA models

| Model | Explicit knowledge (EK) | Implicit knowledge (IK) | Language knowledge (LK) |
|---|---|---|--|
| 1a (R. Ellis & Loewen, 2007) | Un-GJT-U, MKT | T-GJT, WMT, SPRT | |
| 1b (1a with correlated errors) | Un-GJT-U, MKT <i>Un-GJT-U/MKT</i> | T-GJT, WMT, SPRT <i>WMT/SPRT</i> | |
| 2a (Bowles, 2011) | Un-GJT, MKT | T-GJT, WMT, SPRT | |
| 2b (2a with correlated errors) | Un-GJT, MKT <i>Un-GJT/MKT</i> | T-GJT, WMT, SPRT <i>WMT/SPRT</i> | |
| 3a (Gutiérrez, 2013) | T-GJT-U, Un-GJT-U, MKT | T-GJT-G, Un-GJT-G, WMT, SPRT | |
| 3b (3a with correlated errors) | T-GJT-U, Un-GJT-U, MKT <i>T-GJT-U/Un-GJT-U</i> | T-GJT-G, Un-GJT-G, WMT, SPRT <i>T-GJT-G/Un-GJT-G / WMT/SPRT</i> | |
| 3c (Gutiérrez, 2013) | Un-GJT-G, Un-GJT-U, MKT | T-GJT-G, T-GJT-U, WMT, SPRT | |
| 3d (3c with correlated errors) | Un-GJT-G, Un-GJT-U, MKT <i>Un-GJT-G/Un-GJT-U</i> | T-GJT-G, T-GJT-U, WMT, SPRT <i>T-GJT-G/T-GJT-U and WMT/SPRT</i> | |
| 4a (our hypothesized model) | Un-GJT-U, T-GJT-U, MKT | WMT, SPRT | |
| 4b (T-GJT-U cross-loading) | Un-GJT-U, T-GJT-U, MKT | T-GJT-U, WMT, SPRT | |
| 4c (Un-GJT-U cross-loading) | Un-GJT-U, T-GJT-U, MKT | Un-GJT-U, WMT, SPRT | |
| 5 | T-GJT-G, T-GJT-U, Un-GJT-G, Un-GJT-U, MKT | WMT, SPRT | |
| 6 | T-GJT, Un-GJT, MKT | WMT, SPRT | |
| 7 (one-factor) | | | T-GJT-U, Un-GJT-U, MKT, WMT, SPRT ^a |
| 8a (one-factor with correlated errors) | | | T-GJT-U, Un-GJT-U, MKT, WMT, SPRT <i>T-GJT-U/Un-GJT-U and WMT/SPRT</i> |

Continued

Table 10. continued

| Model | Explicit knowledge (EK) | Implicit knowledge (IK) | Language knowledge (LK) |
|---|-------------------------|-------------------------|---|
| 8b (one-factor with correlated errors) | | | T-GJT-U, Un-GJT-U, MKT, WMT, SPRT <i>T-GJT-U/Un-GJT-U, MKT/Un-GJT-U and WMT/SPRT</i> |
| 8c (one-factor with correlated errors) | | | T-GJT-U, Un-GJT-U, MKT, WMT, SPRT <i>T-GJT-U/Un-GJT-U, MKT/Un-GJT-U, MKT/T-GJT-G and WMT/SPRT</i> |
| 9a (one-factor) | | | T-GJT-G, T-GJT-U, Un-GJT-G, Un-GJT-U, MKT, WMT, SPRT |
| 9b (one-factor with correlated errors) | | | T-GJT-G, T-GJT-U, Un-GJT-G, Un-GJT-U, MKT, WMT, SPRT <i>T-GJT-G/T-GJT-U, Un-GJT-G/Un-GJT-U and WMT/SPRT</i> |
| 9c (one-factor with correlated errors) | | | T-GJT-G, T-GJT-U, Un-GJT-G, Un-GJT-U, MKT, WMT, SPRT <i>T-GJT-G/Un-GJT-G, T-GJT-U/Un-GJT-U and WMT/SPRT</i> |

Note. Correlated errors are italicized.
^a No correlated errors.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

1 shows which measured variables loaded on which factors, as well as
2 the relationship between error terms for models in which method effect
3 was accounted for. LISREL Version 9.1 (Jöreskog & Sorbom, 2012) was
4 used for running CFA.

5 The models labeled 1a to 3d tested the models adopted in previous
6 studies (i.e., Bowles, 2011; R. Ellis & Loewen, 2007; Gutiérrez, 2013;
7 Zhang, 2015). The details of these models are not explained here but are
8 available in Table 10. The table provides all immediately necessary
9 information about each of these models. Regarding models from pre-
10 vious studies, a version of each with correlated error terms for the measured
11 variables was also tested. In CFA, the specification of correlated
12 errors is made on the basis of source or method effects, which explains
13 the additional indicator covariation that resulted from common assessment
14 methods (Brown, 2006). In other words, in correlating errors, the
15 aim was to account for the method effect and to improve the model fit
16 indices. The models from previous studies were alternative models to
17 the model hypothesized in the present study.

18 The next stage of the analysis involved testing Model 4a, developed
19 as a part of this study. In this case, the T-GJT-U, Un-GJT-U, and MKT
20 loaded on the EK factor, and the SPRT and WMT loaded on the IK
21 factor. Next, rebuttals to Model 4a were examined by testing rival
22 models (Models 4b–9c). In Models 4b and 4c, the T-GJT-U and Un-GJT-
23 U cross-loaded on both factors, respectively. In Model 5, all combina-
24 tions of GJTs loaded on the EK factor, and in Model 6, the total score
25 from both the grammatical and ungrammatical sentences of both GJTs
26 loaded on the EK factor. In Models 5 and 6, the SPRT and WMT loaded
27 on the IK factor.

