

AGE MATTERS, AND SO MAY RATERS

Rater Differences in the Assessment of Foreign Accents

1
2
3
4
5
6
7
8
9
10
11
12
13
14 **Becky H. Huang**

15 *University of Texas at San Antonio*

16
17
18 **Sun-Ah Jun**

19 *University of California Los Angeles*

20
21
22
23
24 Research on the age of learning effect on second language learners' foreign accents utilizes human judgments to determine speech production outcomes. Inferences drawn from analyses of these ratings are then used to inform theories. The present study focuses on rater differences in the age of learning effect research. Three groups of raters who differed in their native language background and language experience participated in the study: inexperienced native English speaker (NES) raters, experienced NES raters, and advanced nonnative English speaker (NNES) raters. All raters evaluated 64 speech samples taken from both NESs and NNESs who varied in their age of arrival in the second-language-speaking country. Results from the study showed that experienced NES raters were better than the other two rater groups at distinguishing NESs from NNESs. Although no rater group differences were found in scaling speakers' foreign accents or in

25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41 This study was supported by a postdoctoral research grant from the National Science Council of Taiwan and a faculty seed grant from the University of Texas at San Antonio to the first author, and a UCLA Faculty Research Grant to the second author. We thank the four anonymous *SSLA* reviewers and Jason Bishop for their constructive and detailed comments on the drafts of the study. We are also grateful to all the participants in this study for their time and to Daniel Sass for his assistance with data analysis. All remaining errors are our own.

42
43
44
45
46
47 Correspondence concerning this article should be addressed to Becky Huang, Department of Bicultural-Bilingual Studies, The University of Texas at San Antonio, One UTSA Circle, San Antonio, TX 78249. E-mail: becky.huang@utsa.edu

1 their interrater reliability, inexperienced NES raters were stricter in their
2 ratings than both experienced NES raters and advanced NNES raters.
3 Implications for research on age of learning effects are discussed.
4

5
6
7
8 Research on the effect of age on learning a second language (L2) has
9 consistently shown that L2 learners who are immersed in a L2 at
10 a younger age produce more targetlike phonology—in particular, more
11 nativelike accents—than those immersed in the L2 later in their life
12 (Asher & García, 1969; Flege, Munro, & MacKay, 1995; Flege, Yeni-
13 Komshian, & Liu, 1999; Huang & Jun, 2011; Moyer, 1999, 2004; Munro &
14 Mann, 2005; Oyama, 1976). The age of learning variable remains the
15 most powerful predictor of a nativelike accent, even after controlling
16 for potential confounding variables, such as L2 learners' educational
17 attainment and their amount of exposure to the L2 (Flege et al., 1999).
18 Researchers have argued that L2 phonological production is strongly
19 constrained by the age of learning variable because of the neuromus-
20 cular demands involved in producing targetlike phonology (Scovel,
21 2000). To produce nativelike phonology, it is hypothesized that the cutoff
22 age is around 6 for segments and 12 for suprasegments (Long, 2005).
23

24 Empirical studies that examine the age effect on L2 phonological produc-
25 tion generally use human judgments to determine L2 learners' outcomes,
26 particularly in the domain of the degree to which learners have foreign
27 accents (FAs). On the basis of the inferences drawn from analyses of
28 raters' ratings, researchers then test or formulate theories of language
29 learning. However, although research in the fields of speech processing
30 and oral language assessment has revealed robust differences among
31 raters, rater differences have received little attention in research on the
32 age effect. In the following sections we first summarize the speech pro-
33 cessing and language assessment literature before reviewing in more
34 detail studies relating specifically to the focus of the present study—
35 namely, age of learning effects and the effects of rater differences. We then
36 describe the methodology of the current study and report the results.
37 We conclude with a discussion of the implications of rater differences for
38 the age of learning research and provide suggestions for future research.
39
40

41 42 **BACKGROUND**

43 44 **Rater Differences in Speech Processing and Oral Language 45 Assessment**

46
47
48 Research on speech processing and oral language assessment has gen-
49 erally shown that raters' or listeners' native language background and

1 language experience can influence their perception of FAs, their ratings
2 of the intelligibility and comprehensibility of nonnative speech, and
3 their judgments of nonnative speakers' (NNSs') language proficiency
4 (e.g., Bradlow & Bent, 2003; Brodkey, 1972; Gass & Varonis, 1984; Huang,
5 2014; Winke & Gass, 2013; Winke, Gass, & Myford, 2013; Xi & Mollaun,
6 2011). The influence of rater¹ differences has been an important topic
7 in the study of speech processing because this topic addresses funda-
8 mental questions about human perception and learning (see Carey,
9 Mannell, & Dunn, 2011, for a thorough review of dominant theories). It is
10 also one of the most heavily surveyed areas in oral language assessment
11 due to its direct implications for the reliability, validity, and fairness of
12 assessments (Bachman, Lynch, & Mason, 1995; Kunnan, 2000).

14 Turning first to the effects of raters' native language background
15 (native speaker [NS] vs. NNS) on ratings of nonnative speech, results
16 from empirical research are quite mixed. In terms of rank ordering
17 nonnative speech, Derwing and Munro (2013), for example, found that
18 the rank-ordering patterns of highly proficient NNS raters closely paral-
19 leled those of NS raters. Ratings by highly proficient NNS raters also
20 significantly correlated with those of NS raters. In terms of the numer-
21 ical ratings assigned to nonnative speech, some studies have found NS
22 raters to be more lenient than their NNS counterparts (e.g., Caban,
23 2003; Fayer & Krasinski, 1987), whereas others have observed no signif-
24 icant differences between the two groups (Kim, 2009; Zhang & Elder,
25 2011). However, analyses of interviews and raters' written comments
26 have revealed that NS and NNS raters may weigh various features of
27 speech samples differently in the justifications of their ratings (Kim,
28 2009; Zhang & Elder, 2011). To illustrate, Kim (2009) conducted a mixed-
29 methods study to assess how NS and NNS teachers evaluate students'
30 oral English performance. Both NS and NNS teachers judged 10 Korean
31 students' English speech samples on a 4-point scale and provided
32 written comments to justify their ratings. Although quantitative analyses
33 of teacher ratings revealed no significant group differences, NS teacher
34 raters commented most frequently on overall language use, pronuncia-
35 tion, and vocabulary, whereas NNS teachers emphasized pronuncia-
36 tion, vocabulary, and intelligibility. Native speaker raters also provided
37 more elaborate comments on speakers' pronunciation and grammar
38 than their NNS peers, suggesting potential qualitative differences in the
39 two groups' rating behavior.

42 An additional problem in investigating the influence of rater experi-
43 ence concerns how experience is defined and operationalized, which
44 varies substantially across studies (Isaacs & Thomson, 2013). Some
45 studies have defined it as the exposure to foreign-accented speech in
46 controlled experiments (Bradlow & Bent, 2003, 2008; Clarke & Garrett,
47 2004), some have appealed to raters' professions (e.g., second or foreign
48 language teachers or speech therapists vs. undergraduate students
49

1 without any linguistics training) as a proxy (Carey et al., 2011; Kennedy &
2 Trofimovich, 2008), some have relied on participants' self-reports of
3 familiarity with foreign accents (Munro, Derwing, & Morton, 2006; Winke
4 et al., 2013), and still others have used some combination of these
5 different criteria (Huang, 2014). There is also much variation in basic
6 research design across studies. Whereas some have utilized experimental
7 methodology to test the effects of perceptual training on the intelligibility
8 and comprehension of foreign-accented speech (e.g., Bradlow & Bent,
9 2003, 2008; Clarke & Garrett, 2004), others have used quasi-experimental
10 designs to examine the effects of language experience by studying dif-
11 ferent rater groups (e.g., Brodkey, 1972; Huang, 2014; Isaacs & Thomson,
12 2013; Kennedy & Trofimovich, 2008; Major, 2007; Munro et al., 2006;
13 Xi & Mollaun, 2011).

15 Despite differences in design and in the definition and measurement
16 of language experience, the overall picture emerging from this body of
17 work suggests that raters' language experience facilitates their compre-
18 hension of nonnative speech. Language experience has also contributed
19 to modifying raters' perceptions of FAs and judgments of L2 speakers'
20 proficiency, though the relationship between the degree of severity in
21 the rating and the raters' experience remains unclear. Some studies
22 have found that experienced raters tend to be more lenient than lingu-
23 stically naïve raters in their ratings (Barnwell, 1989; Carey et al., 2011;
24 Hsieh, 2011; Winke et al., 2013), but others have found the opposite to be
25 true (Galloway, 1980; Hadden, 1991). Furthermore, studies comparing
26 specific features used by raters (e.g., global proficiency, pronunciation,
27 or fluency) have also shown that language experience may play a signif-
28 icant role in their decision-making processes (Chalhoub-Deville, 1995;
29 Hsieh, 2011; Winke & Gass, 2013; Xi & Mollaun, 2011).

34 Rater Differences in the Age of Learning Effect Research

AQ1

37 We found a total of nine studies that examined rater differences in the
38 assessment of the age effect,² and all of them included rater differences
39 as a secondary research question rather than as the primary focus.³
40 There was also much variation across studies in the specific background
41 variables examined, which included the investigation of NS raters
42 with mainstream versus regional accents (Bongaerts, 1999b; Bongaerts,
43 Planken, & Schils, 1995), NS raters with and without experience in
44 teaching and linguistic training (Bongaerts, Mennen, & van der Slik,
45 2000; Bongaerts, van Summeren, Planken, & Schils, 1997; Nikolov, 2000,
46 Study 1; Thompson, 1991), NNSs with limited versus substantial expo-
47 sure to the target language (Flege, 1988; MacKay, Flege, & Imai, 2006),
48 and children versus adults (Nikolov, 2000).

