



## Original Articles

# A critical period for second language acquisition: Evidence from 2/3 million English speakers



Joshua K. Hartshorne<sup>a,b,\*</sup>, Joshua B. Tenenbaum<sup>a</sup>, Steven Pinker<sup>c</sup>

<sup>a</sup> Department of Brain & Cognitive Sciences, Massachusetts Institute of Technology, Building 46, 77 Massachusetts Avenue, MIT, Cambridge, MA 02139, United States

<sup>b</sup> Department of Psychology, Boston College, McGuinn Hall 527, Chestnut Hill, MA 02467, United States

<sup>c</sup> Department of Psychology, Harvard University, William James Hall 970, 33 Kirkland St., Cambridge, MA 02138, United States

## ARTICLE INFO

## Keywords:

Language acquisition  
Critical period  
L2 acquisition

## ABSTRACT

Children learn language more easily than adults, though when and why this ability declines have been obscure for both empirical reasons (underpowered studies) and conceptual reasons (measuring the ultimate attainment of learners who started at different ages cannot by itself reveal changes in underlying learning ability). We address both limitations with a dataset of unprecedented size (669,498 native and non-native English speakers) and a computational model that estimates the trajectory of underlying learning ability by disentangling current age, age at first exposure, and years of experience. This allows us to provide the first direct estimate of how grammar-learning ability changes with age, finding that it is preserved almost to the crux of adulthood (17.4 years old) and then declines steadily. This finding held not only for “difficult” syntactic phenomena but also for “easy” syntactic phenomena that are normally mastered early in acquisition. The results support the existence of a sharply-defined critical period for language acquisition, but the age of offset is much later than previously speculated. The size of the dataset also provides novel insight into several other outstanding questions in language acquisition.

## 1. Introduction

People who learned a second language in childhood are difficult to distinguish from native speakers, whereas those who began in adulthood are often saddled with an accent and conspicuous grammatical errors. This fact has influenced many areas of science, including theories about the plasticity of the young brain, the role of neural maturation in learning, and the modularity of linguistic abilities (Johnson & Newport, 1989; Lenneberg, 1967; Morgan-Short & Ullman, 2012; Newport, 1988; Pinker, 1994). It has also affected policy, driving debates about early childhood stimulation, bilingual education, and foreign language instruction (Bruer, 1999).

However, neither the nature nor the causes of this “critical period” for second language acquisition are well understood. (Here, we use the term “critical period” as a theory-neutral descriptor of diminished achievement by adult learners, whatever its cause.) There is little consensus as to whether children’s advantage comes from superior neural plasticity, an earlier start that gives them additional years of learning, limitations in cognitive processing that prevent them from being distracted by irrelevant information, a lack of interference from a well-learned first language, a greater willingness to experiment and

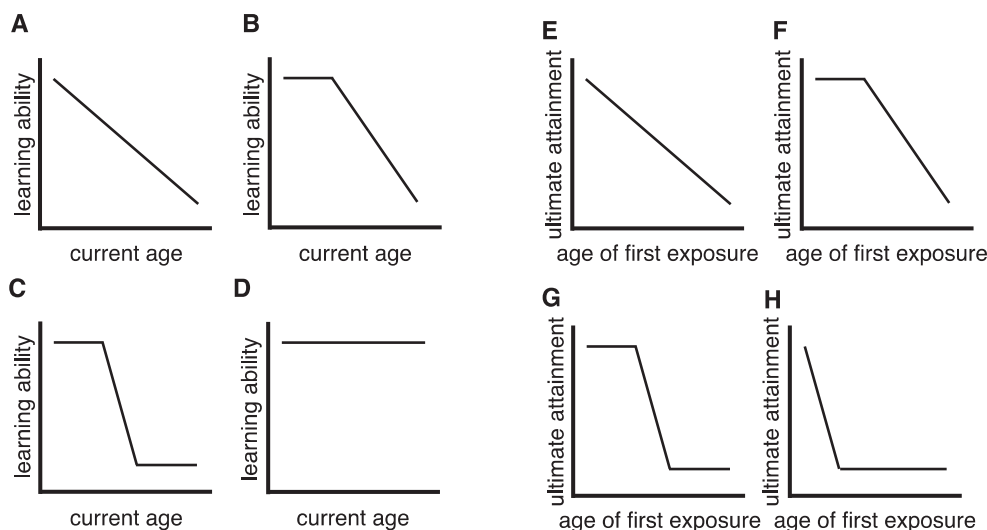
make errors, a greater desire to conform to their peers, or a greater likelihood of learning through immersion in a community of native speakers (Birdsong, 2017; Birdsong & Molis, 2001; Hakuta, Bialystok, & Wiley, 2003; Hernandez, Li, & MacWhinney, 2005; Johnson & Newport, 1989; Newport, 1990; Pinker, 1994). We do not even know how long the critical period lasts, whether learning ability declines gradually or precipitously once it is over, or whether the ability continues to decline throughout adulthood or instead reaches a floor (Birdsong & Molis, 2001; Guion, Flege, Liu, & Yeni-Komshian, 2000; Hakuta et al., 2003; Jia, Aaronson, & Wu, 2002; Johnson & Newport, 1989; McDonald, 2000; Sebastián-Gallés, Echeverría, & Bosch, 2005; Vanhove, 2013).

### 1.1. Learning ability vs. ultimate attainment

As noted by Patkowski (1980), researchers interested in critical periods focus on two interrelated yet distinct questions:

- (1) How does learning ability change with age?
- (2) How proficient can someone be if they began learning at a particular age?

\* Corresponding author at: Department of Psychology, Boston College, McGuinn Hall 527, Chestnut Hill, MA 02467, United States.  
E-mail address: [Joshua.hartshorne@bc.edu](mailto:Joshua.hartshorne@bc.edu) (J.K. Hartshorne).



**Fig. 1.** (A–D) Schematic depictions of four theories of how *language learning ability* might change with age. (E–H) Schematic depictions of four theories of how *ultimate attainment* might vary with age of first exposure to the language. Note: While the curves hypothesized for learning ability and ultimate attainment resemble one another, there is little systematic relationship between the two; see the main text.

The questions are different because language acquisition is not instantaneous. For example, an older learner who (hypothetically) acquired language at a slower rate could, in theory, still attain perfect proficiency if he or she persisted at the learning long enough.

The question of ultimate attainment (2) captures the most public attention because it directly applies to people's lives, but the question of learning ability (1) is more theoretically central. Does learning ability decline gradually from birth (Guion et al., 2000; Hernandez et al., 2005), whether from neural maturation, interference from the first language, or other causes (Fig. 1A)? Alternatively, is there an initial period of high ability, followed by a continuous decline (Fig. 1B), or a decline that reaches a floor (Fig. 1C) (Johnson & Newport, 1989)? Or does ability remain relatively constant (Fig. 1D), with adults failing to learn for some other reason such as less time and interest (Hakuta et al., 2003; Hernandez et al., 2005)?

Unfortunately, learning ability is a hidden variable that is difficult to measure directly. Studies that compare children and adults exposed to comparable material in the lab or during the initial months of an immersion program show that adults perform better, not worse, than children (Huang, 2015; Krashen, Long, & Scarcella, 1979; Snow & Hoefnagel-Höhle, 1978), perhaps because they deploy conscious strategies and transfer what they know about their first language. Thus, studies that are confined to the initial stages of learning cannot easily measure whatever it is that gives children their long-term advantage. (Note that strictly speaking, these studies measure learning *rate*, not learning *ability*. While these are conceptually distinct, in practice they are difficult to disentangle, and the distinction has played little role in the literature. In the present paper, we will use the terms interchangeably.)