28 Models 7–9c allowed testing rival models to Model 4a. In these models,
29 all the measured variables loaded on a single factor, labeled language
30 knowledge (LK). As can be seen in Table 10, in order to account for the
31 method effect, the error terms of various measured variables were also
32 correlated. Several versions of Model 4, in which the error terms of var-
33 ious measures were correlated, were subsequently tested. However, as
34 these models did not improve Model 4 significantly, the related findings
35 are not reported here (see Appendix C for the covariance matrix used
36 for the preceding analyses).

37 To evaluate and compare the plausibility of the CFA models, a profile
38 of model fit indices recommended by Hu and Bentler (1999) and Mueller
39 and Hancock (2008) was used and reported: (a) the chi squared (χ^2),
40 with its degrees of freedom and *p*-value; (b) the standardized root mean
41 square residual (SRMR); (c) the root mean square error of approxima-
42 tion (RMSEA); the comparative fit index (CFI); (d) the normal fit index
43 (NFI); (e) the nonnormed fit index (NNFI); and (f) the goodness-of-fit
44 statistic (GFI). For a model to be deemed a good fit to the data, it had to
45 meet the following criteria: the chi-squared should not be statistically
46
47
48
49

1 significant at a .05 level; the SRMR and RMSEA should be lower than .08
 2 and .06, respectively; and values greater than .95 for the CFI and NNFI
 3 and .90 for the NFI and GFI indicate a good model fit. Table 11 summa-
 4 rizes the model fit indices for all CFA models examined.

5 As can be seen in Table 11, in terms of model fit indices, only six of the
 6 models (4a, 4b, 4c, 6, 8a, and 8b) fit the data well. Fit indices for the
 7 remaining models were outside the acceptable range.

8 Model 4a was the model hypothesized in this study and was com-
 9 pared to the other models with acceptable fit indices. These models are
 10 nested models (with the exception of Model 6, which employs different
 11 measures from the rest), and, therefore, a formal chi-squared difference
 12 test was conducted by $\Delta\chi^2_{(df_1 - df_2)} = \chi^2_{df_1} - \chi^2_{df_2}$ and was distributed as a chi-
 13 squared distribution with $df = df_1 - df_2$ (Mueller & Hancock, 2008). Table 12
 14 presents the results yielded by the chi-squared difference tests.

15 As can be seen in Table 12, the chi-squared difference test was not
 16 statistically significant for any of the comparisons. These results sug-
 17 gest that no model is statistically different from any other. It should be
 18
 19
 20
 21

22 **Table 11.** Summary of the model fit indices for the tested CFA
 23 models

| Index | CFI | NFI | NNFI | GFI | RMSEA | SRMR | Chi-squared |
|-----------|-------|--------|--------|-------|-------|-------|-------------------------------|
| Criterion | ≤ .95 | ≤ .90 | ≤ .95 | ≤ .90 | ≥ .06 | ≥ .08 | Nonsignificant |
| Model 1a | 0 | -30.41 | -116.1 | .24 | 1.73 | 0.53 | * $\chi^2 = 945.12, df = 4$ |
| Model 1b | .5 | .6 | -1.48 | .95 | .25 | .09 | * $\chi^2 = 11.94, df = 2$ |
| Model 2a | 0 | -37.61 | -144.9 | .2 | 1.89 | 2.07 | * $\chi^2 = 1,132.75, df = 4$ |
| Model 2b | .37 | .51 | -2.16 | .94 | .28 | .11 | * $\chi^2 = 14.24, df = 2$ |
| Model 3a | .77 | .65 | .63 | .92 | .1 | 0.08 | * $\chi^2 = 23.317, df = 13$ |
| Model 3b | .86 | .75 | .7 | .95 | .09 | .07 | $\chi^2 = 16.35, df = 10$ |
| Model 3c | .22 | .27 | -.26 | .88 | .18 | .14 | * $\chi^2 = 47.74, df = 13$ |
| Model 3d | .56 | .55 | .07 | .91 | .16 | .11 | * $\chi^2 = 29.75, df = 10$ |
| Model 4a | .99 | .90 | .96 | .98 | .05 | .06 | $\chi^2 = 4.66, df = 4$ |
| Model 4b | 1 | .98 | 1.17 | .99 | 0 | .03 | $\chi^2 = 1.09, df = 3$ |
| Model 4c | 1 | .97 | 1.13 | .99 | 0 | .04 | $\chi^2 = 1.46, df = 3$ |
| Model 5 | .76 | .64 | .61 | .76 | .1 | .09 | * $\chi^2 = 23.864, df = 13$ |
| Model 6 | 1 | .91 | 1.19 | .99 | 0 | .05 | $\chi^2 = 2.53, df = 4$ |
| Model 7 | .9 | .82 | .8 | .96 | .1 | .09 | $\chi^2 = 8.93, df = 5$ |
| Model 8a | .96 | .91 | .88 | .98 | .08 | .06 | $\chi^2 = 4.45, df = 3$ |
| Model 8b | .99 | .95 | .95 | .99 | .05 | .05 | $\chi^2 = 2.38, df = 2$ |
| Model 8c | .86 | .86 | -.45 | .97 | .27 | .07 | * $\chi^2 = 6.6, df = 1$ |
| Model 9a | .68 | .57 | .53 | .92 | .11 | .1 | * $\chi^2 = 28.1, df = 14$ |
| Model 9b | .75 | .66 | .51 | .94 | .11 | .09 | * $\chi^2 = 22.4, df = 11$ |
| Model 9c | .88 | .75 | .78 | .95 | .08 | .08 | $\chi^2 = 16.24, df = 11$ |

49 * χ^2 is significant at the .05 level.