1 Bongaerts and colleagues (Bongaerts, 1999b; Bongaerts et al., 2000;
2 Bongaerts et al., 1995; Bongaerts et al., 1997) conducted a series of
3 empirical studies to identify highly advanced late L2 learners with
4 nativelike proficiency. The main purpose of their work, a test of the critical
5 period hypothesis, was to determine whether late L2 learners can
6 “beat the odds” of a late start and achieve nativelike proficiency in their
7 L2. In an initial study, Bongaerts et al. (1995) found that NS raters from
8 northern England who were unfamiliar with regional British accents
9 judged the speech of NSs of British English with regional accents to be
10 less nativelike than that of Dutch learners of British English with training
11 in Received Pronunciation. These curious results led the authors to
12 conduct a subsequent study, published in the same paper, in which they
13 substituted NSs of regional British accents with NSs of neutral accents
14 and included NS raters both with and without experience in language
15 teaching and linguistics. Bongaerts and colleagues compared ratings by
16 two NS rater groups (experienced vs. inexperienced) and found that
17 they were not reliably different from each other. Using a widely adopted
18 *nativelikeness criterion* (Abrahamsson & Hyltenstam, 2009)—namely,
19 the criterion of two standard deviations within the mean of the NS control
20 group—the authors identified five late Dutch learners of British
21 English who were perceived as speaking with a nativelike accent. With
22 evidence of late learners having achieved nativelike proficiency, the
23 authors argued against the critical period hypothesis. In another study,
24 Bongaerts (1999b) examined a different first language (L1)–L2 pairing
25 (Dutch learners of French) and again included both experienced and
26 inexperienced NSs of French as raters. The researchers replicated their
27 previous finding, showing no reliable differences between the judgments
28 of the experienced and inexperienced raters. The results also
29 revealed three adult learners who demonstrated a nativelike accent in
30 their L2 French.
31
32
33

34 In contrast to these findings, Bongaerts and colleagues found discrepancies
35 between the judgments of raters with and without linguistic
36 experience in a fourth study (Bongaerts et al., 2000). In this study, two
37 groups of Dutch NSs, experienced and inexperienced, rated the speech
38 production of L2 learners of Dutch and of Dutch NS control participants.
39 Experienced raters were, again, trained L2 teachers, whereas inexperienced
40 raters lacked any experience in language teaching or linguistics.
41 Using the nativelikeness criterion derived from experienced raters’
42 average ratings, the researchers identified four late L2 learners with
43 nativelike accents, but only two of them were also perceived as native-
44 like based on inexperienced raters’ ratings. The researchers combined
45 the results and made inferences from the two, instead of four, nativelike
46 cases as evidence against the critical period hypothesis.
47

48 In an attempt to replicate the work by Bongaerts and colleagues,
49 Nikolov (2000) also found raters’ behavior to vary as a function of

1 linguistic experience. She included L2 learners of Hungarian (Study 1)
2 and English (Study 2) as well as NSs of each target language as control
3 participants. For both studies, Nikolov included three groups of NS raters
4 in the respective L2, one adolescent and two adult groups. In Study 1, the
5 adults were either university students or language teachers, whereas
6 both groups of adults were postgraduate students in Study 2. Results
7 from the two studies showed that adult raters rank ordered the speakers
8 in a similar manner. Adolescents' ratings, in contrast, were not as highly
9 correlated with the adults' ratings. Crucially, in Study 1, Nikolov found
10 that teacher raters, compared with university students, were better
11 able to distinguish NSs and NNSs.

12
13 In line with these studies showing positive rater differences, Thompson
14 (1991) investigated factors predicting Russian immigrants' degrees of
15 FA in English by assessing the judgments of two groups of English NSs
16 varying in their experience in language teaching and their contact with
17 Russian immigrants in the United States (experienced vs. inexperienced).
18 The results revealed that, compared to inexperienced NS raters,
19 experienced NS raters were more lenient and more reliable in their assess-
20 ment of FAs. It is worth noting that Thompson did not use the commonly
21 adopted nativelikeness criterion to identify nativelike L2 learners. Instead,
22 she examined mean ratings and found two nativelike speaker participants
23 with perfect scores.

24
25 In a study examining the effects of the amount of L1 use on Italian
26 immigrants' L2 pronunciation outcomes, Flege, Frieda, and Nozawa
27 (1997) also included two groups of English NSs as raters. One rater group
28 was from Alabama, United States, and self-reported having limited expo-
29 sure to FAs; the other rater group was from Ottawa, Canada, and reported
30 having regular exposure to a variety of FAs, including Italian-accented
31 English. To our knowledge, this is the only study in the age effect litera-
32 ture that includes detailed background information about the raters
33 employed. The speaker participants in the study included two groups of
34 Italian L1 immigrants who differed in the amount of their L1 use (high vs.
35 low), and one control group of Canadian English NSs. Flege and his col-
36 leagues found that the two groups of raters were not reliably different in
37 rank ordering the talkers (NS > low L1 use > high L1 use). However,
38 raters from Ottawa were better at detecting the FAs of Italian immigrants
39 to Canada and at distinguishing NSs and NNSs. Flege and his colleagues
40 attributed this advantage to the Canadian raters' experience and famil-
41 iarity with the specific type of FAs under investigation.

42
43 All seven of the studies just described used NS raters. The only two
44 studies that included NNS raters were those reported by Flege (1988)
45 and MacKay et al. (2006), who examined NNS raters' ability to gauge FAs
46 in their L2. Flege (1988) included four speaker groups of Chinese
47 learners of English who varied in their age of arrival (AoA) and length of
48 residence in the United States and a group of English NSs as controls.
49

1 Three groups of raters were recruited to judge the FAs of these speakers:
2 one English NS group and two NNS groups with 1 and 5 years of resi-
3 dence in the United States, respectively. Flege found that the two groups
4 of NNS raters rank ordered the speakers in a similar manner as the NS
5 raters. However, NNS raters with longer residence in the United States
6 were better at distinguishing between NSs and NNSs than were NNS
7 raters with shorter residence. The ratings of NNS raters with longer resi-
8 dence also did not reliably differ from those of the NS raters. The com-
9 bined results suggested that experience (from longer residence) may
10 afford the NNS raters perceptions of the L2 sounds that are more similar
11 to those of NS raters. A major limitation of the study was that most of the
12 raters were drawn from the speaker participants and thus rated their
13 own speech, which raised concerns about the validity of their ratings.

14
15 Similar to Flege's (1998) study, MacKay et al. (2006) evaluated NNSs'
16 ability to scale L2 speakers' FAs.⁴ However, unlike the NNS raters in
17 Flege's study, who shared the L1 (Chinese) of the NNS speakers, the
18 NNS raters in MacKay et al.'s study spoke a different L1 (Arabic) than
19 the NNS speakers (Italian). The experiment involved English NS control
20 speakers and four groups of Italian learners of English, differing in their
21 AoAs and L1 usage (early AoA–low L1 use; early AoA–high L1 use; late
22 AoA–low L1 use; and late AoA–high L1 use). Two groups of raters, one
23 English NS group and one NNS group with an Arabic L1 background,
24 evaluated all speakers' speech samples. The NNS raters themselves had
25 AoAs of 15 and older, and their average length of residence was 3 years.
26 Results from the study revealed that NNS raters' ratings were highly cor-
27 related with those of NS raters, suggesting that even NNS raters who do
28 not share the L1 of the NNS speakers can reliably scale the speakers' FAs.

29
30 In summary, existing studies paint a mixed picture about the role
31 of language experience in rating behavior;⁵ some have suggested
32 that NS raters' language experience may lead them to be more lenient
33 (Bongaerts et al., 2000; Nikolov, 2000; Thompson, 1991) or may help
34 them distinguish NSs from NNSs (Flege et al., 1997; Nikolov, 2000). Other
35 studies, however, have found no such effect (Bongaerts 1999b; Bongaerts
36 et al., 1995; Bongaerts et al., 1997). Two studies that compare NNS
37 raters with NS raters have found that NNS raters can scale the FAs of
38 L2 speakers as well as their NS peers (Flege, 1988; MacKay et al., 2006),
39 but these findings await replication. It is also unknown how NNS raters
40 compare to NS raters in terms of the degree of severity in their judg-
41 ments and their ability to distinguish between NSs and NNSs.

42
43 By systematically investigating the effect of rater differences on
44 the assessment of FAs, the present study aims to clarify how raters'
45 language experiences affect their judgments of foreign accents. In par-
46 ticular, we focus on raters' L1 background and their experience with FAs
47 and language teaching, using a quasi-experimental design. This included
48 comparison of three groups of raters: advanced nonnative English
49

1 speakers (advanced NNEs), inexperienced native English speakers
2 (inexperienced NESs), and experienced native English speakers (experi-
3 enced NESs). Our research questions focused on whether raters' L1 and
4 linguistic experience influenced any of the following:

- 5
- 6
- 7 1. The reliability of their foreign accent ratings.
- 8 2. Their scaling of L2 speakers' degree of foreign accents.
- 9 3. The degree of severity in their ratings of L2 speakers' FAs.
- 10 4. Their ability to reliably distinguish NSs from NNSs.