Thus, although the question of learning ability (1) is more theoretically central, empirical studies have largely probed the more tractable question of how ultimate attainment changes as a function of age of first exposure (2). Here, too, there are a number of theoretically interesting possibilities (Fig. 1E–H). The hope has been that identifying the shape of the ultimate attainment curve might tell us something about the shape of the learning ability curve (cf. Birdsong, 2006; Hakuta et al., 2003; Johnson & Newport, 1989). Unfortunately, this turns out not to be the case. Despite the similarities between the two sets of hypothesized curves (e.g., compare Fig. 1A and E), they bear little relationship to one another: The same ultimate attainment curve (e.g., Fig. 1E) is consistent with many different learning ability curves (Fig. 1A–D).

Here is why learning ability curves (Fig. 1A–D) and ultimate attainment curves (Fig. 1E–H) should not be conflated: If, hypothetically, learning ability plummeted at age 15 but it took 10 years of experience to master a language completely, then ultimate attainment would

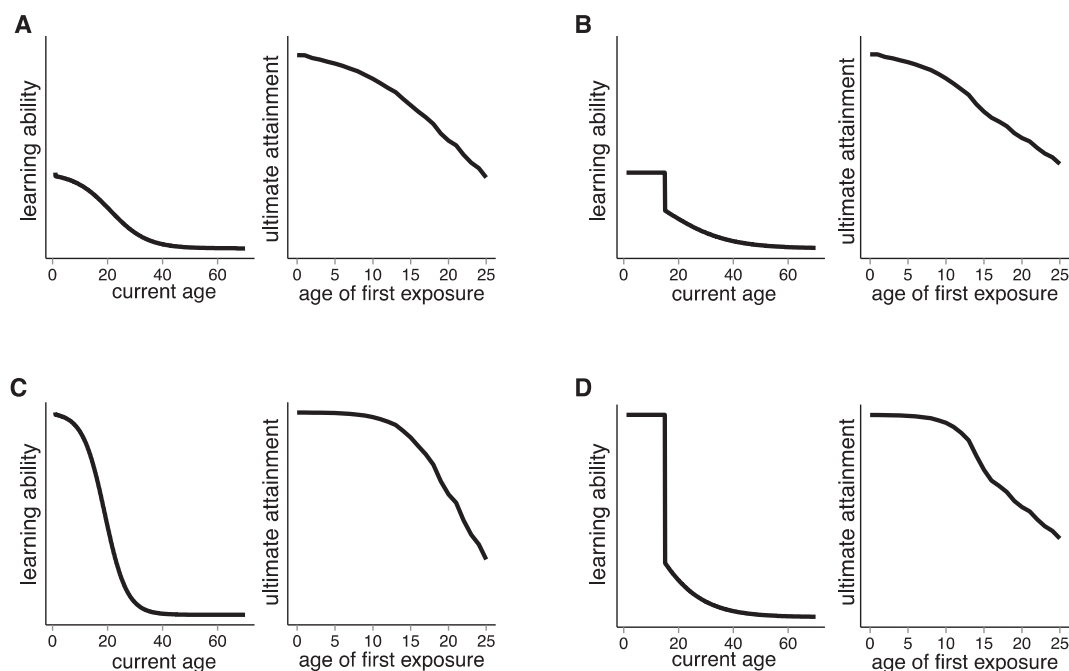
decline starting at an age of exposure of 5 (since someone who began at 6 years old would learn at peak capacity for only 9 of the 10 years required, someone who began at 7 years old would learn for only 8 of those years, and so on). It would be erroneous, in that case, to conclude that a decline in *ultimate attainment* starting at age 5 implied that children's *learning ability* declines starting at age 5. Conversely, showing that people who began learning at a certain age reached native-like proficiency merely indicates that they learned fast enough, not that they learned as fast as a native speaker, just as the fact that two runners both finished a race indicates only that they both started early enough and ran fast enough, not that they ran at the exact same speed.

As a result, it is impossible to directly infer developmental changes in underlying ability (the theoretical construct of interest) from age-related changes in ultimate attainment (the empirically available measurements). Fig. 2 shows that two very distinct ability curves, one with a steady decline from infancy (2A), the other with a sudden drop in late adolescence (2B), can give rise to indistinguishable ultimate attainment curves. (The curves are generated by our ELSD model, described below, but the point is model-independent.) Conversely, a rapid drop in ultimate attainment beginning at age 10 could be explained by a continuous decline in learning ability beginning in infancy (Fig. 2C) or by a discontinuous drop in learning rate at 15 years old (Fig. 2D). Moreover, *quantitative* differences in the magnitude of a hypothetical decline in underlying learning ability (which are not specified in existing theories) can give rise to *qualitative* differences in the empirically measured ultimate attainment curves, such as a gentle decline versus a sudden drop-off: compare Fig. 2A with 2C, and Fig. 2B with 2D.

## 1.2. The present study

As we have seen, to understand how language-learning ability changes with age, we must disentangle it from age of exposure, years of experience, and age at testing. Unfortunately, this challenge is insuperable with any study that fails to use sufficiently large samples and ranges, because any imprecision in measuring the effects of amount of exposure on attainment, the effects of age of first exposure on attainment, or both, will render the results ambiguous or even uninterpretable.

Moreover, an underlying ability curve can be ascertained only if the measure of language attainment is sufficiently sensitive: If learners hit an artificial ceiling, any gains from an earlier age of exposure or a greater amount of exposure will be concealed. Indeed, the concept of native proficiency entails *extreme* levels of accuracy. An error rate that would be considered excellent in other academic or psychological settings, such as 0.75%, represents a conspicuous immaturity in the



**Fig. 2.** Simulation results showing how the mapping between hypothetical changes in underlying learning rate (the left graph in each pair) and empirically measured changes in ultimate attainment is many-to-many. These quantitative predictions were derived from the ELSD model, described below, but the basic point is model-independent.

context of language. For example, over-regularizations of irregular verbs, such as *runned* and *breaked*, are among the most frequently noted errors in preschoolers' speech (Pinker, 1999), despite occurring in only 0.75% of utterances (and on 2.5% of past-marked irregular verbs; Marcus et al., 1992).

These basic mathematical facts raise a significant practical problem: Detecting an error that occurs as little as 0.75% of the time requires a lot of data: A preschooler has to produce 92 utterances to have a better than even chance of producing an over-regularization. Thus, to detect even “conspicuous” errors, such as childhood over-regularization, we need to test many subjects on many items.

Below, we describe a study of syntax that attempts to meet these challenges using novel experimental and analytical techniques. To foreshadow, the age at which syntax-learning ability begins to decline is much later than usually suspected, and it takes both native and non-native speakers longer to reach their ultimate level of attainment than has been previously assumed. While both findings are unexpected, we show that the apparent inconsistencies with prior findings can be explained by the much higher precision afforded by our methods. Indeed, the findings below should not be surprising in retrospect. More importantly, these findings appear robust and emerge in a variety of different analyses.