AQ38

Table 12. Chi-squared difference formal test results

| | | |
|---|--|---|
| Model 4a: $\chi^2 = 4.66,$ $df = 4$ | Model 4b: $\chi^2 = 1.09, df = 3$ | $\Delta \chi^2 = 3.57, df = 1, p \text{ value} = .06$ |
| | Model 4c: $\chi^2 = 1.46, df = 3$ | $\Delta \chi^2 = 3.2, df = 1, p \text{ value} = .07$ |
| | Model 8a: $\chi^2 = 4.45, df = 3$ | $\Delta \chi^2 = .21, df = 1, p \text{ value} = .65$ |
| | Model 8b: $\chi^2 = 2.38, df = 2$ | $\Delta \chi^2 = 2.28, df = 2, p \text{ value} = .32$ |

noted that Model 6 could not be statistically compared to the other models because it utilizes different sets of measured variables.

In the next step, the factor loadings in Model 4a were compared to those in all the other models (except for Model 6) that provided a good fit to the data. The model with higher and significant factor loadings is considered to fit the data better. Table 13 summarizes the factor loadings for these six models.

The examination of the factor loadings revealed that only Model 4a had significant loadings for all the measured variables. In Models 4b and 4c, the T-GJT-U and Un-GJT-U did not load on the IK factor significantly. Therefore, Model 4a, a more parsimonious model, is superior to the others. In Model 8a and 8b, SPRT and WMT did not load on the LK factor significantly. In Model 8a, the loading of the MKT was not statistically significant either. Finally, in Model 6, the T-GJT, MKT, and WMT did not load on their corresponding factors significantly. The only issue with Model 4a that needs further explanation is that, at 1.34, the loading of the Un-GJT-U exceeded 1, indicating that the error variance for this measured variable was negative. However, according to Jöreskog (1999), standardized coefficients (loadings) that are greater than 1, especially if they are not significant and are smaller than 2.8, are not an issue. In addition, in Model 4a, the correlation between the two factors was not statistically significant. These results lead to the conclusion that Model 4a, which was the hypothesized

Table 13. Factor loadings for Models 4a, 4b, 4c, 8a, 8b, and 6

| Model | T-GJT | T-GJT-U | Un-GJT | Un-GJT-U | MKT | SPRT | WMT |
|----------|-------|--------------------------------|--------|--------------------------------|--------|------|--------|
| Model 4a | | *.36 | | **1.34 | *.33 | *.42 | *.58 |
| Model 4b | | *.53 for EK and -.33 for IK | | **1.07 | ** .43 | *.42 | ** .57 |
| Model 4c | | ** .5 | | **1.1 for EK and .47 for IK | ** .44 | *.37 | ** .65 |
| Model 8a | | *.71 | | *2.28 | .20 | .05 | .08 |
| Model 8b | | ** .25 | | **3.14 | ** .91 | .05 | .07 |
| Model 6 | .08 | | *2.63 | | .17 | *.76 | .32 |

** Loading is significant at the .01 level.

* Loading is significant at the .05 level.

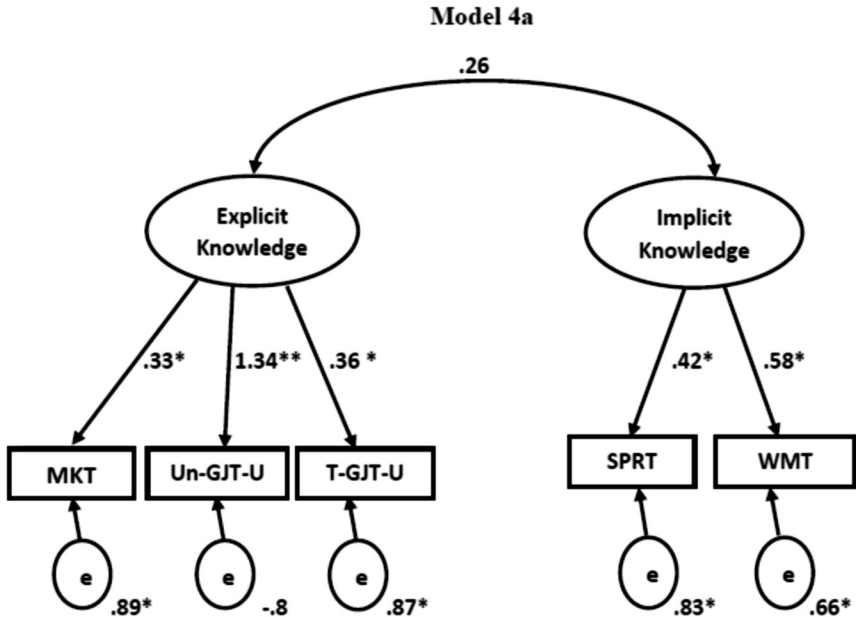
1 model of the study, provides the best fit to the data. Figure 1 illus-
2 trates the best fitting Model 4a.