11 12 **METHOD**

13 **Participants: Speakers**

14
15
16
17 Four groups of speakers participated in the study: three L2 learner
18 groups, varying in AoA (child arrivals, adolescent arrivals, and adult
19 arrivals), and one NES control group ($n = 14$). The L2 learner groups
20 immigrated to the United States either in childhood (ages 5–11; $n = 22$),
21 adolescence (ages 12–16; $n = 14$), or adulthood (ages 17–25; $n = 14$) and
22 had lived in the United States for at least 7 years. They all spoke Mandarin
23 Chinese as their L1 and were at least college educated or were under-
24 graduate students at the time of testing. To control for amount of
25 L2 input, we matched the three groups on their length of residence
26 (LoR) in the United States and their self-reported English input in the
27 past 5 years (estimated as a percentage) in three different domains
28 (oral, media, and literacy).
29

30 All of the NES control speaker participants were NSs of American
31 English and had no significant exposure to a L2 other than formal foreign
32 language instruction in high school. To avoid any dialectal differences,
33 we screened and selected participants without a noticeable regional
34 accent. See Table 1 for a summary of the demographic information for
35 all speaker participants.
36
37
38
39

40 **Participants: Raters**

41
42 Three groups of raters participated in the current study: advanced
43 NNEs, inexperienced NESs, and experienced NES raters. All raters were
44 at least college educated or were current undergraduate students at the
45 time of testing. They reported normal hearing and no history of hearing
46 or speech problems.
47

48 We carefully screened all potential rater participants via a survey
49 instrument before inviting them to participate in the study. We used

Table 1. Means and Standard Deviations (in Parentheses) of Demographic Information by Speaker Group

Demographic variable	NES controls (<i>n</i> = 14)	L2 child arrivals (<i>n</i> = 22)	L2 adolescent arrivals (<i>n</i> = 14)	L2 adult arrivals (<i>n</i> = 14)
Chronological age	27 (7)	24 (5)	29 (6)	36 (5)
Age of arrival	NA	8.55 (1.77)	13.93 (1.64)	21.79 (2.83)
Gender	7f, 7m	15f, 7m	4f, 10m	11f, 3m
Length of residence	NA	15.18 (4.65)	14.74 (6.04)	13.58 (4.87)
English input-oral	NA	71 (15)	75 (14)	68 (14)
English input-media	NA	81 (17)	83 (20)	74 (24)
English input-literacy	NA	90 (10)	87 (10)	81 (23)

Note. NES = native English speaker; f = female; m = male; NA = not applicable. Chronological age = current age in years; age of arrival = age at which the participant arrived in the United States; length of residence = length of stay in the United States in years; English input-oral = self-estimated English exposure in oral domain (i.e., listening or speaking) in the past 5 years in percent form; English input-media = self-estimated English exposure via media (i.e., TV, radio, movies) in the past 5 years in percent form; English input-literacy = self-estimated English exposure in literacy domain (i.e., reading or writing) in the past 5 years in percent form. One-way ANOVAs revealed no significant differences among the three groups in LoR, $F(2, 47) = 0.42, p = .657$, and English input in the oral domain, $F(2, 47) = 2.98, p = .060$, the media domain, $F(2, 47) = 1.09, p = .346$, or in literacy domain, $F(2, 47) = 1.88, p = .165$.

different screening criteria for each of the three groups. Advanced NNEs were required to (1) speak Mandarin Chinese as their L1, (b) have come to the United States after college and have lived in the United States for at least 7 years, (c) self-report high levels of oral English proficiency (at least a rating of 7 on a scale of 1–9), and (d) self-report high levels of familiarity with FAs (at least a rating of 3 on a scale of 1–4). Almost half of the advanced NNE raters had also taught English as a foreign language in their home countries. Both inexperienced NES and experienced NES rater groups were required to be NSs of American English and to speak without a noticeable regional accent and to have had no significant exposure to a L2 prior to high school. Inexperienced NES raters also met the following criteria: (a) They did not speak any L2 or foreign language fluently at the time of participation, (b) they had not taken any class related to speech or phonetics, (c) they lacked any experience in language teaching or the rating of nonnative speech, and (4) they self-reported little familiarity with FAs. To qualify for the study, experienced NES raters were required to (a) be current or past English as a second language (ESL) teachers or speech therapists,⁶ (b) have worked in that profession for at least one full year, and (c) self-report high levels of familiarity with FAs (at least a rating of 3 on a scale of 1–4). Three of the experienced NES raters also had knowledge of Mandarin Chinese via either formal language classes or self-study. See Table 2 for a summary of the demographic information of all rater participants.

Table 2. Means and Standard Deviations (in Parentheses) of Demographic Information by Rater Group^a

Demographic variable	Advanced NNES (<i>n</i> = 10)	Inexperienced NES (<i>n</i> = 10)	Experienced NES (<i>n</i> = 10)
Chronological age	32 (3)	24 (3)	31 (4)
Gender	5f, 5m	5f, 5m	5f, 5m
Length of residence	9 (3)	NA	NA
English oral proficiency	7.10 (0.74)	NA	NA
Familiarity with foreign accents	3.10 (0.32)	2 (0)	3.50 (0.53)
TESL/TEFL experience	0.50 (0.85)	0	80.80 (1.95)

Note. NNES = nonnative English speaker; NES = native English speaker. English oral proficiency = self-rated English oral proficiency on a 1-9 scale (1 = poor, 9 = nativelike); familiarity with foreign accents = self-rated familiarity with foreign accents on a 1-4 scale (1 = not familiar at all. I am not able to tell whether the speaker's first language is English or other languages, 2 = somewhat familiar. I can sometimes tell whether the speaker's first language is English or other languages, 3 = moderately familiar. I can often tell whether the speaker's first language is English or other languages, 4 = very familiar. I can always tell whether the speaker's first language is English or other languages); TESL/TEFL experience = years of experience in teaching learners of English as a second/foreign language. One-way ANOVA revealed a significant group effect, $F(2, 27) = 47.90$; $p < .001$, and Bonferroni post hoc tests suggested that both advanced NNES and experienced NES groups reported significantly higher familiarity with foreign accents than inexperienced NES group (both $p < .001$). Experienced NES raters reported slightly higher familiarity than advanced NNES raters, though the difference was only marginally significant ($p = .054$).

Procedure

Collection of Speech Samples. All speakers were recorded reading a paragraph in a quiet room at a university laboratory or their private residences (Weinberger, 2014; see the Appendix for the paragraph). They were given 1 min to review the paragraph before recording and were instructed to read at their own pace. Recordings were completed using a high-quality head-mounted microphone (Shure SM 10A) in the program Audacity (Audacity Team, 2000, Version 1.2.5) on a laptop (IBM Thinkpad x60s).

Collection of Ratings. We normalized the intensity of all speech files (75 dB) using PRAAT (Version 5.0, Boersma & Weenink, 2009) and assigned an identification number by the order of speakers' AoAs. We then generated two lists of randomized orders using an online number randomizer (Urbaniak & Plous, 2013). The two orders were then balanced within each rater group. The average length of the speech files was approximately 28 s (range = 22–40).

Each rater participant rated all 64 speech files in a single 30–40 min session. There were two data-collection sites: one was a university laboratory and the other was a quiet office in a research institution. Following previous research (Bongaerts, 1999b; Bongaerts et al., 1997),

1 rater participants were not informed of the purpose of the study but
2 were told that they would listen to and rate speech samples produced
3 by unspecified proportions of NSs and NNSs of American English. No
4 practice or training was provided. Rater participants listened to each
5 speech file via headphones (Sennheiser HD 203) on a laptop (IBM Think-
6 pad x60s). Speech files were played via the MATLAB program (Version
7 6.5.1). Rater participants were instructed to judge the degree of FAs
8 using a 1–9 scale (*1* = strong FA, *9* = native English speaker).

9
10 To standardize the rating procedure and condition across the two
11 sites, we used the same models of headphones and laptops and the
12 same presentation program. The two researchers working on the pro-
13 ject also adhered to the same script when giving instructions to the
14 rater participants. The distribution of raters from the three rater groups
15 was roughly equal between the two sites.

16 The current study presents several methodological improvements
17 from previous research. First, to understand the comprehensive pic-
18 ture of rater effects, the study included speakers with a range of AoAs;
19 second, it also controlled for speakers' input, which was operational-
20 ized as a combination of LoR and amount of L2 input. Previous studies
21 either focused only on adult learners (Bongaerts, 1999b; Bongaerts
22 et al., 2000; Bongaerts et al., 1995; Bongaerts et al., 1997; Nikolov, 2000)
23 or included child and adult learners, but not adolescent learners (Flege,
24 1988; MacKay et al., 2006). In an earlier study that covered a range of
25 AoAs (Thompson, 1991), the researcher did not control for learners'
26 AoAs (Thompson, 1991), the researcher did not control for learners'
27 input, making the interpretations of the results difficult. The current
28 design permits investigations of the potential interactions between
29 speakers' AoAs and differences in behavior among raters. Additionally,
30 we carefully screened the raters to ensure they met our selection crite-
31 ria, and we used a larger sample of raters than that of most previous
32 studies. We also chose a paragraph that contained the full inventory of
33 American English phones to elicit read-aloud production. Compared to
34 the sentence-level productions used in previous research, the produc-
35 tion of paragraph-length samples provided more speech materials for
36 raters to evaluate, which in turn afforded us better chances to identify
37 possible differences among raters' perceptions. Finally, the current study
38 adopted a broader definition of language experience and used a combi-
39 nation of both objective (raters' profession) and subjective (self-reports
40 of familiarity with FAs) measures of language experience.

41 42 43 44 **RESULTS**

45
46
47 To examine interrater reliability and determine whether it was statisti-
48 cally appropriate to combine rater participants' FA ratings, we calcu-
49 lated average absolute agreement using a two-way mixed intraclass

1 correlation (ICC) for each rater group. The ICC measures were .92 (95%
2 CI [.88, .95]) for the advanced NNES raters, .95 (95% CI [.92, .97]) for the
3 experienced NES raters, and .97 (95% CI [.95, .98]) for the inexperienced
4 NES raters. The ICCs were all very high across rater groups and were
5 comparable to the values reported in previous studies that used ICCs
6 as an interrater reliability index (e.g., Derwing & Munro, 2013; MacKay
7 et al., 2006; Munro et al., 2006; Trofimovich, Lightbown, Halter, & Song,
8 2009). We thus computed the average of these ratings within each group
9 to use as dependent variables in the statistical models. The confidence
10 intervals suggest a narrow interval around the ICC point estimates and
11 indicate that the only difference (based on the confidence intervals)
12 was between advanced NNES and inexperienced NES raters.⁷ We also
13 examined the correlation matrix to substantiate the ICCs, and we found
14 that the coefficients were consistent and large (smallest $r = .49$) within
15 each rating group. However, the average Spearman's rho correlation
16 coefficients were smaller for the advanced NNES raters ($\rho_{\text{mean}} = .61$, $s = .01$)
17 than for the experienced NES raters ($\rho_{\text{mean}} = .86$, $s = .01$) and the inexpe-
18 rienced NES raters ($\rho_{\text{mean}} = .87$, $s = .01$).
19