## 2. Method

### 2.1. Overview

Initial power calculations suggested that several hundred thousand subjects of diverse ages and linguistic backgrounds would be required to disentangle age of first exposure, age at testing, and years of exposure (we return to issues of power in the discussion, below). The standard undergraduate subject pools are not nearly large or diverse enough to achieve this, nor are crowdsourcing platforms like Amazon Mechanical Turk (Stewart et al., 2015). Inspired partly by Josh Katz's Dialect Quiz for the *New York Times*, we developed an Internet quiz we hoped would be sufficiently appealing as to attract large numbers of participants. In order to go viral, the quiz needed to be entertaining and

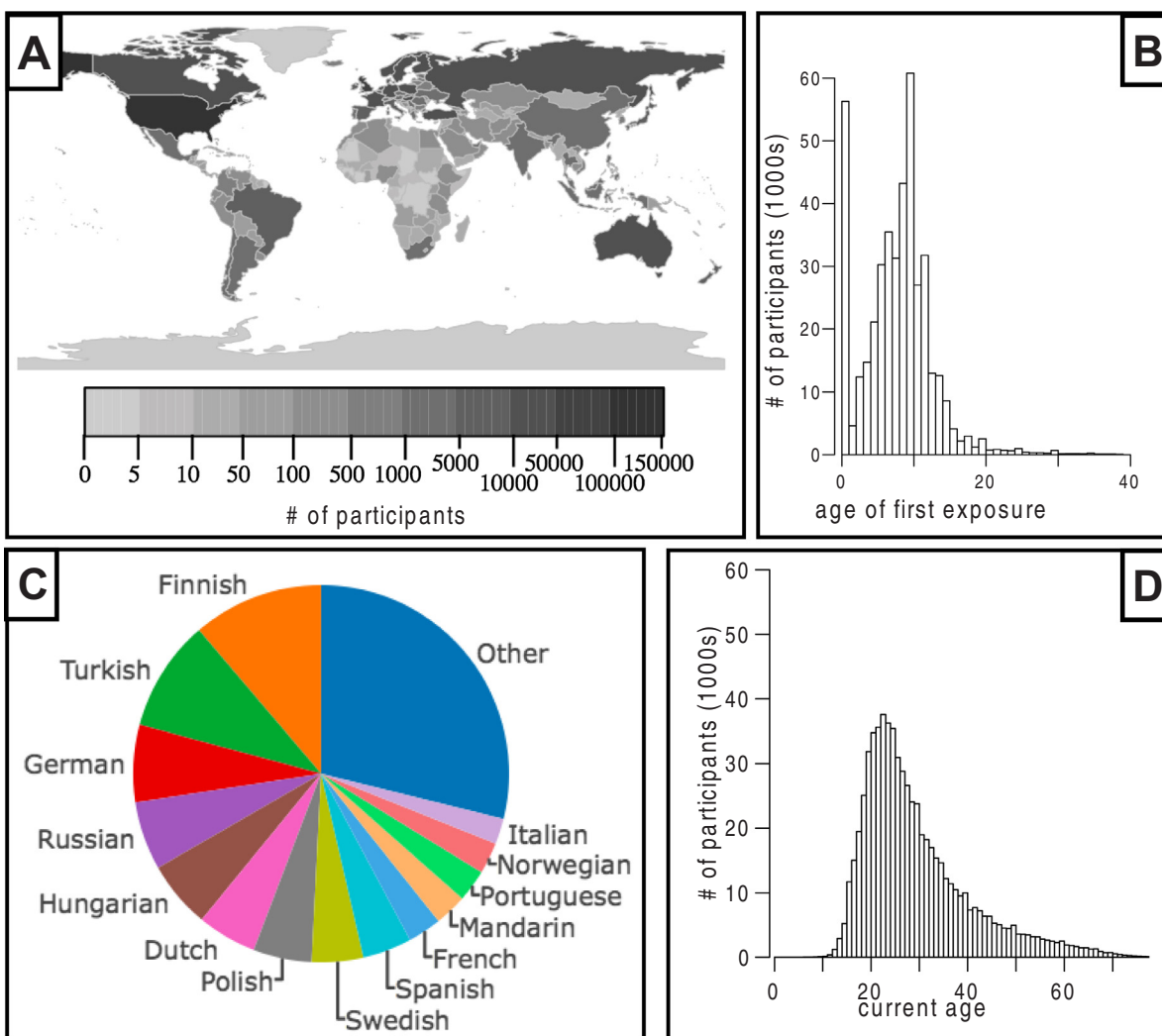
intrinsically motivating while also quick to complete, since Internet volunteers rarely spend more than 10 min on a quiz. At the same time, to yield useful data the quiz had to include a robust, comprehensive measure of syntactic knowledge without an artificial ceiling, as well as elicit demographic data about age and linguistic background. Below, we describe how we addressed these desiderata. Procedures were approved by the Committee on the Use of Humans as Experimental Subjects at Massachusetts Institute of Technology.

### 2.2. Procedure

Potential subjects were invited to take a grammar quiz ([www.gameswithwords.org/WhichEnglish](http://www.gameswithwords.org/WhichEnglish)), the results of which would allow a computer algorithm to guess their native language and their dialect of English. After providing informed consent, subjects provided basic demographic details (age, gender, education, learning disability) and indicated whether they had taken the quiz before. They then completed the quiz and were presented with the algorithm's top three guesses of their native language and their dialect, which was based on the Euclidean distance between the vector of the subject's responses and the vector of mean responses for each language and dialect. Participants found this aspect of the quiz highly engaging, and the quiz was widely shared on social media. For instance, it was shared more than 300,000 times on Facebook.

After seeing the guesses, subjects were invited to help us improve the algorithm by filling out a demographic questionnaire. (Although early answers were used to tune the algorithm, the algorithm's accuracy quickly plateaued and was not tuned further.) This included all the countries they had lived in for at least 6 months, and all the languages they spoke from birth.<sup>1</sup> Participants who listed multiple countries were asked to indicate their current country. For some countries (such as the USA), additional localizing information was collected. Participants who did not report speaking English from birth were asked at what age they

<sup>1</sup> The first several thousand participants were asked to list their “native languages.” Based on participant feedback, this was adjusted to “native languages (learned from birth).”



**Fig. 3.** (A) Current country of residence of participants (excluding participants with multiple residences). (B) Histogram of participants by age of first exposure to English. (C) Native languages of the bilinguals (excluding English). (D) Histogram of participants by current age.

began learning English, how many years they had lived in an English-speaking country, and whether any immediate family members were native speakers of English. Approximately 80% of subjects who completed the syntax questions also completed this demographic questionnaire. The data reported here come from those subjects.

### 2.3. Participants

All participants gave informed consent. 680,333 participants completed the experiment, excluding repeats. We further excluded participants who gave inconsistent or implausible responses to the demographic questions (listing a current age less than the age of first exposure to English; listing a current age that is less than the number of years spent in an English-speaking country; reporting college attendance and a current age of less than 16, or reporting graduate school attendance and a current age of less than 19), resulting in 669,800 participants. Finally, based on the histogram of ages, we excluded participants younger than 7 and older than 89 as implausible. Note: a number of participants ages 7–10 reported in the comments that their parents helped by reading the quiz to them, adding credibility to those data. The resulting number of participants for the analyses was 669,498.

The sample was demographically diverse (Fig. 3). Thirty-eight languages were represented by at least 1000 native speakers, not counting

individuals who had multiple native languages. The most common native languages other than English were Finnish (N = 39,962), Turkish (N = 36,239), German (N = 24,995), Russian (N = 22,834), and Hungarian (N = 22,108).