3
4
5 **DISCUSSION**

6
7 **Construct Validity of the GJTs**

8
9
10 The current study aimed at investigating whether manipulating the
11 time condition and/or stimulus type transforms GJTs into distinct
12 measures of EK and IK. Unlike the previous validation studies (Bowles,
13 2011; R. Ellis, 2005; R. Ellis & Loewen, 2007; Zhang, 2015), in this work it
14 was hypothesized that GJTs of any kind are too coarse to be measures
15 of IK. Grammaticality judgment tests draw attention to form, and applying
16 time pressure does not necessarily prevent L2 learners from accessing
17 EK (Suzuki & DeKeyser, in press). In addition, contrary to what Gutiérrez
18 (2013) proposed, here it was hypothesized that the ungrammatical
19 sentences of GJTs would provide a more valid measure of (explicit)
20 L2 knowledge of the target structures.
21

AQ14



22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45 **Figure 1.** The best fitting CFA model (Model 4a).
46 *Note.* MKT = metalinguistic knowledge test, Un-GJT-U = ungrammatical
47 sentences of the untimed GJT, T-GJT-U = ungrammatical sentences
48 of the timed GJT, SPRT = self-paced reading task, and WMT = word-
49 monitoring task.

AQ39

1 To test these hypotheses, four different types of GJTs were devel-
2 oped: grammatical sentences in a timed GJT, ungrammatical sentences
3 in a timed GJT, grammatical sentences in an untimed GJT, and ungram-
4 matical sentences in an untimed GJT. In the tests, participants' total
5 scores for the timed and untimed GJTs were calculated as well. Next,
6 the learners' performance on the GJTs was compared with their perfor-
7 mance on two other online processing measures (the WMT and SPRT),
8 on one hand, and a well-established measure of EK (the MKT) on the
9 other. Several CFA models were tested, including the ones adopted in
10 previous studies, in which good fit indices were reported.

11 The CFA produced the best fitting two-factor model consisting of
12 EK loaded onto the ungrammatical sentences of both the timed and
13 untimed GJTs and the MKT and of IK loaded onto the WMT and SPRT.
14 All models similar to those adopted in previous studies failed to achieve
15 acceptable fit to the data, suggesting that the prior findings are of ques-
16 tionable validity. The bivariate correlation coefficients also showed
17 that the grammatical sentences of both GJTs only correlated with each
18 other, and not with any other measures. This finding may suggest that
19 grammatical sentences in GJTs behave differently from the remaining
20 measures employed in this work.

21 The correlation between the two factors in the best fitting model
22 hypothesized in the current study (Model 4a) was small in magnitude
23 and not statistically significant ($r = .26$). As noted by Brown (2006), "The
24 size of the factor correlations in multifactorial CFA solutions should be
25 interpreted with regard to the *discriminant validity* of the latent con-
26 structs. Small, or statistically non-significant, factor covariances are not
27 usually problematic and are typically retained in the solution" (p. 131).
28 The small, nonsignificant correlation between the factors in the current
29 study can thus serve as evidence that the measures employed tapped
30 into two distinct constructs. Previous studies, such as those conducted
31 by Bowles (2011) and Zhang (2015), yielded large and significant corre-
32 lations between the two factors in the models reported to have the best
33 fit to the data. More specifically, Bowles (2011) reported a correlation of
34 .87 between the two factors, whereas Zhang (2015) obtained .86. These
35 high correlations weaken their argument that timed and untimed GJTs
36 are distinct measures for EK and IK. This is consistent with Kachisnke
37 and Vafaei's (2014) conclusion that the data from the previous valida-
38 tion studies fit well in the one-factor model accounting for the method
39 effect as well as in the two-factor models.

40 Given the small sample size in the current study, further post hoc
41 analysis was conducted to assess whether the two factors in the best
42 fitting model (Model 4a) can be recovered with the current sample size.
43 MacCallum, Widaman, Zhang, and Hong (1999) and Velicer and Fava
44 (1998) provide a thorough review of studies investigating the role of
45 sample size in factor analysis. They demonstrated that the minimum
46
47
48
49

1 sample size required largely depends on several factors, among which
2 communalities were the most influential. According to MacCallum,
3 Widaman, Preacher, and Hong (2001), “The level of communality has an
4 especially strong interaction with N such that when communalities are
5 high, good recovery of population factors can be achieved with rela-
6 tively small samples” (p. 612). It has been suggested that low commu-
7 nalities in particular pose a serious problem in small sample sizes. The
8 average communality was thus computed⁸ for Model 4a, and it was .51.
9 The value is above the low communality level (.35), as suggested by
10 MacCallum et al. (2001), and this lends some support for the adequacy
11 of this sample.
12

13 In sum, the present findings challenge the prior claims that timed and
14 untimed GJTs are distinct measures of EK and IK constructs. Thus,
15 based on the findings of the extant studies, it cannot be concluded that
16 timed GJTs are measures of IK. If behavioral measures are considered to
17 lie on a continuum from more explicit to more implicit, GJTs are prob-
18 ably considered closer to the explicit end of the continuum.
19
20

21 **Using Online Processing Measures to Capture IK**

22
23
24

25 The current study demonstrated that the online psycholinguistic
26 measures (the WMT and SPRT) are tapping into a different construct
27 than the one the GJTs are drawing on. This suggests that those newer
28 tasks probably lie closer to the implicit end of the continuum relative to
29 the GJTs. As delineated previously, online sentence-processing tasks
30 can minimize the involvement of EK by capturing the online sensitivity
31 to violations while attention is directed to meaning. Self-paced reading
32 tasks and WMTs have been primarily utilized in psycholinguistic inves-
33 tigation (e.g., Clahsen & Felser, 2006) and have not been explicitly
34 employed for addressing the issues of measurements of EK and IK until
35 recently (Suzuki, 2015). The current findings corroborate those reported
36 by Suzuki and DeKeyser (in press), who noted that real-time grammat-
37 ical processing can index IK. In the work of Suzuki and DeKeyser (in
38 press), the WMT was the only linguistic measure for IK and was not
39 contrasted with the results of any other processing measure, such as a
40 SPRT. The WMT adopted by Suzuki and DeKeyser (in press) was admin-
41 istered in the auditory modality,⁹ whereas the WMT and the SPRT in the
42 current study were given in the visual modality. Despite the modality
43 difference, the critical design of the task was shared, and their perfor-
44 mances converged. Combined, these results suggest that online com-
45 prehension tasks are good candidates for IK measures.
46
47