20 We then conducted the traditional 4 (between, or speaker group,
21 factor) \times 3 (within, or rater group, factor) MANOVA to assess the effect
22 of raters on the dependent variable, rater FA ratings. This multivariate
23 model rests on the following assumptions: homogeneity of variance,
24 equality of covariance matrices, multivariate normality, and independ-
25 ent errors. Although preliminary analyses revealed these assumptions
26 were met when excluding the NES speaker group from the analyses, the
27 inclusion of this group violated the equality of covariance matrices and
28 the homogeneity of variance assumptions, as the estimated standard
29 deviations were much smaller for this group. These findings were expected,
30 given that the NES speaker group should have much higher ratings and less
31 variability among raters than the other speaker groups. Consequently,
32 including this group in the model biases the test statistics (i.e., it underesti-
33 mated the standard errors and, therefore, the test statistics) and statistical
34 inferences by violating the statistical assumptions of the model (i.e., overall
35 MANOVA and pairwise comparisons). Although one conventional approach
36 is to transform the dependent variable, these assumptions remained
37 violated after conducting Box-Cox transformations. To resolve these
38 concerns, we elected to employ robust analyses that use nonpooled
39 standard errors and did not rely on traditional model assumptions.
40

41 Empirical evidence suggests that the approximate degrees of freedom (ADF)
42 tests are relatively robust for tests of mean equality when the aforemen-
43 tioned assumptions are violated (for a review, see Keselman, Algina, & Kowalchuk,
44 2001; Lix & Keselman, 1995). Moreover, the statistical power of these ADF
45 tests is generally comparable to traditional multivariate and univariate mixed
46 ANOVA approaches, even when model assumptions are satisfied. Especially important for this study is
47
48
49

1 the observation that the ADF tests are also more powerful than tradi-
 2 tional tests when both variances and group cell sizes are unequal. For
 3 these reasons, we conducted a 4 (speaker group) \times 3 (rater experience)
 4 mixed ADF ANOVA followed by robust Welch-James ADF pairwise mean
 5 comparisons. Recall that the Welch-James ADF pairwise comparison
 6 procedure does not pool the variances; thus the smaller variances for
 7 the NES speakers will not influence the standard errors when comparing
 8 means across the other groups.
 9

10 Effect sizes (i.e., Cohen's d) were calculated differently depending on
 11 whether the pairwise comparison was a between- or within-groups
 12 comparison, as within-groups comparisons should take into account
 13 the correlation between variables (Lipsey & Wilson, 2001). Therefore,
 14 the between-group effect sizes were calculated by taking the mean dif-
 15 ference ($M_1 - M_2$) and dividing by the pooled standard deviation (s_p).
 16 The within-group effect sizes were computed using the following equa-
 17 tion: $t_c[2(1 - r)/n]^{1/2}$, where t_c is the t statistic (or square root of the
 18 F statistic in Tables 3 and 4), r is the correlation coefficient, and n is
 19 the sample size (see Dunlap, Cortina, Vaslow, & Burke, 1996).
 20

21 Approximate degrees of freedom mixed ANOVA analyses revealed a
 22 statistically significant main effect of rater group, $F(2, 38.63) = 225.19$,
 23 $p < .001$, $\eta_p^2 = .87$, and speaker group, $F(3, 116.20) = 116.20$, $p < .001$,
 24 $\eta_p^2 = .69$, and a rater experience by speaker group interaction, $F(6, 32.45) =$
 25 45.09 , $p < .001$, $\eta_p^2 = .40$. Based on Cohen's (1988) tentative standards
 26 (small $\eta_p^2 = .01$; medium $\eta_p^2 = .06$, large $\eta_p^2 = .14$), the effect size was large
 27 for each of the main effects and interaction. To better understand these
 28 differences, in particular the significant interaction, Welch-James ADF
 29 pairwise comparisons using a Bonferroni adjustment ($\alpha = .05/21 = .002$)
 30 were conducted. Note that, to reduce the probability of a Type I error,
 31 not all mean comparisons were conducted. Instead, only means of in-
 32 terest were tested for statistical significance (see Tables 3 and 4).
 33

34 Pairwise analyses comparing differences among speaker groups at
 35 each rater group (see Figure 1 and Table 3) revealed that NES speakers
 36 received higher ratings than child arrivals, and child arrivals received
 37 higher ratings than adolescent arrivals, regardless of rater groups.
 38 However, no difference between adolescent arrivals and adult arrivals
 39 emerged, regardless of rater background. When comparing rater groups
 40 at each speaker group level (see Figure 1 and Table 4), analyses indi-
 41 cated that advanced NNES and experienced NES raters overall assigned
 42 higher scores than inexperienced NES raters, with a much smaller dif-
 43 ference between advanced NNES and experienced NES raters. It is worth
 44 noting that experienced NES raters assigned significantly higher ratings
 45 to NESs than both inexperienced NES and advanced NNES raters, but the
 46 ratings for NESs did not significantly differ between inexperienced NES
 47 and advanced NNES raters. These results suggest that experienced NES
 48 raters were better at identifying NS norms than the other two rater groups.
 49

Table 3. Mean, Standard Deviation, Effect Size (d), the Robust Welch-James ADF Pairwise Comparison (F statistic), and the Corresponding p Value by Rater Group

Rater group	Comparison	M_1	SD_1	M_2	SD_2	d	F statistic	p value
Inexperienced NES	Native English speaker vs. child arrivals	8.56	0.28	6.01	1.62	1.99	52.33	< .001
Inexperienced NES	Child arrivals vs. adolescent arrivals	6.01	1.62	4.02	1.32	1.32	16.31	< .001
Inexperienced NES	Adolescent arrivals vs. adult arrivals	4.02	1.32	3.24	1.06	0.65	2.96	.098
Experienced NES	Native English speaker vs. child arrivals	8.96	0.06	7.35	1.26	1.64	36.37	< .001
Experienced NES	Child arrivals vs. adolescent arrivals	7.35	1.26	5.70	1.03	1.40	18.41	< .001
Experienced NES	Adolescent arrivals vs. adult arrivals	5.70	1.03	4.85	1.24	0.75	3.92	.059
Advanced NES	Native English speaker vs. child arrivals	8.51	0.44	6.89	1.23	1.63	32.24	< .001
Advanced NNEs	Child arrivals vs. adolescent arrivals	6.89	1.23	5.57	1.10	1.11	11.16	.002
Advanced NNEs	Adolescent arrivals vs. adult arrivals	5.57	1.10	5.14	1.11	0.39	1.08	.307

Note. M_1 and SD_1 as compared to M_2 and SD_2 are the means and standard deviations of ratings for the first and second speaker group in each comparison. Nonbold comparisons were not statistically significant after the Bonferroni adjustment ($\alpha = .05/21 = .002$) for Type I error.

Table 4. Mean and Standard Deviation (in parentheses) of the Three Rater Groups' Ratings for Each Speaker Group, along with the Robust Welch-James ADF Pairwise Comparison Results by Speaker Group

Speaker group	Advanced NNES	Inexperienced NES	Experienced NES	Comparison
Native English speaker	8.51 (0.44)	8.56 (0.28)	8.96 (0.06)	Inexperienced NES vs. experienced NES ($F = 37.28$; $p < .001$; $d = -2.03$) Experienced NES vs. advanced NNES ($F = 16.06$; $p < .001$; $d = 1.19$) Inexperienced NES vs. advanced NNES ($F = 0.39$; $p = .542$; $d = 0.09$)
Child arrivals	6.89 (1.23)	6.01 (1.62)	7.35 (1.26)	Inexperienced NES vs. experienced NES ($F = 125.95$; $p < .001$; $d = -0.72$) Experienced NES vs. advanced NNES ($F = 12.82$; $p < .001$; $d = 0.37$) Inexperienced NES vs. advanced NNES ($F = 33.61$; $p < .001$; $d = -0.52$)
Adolescent arrivals	5.57 (1.10)	4.02 (1.32)	5.70 (1.03)	Inexperienced NES vs. experienced NES ($F = 99.70$; $p < .001$; $d = -1.28$) Experienced NES vs. advanced NES ($F = 0.96$; $p = .345$; $d = 0.12$)
Adult arrivals	5.14 (1.11)	3.24 (1.06)	4.85 (1.24)	Inexperienced NES vs. advanced NNES ($F = 58.18$; $p < .001$; $d = -1.24$) Inexperienced NES vs. experienced NES ($F = 233.96$; $p < .001$; $d = -1.26$) Experienced NES vs. advanced NNES ($F = 1.88$; $p = .193$; $d = -0.24$)
				Inexperienced NES vs. advanced NNES ($F = 115.02$; $p < .001$; $d = -1.73$)

Note. Nonbold comparisons were not statistically significant after the Bonferroni adjustment ($\alpha = .05/21 = .002$) for Type I error.

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

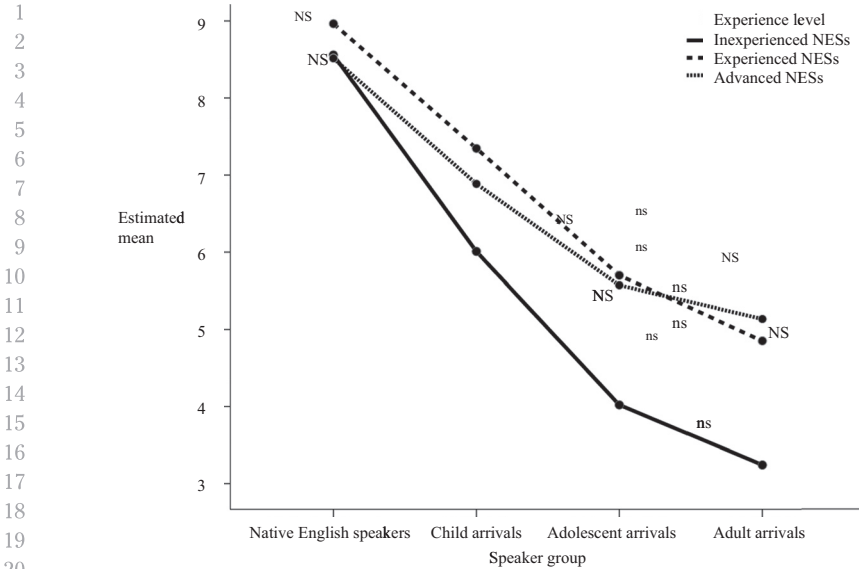


Figure 1. Estimated means from the ADF mixed ANOVA by rater group by speaker group. NS denotes nonsignificant within-group (rater group) mean differences, whereas ns represents nonsignificant between-groups (speaker group) mean differences, for those means nearest to the corresponding labels. All results are based on the Welch-James ADF pairwise comparison tests.