Analyses focused on three subject groups. *Monolinguals* (N = 246,497) grew up speaking English only; their age of first exposure was coded as 0. *Immersion learners* (N = 45,067) were either simultaneous bilinguals who grew up learning English simultaneously with another language (age of first exposure = 0), or later learners who learned English primarily in an English-speaking setting (defined as spending at least 90% of their life since age of first exposure in an English-speaking country). *Non-immersion learners* (N = 266,701) had spent at most 10% of post-exposure life in an English-speaking country and no more than 1 year in total.<sup>2</sup> Subjects with intermediate amounts

<sup>2</sup> A small proportion of the non-immersion learners (2.7%) reported ages of first exposure between 1 and 3 years. These learners scored quite poorly (the ultimate attainment of those with ages of exposure of 1 year was as poor as those with ages of exposure in their 20s) and exhibited noisy performance curves that, unlike those of all other learners, failed to show any improvement with age (Fig. S1). While this might be a genuine and surprising finding, it more likely reflects the idiosyncratic histories or questionnaire responses of these learners. Unlike the later non-immersion learners, many of whom cited school instruction as their initial source of their exposure, the early non-immersion learners gave little indication about the nature of their first exposure, and it is possible that they had little formal instruction and had learned primarily through television and movies (frequently cited by non-immersion learners as significant sources of English

of immersion (N = 122,068) were not analyzed further.

## 2.4. Materials

We took a shotgun approach to assessing syntax, using as diverse a set of items as we could fit into a short quiz, addressing such phenomena as passivization, clefting, agreement, relative clauses, preposition use, verb syntactic subcategorization, pronoun gender and case, modals, determiners, subject-dropping, aspect, sequence of tenses, and *wh*-movement. This broad approach has two advantages. First, it provides a more comprehensive assessment of syntactic phenomena than many prior studies, which focused on a smaller number of phenomena (Flege, Yeni-Komshian, & Liu, 1999; Johnson & Newport, 1989; Mayberry & Lock, 2003). Second, this diversity provides some robustness to transfer from the first language. That is, while native speakers of some languages may find certain phenomena easier to master than others (e.g., Spanish-speakers may find tense reasonably natural while Mandarin-speakers may find word-order restrictions intuitive), the diversity of items should help wash out these differences (see also discussion below).

### 2.4.1. Item selection

Items were subjected to several rounds of pilot testing to select a sufficient number of critical items that were diagnostic of proficiency (neither too easy nor too hard) and that represented a wide range of grammatical phenomena, while requiring less than 10 min to complete. These included phenomena known to present difficulties for children, such as passives and clefts, and for non-native speakers, such as tenses and articles. We focused particularly on items known to be difficult for speakers of a variety of first languages: in particular, Arabic, French, German, Hindi, Japanese, Korean, Mandarin, Russian, Spanish, or Vietnamese. Based on previous experiments on [gameswithwords.org](http://gameswithwords.org), we expected these to be among the most common native languages.

In addition to the critical items, we included items designed to distinguish among English dialects drawn from websites describing “Irishisms,” “Canadianisms”, and so on. These items were not used for assessing language proficiency and were not used in the data analyses below, but were important for recruiting subjects (see above). Several rounds of pilot-testing reduced this set to the smallest number of items that could reliably distinguish major English dialects.

As in most previous studies, we solicited grammaticality judgments (e.g., “Is the following grammatical: *Who whom kissed?*”). In order to shorten the test and improve the subject experience, where possible we grouped multiple grammaticality judgments into a single multiple-choice question. Because the grammaticality judgment task is time-consuming and unsuitable for probing certain grammatical phenomena, we also included items that required matching a sentence to a picture (e.g., to probe topicalization and the application of linking rules). Several rounds of piloting were used to construct a test that involved items of a range of difficulty.

The final set of 132 items is provided in the [Supplementary Materials](#). Of these, 95 were critical items, defined as items for which the same response was selected by at least 70% of the native English speaking adults 18–70 years old in our full dataset in each of thirteen broadly-defined English dialects (Standard American, African American Vernacular English, Canadian, English, Scottish, Irish, North Irish, Welsh, South African, Australian, New Zealand, Indian, and Singaporean). (For obvious reasons, the exact number of critical items was not known until after the data was collected.) All analyses below are restricted to this set.

Many prior studies classify items according to the syntactic phenomenon they test. While this is straightforward for certain types of tests, such as our sentence-picture matching items, the accuracy of these categorizations for grammaticality judgments is unclear. For

(footnote continued)

input). Given this uncertainty, we excluded these participants from the main analyses.

instance, in judging a sentence to be grammatical, subjects can hardly be expected to know which syntactic rule the experimenter deliberately did *not* violate. Likewise, ungrammatical sentences may implicate different rules depending on what the intended message was: *I eats dinner* could involve an agreement error on the verb or a failure of pronoun selection. Thus, the syntactic violation that catches the subject’s eye may not be the one the experimenter had in mind. Because our goal was merely to have a diverse set of items, an exact count of syntactic phenomena is less important than demonstrating diversity. Thus, we have bypassed these theoretically thorny issues by avoiding categorization and simply providing the entire stimulus set in the [Supplementary Materials](#). As a result, readers can judge for themselves whether the items are sufficiently diverse.

### 2.4.2. Test reliability

Reliability for the critical items was high across the entire dataset (Cronbach’s alpha = 0.86). Because monolingual subjects were close to ceiling, reliability is expected to be lower for that subset. Reliability is a measure of covariation, and the monolinguals exhibited very little variation (the majority missed fewer than 3 items), exactly as one would expect for a valid test. However, reliability for monolinguals was still well above chance (0.66), indicating that what few errors they made were not randomly distributed (as would be expected from mere sloppiness) nor concentrated on a few “bad” items (in which case, there would be little variance). Thus, our test was sensitive to differences in grammatical knowledge even for monolinguals who were close to ceiling. It is difficult to compare these numbers to prior studies, since most did not report reliability (but see DeKeyser, 2000; DeKeyser, Alfi-Shabtay, & Ravid, 2010; Granena & Long, 2013).

### 2.4.3. Data

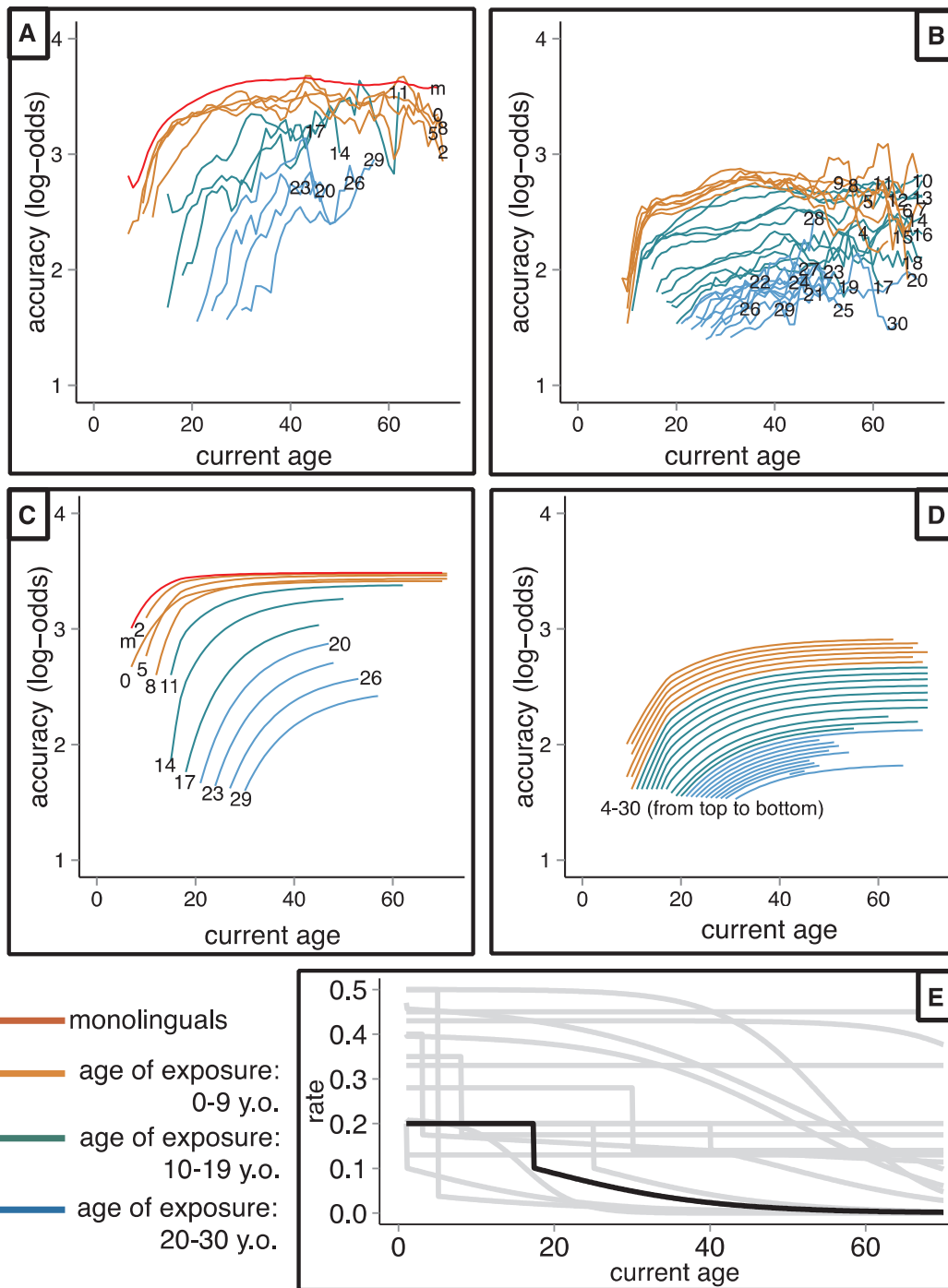
The resulting dataset is available at <http://osf.io/pyb8s>.