48 Recently, in SLA studies on implicit “learning,” researchers have
49 also started to apply RT-based online measures to assess to what extent

1 implicit learning took place (e.g., Leung & Williams, 2011; Paciorek &
2 Williams, 2015). Processing measures are more advantageous for capturing
3 how L2 learners access linguistic knowledge in real time and can
4 potentially reveal implicit learning processes. The recent applications
5 of online measures in implicit learning research support the current
6 finding that online comprehension measures draw on IK. Online processing
7 measures offer great promise for further validation of EK and IK
8 measures, as this study demonstrated that they can be employed to
9 further scrutinize the validity evidence for the GJTs.
10

11 12 13 14 **Suggestions for Further Research** 15

16 The current study is not without limitations and offers several venues
17 for further research. First, in order to investigate the construct validity
18 of behavioral measures hypothesized to measure IK, the issue of awareness
19 remains central. If the key definition of IK is that it definitively *does*
20 *not involve* awareness, measures of awareness should be included in
21 the validation studies. The current study did not employ any such
22 measure of awareness.
23

24 Second, a test may be valid for certain purposes, such as a particular
25 learning context or population of test takers, but not for others (Henning,
26 1987). The participants in the current study were learners with both
27 classroom-based and naturalistic learning experience. The validity of
28 GJTs may be assessed in different ways, depending on the specific context
29 of the study in which GJTs are used. For instance, heritage learners
30 with less formal instruction may perform on GJTs differently from the
31 sample recruited for the present study, as they are posited to be less
32 likely to possess EK compared to classroom learners.
33

34 Third, the current test battery only targeted four target structures.
35 The original research by Rod Ellis targeted seventeen structures, with a
36 similar number used in the subsequent studies. The smaller number of
37 types of target structures tested in the current study may limit the generalizability
38 of the present findings.

39 In addition, the current study focused only on Chinese L2 learners,
40 whose L1 is typologically different from English. If the same test battery
41 were administered to another L1 group, the results would likely be different.
42 Because the explicit and implicit learning processes may be interactively
43 influenced by the target structures to be acquired and learners' prior
44 knowledge, including their L1 (e.g., DeKeyser, 2003; Leung & Williams,
45 2014; Williams, 2005), it may be worth expanding the current research
46 to a different population and/or testing other linguistic structures.
47

48 Finally, in the model (Model 4a) that provided the best fit to the data,
49 as well as in the remaining models, the factor loadings of the measured

1 variables were relatively low. However, for stable latent variables to be
 2 measured, models are usually required to achieve higher factor loadings.
 3 Obviously, more rigorous further studies on the validity of GJTs are
 4 urgently needed.
 5

7 CONCLUSION

10 The current study set out to investigate the validity of GJTs as measures
 11 of EK and IK. The data analyses were extremely stringent and included
 12 a thorough execution of CFA. Thus it was possible to reveal the limita-
 13 tions of the previous research. Specifically, the study provided evidence
 14 that challenged conclusions from previous studies that time pressure
 15 renders GJTs measures of IK. The claim that timed GJTs are measures of
 16 IK no longer holds when online comprehension tasks that can capture
 17 grammatical sensitivity are employed. Given the nature of GJTs—they
 18 involve focus on form—GJTs may be considered to be located closer to
 19 the explicit, rather than the implicit, end of the continuum (DeKeyser,
 20 2003). Endorsing the recent call for the use of processing measures for
 21 capturing implicit processes (Andringa & Curcic, 2015; Suzuki & DeKeyser,
 22 in press), the current study demonstrated the potential of online com-
 23 prehension tasks as measures of IK.
 24

26 *Received 31 March 2015*

27 *Accepted 10 November 2015*

28 *Final Version Received 3 September 2015*

31 NOTES

34 1. Although EK can be verbalized, obviously not everyone has the metalinguistic
 35 means to articulate the rules clearly and completely (DeKeyser, 2009). Therefore, lack of
 36 verbalization ability is not necessarily evidence that learners do not possess EK.

37 2. The idea presented previously is controversial, and several other SLA scholars
 38 (e.g., DeKeyser, 1997, 2009) believe that the ultimate goal of SLA—fluent and accurate use
 39 of a L2—can also be accomplished by using automatized EK.

40 3. In a strong, noninterface position, the possibility that EK can be transformed into
 41 IK and vice versa is completely ruled out (Bowles, 2011; Hulstijn, 2002). However, in a
 42 weaker version of this position, the possibility of the transformation of IK into EK is rec-
 43 ognized (Bialystok, 1994).

44 4. The participants' responses to the question about the age at which they started
 45 learning English at school were variable. The minimum reported age was 1, and the
 46 maximum was 18. However, the average reported age was 9.75 years old, and the median
 47 was 10.

48 5. The corpus used in the current study was the Corpus of Contemporary American
 49 English (COCA): <http://corpus.byu.edu/coca/>. All words were taken from the 1,000 most
 frequent word families of English.

6. In some dialects of English, this sentence may not be considered unambiguously
 ungrammatical. However, these sentences were piloted with NSs of American English,
 and all NS participants agreed that "a lot of sugars" is wrong.

AQ40

AQ15

AQ41

7. Following previous studies of this type (e.g., Gutiérrez, 2013), this time limit was set to 3–6 s for NSs in the pilot.