As seen in Figure 1 and Table 4, the following three nonsignificant rater group differences emerged: inexperienced NES versus advanced NNES raters at the NES speaker group and experienced NES versus advanced NNES raters at the adolescent and adult arrivals groups. For rating L2 speakers, whereas the experienced NES and advanced NNES raters had relatively consistent downward rating trends with the increase of L2 speakers' AoAs, these downward changes were much steeper for inexperienced NES raters. This implies that inexperienced NES raters tend to provide lower ratings, with these ratings decreasing at a greater rate than for the other two rater groups. Although slight differences existed between experienced NES and advanced NNES raters at times (see NES and child arrival comparisons), they were relatively small, especially for adolescent and adult arrivals.

To determine whether the three groups of raters identified the same L2 speakers as having achieved nativelike proficiency, we adopted the common nativelikeness criterion in the literature—that is, within two standard deviations of the average ratings assigned to NES controls—and we called the value derived from the criterion a *criterion score*. Judging from the ratings assigned to NES control speakers, experienced

1 NES raters appeared to be better able to distinguish between NSs and
 2 NNSs and assigned the highest ratings to NES control speakers of the
 3 three rater groups. Experienced NES raters also exhibited the least varia-
 4 tion ($M = 8.96, SD = 0.06$); thus their criterion score was the highest
 5 among the three groups (8.84). Advanced NNES raters' mean ratings of
 6 NES control speakers ($M = 8.51$) were similar to those of inexperienced
 7 NES raters ($M = 8.56$) despite their nonnative background. However,
 8 there was less variation among inexperienced NES raters ($SD = 0.28$)
 9 than among advanced NNES raters ($SD = 0.44$). Advanced NNES raters'
 10 criterion score was thus the lowest (7.63).
 11

12 We identified the most nativelike cases using advanced NNES raters'
 13 criterion score and the fewest cases by experienced NES raters' crite-
 14 rion score (see Table 5). Critically, no L2 speakers with an AoA of more
 15 than 5 were identified as nativelike by the experienced NES raters' crite-
 16 rion score, suggesting a potential terminus of plasticity at the age of 5.
 17 However, by the criterion score of the inexperienced NES rater group,
 18 two additional L2 speakers beyond the AoA of 5—one with an AoA of 6
 19
 20

21 **Table 5.** Nativelike Criteria and the Number of Nativelike L2 Speakers
 22 by Speaker Group by Rater Group
 23

24 Nativelike criteria and 25 number of nativelike 26 L2 speakers	Advanced NNES Raters	Inexperienced NES Raters	Experienced NES Raters
27 Nativelike 28 criterion score	7.63	8	8.84
29 Number of child arrivals 30 identified as nativelike 31 with AoAs ≤ 5 32 (initials and AoAs)	$n = 1$ YL, AoA = 5**	$n = 1$ YL, AoA = 5**	$n = 1$ YL, AoA = 5**
33 Number of child arrivals 34 identified as nativelike 35 with AoAs > 5 36 (initials and AoAs)	$n = 7$ LI, AoA = 6* SC, AoA = 8* MY, AoA = 7 AF, AoA = 7 KS, AoA = 8 AW, AoA = 9 SL, AoA = 10	$n = 2$ LI, AoA = 6* SC, AoA = 8*	$n = 0$
37 Number of adolescent 38 arrivals speakers 39 identified as nativelike	$n = 0$	$n = 0$	$n = 0$
40 Number of adult arrivals 41 speakers identified 42 as nativelike	$n = 0$	$n = 0$	$n = 0$

43
 44
 45
 46
 47
 48 Note. ** = nativelike cases identified by all three rater groups; * = nativelike cases identified by two of
 49 the three rater group.

1 and the other of 8—met the nativelike criterion. According to the ad-
2 vanced NNES raters' criterion score, five L2 speakers with an AoA of
3 more than 5 (including one child arrival with an AoA of 10) qualified as
4 having obtained a nativelike accent.
5

7 DISCUSSION

10 The current study set out to address a methodological gap in the litera-
11 ture on age of learning effects on L2 production outcomes, specifically
12 focusing on rater differences in the assessment of FAs. The study pro-
13 vided a systematic and detailed comparison of the judgments of three
14 groups of raters with different L1 backgrounds and language experi-
15 ence. Using a 1–9 numerical scale, all raters judged the degrees of FAs
16 in the paragraph read by 14 NES control participants and 50 L2 speakers
17 with varying AoAs. To examine differences across rater groups in inter-
18 rater reliability, two reliability measures were calculated for each rater
19 group and were compared across groups. Approximate degrees of free-
20 dom mixed ANOVA analyses were conducted to examine the differences
21 in raters' scaling of the speakers' FAs, the severity of ratings, and their
22 ability to distinguish NESs and NNEs.
23

24 Results from the interrater reliability analysis showed that advanced
25 NNES raters were the least reliable among the three rater groups, but
26 the two NES rater groups exhibited similar degrees of rating consis-
27 tency. These findings contradicted the results of Thompson's (1991)
28 study, in which experienced NES raters were much more reliable than
29 their inexperienced NES peers. Although a direct comparison of the cur-
30 rent study with that of Thompson is not possible, we speculate that the
31 discrepancies between the two may be attributed to the differences in
32 the AoA ranges of the speakers, their level of proficiency in the L1, or the
33 rating scales. The AoAs of the L2 speakers in Thompson's study ranged
34 from 4 to 42, but the distribution of the AoA range was not specified. It
35 could be that the majority of the speakers in Thompson's study were in
36 a similar AoA range, which made it easy for raters to reach consensus on
37 the speakers' FAs. In contrast, the current study included approximately
38 equal numbers of participants in each AoA bracket. Furthermore,
39 Thompson included only Russian learners of English who still main-
40 tained a high level of proficiency in their L1 and suggested that this may
41 have resulted in heavier FAs. In contrast, the speakers in the present
42 study self-reported substantial exposure to their L2 and may thus have
43 lighter, and therefore harder-to-discern, FAs. Finally, the wider range of
44 the scale (1–9) used in the current study may also lower the interrater
45 reliability more so than the 1–5 scale in Thompson's study.
46

48 The results revealed no significant differences in raters' rank orders
49 of the speakers' FAs; the three rater groups rank ordered the speakers

1 in a similar manner (NES > L2 child; L2 child > L2 adolescent; L2 adoles-
2 cent = L2 adult). This finding replicated previous research (Derwing &
3 Munro, 2013; Flege, 1988; Flege et al., 1997; MacKay et al., 2006) and
4 suggests that raters can reliably evaluate the relative degrees of FA
5 among speakers regardless of their own L1 backgrounds and language
6 experience, even when the speech material is a paragraph rather than a
7 single sentence.
8

9 Despite the similarities in their rank orders, the raters diverged in the
10 severity of their ratings of L2 speakers' FAs. The analyses of their aggre-
11 gated ratings revealed that, in general, inexperienced NES raters were
12 consistently stricter (i.e., assigned lower ratings) than experienced NES
13 raters and advanced NNES raters when judging L2 speakers' FAs. The
14 discrepancies between experienced and inexperienced NES raters'
15 severity corroborated previous research by Thompson (1991) but con-
16 tradicted Bongaerts (1999b) and Bongaerts et al. (1997). The results
17 were also consistent with previous research on oral language assess-
18 ments showing that experienced NS raters were more lenient than inex-
19 periented raters (Fayer & Krasinski, 1987; Hsieh, 2011) but contradicted
20 the findings of Hadden (1991), which found the opposite effect.
21

22 The discrepancies between the three groups' severity in ratings
23 may be attributed to a cognitive or exposure-based account, a social
24 account, or a combination of both. Usage-based models emphasize the
25 importance of experience and exposure in language acquisition (Kuhl,
26 2004; Tomasello, 1992); learners start with extracting phonemic proto-
27 types or lexically specific schemas from the linguistic environment.
28 Repeated exposure to the same prototypes or schemas then facilitates
29 the development of phonetic categories and eventually mental grammar.
30 Speech processing research (Bradlow & Bent, 2003, 2008) also provides
31 some support for this account by demonstrating that, with exposure to
32 foreign-accented speech in a controlled experimental session, listeners
33 can adapt to nonnative speech and can improve intelligibility and com-
34 prehension of the accented speech. Due to prior language experience or
35 exposure to FAs, experienced NES and advanced NNES raters in the
36 current study may thus have developed phonetic categories for the
37 accented speech, making it easier for them to process it. It is therefore
38 possible that the higher ratings assigned by some raters reflected this
39 ease of processing.
40

41 Furthermore, experienced NES and advanced NNES raters may also
42 have been more lenient than inexperienced NES raters because they were
43 familiar with a fuller range of FAs, having encountered L2 speakers with
44 stronger FAs and more limited proficiency than those of the advanced L2
45 speakers in the current study (who were all long-term U.S. residents).
46 Experienced NES and advanced NNES raters may have thus used the
47 relatively lower proficiency L2 speakers from their prior experience
48 as a baseline for comparison, a rating strategy documented in prior
49

1 studies (Isaacs & Thomson, 2013). It is worth noting that two recent
2 studies have found that NSs' level of education and professional training
3 influence their linguistic skills (specifically lexical knowledge and speaking
4 proficiency; Mulder & Hulstijn, 2011) as well as their judgments about
5 the structure of long-distance dependencies (Dąbrowska, 2010). These
6 findings seem to corroborate the cognitive or exposure-based account
7 beyond the domain of speech perception.
8