## 3. Results

### 3.1. Learning rate

We focus first on the difficult but theoretically important question of the underlying learning rate. We defer the traditional question of level of ultimate attainment to a later section. Note that all analyses are conducted in terms of log-odds (the log-transformed odds of a correct answer, using the empirical logit method to avoid division by zero) rather than percent correct. Although prior work on critical periods has tended to use percent correct, this is problematic. Specifically, percentage points are not all of equal value, being more meaningful closer to 0% or 100% than when near 50% (Jaeger, 2008). That is, the difference between 95% and 96% is “larger” than the difference between 55% and 56%. Thus, the use of percentages artificially imposes ceiling effects, inflating both Type I and Type II error rates, particularly for interactions. Similarly, graphing results in terms of percentage correct distorts the results (particularly the shapes of curves), and so we have graphed in terms of log odds. For reference, we have included percent correct on the right-hand side of many of the graphs.

Fig. 4 plots the level of performance against current age in separate curves for participants with different ranges of age of first exposure. It simultaneously reveals the effects of age of first exposure (the differences among the curves) and total years of exposure (the left-to-right position along each curve). Immersion learners—who were less numerous than the other groups—were aggregated into three-year bins for age of exposure, except for the simultaneous bilinguals (age of exposure = 0), who constituted their own bin. Curves were smoothed with a five-year floating window (analyses on non-smoothed data are discussed in the next subsection), and each of the estimated performance curves (described below) was restricted to consecutive ages for which there were at least ten participants in the five-year window, leaving 244,840 monolinguals, 44,412 immersion learners, and



**Fig. 4.** (A and B) Performance curves for monolinguals and immersion learners (A) and non-immersion learners (B) under 70 years old, smoothed with five-year floating windows. (C and D) Corresponding curves for the best-fitting model. (E) Learning rate for the best-fitting model (black), with examples of the many hypotheses for how learning rate changes with age that were considered in model fitting (grey). For additional detail, see Fig. 7, S3, and S6.

257,998 non-immersion learners.

In order to estimate how underlying learning ability changes with age, we used a novel computational model to disentangle current age, age of first exposure, and amount of experience. Specifically, we modeled syntax acquisition as a simple exponential learning process:

$$g(t) = 1 - e^{-\int_{t_c}^t E r dt} \tag{1}$$

where  $g$  is grammatical proficiency,  $t$  is current age,  $t_c$  is age of first exposure,  $r$  is the learning rate, and  $E$  is an experience discount factor, modeled separately for simultaneous bilinguals, immigrants, and non-

immersion learners, reflecting the fact that they may receive less English input than monolinguals. We modeled a possible developmental change in the learning rate  $r$  as a piecewise function in which  $r$  is constant from birth to age  $t_c$ , whereupon it declines according to a sigmoid with shape parameters  $\alpha$  and  $\delta$  ( $\alpha$  controls the steepness of the sigmoid, and  $\delta$  moves its center left or right):

$$r(t) = \begin{cases} r_0, & t \leq t_c \\ r_0 \left( 1 - \frac{1}{1 + e^{-\alpha(t-t_c-\delta)}} \right), & t > t_c \end{cases} \tag{2}$$

The piecewise structure of this Exponential Learning with Sigmoidal

Decay (ELSD) model, and the fact that sigmoid functions can accommodate both flat and steep declines, allows it to capture a very wide range of developmental trajectories, including all of those discussed in the literature. Learning rate may be initially high or low, begin declining at any point in the lifespan (or not at all), decline rapidly or gradually, decline continuously or discontinuously, etc. Examples of the many possibilities encompassed by the model include the different curves shown in Figs. 2 and S2, as well as the gray lines in Fig. 4E.

The model was fitted simultaneously to the performance curves for monolinguals, immersion learners, and non-immersion learners (cf. Fig. 4A and B). Parameters were fit with Differential Evolution (Mullen, Aridia, Gil, Windover, & Cline, 2011) and compared using Monte Carlo split-half cross-validated  $R^2$ , which avoids over-fitting. The best-fitting model ( $R^2 = 0.89$ ) involved a rate change beginning at 17.4 years (Fig. 4E). The fit was significantly better than the best fit for alternative models in which learning rate did not change ( $R^2 = 0.66$ ) or changed according to a step function with no further decline in the learning rate after the initial drop ( $R^2 = 0.70$ ). Details on these and related models can be found in the [supplementary materials](#).

### 3.2. Interim discussion

Though the ELSD model is necessarily simplified, the good fit between model and data, and the poorer fit by reasonable alternatives, offers good support for the existence of a critical period for language acquisition, and suggests that our estimate of when the learning rate declines (17.4 years old) is likely to be reasonably accurate.

This age is much later than what is usually found for the offset of the critical period for native-like *ultimate* attainment of syntax. However, as discussed in the Introduction, because language acquisition takes time, there is no reason to suppose that the last age at which native-like ultimate attainment can be achieved is the same as the age at which underlying ability declines (see also Patkowski, 1980). Instead, the relationship between ultimate attainment and critical periods is complex, depending also on how long it takes to learn a language. The ELSD model disentangles these factors. In order to better understand the results of the above analyses, we look at these issues in turn.

### 3.3. The duration of learning

Little is known about how long it takes learners to reach asymptotic performance. On the one hand, developmentalists have observed that by 3–5 years of age, most children show above-chance sensitivity to many syntactic phenomena (Crain & Thornton, 2011; Pinker, 1994). Indeed, our youngest native speakers (~7 years old) were already scoring very well on our quiz (Fig. 5B).

While certainly an important fact about acquisition, this is the wrong standard for research into critical periods. The question has never been “why do non-native speakers not match the competency level of preschooler?” Many of them do. In fact, in our dataset, even non-native immersion learners who began learning in their late 20s eventually surpassed the youngest native speakers in our dataset (Fig. 4A).