8. The value was computed by the average of the squared factor loadings (Brown, 2006).

9. The target word was presented visually on the computer screen, but the carrier sentence was presented auditorily.

REFERENCES

- Andringa, S., & Curcic, M. (2015). How explicit knowledge affects online L2 processing [Special issue]. *Studies in Second Language Acquisition*, 37(2), 237–268. doi:10.1017/S0272263115000017
- Bley-Vroman, R. (1991). The logical problem of foreign language learning. *Linguistic Analysis*, 20(1–2), 3–49.
- Bowles, M. A. (2011). Measuring implicit and explicit linguistic knowledge. *Studies in Second Language Acquisition*, 33(2), 247–271.
- Brown, T. A. (2006). *Confirmatory factor analysis for applied research*. New York: Guilford Press.
- Clahsen, H., & Felser, C. (2006). Grammatical processing in language learners. *Applied Psycholinguistics*, 27(1), 3–42.
- Cronbach, L. J. (1988). Five perspectives on validation argument. In H. Wainer & H. Braun (Eds.), *Test validity* (pp. 3–17). Hillsdale, NJ: Erlbaum. **AQ16**
- Cronbach, L. J., & Meehl, P. E. (1955). Construct validity in psychological tests. *Psychological Bulletin*, 52(4), 281–302.
- DeKeyser, R. M. (2003). Implicit and explicit learning. In C. J. Doughty & M. H. Long (Eds.), *The Handbook of Second Language Acquisition* (pp. 312–348). Oxford, UK: Blackwell.
- DeKeyser, R. M. (2007). Skill acquisition theory. In B. V. J. Williams (Ed.), *Theories in second language acquisition* (pp. 97–114). Mahwah, NJ: Erlbaum. **AQ18**
- DeKeyser, R. M. (2009). Cognitive-psychological processes in second language learning. In M. H. Long & C. J. Doughty (Eds.), *The handbook of language teaching* (pp. 119–138). Oxford, UK: Wiley-Blackwell.
- DeKeyser, R. (2012). Interactions between individual differences, treatments, and structures in SLA. *Language Learning*, 62(s2), 189–200. **AQ19**
- DeKeyser, R., & Criado, R. (2012). Automatization, skill acquisition, and practice in second language acquisition. In C. Chapelle (Ed.), *The encyclopedia of applied linguistics* (pp. 4501–4504). Oxford, UK: Wiley-Blackwell. **AQ20**
- Ellis, N. C. (1993). Rules and instances in foreign language learning: Interactions of explicit and implicit knowledge. *European Journal of Cognitive Psychology*, 5(3), 289–318.
- Ellis, N. C. (Ed.). (1994). *Implicit and explicit learning of languages*. Academic Pr. **AQ21**
- Ellis, N. C. (2005). At the interface: Dynamic interactions of explicit and implicit language knowledge. *Studies in Second Language Acquisition*, 27(2), 305–352.
- Ellis, N. C. (2008). The dynamics of second language emergence: Cycles of language use, language change, and language acquisition. *The Modern Language Journal*, 92(2), 232–249.
- Ellis, N. C. (2011). Implicit and explicit SLA and their interface. *Implicit and Explicit Language Learning: Conditions, Processes, and Knowledge in SLA and Bilingualism*, 35–47.
- Ellis, R. (1991). Grammaticality judgments and learner variability. In H. Burmeister & P. Rounds (Eds.), *Variability in second language acquisition: Proceedings of the tenth meeting of the Second Language Acquisition Forum* (pp. 25–60). Eugene: Department of Linguistics, University of Oregon.
- Ellis, R. (1993). Second language acquisition and the structural syllabus. *TESOL Quarterly*, 27, 91–113.
- Ellis, R. (2004). The definition and measurement of L2 explicit knowledge. *Language Learning*, 54(2), 227–275.
- Ellis, R. (2005). Measuring implicit and explicit knowledge of a second language. *Studies in Second Language Acquisition*, 27(2), 141–172.
- Ellis, R. (2009). *Implicit and explicit knowledge in second language learning, testing and teaching* (Vol. 42). Bristol, UK: Multilingual Matters.