9 Alternatively, even if the required cognitive effort were equal, experi-
10 enced NES and advanced NNES raters may have been more sympathetic
11 with and favorable toward NNEs due to their familiarity with the
12 NNEs' FAs (the advanced NNES raters were themselves NNEs). In
13 other words, teaching experience and accent familiarity may have
14 helped these raters recognize the NNEs' country of origin and ethnicity,
15 or it may have even directly activated raters' favorable attitudes toward
16 NNEs, leading to more lenient ratings. This can be illustrated in a quali-
17 tative study that investigated the effect of raters' familiarity with speakers'
18 FAs on their rating process. Winke et al. (2013) asked heritage language
19 speakers of Mandarin Chinese, Korean, and Spanish to rate ESL speech
20 samples from speakers who had one of these three languages as their
21 L1. The rating session was video recorded, and raters later watched
22 these videos and discussed their decision-making processes. The
23 heritage language speaker raters often claimed to recognize the L1s of
24 some ESL speakers and reported that their ratings of them were more
25 lenient as a result. To tease these two accounts apart, and to under-
26 stand raters' cognitive processing and reactions to nonnative speech,
27 future studies may utilize the sort of think-aloud protocols that are com-
28 monly employed to examine decision-making processes (e.g., Ericsson &
29 Simon, 1980; Pressley & Afflerbach, 1995).
30

31 Although there were no significant differences between the aggre-
32 gated ratings of experienced NES and advanced NNES raters in most
33 pairwise comparisons in the current study, experienced NES raters were
34 significantly more lenient when rating child L2 speakers compared to
35 advanced NNES raters, but not reliably so when rating the adolescent
36 and adult L2 speakers. The results suggest that raters' behavior varied
37 by speakers' proficiency. We argue that our experienced NES raters'
38 leniency may be attributed to their language experience as well. As pre-
39 vious research comparing teacher and nonteacher raters has shown,
40 teachers are generally more lenient toward L2 learners than nonteach-
41 ers (Fayer & Krasinski, 1987; Hsieh, 2011). Although the experienced
42 NES raters in the current study did recognize the slight FAs in the child
43 L2 speakers, their ratings may reflect knowledge of these speakers'
44 likely improvement from some previous state.
45

46 Turning now to the results regarding the ability to distinguish
47 between NESs and NNEs, experienced NES raters were, as a group, better
48 at distinguishing native and nonnative speech, consistently assigning
49

1 the highest ratings to NESs. This particular finding is in line with what
2 Flege et al.'s (1997) study revealed about the perceptual differences
3 between NES raters in Alabama and those in Ottawa. It is likely that
4 experienced NES raters' ESL or English as a foreign language teaching
5 experience, combined with their familiarity with FAs and NES speech,
6 facilitated the identification of NS norms. In line with the cognitive or
7 exposure-based account, their experiences could have provided them
8 with useful knowledge about variation in both NES and NNES speech,
9 resulting in increased accuracy in distinguishing the two. In contrast,
10 inexperienced NES raters have little experience with the speech varia-
11 tion in L2 speech, possibly leading them to interpret L1 variants as FAs
12 (Markham, 1997). Therefore, advanced NNES raters and inexperienced
13 NES raters did not differ reliably in their ratings of the NES speakers. We
14 interpret these results as evidence supporting the hypothesis that
15 NNSs, particularly those with substantial experience and high profi-
16 ciency in the target language, can develop phonetic representations
17 that are identical to *some* of those possessed by NSs of the target
18 language. It is important to note that we are not arguing that NNSs can
19 achieve nativelike mental representations of *all* of the segmental and
20 suprasegmental features in the target language. In fact, the perceptions
21 of FAs may involve potential paralinguistic cues such as speakers' con-
22 fidence and both segmental and suprasegmental deviations (Piller,
23 2002). Research on crosslinguistic speech perception has shown that
24 even early L2 learners with near-native proficiency in the target language
25 may not achieve nativelike segmental perception (Navarra, Sebastián-
26 Gallés, & Soto-Faraco, 2005; Pallier, Bosch, & Sebastián-Gallés, 1997;
27 Sebastián-Gallés & Soto-Faraco, 1999). Nonetheless, corroborating findings
28 on the perceptions of regional dialects and FAs (Baker, Eddington, &
29 Nay, 2009; Flege et al., 1997), our results also demonstrate that NNSs
30 can develop global sound representations of the target language resem-
31 bling those of linguistically inexperienced NSs before developing the
32 ability to produce these sounds in a targetlike way.

36 The discrepancies in the raters' ability to distinguish NS and NNS
37 speakers resulted in different criteria for determining cases of native-
38 like L2 speakers. Because the experienced NES raters were better and
39 more consistent as a group in identifying NSs, their nativelike criterion
40 score was the highest among the three rater groups. Using their crite-
41 rion, no L2 speakers with an AoA of more than 5 were identified as
42 speaking with a nativelike accent. In contrast, several L2 speakers with
43 an AoA of more than 5 met the nativelike criteria by both inexperienced
44 NES raters and advanced NNES raters. One child L2 speaker with an
45 AoA of 10 also qualified as a nativelike learner based on the advanced
46 NNES raters' criterion score, though not by the other two NES rater
47 groups' standards. The discrepancies in the number of nativelike cases
48 can have important implications for the theories about age effects in L2
49

1 acquisition. Such discrepancies, particularly those between the two
2 NES rater groups, would lead to different conclusions about the terminus
3 of the critical period for L2 speech production. If we adopt experienced
4 NES raters' criterion, the cutoff age would be around 5, which is
5 close to Michael Long's (2005) proposal of age 6. However, using the
6 other two rater groups' criteria would result in different conclusions
7 because the upper limit of nativelike cases increased to 8 for inexperienced
8 NES raters and 10 for advanced NNS raters.

10 The implications of the current study are thus multifold. First, the
11 results suggest that, for future assessment tasks involving only rank
12 ordering of speakers, advanced NNSs can be just as competent as NSs.
13 Native speakers have long served as the norm for language instruction
14 as well as assessment, and their evaluations have been regarded as
15 more valid and informative than those of NNSs (Munro et al., 2006;
16 Zhang & Elder, 2011). However, the current results show that the evaluations
17 of NNSs with advanced proficiency are not reliably different from
18 those of NSs in the scaling of L2 speakers' speech production. The NNS
19 raters' degrees of severity in rating also closely paralleled those of experienced
20 NS raters, except for the child L2 speaker group, which was
21 judged less favorably by NNS raters than by the NS raters. In other
22 words, the evaluations of advanced NNS raters resembled those of
23 experienced NS raters in both rank ordering and overall severity. The
24 results spell good news for the large number of NNS teachers world-
25 wide and reassure that their evaluations and feedback to their L2
26 students are valid and meaningful.

28 Additionally, the discrepancies in the raters' severity, particularly
29 those between experienced and inexperienced NS raters, have implications
30 for the validity of research comparing L2 speakers' outcomes
31 across studies. Recently, meta-analysis research has been gaining popularity
32 in L2 acquisition research (Norris & Ortega, 2000), as it illuminates
33 and advances the field by systematically synthesizing the evidence
34 across empirical investigations. However, given the rater differences
35 evident in the current study, future metasynthesis of age effect research
36 would need to take into consideration the fact that such differences may
37 constitute a confound. Inferences drawn from cross-study comparisons
38 that fail to take rater backgrounds into account are therefore tenuous.

40 Finally, the existence of divergent rating outcomes between experienced
41 and inexperienced NS raters also raises concerns about the
42 *native speaker construct*. Native speakers have traditionally been defined
43 as speakers of the standard variety of the language (Davies, 1999), but
44 this definition lacks specificity. Researchers have challenged the absolute
45 distinction between NSs and NNSs, arguing against such a simple
46 dichotomy (Davies, 2003; Faez, 2011). As revealed in the current study,
47 NSs' language and teaching experience, or lack thereof, may also influence
48 their speech perception and judgments of NNSs' L2 proficiency.

1 The current conceptualization of NSs thus fails to appreciate the variation
2 in L1 proficiency among even NSs themselves (Cook, 1999; Davies, 2003).

3 In addition to serving as the benchmark for comparison in L2 acquisition
4 research, NS norms have also been used widely in language teaching
5 and testing. The lack of detail in the current definition, however, raises
6 questions about who the ideal NSs actually are. Would simultaneous
7 bilinguals (those who have been exposed to two languages since birth)
8 of the standard variety of both of their languages be regarded as NSs of
9 both languages? Would second-generation immigrants whose exposure
10 to the target language began at school age but who have become dominant
11 in the target language be considered NSs? If not, would the construct
12 of NSs and NNSs become one of monolinguals and bilinguals?
13 Although we agree on the necessity of concrete and well-defined benchmarks
14 in empirical research, we urge researchers to acknowledge the complexities
15 of the NS construct, to provide justifications for their decisions, and to
16 include detailed descriptions of any raters (NSs or NNSs) their study
17 employs. In a review of factors that affect L2 FAs, Piske, MacKay, and
18 Flege (2001) discussed the potential rater effect in the literature and
19 proposed that researchers include a representative sample of raters from
20 different backgrounds rather than one particular type. The proposal
21 appeared to be a feasible and balanced solution. However, we would like
22 to suggest that researchers also evaluate the purpose of their own study
23 before making a decision on the specific type of raters to include. For
24 example, if the purpose of the study is to evaluate the L2 speakers' FAs
25 as perceived by NSs living in metropolitan cities (who are presumably
26 familiar with FAs), then it would make sense to include judgments by
27 experienced NSs rather than those of inexperienced NSs. Conversely,
28 if the intent is to assess L2 speakers' FAs as perceived by NSs with
29 limited exposure to L2 speakers, using inexperienced NSs would provide
30 more informative assessments than would using experienced NSs.