Instead, the puzzle driving this entire research domain is why later learners do not reach the same proficiency level of mature native speakers. That is a much higher standard. Many other aspects of syntax continue to develop in the school-age years (Berman, 2004, 2007; Nippold, 2007), and prior studies have not been able to determine the age at which syntactic development concludes. Even for those aspects of syntax that preschoolers are sensitive to, they are rarely at ceiling, and they typically do worse than college-age adults, whether assessed through comprehension, elicited production, or spontaneous production (e.g., Kidd & Bavin, 2002; Kidd & Lum, 2008; Marcus et al., 1992; Messenger, Branigan, McLean, & Sorace, 2012; Rowland & Pine, 2000). However, while we know that performance continues to improve into the school ages, the literature has little to say about when children

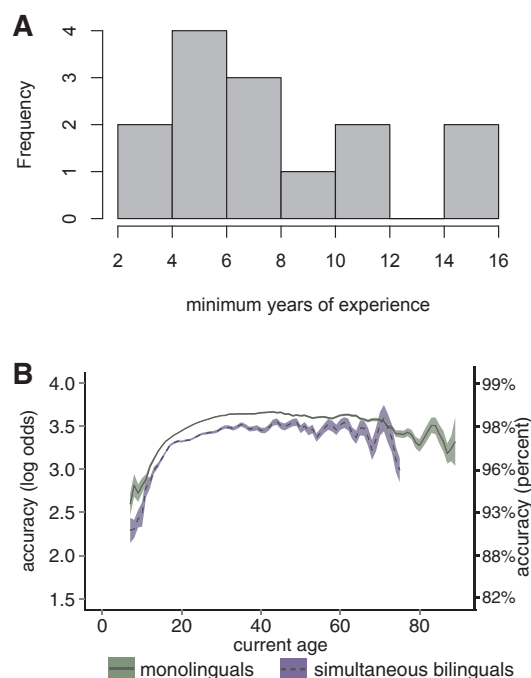


Fig. 5. (A) Histogram of cutoffs used for minimum years of experience to asymptotic learning in previous studies of syntax (Abrahamsson, 2012; Birdsong & Molis, 2001; DeKeyser, 2000; DeKeyser et al., 2010; Flege et al., 1999; Granena & Long, 2013; Jia et al., 2002; Johnson & Newport, 1989, 1991; Mayberry & Lock, 2003; Mayberry, Lock, & Kazmi, 2002; McDonald, 2000; Weber-Fox & Neville, 1996). Papers with multiple studies are included only once, except for McDonald (2000), which used different cutoffs in two different studies. (B) Accuracy for monolinguals ( $N = 246,497$ ) and simultaneous bilinguals ( $N = 30,397$ ). Shaded area represents  $\pm 1$  SE. This highlights information also available in Fig. 4A.

attain adult levels of accuracy. Moreover, the common practice of comparing children to college-aged adults necessarily renders undetectable any post-college development.

Even less is known about how long non-native speakers continue to improve on the target language. While a few studies found limited continued improvement for immersion learners after the first five years (Johnson & Newport, 1989; Patkowski, 1980), these studies had minimal power to detect continued improvement (see below). Specifically, looking at samples of non-native learners who were selected to have at least three years (Johnson & Newport, 1989) or five years (Patkowski, 1980) of experience, these authors found that while age of first exposure predicted performance, length of experience did not. In contrast, analysis of US Census data suggests that learning continues for decades (Stevens, 1999), though the validity of this self-report data is uncertain. Analysis of foreign language education suggests learning in that context may continue for a couple of decades, though this may merely reflect the slower pace of non-immersion learning (Huang, 2015).

This empirical uncertainty is reflected directly in the ultimate attainment literature. Ultimate attainment analyses require restricting analysis to those subjects who have been learning the target language long enough to have reached asymptote (e.g., Johnson & Newport, 1989). In the absence of any clear evidence, researchers have chosen a diverse set of cut-offs, ranging anywhere from three (Birdsong & Molis, 2001; McDonald, 2000) to fifteen years (Abrahamsson, 2012) (Fig. 5A).

Inspection of Fig. 5B suggests that native speakers did not reach asymptote until around 30 years old, though most of the learning takes place in the first 10–20 years. The results for later learners shown in Fig. 4 similarly suggest a protracted period of learning (for detailed results, see Figs. S21 and S22 in the [Supplementary Materials](#), and

surrounding discussion). Note that the increases in performance after the first 15–20 years are modest, which accords with the fact that they are not routinely noticed.

While this prolonged learning trajectory was not anticipated in the language learning literature, it joins mounting evidence that many cognitive abilities continue to develop through adolescence and even adulthood, including working memory, face recognition, magnitude estimation, and various measures of crystallized intelligence (Germine, Duchaine, & Nakayama, 2011; Halberda, Ly, Wilmer, Naiman, & Germine, 2012; Hartshorne & Germine, 2015).

Thus, even native speakers—who are able to make full use of the critical period—take a very long time to reach mature, native-like proficiency. By implication, someone who started relatively late in the critical period—that is, someone who had limited time to learn at the high rate the critical period provides—would simply run out of time. In order to follow up on this issue and test this implication, we turn to analysis of ultimate attainment.

### 3.4. Ultimate attainment

Based on the results above, we expect that the last age of first exposure at which native-like attainment is still within reach is likely well prior to 17. Below, we first estimate this age from our own data and then compare that against previous estimates.

Following the usual practice, we first restrict the analysis to those subjects who have been learning English long enough to have reached asymptote (e.g., Johnson & Newport, 1989). As described in the previous section, there is no consensus as to how long “long enough” is (see Fig. 5A). This stems from the fact that, prior to our own study, there was little data to constrain hypotheses (see previous section). Inspection of Figs. 4 and 5 suggests 30 years old as a reasonable cutoff.

Thus, to estimate the age at which mastery of a second language is no longer attainable, we analyzed ultimate attainment curves by focusing on the 11,371 immersion learners and 29,708 non-immersion learners who had at least 30 years of experience (ensuring asymptotic learning) and who were at most 70 years old (avoiding age-related decline) (Fig. 6). We fitted these curves using multivariate adaptive regression splines (Friedman, 1991; Milborrow, 2014). Immersion learners showed only a minimal decline in ultimate attainment until an age of first exposure of 12 years ( $B = -0.009$ ; 0.01 SDs/year), after which the decline became significantly steeper ( $B = -0.06$ ; 0.07 SDs/year). Non-immersion learners showed similar results: From 4 years to 9 years, proficiency showed no decline (in fact it increased slightly;  $B = 0.01$ ; 0.01 SDs/year), followed by a steep decline ( $B = -0.06$ ; 0.07 SDs/year). Two other methods of estimating changes in slope provided similar results (see Supplementary Materials).

While these analyses employ the standard method of analyzing subjects who have (presumably) already reached ultimate attainment, the density of our data allows a more direct analysis. Fig. 7 re-plots the data in Fig. 4 against years of experience, aligning the curves for the learners who began at different ages at the onset of learning. Inspection reveals that the learning trajectories for immersion learners who began in the first decade of life (the orange curves) are almost indistinguishable (Fig. 7A). We see a similar trend for the non-immersion learners (Fig. 7B).