AQ17

AQ22

- 1 Ellis, R., & Loewen, S. (2007). Confirming the operational definitions of explicit and
2 implicit knowledge in Ellis (2005): Responding to Isemonger. *Studies in Second*
3 *Language Acquisition*, 29(1), 119–126. doi:10.1017/S0272263107070052
- 4 Forster, K. I., & Forster, J. C. (2003). DMDX: A Windows display program with millisecond
5 accuracy. *Behavior Research Methods, Instruments, & Computers*, 35(1), 116–124.
- 6 Granena, G. (2013). Individual differences in sequence learning ability and second
7 language acquisition in early childhood and adulthood. *Language Learning*, 63(4),
8 665–703.
- 9 Gutiérrez, X. (2013). The construct validity of grammaticality judgment tests as measures
10 of implicit and explicit knowledge. *Studies in Second Language Acquisition*, 35(3),
11 423–449. doi:10.1017/S0272263113000041
- 12 Haig, J. (1991). *Universal grammar and second language acquisition: The influence of task*
13 *type on late learner's access to the subadjacency principle* (TESL monograph). McGill
14 University, Montreal, Quebec, Canada.
- 15 Han, Y., & Ellis, R. (1998). Implicit knowledge, explicit knowledge and general language
16 proficiency. *Language Teaching Research*, 2(1), 1–23.
- 17 Henning, G. (1987). *A guide to language testing: Development, evaluation, research*. Cam-
18 bridge, MA: Newbury House.
- 19 Hulstijn, J. H. (2001). Intentional and incidental second language vocabulary learning:
20 A reappraisal of elaboration, rehearsal and automaticity. In P. Robinson (Ed.), *Cogni-*
21 *tion and second language instruction* (pp. 258–286). Cambridge, UK: Cambridge Univer-
22 sity Press.
- 23 Hulstijn, J. H. (2002). Towards a unified account of the representation, processing and
24 acquisition of second language knowledge. *Second Language Research*, 18(3), 193–223.
25 doi:10.1191/0267658302sr207oa
- 26 Hulstijn, J. H. (2005). Theoretical and empirical issues in the study of implicit and explicit
27 second-language learning. *Studies in Second Language Acquisition*, 27(2), 129–140.
- 28 Isemonger, I. M. (2007). Operational definitions of explicit and implicit knowledge:
29 Response to R. Ellis (2005) and some recommendations for future research in this
30 area. *Studies in Second Language Acquisition*, 29(1), 101–118.
- 31 Jiang, N. (2004). Morphological insensitivity in second language processing. *Applied Psy-*
32 *cholinguistics*, 25(4), 603–634.
- 33 Jiang, N. (2007). Selective integration of linguistic knowledge in adult second language
34 learning. *Language Learning*, 57(1), 1–33.
- 35 Jiang, N., Hu, G., Lukyanchenko, A., & Cao, Y. (2010). *Insensitivity to morphological errors in*
36 *L2: Evidence from word monitoring*. Paper presented at the SLRF 2010, October 14–17,
37 2010, College Park, MD.
- 38 Jiang, N., Novokshanova, E., Masuda, K., & Wang, X. (2011). Morphological congruency
39 and the acquisition of L2 morphemes. *Language Learning*, 61(3), 940–967.
- 40 Jöreskog, K. G. (1999). How large can a standardized coefficient be? Unpublished report.
41 *SSI Central, Inc*. Retrieved from <http://www.ssicentral.com/lisrel/techdocs/HowLargeCanaStandardizedCoefficientbe.pdf>
- 42 Jöreskog, K. G., & Sorbom, D. (2012). LISREL 9.1 [Computer software]. Lincolnwood, IL:
43 Scientific Software International.
- 44 Kaan, E. (2014). Predictive sentence processing in L2 and L1: What is different? *Linguistic*
45 *Approaches to Bilingualism*, 4(2), 257–282.
- 46 Kachinske, E., & Vafaei, P. (2014). *Reexamining the validity of GJTs as measures of implicit*
47 *knowledge: Reanalysis of Ellis & Loewen (2007), Bowles (2011), and Gutierrez (2013)*.
48 Paper presented at the Second Language Research Forum (SLRF), University of South
49 Carolina, SC.
- 50 Kane, M. T. (2001). Current concerns in validity theory. *Journal of Educational Measure-*
51 *ment*, 38(4), 319–342.
- 52 Kaufman, S. B., DeYoung, C. G., Gray, J. R., Jiménez, L., Brown, J., & Mackintosh, N. (2010).
53 Implicit learning as an ability. *Cognition*, 116(3), 321–340.
- 54 Kline, P. (1994). *An easy guide to factor analysis*. London, UK: Rutledge.
- 55 Kline, R. B. (2005). *Principles and practice of structural equation modeling* (2nd ed.).
56 New York: Guilford Press.
- 57 Leung, J. H. C., & Williams, J. N. (2011). The implicit learning of mappings between forms
58 and contextually derived meanings. *Studies in Second Language Acquisition*, 33(1),
59 33–55.