31 We suggest several directions for future research. First, although the
32 size of the current rater sample ($n = 30$) is reasonable and relatively
33 large for studies on age effects and experimental or quasi-experimental
34 L2 acquisition research in general, it is not sufficient for more sophisticated
35 statistical techniques that also take into account variation on the
36 speaker level, such as multifaceted Rasch models (Hsieh, 2011; Isaacs &
37 Thomson, 2013; Winke et al., 2013) or G-theory (Xi & Mollaun, 2011).
38 Future researchers may want to include a larger sample and to utilize
39 such advanced analyses. Additionally, in the present study, we included
40 only one type of speech stimuli—namely, the production of a read-aloud
41 passage. Previous research has also shown that speakers' proficiency
42 may vary depending on the type of L2 tasks under study (Moyer, 2004;
43 Thompson, 1991); future work investigating the effect of task types
44 will advance our understanding of the matter. Finally, further
45
46
47
48
49

1 research is also needed to better understand the roles of other rater
2 background variables, such as gender and attitudes toward NNSs, and
3 to assess the interaction between such variables and task types.
4

6 CONCLUSIONS

7
8
9 To conclude, the current study found that L1 background and language
10 experience significantly influenced raters' rating scores and ability to
11 distinguish NSs from near-native early L2 learners. However, such experi-
12 ence was not found to predict raters' scaling of L2 speakers' proficiency.
13 Results from the current study therefore demonstrate that differences
14 among raters have important implications for studying age of learning
15 effects and should be controlled for in this line of research. The results
16 also indicate that rater differences can affect the number of natively-like L2
17 learners identified. Although no prior research has yet adapted NNSs'
18 ratings as criterion scores to determine the cutoff age for the critical
19 period, NNSs have been included as raters in several studies investigating
20 age effects on L2 acquisition (Flege, 1988; MacKay et al., 2006). Addition-
21 ally, as discussed in the literature review section, information about the
22 experience level of NS raters included in age effect research is generally
23 very limited. The current study revealed variation among NS raters' judg-
24 ments, which can potentially lead to different conclusions about the
25 terminus of the critical period. This particular finding thus highlights
26 critical implications for theories of L2 acquisition.
27

28 Rater differences also impact the severity of ratings, rendering compari-
29 sons or synthesis across studies problematic. In addition, the findings
30 raise questions about the general reliability and validity of instruments
31 used in measuring L2 outcomes. We hope that this study will draw
32 researchers' attention to the methodological issues in L2 acquisition
33 research and that more research will be conducted to help ensure the reli-
34 ability and validity of the inferences informing L2 acquisition theories.
35
36

37 *Received 23 July 2013*

38 *Accepted 10 April 2014*

39 *Final Version Received 20 June 2014*
40

42 NOTES

43
44 1. Previous research used either *listeners* or *raters* depending on the purpose of their
45 studies. We chose to use *raters* instead of *listeners* to denote active evaluation, rather than
46 passive listening, on the part of the participants.

47 2. We are aware of debates on the critical period hypothesis, which is not accepted
48 by all researchers as a useful account for the age of learning effect. For example, rather
49 than attributing the age effect to a critical period, Flege (1995) proposed the speech
learning model to explain age effects on the acquisition of L2 segments. He argued that L1

1 and L2 sounds share the same phonological space and mutually influence each other.
 2 Because adults have more robust L1 categories compared to children, it is more chal-
 3 lenging for adults to form new L2 categories. However, because the focus of this study is
 4 on rater effects (a measurement issue) in age of learning research rather than on the
 5 critical period hypothesis per se, we only discuss empirical studies that are relevant to
 6 the current study.

7 3. One reviewer brought to our attention an additional study by Bongaerts (Bon-
 8 gaerts, 1999a) that involved nonnative speaker judges. This particular study, which was
 9 written in Dutch, is not included in the current article because we limit our literature re-
 10 view to publications in English only.

11 4. We use *scale* and *rank* interchangeably throughout the article.

12 5. Note that here we limit our discussion to empirical studies that examined the age
 13 of learning effect. We understand that the distinction we made between rater studies in
 14 age effect research and in other research may be arbitrary. However, the distinction is
 15 essential to the aim of the current article, which is to raise awareness of an important
 16 methodological issue in age of learning effect research and in L2 acquisition research
 17 generally.

18 6. We decided to include speech therapists because speech therapists, along with
 19 phoneticians and ESL teachers, have been referred to as experienced or “expert” raters in
 20 L2 pronunciation research (Isaacs & Thomson, 2013). Speech therapists may develop
 21 familiarity with foreign accents through their work in accent modification or reduction
 22 (see <http://www.asha.org/public/speech/development/accent-modification/>). There was
 23 only one speech therapist in the experienced NES rater group. The other nine raters in the
 24 group were current or past ESL teachers. The speech therapist met the criteria for inclu-
 25 sion, as did the other nine raters in the group.

26 7. We conducted a *t* test comparing the aggregated ratings of the three experienced
 27 NES raters who speak Mandarin and the seven who do not, and the differences were not
 28 significant, $t(8) = -.433$, $p = .66$. On the basis of these results, we do not differentiate the
 29 two subgroups in our analyses.

30 REFERENCES

- 31 Abrahamsson, N., & Hylténstam, K. (2009). Age of onset and nativelikeness in a
 32 second language: Listener perception versus linguistic scrutiny. *Language Learning*,
 33 59, 249–306.
- 34 Asher, J. J., & García, R. (1969). The optimal age to learn a foreign language. *Modern*
 35 *Language Journal*, 53, 334–341.
- 36 Audacity Team. (2000). Audacity (Version 1.2.5) [Computer software]. Retrieved from
 37 <http://audacity.sourceforge.net/>
- 38 Bachman, L. F., Lynch, B. K., & Mason, M. (1995). Investigating variability in tasks and
 39 rater judgements in a performance test of foreign language speaking. *Language*
 40 *Testing*, 12, 238–257.
- 41 Baker, W., Eddington, D., & Nay, L. (2009). Dialect identification: The effects of region of
 42 origin and amount of experience. *American Speech*, 84, 48–71.
- 43 Barnwell, D. (1989). “Naïve” native speakers and judgments of oral proficiency in Spanish.
 44 *Language Testing*, 6, 152–163.
- 45 Boersma, P., & Weenink, D. (2009). Praat: Doing phonetics by computer (Version 5.1.05)
 46 [Computer program]. Retrieved from <http://www.fon.hum.uva.nl/praat/>
- 47 Bongaerts, T. (1999a). De keuze van beoordelaars in onderzoek naar uitspreekvaardigheid
 48 in een vreemde taal [Selecting raters for research on L2 pronunciation]. In M. Gerrit-
 49 sen & D. Springorum (Eds.), *Een bundel bedrijfscommunicatie voor Ger Peerbooms bij*
gelegenheid van zijn 65e verjaardag (pp. 1–10). Nijmegen, the Netherlands: Nijmegen
 University Press.
- Bongaerts, T. (1999b). Ultimate attainment in foreign language pronunciation: The case of
 very advanced late foreign language learners. In D. Birdsong (Ed.), *Second language*
acquisition and the Critical Period Hypothesis (pp. 133–159). Mahwah, NJ: Erlbaum.
- Bongaerts, T., Mennen, S., & van der Slik, F. (2000). Authenticity of pronunciation in natu-
 ralistic second language acquisition: The case of very advanced late learners of Dutch
 as a second language. *Studia linguistica*, 54, 298–308.

- 1 Bongaerts, T., Planken, B., & Schils, E. (1995). Can late starters attain a native accent in a
 2 foreign language? A test of the Critical Period Hypothesis. In D. Singleton & Z. Lengyel
 3 (Eds.), *The age factor in second language acquisition* (pp. 30–50). Clevedon, UK: Multi-
 4 lingual Matters.
- 5 Bongaerts, T., van Summeren, C., Planken, B., & Schils, E. (1997). Age and ultimate attain-
 6 ment in the pronunciation of a foreign language. *Studies in Second Language Acquisition*,
 7 *19*, 447–465.
- 8 Bradlow, A., & Bent, T. (2003). Listener adaptation to foreign-accented English. In
 9 M. J. Solé, D. Recasens, & J. Romero (Eds.), *Proceedings of the 15th International*
 10 *Congress of Phonetic Sciences* (pp. 1581–1583). Barcelona, Spain: Universitat Autònoma
 11 de Barcelona.
- 12 Bradlow, A. R., & Bent, T. (2008). Perceptual adaptation to non-native speech. *Cognition*,
 13 *106*, 707–729.
- 14 Brodkey, D. (1972). Dictation as a measure of mutual intelligibility. *Language Learning*, *22*,
 15 203–220.
- 16 Caban, H. L. (2003). Rater group bias in the speaking assessment of four L1 Japanese ESL
 17 students. *Second Language Studies*, *21*, 1–43.
- 18 Carey, M. D., Mannell, R. H., & Dunn, P. K. (2011). Does a rater's familiarity with a candi-
 19 date's pronunciation affect the rating in oral proficiency interviews? *Language*
 20 *Testing*, *28*, 201–219.
- 21 Chalhoub-Deville, M. (1995). A contextualized approach to describing oral language profi-
 22 ciency. *Language Learning*, *45*, 251–281.
- 23 Clarke, C. M., & Garrett, M. F. (2004). Rapid adaptation to foreign-accented English. *Journal*
 24 *of the Acoustical Society of America*, *116*, 3647–3658.
- 25 Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale,
 26 NJ: Erlbaum.
- 27 Cook, V. (1999). Going beyond the native speaker in language teaching. *TESOL Quarterly*,
 28 *33*, 185–209.
- 29 Dąbrowska, E. (2010). Naive vs. expert intuitions: An empirical study of acceptability judg-
 30 ments. *The Linguistic Review*, *27*, 1–23.
- 31 Davies, A. (1999). Standard English: Discordant voices. *World Englishes*, *18*, 171–186.
- 32 Davies, A. (2003). *The native speaker: Myth and reality*. Clevedon, UK: Multilingual Matters.
- 33 Derwing, T. M., & Munro, M. J. (2013). The development of L2 oral language skills in two
 34 L1 groups: A 7-year study. *Language Learning*, *63*, 163–185.
- 35 Ericsson, K. A., & Simon, H. A. (1980). Verbal reports as data. *Psychological Review*, *87*,
 36 215–251.
- 37 Faez, F. (2011). Reconceptualizing the native/nonnative speaker dichotomy. *Journal of*
 38 *Language, Identity & Education*, *10*, 231–249.
- 39 Fayer, J. M., & Krasinski, E. (1987). Native and nonnative judgments of intelligibility and
 40 irritation. *Language Learning*, *37*, 313–326.
- 41 Flege, J. E. (1988). Factors affecting degree of perceived foreign accent in English sen-
 42 tences. *Journal of the Acoustical Society of America*, *84*, 70–79.
- 43 Flege, J. E., Frieda, E. M., & Nozawa, T. (1997). Amount of native-language (L1) use affects
 44 the pronunciation of an L2. *Journal of Phonetics*, *25*, 169–186.
- 45 Flege, J. E., Munro, M. J., & MacKay, I. R. (1995). Factors affecting strength of perceived
 46 foreign accent in a second language. *Journal of the Acoustical Society of America*, *97*,
 47 3125–3135.
- 48 Flege, J. E., Yeni-Komshian, G. H., & Liu, S. (1999). Age constraints on second-language
 49 acquisition. *Journal of Memory and Language*, *41*, 78–104.
- Galloway, V. B. (1980). Perceptions of the communicative efforts of American students of
 Spanish. *Modern Language Journal*, *64*, 428–433.
- Gass, S., & Varonis, E. (1984). The effect of familiarity on the comprehensibility of nonna-
 tive speech. *Language Learning*, *34*, 65–89.
- Hadden, B. L. (1991). Teacher and nonteacher perceptions of second-language communi-
 cation. *Language Learning*, *41*, 1–24.
- Hsieh, C.-N. (2011). Rater effects in ITA testing: ESL teachers' versus American undergrad-
 uates' judgments of accentedness, comprehensibility, and oral proficiency. *Spain*
Fellow Working Papers in Second or Foreign Language Assessment, *9*, 47–74.
- Huang, B. H. (2014). The effects of age on second language grammar and speech produc-
 tion. *Journal of Psycholinguistic Research*, *43*, 397–420.