We confirmed these observations with permutation analysis. Specifically, we calculated the average difference between each performance curve and the performance curve for the youngest learners of that type (the simultaneous bilinguals for immersion learners, the learners with an age of first exposure of 4 years for the non-immersion learners). A positive score indicated that the performance curve was, on average, below the curve for the earliest learners. We then constructed an empirical distribution by randomly permuting the age of exposure across participants at a given number of years of experience. The curves were again smoothed with five-year floating windows and the difference scores were again calculated. This was repeated 1000 times. The

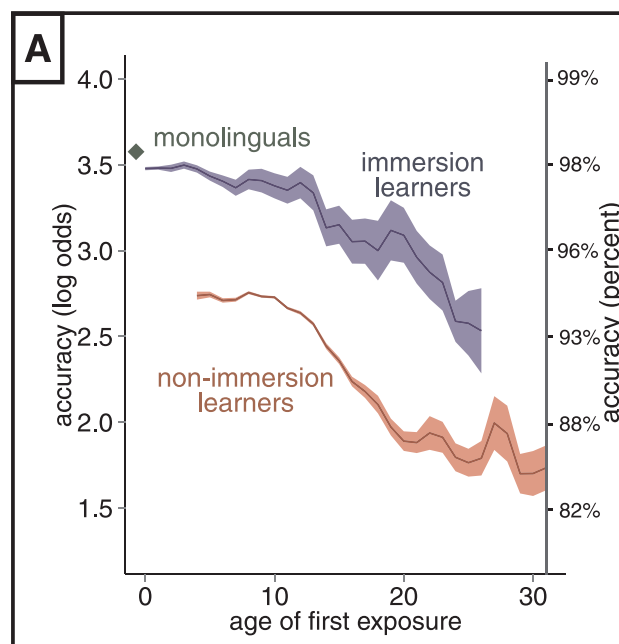


Fig. 6. Ultimate attainment for monolinguals, immersion learners, and non-immersion learners, smoothed with a three-year floating window. Shaded areas represent  $\pm 1$  SE. Attainment for monolinguals was significantly higher than that of simultaneous bilinguals (immersion learners with exposure age = 0) ( $p < .01$ ).

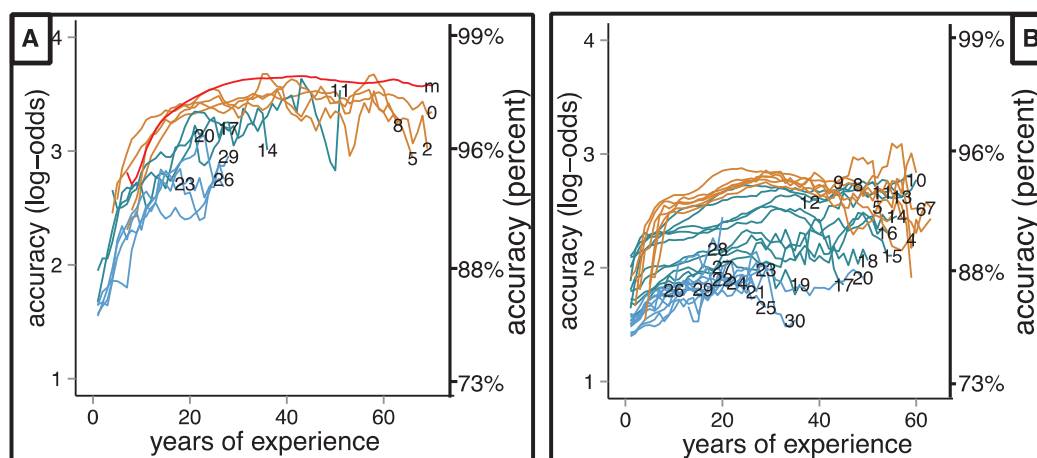
percentage of cases in this distribution in which the difference score for a given performance curve is larger than the actual difference score for that performance curve serves as a one-tailed  $p$ -value (all comparisons reported as significant are also significant as two-tailed tests). These analyses revealed that the performance curves for immersion learners with average exposure ages of 2, 5, and 8 years were not significantly different from those of simultaneous bilinguals (exposure age = 0;  $ps > 0.31$ ), while the curves for later learners were significantly lower ( $ps < 0.01$ ). Similarly, non-immersion learners with ages of exposure of 5–11 years were indistinguishable from our earliest non-immersion learners (4 years;  $ps > 0.31$ ), whereas later learners learned significantly more slowly ( $ps < 0.01$ ).

#### 3.4.1. Comparison with previous ultimate attainment results

Both traditional ultimate attainment analyses and permutation analyses indicated that learners must start by 10–12 years of age to reach native-level proficiency. Those who begin later literally run out of time before the sharp drop in learning rate at around 17–18 years of age. For non-immersion learners, the ceiling was lower but the overall story was the same: little difference between learners who start within the first decade of life, with a ceiling that noticeably drops for later learners. These findings are consistent with the protracted trajectory of learning that we observe in our data (see previous section).

However, our results for immersion learners diverge from those of some previous studies (there are no similar studies of non-immersion learners). For instance, Johnson and Newport's (1989) study of immersion learners found no correlation between ultimate attainment and age of first exposure after an onset age of 16, whereas we see a strong relationship (for review, see Qureshi, 2016). In principle, this could be due to differences in subject population or the types of grammar rules tested. Indeed, researchers frequently argue that such differences have large effects on ultimate attainment, based on the fact that studies of different populations or stimuli have produced different results (Abrahamsson, 2012; Birdsong & Molis, 2001; DeKeyser, 2000; DeKeyser et al., 2010; Flege et al., 1999; Granena & Long, 2013; Hakuta et al., 2003; Jia et al., 2002; Johnson & Newport, 1989; Vanhove, 2013;





**Fig. 7.** Accuracy as a function of years of experience, by age of first exposure for immersion learners (A) and non-immersion learners (B). Color scheme is same as in Fig. 4. Red: monolinguals. Orange: AoFE < 11. Green: 10 < AoFE < 21. Blue: AoFE > 20. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Weber-Fox & Neville, 1996).

However, a recent analysis by Vanhove (2013) raised questions about whether these differences are statistically meaningful. Whereas most prior studies had between 50 and 250 subjects, Vanhove demonstrates that precisely measuring how ultimate attainment changes as a function of age of first exposure requires thousands. Only one previous dataset, based on US Census data, reaches sufficient sample size (Hakuta et al., 2003; Stevens, 1999). However, this study was based on a self-report of proficiency on a four-point scale, which is unlikely to have much precision. Thus, differences across findings in the literature could reflect nothing more than random noise.

Thus, in order to better understand whether the differences in our findings and those of prior studies are meaningful, we need to consider the precision of these findings. We estimated precision using bootstrapping, simulating running many different studies by resampling with replacement from our own data (Efron & Tibshirani, 1993). The results of each simulation will be slightly different, and so the range of results across simulations simulates the variability we would expect from statistical noise alone. Crucially, we can simulate running studies with different sample sizes. Thus, we can ask whether Johnson and Newport's (1989) findings are within what we might have found had we used our own methods but tested the same number of subjects ( $N = 69$ ).

For our simulations, we considered two different sample sizes:  $N = 69$ , the size of the classic Johnson and Newport (1989) study, and  $N = 275$ , larger than the largest prior study, with the exception of the aforementioned Census studies. For comparison, we also simulated studies with  $N = 11,371$ , the number of subjects in our own ultimate attainment results described in the previous section.