AQ24

AQ27
AQ28

AQ23

AQ25

AQ26

AQ29

- 1 Leung, J. H. C., & Williams, J. N. (2014). Crosslinguistic differences in implicit language
2 learning. *Studies in Second Language Acquisition*, 36(4), 733–755.
- 3 Loewen, S. (2009). Grammaticality judgment tests and the measurement of implicit and
4 explicit L2 knowledge. In R. Ellis, S. Loewen, C. Elder, R. Erlam, J. Philip, & H. Reinders
5 (Eds.), *Implicit and explicit knowledge in second language learning, testing and teaching*
6 (pp. 94–112). Tonawanda, NY: Multilingual Matters. **AQ30**
- 7 MacCallum, R. C., Widaman, K. F., Preacher, K. J., & Hong, S. (2001). Sample size in factor
8 analysis: The role of model error. *Multivariate Behavioral Research*, 36(4), 611–637.
- 9 MacCallum, R. C., Widaman, K. F., Zhang, S., & Hong, S. (1999). Sample size in factor
10 analysis. *Psychological Methods*, 4(1), 84–99.
- 11 Meara, P. M. (2005). *LLAMA language aptitude tests*. Swansea, UK: Lognostics. **AQ31**
- 12 Mueller, R. O., & Hancock, G. R. (2008). Best practices in structural equation modeling. In
13 J. Osborne. (Ed.), *Best practices in quantitative methods* (pp. 488–508). Thousand Oaks,
14 CA: Sage.
- 15 Murphy, V. A. (1997). The effect of modality on a grammaticality judgment task. *Second*
16 *Language Research*, 13(1), 34–65. **AQ32**
- 17 Nunnally, J. (1978). *Psychometric theory* (2nd ed.). New York: McGraw-Hill.
- 18 Paciorek, A., & Williams, J. N. (2015). Implicit learning of semantic preferences of verbs
19 [Special issue]. *Studies in Second Language Acquisition*, 37(2), 359–382. doi:10.1017/
20 S0272263115000108
- 21 Paradis, M. (2009). *Declarative and procedural determinants of second languages*.
22 Philadelphia, PA: John Benjamins.
- 23 Pett, M. A., Lackey, N. R., & Sullivan, J. J. (2003). *Making sense of factor analysis: The use of*
24 *factor analysis for instrument development in health care research*. Thousand Oaks, CA:
25 Sage. **AQ33**
- 26 Reed, J., & Johnson, P. (1994). Assessing implicit learning with indirect tests: Determining
27 what is learned about sequence structure. *Journal of Experimental Psychology:*
28 *Learning, Memory, and Cognition*, 20(3), 585–594. **AQ34**
- 29 Roberts, L., & Liszka, S. A. (2013). Processing tense/aspect-agreement violations on-line in
30 the second language: A self-paced reading study with French and German L2 learners
31 of English. *Second Language Research*, 29(4), 413–439. doi:10.1177/0267658313503171
- 32 Segalowitz, N., & Hulstijn, J. (2005). Automaticity in bilingualism and second language
33 learning. In J. F. Kroll & A. M. B. de Groot (Eds.), *Handbook of bilingualism: Psycholin-*
34 *guistic approaches* (pp. 371–388). New York: Oxford University Press.
- 35 Suzuki, Y. (2015). *Using new measures of implicit L2 knowledge to study the interface of*
36 *explicit and implicit knowledge* (Unpublished doctoral dissertation). University of
37 Maryland, College Park.
- 38 Suzuki, Y., & DeKeyser, R. M. (in press). Comparing elicited imitation and word moni-
39 toring as measures of implicit knowledge. *Language Learning*, 65(4).
- 40 Tinsley, H. E., & Tinsley, D. J. (1987). Uses of factor analysis in counseling psychology
41 research. *Journal of Counseling Psychology*, 34(4), 414–424. **AQ35**
- 42 Ullman, M. T. (2001). A neurocognitive perspective on language: The declarative/procedural
43 model. *Nature Reviews Neuroscience*, 2(10), 717–726.
- 44 Velicer, W. F., & Fava, J. L. (1998). Affects of variable and subject sampling on factor pat-
45 tern recovery. *Psychological Methods*, 3(2), 231.
- 46 Williams, J. N. (2005). Learning without awareness. *Studies in Second Language Acquisition*,
47 27(2), 269–304. doi:10.1017/S0272263105050138
- 48 Williams, J. N. (2009). Implicit learning. In W. C. Ritchie & T. K. Bhatia (Eds.), *New handbook*
49 *of second language acquisition* (pp. 319–353). Bingley, UK: Emerald Group.
- Zhang, R. (2015). Measuring university-level L2 learners' implicit and explicit linguistic
knowledge. *Studies in Second Language Acquisition*, 37, 457–486.

APPENDIX A

MKT RUBRIC

1. Third-person -s

- a. Full explanation should consist of = because the noun/subject/“the name of the word appearing in the sentence” is singular, you should use the verb + s/’s should be added to the verb/not plural.
- b. Participants will have to mention the singular noun or third-person singular and verb form.

2. Present perfect/simple past

- a. Full explanation should consist of = mention of simple past/an event occurring at a specific time in the past.

3. Countable/uncountable

- a. Full explanation should consist of = mention of specific terminology, such as countable/uncountable or can be plural/cannot be plural.

4. Hypothetical/second conditional

- a. Full explanation should consist of = the modal/the verb should be in the past tense to agree with the tense of the verb in the first sentence, or because the sentence describes a hypothetical situation/unreal past/hypothesis/assumption/supposition, or subjunctive mood.

APPENDIX B

CORRELATIONAL MATRIX

| Task | T-GJT | T-GJT-G | T-GJT-U | Un-GJT | Un-GJT-G | Un-GJT-U | MKT | SPRT | WMT |
|----------|-------|---------|---------|--------|----------|----------|-------|------|------|
| T-GJT | — | .65** | .76** | .22 | .07 | .22 | .06 | .02 | -.13 |
| T-GJT-G | | — | .01 | -.02 | .29* | -.21 | -.14 | .1 | -.09 |
| T-GJT-U | | | — | .33** | -.16 | .49** | .22* | -.07 | -.11 |
| UN-GJT | | | | — | .48** | .85** | .45** | .12 | .16 |
| UN-GJTG | | | | | — | -.01 | .17 | .02 | .01 |
| UN-GJT-U | | | | | | — | .47** | .13 | .14 |
| MKT | | | | | | | — | .08 | -.03 |
| SPRT | | | | | | | | — | .24* |
| WMT | | | | | | | | | — |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

APPENDIX C

COVARIANCE MATRIX

| | T-GJT | T-GJT-G | T-GJT-U | Un-GJT | Un-GJT-G | Un-GJT-U | MKT | SPRT | WMT |
|----------|----------|---------|----------|--------|----------|----------|--------|------------|------------|
| T-GJT | 42.448 | | | | | | | | |
| T-GJT-G | 3.216 | 0.580 | | | | | | | |
| T-GJT-U | 23.896 | 0.025 | 23.496 | | | | | | |
| Un-GJT | 1.343 | -0.019 | 1.494 | 0.887 | | | | | |
| Un-GJT-G | 0.338 | 0.156 | -0.550 | 0.321 | 0.511 | | | | |
| Un-GJT-U | 8.611 | -0.981 | 14.500 | 4.950 | -0.051 | 37.810 | | | |
| MKT | 0.142 | -0.038 | 0.385 | 0.152 | 0.044 | 1.028 | 0.128 | | |
| SPRT | 42.648 | 22.157 | -102.581 | 33.229 | 3.728 | 232.541 | 8.417 | 89,347.691 | |
| WMT | -157.921 | -13.305 | -103.949 | 28.314 | 1.201 | 163.838 | -2.103 | 13,540.185 | 35,451.312 |