- 1 Huang, B. H., & Jun, S.-A. (2011). Specifying the age-related effect on the acquisition of
2 second language prosody. *Language and Speech, 54*, 387–414.
- 3 Isaacs, T., & Thomson, R. I. (2013). Rater experience, rating scale length, and judgments of
4 L2 pronunciation: Revisiting research conventions. *Language Assessment Quarterly,*
5 *10*, 135–159.
- 6 Kennedy, S., & Trofimovich, P. (2008). Intelligibility, comprehensibility, and accentedness
7 of L2 speech: The role of listener experience and semantic context. *Canadian Modern*
8 *Language Review, 64*, 459–489.
- 9 Keselman, H. J., Algina, J., & Kowalchuk, R. K. (2001). The analysis of repeated measures
10 designs: A review. *British Journal of Mathematical and Statistical Psychology, 54*, 1–20.
- 11 Kim, Y. H. (2009). An investigation into native and non-native teachers' judgments of oral
12 English performance: A mixed methods approach. *Language Testing, 26*, 187–217.
- 13 Kuhl, P. K. (2004). Early language acquisition: Cracking the speech code. *Nature Reviews*
14 *Neuroscience, 5*, 831–843.
- 15 Kunnan, A. J. (Ed.). (2000). *Fairness and validation in language assessment: Selected papers*
16 *from the 19th Language Testing Research Colloquium, Orlando, Florida* (Vol. 9). Cambridge,
17 UK: Cambridge University Press.
- 18 Lipsey, M. W., & Wilson, D. B. (2001). *Practical meta-analysis*. Thousand Oaks, CA: Sage.
- 19 Lix, L. M., & Keselman, H. J. (1995). Approximate degrees of freedom tests: A unified
20 perspective on testing for mean equality. *Psychological Bulletin, 117*, 547–560.
- 21 Long, M. H. (2005). Problems with supposed counter-evidence to the Critical Period Hypo-
22 thesis. *International Review of Applied Linguistics, 43*, 287–317.
- 23 MacKay, I. R., Flege, J. E., & Imai, S. (2006). Evaluating the effects of chronological age and
24 sentence duration on degree of perceived foreign accent. *Applied Psycholinguistics,*
25 *27*, 157–183.
- 26 Major, R. C. (2007). Identifying a foreign accent in an unfamiliar language. *Studies in*
27 *Second Language Acquisition, 29*, 539–556.
- 28 Markham, D. (1997). *Phonetic imitation, accent, and the learner*. Lund, Sweden: Lund
29 University Press.
- 30 MATLAB (Version 6.5.1) [Computer software]. Natick, MA: MathWorks.
- 31 Moyer, A. (1999). Ultimate attainment in L2 phonology. *Studies in Second Language Acqui-*
32 *sition, 21*, 81–108.
- 33 Moyer, A. (2004). *Age, accent and experience in second language acquisition*. Clevedon, UK:
34 Multilingual Matters.
- 35 Mulder, K., & Hulstijn, J. H. (2011). Linguistic skills of adult native speakers, as a function
36 of age and level of education. *Applied Linguistics, 32*, 475–494.
- 37 Munro, M. J., Derwing, T. M., & Morton, S. L. (2006). The mutual intelligibility of L2 speech.
38 *Studies in Second Language Acquisition, 28*, 111–131.
- 39 Munro, M., & Mann, V. (2005). Age of immersion as a predictor of foreign accent. *Applied*
40 *Psycholinguistics, 26*, 311–341.
- 41 Navarra, J., Sebastián-Gallés, N., & Soto-Faraco, S. (2005). The perception of second
42 language sounds in early bilinguals: New evidence from an implicit measure. *Journal*
43 *of Experimental Psychology: Human Perception and Performance, 31*, 912–918.
- 44 Nikolov, M. (2000). The Critical Period Hypothesis reconsidered: Successful adult learners
45 of Hungarian and English. *International Review of Applied Linguistics in Language*
46 *(IRAL) Teaching, 38*, 109–124.
- 47 Norris, J. M., & Ortega, L. (2000). Effectiveness of L2 instruction: A research synthesis and
48 quantitative meta-analysis. *Language Learning, 50*, 417–528.
- 49 Oyama, S. (1976). A sensitive period for the acquisition of a nonnative phonological
system. *Journal of Psycholinguistic Research, 5*, 261–283.
- Pallier, C., Bosch, L., & Sebastián-Gallés, N. (1997). A limit on behavioral plasticity in
speech perception. *Cognition, 64*, B9–B17.
- Piller, I. (2002). Passing for a native speaker: Identity and success in second language
learning. *Journal of Sociolinguistics, 6*, 179–208.
- Piske, T., MacKay, I. R., & Flege, J. E. (2001). Factors affecting degree of foreign accent in
an L2: A review. *Journal of Phonetics, 29*, 191–215.
- Pressley, M., & Afflerbach, P. P. (1995). *Verbal protocols of reading: The nature of construc-*
tively responsive reading. Hillsdale, NJ: Erlbaum.
- Scovel, T. (2000). A critical review of the critical period research. *Annual Review of*
Applied Linguistics, 20, 213–223.

- 1 Sebastián-Gallés, N., & Soto-Faraco, S. (1999). Online processing of native and non-native
2 phonemic contrasts in early bilinguals. *Cognition*, 72, 111–123.
- 3 Thompson, I. (1991). Foreign accents revisited: The English pronunciation of Russian
4 immigrants. *Language Learning*, 41, 177–204.
- 5 Tomasello, M. (1992). *First verbs: A case study of early grammatical development*. Cambridge,
6 UK: Cambridge University Press.
- 7 Trofimovich, P., Lightbown, P. M., Halter, R., & Song, H. (2009). Comprehension-based
8 practice: The development of L2 pronunciation in a listening and reading program.
9 *Studies in Second Language Acquisition*, 31, 609–639.
- 10 Urbaniak, G. C., & Plous, S. (2013). Research randomizer (Version 4.0) [Computer software].
11 Retrieved from <http://www.randomizer.org/>
- 12 Weinberger, S. (2014). *Speech accent archive* [Database]. Retrieved from <http://accent.gmu.edu>.
- 13 Winke, P., & Gass, S. (2013). The influence of second language experience and accent
14 familiarity on oral proficiency rating: A qualitative investigation. *TESOL Quarterly*, 47,
15 762–789.
- 16 Winke, P., Gass, S., & Myford, C. (2013). Raters' L2 background as a potential source of
17 bias in rating oral performance. *Language Testing*, 30, 231–252.
- 18 Xi, X., & Mollaun, P. (2011). Using raters from India to score a large test. *Language Learning*,
19 61, 1222–1255.
- 20 Zhang, Y., & Elder, C. (2011). Judgments of oral proficiency by non-native and native English
21 speaking teacher raters: Competing or complementary constructs? *Language Testing*,
22 28, 31–50.

APPENDIX

READING PARAGRAPH FOR SPEECH ELICITATION

27 Please call Stella. Ask her to bring these things with her from the store:
28 Six spoons of fresh snow peas, five thick slabs of blue cheese, and
29 maybe a snack for her brother Bob. We also need a small plastic snake
30 and a big toy frog for the kids. She can scoop these things into three red
31 bags, and we will go to meet her Wednesday at the train station.

32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Author Queries

QA	The distinction between surnames can be ambiguous, therefore to ensure accurate tagging for indexing purposes online (eg for PubMed entries), please check that the highlighted surnames have been correctly identified, that all names are in the correct order and spelt correctly.
AQ1	Note that note callouts should not be placed in headings. I have moved the note 2 callout to the sentence below.
AQ2	Dunlap et al. 1996 is not listed in the Refs. Add to Refs.
AQ3	Flege 1995 is not listed in the Refs. Change to "Flege et al. (1995)" or add Flege 1995 to Refs.
AQ4	I have switched notes 5 and 6 b/c the note 5 callout seemed to correspond to the text of (former) note 6, and the note 6 callout matched (former) note 5.