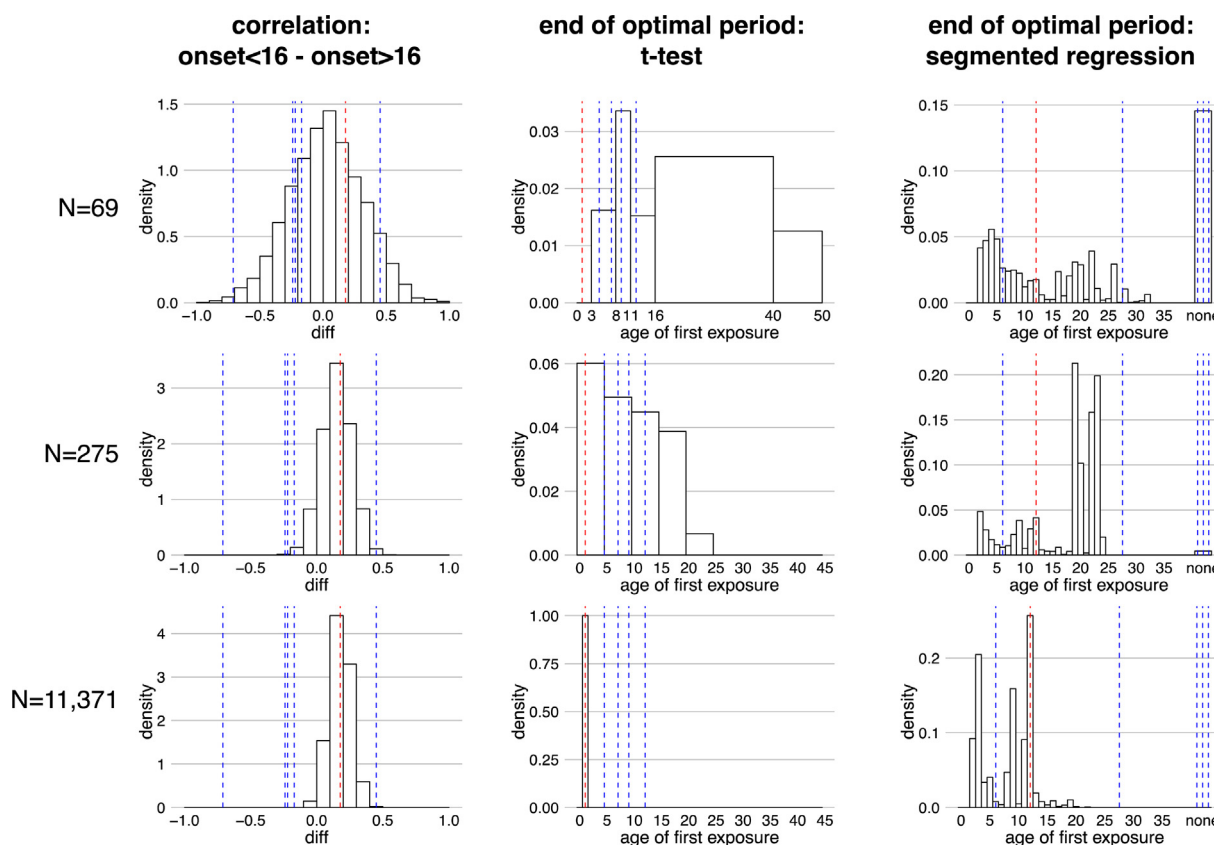
We focused on three different analyses that have been reported in a number of prior studies (Bialystok & Miller, 1999; Birdsong & Molis, 2001; DeKeyser, 2000; DeKeyser et al., 2010; Flege et al., 1999; Johnson & Newport, 1989; Weber-Fox & Neville, 1996). First, we considered Johnson and Newport's finding that the correlation between age of first exposure and ultimate attainment is much stronger before an exposure age of 16 ( $r = -0.87$ ) than after ( $r = -0.16$ ). This finding has proved controversial, with subsequent studies finding much weaker effects or no effect at all (Bialystok & Miller, 1999; Birdsong & Molis, 2001; DeKeyser, 2000; Johnson & Newport, 1989). All these prior findings are well within what one would expect for  $N = 69$  (Fig. 8, upper left). As power increased, the variability in the estimates dropped dramatically, with more highly-powered studies being increasingly unlikely to find any substantial difference in the correlations before and after 16 years old.

Second, Johnson and Newport also reported that individuals who began learning English at 8–10 years old failed to reach monolingual-like ultimate attainment, whereas individuals who began earlier did, suggesting that the “optimal period” for language-learning is 0–7 years old. Once again, there has been considerable variability in subsequent studies, and our own study finds that even simultaneous bilinguals do not quite reach monolingual levels. Vanhove (2013) suggested, based on power calculations, that accurately estimating the end of the optimal period requires thousands of subjects. Although a small study can detect very large effects, the differences between learners who began just within the optimal period and those who began just after are relatively small (Fig. 6) and thus undetectable with a low-power study. Our simulations confirm this analysis (Fig. 8, middle column): in our simulation of Johnson & Newport (Fig. 8, middle column, top), the 95% confidence interval contained almost the entire range. Even with 275 subjects, a wide range of findings would be expected. However, simulations based on our full sample show no variability at all, with learners who began at 1 year of age performing reliably worse than monolinguals (Fig. 8, middle column, bottom).

Third, whereas the previous analysis of the optimal period followed Johnson and Newport's method of using t-tests to compare native speakers to groups of later-learners, subsequent researchers have used instead curve estimation—typically segmented regression with breakpoint estimation—which is argued to be more precise and less prone to false positives (Birdsong & Molis, 2001; Vanhove, 2013; but see DeKeyser et al., 2010). If there is an optimal period, the slope of the ultimate attainment curve should initially be close to 0, followed by a point where it becomes significantly more negative. By this standard of evidence, most studies have failed to find any evidence of an optimal period (Birdsong & Molis, 2001; Flege et al., 1999; Vanhove, 2013). Our simulations suggest these prior findings were false negatives due to low power: Like the majority of prior studies, low-power simulations elicited largely null results, whereas high-power simulations suggested an optimal period ending in early or middle childhood (Fig. 8, right).

### 3.4.2. Interim discussion

Two sets of analyses of our data suggest that learners who begin as late as 10–12 years old reach similar levels of ultimate attainment as native bilinguals. After that age, we find a continuous decline in attainment as a function of age of first exposure, with no evidence that this relationship ceases after a particular age (cf. Johnson & Newport, 1989; Pulvermüller & Schumann, 1994). These findings are consistent with our results for learning rate. Interestingly, these findings held not only for immersion but also non-immersion learners, a population that



**Fig. 8.** We conducted 2500 simulated experiments of monolingual and immersion learners with each of three sample sizes:  $N = 69$  (equivalent to Johnson & Newport, 1989),  $N = 275$  (larger than the largest prior lab-based study), and  $N = 11,371$  (equivalent to the present study). Three analyses were considered. Left: Correlation between age of first exposure and ultimate attainment prior to 16 years old minus after 16 years old. Middle: First subgroup of subjects to be significantly worse than monolinguals in a  $t$ -test (note: the top graph uses the same age bins as Johnson & Newport, 1989). Right: age of first exposure at which performance begins to decline more rapidly, if any. Blue: estimates from Bialystok and Miller (1999), Birdsong and Molis (2001), DeKeyser (2000), DeKeyser et al. (2010), Flege et al. (1999), Johnson and Newport (1989), and Weber-Fox and Neville (1996). While many other papers addressed similar issues, these papers provide the closest analog to Johnson & Newport in that they used a broad-spectrum test of syntax, defined the onset of learning as the age at immigration, and (crucially) report comparable statistics. Red: estimates from current study. Full details available in Supplementary Materials. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

has not been much studied in this regard.

Our findings do contrast with the conclusions of some prior studies of ultimate attainment in immersion learners. However, as our simulations show, these conclusions were probably overfit to point estimates. That is, conclusions depended on the most probable estimate (the optimal period ends at 8 years of age), ignoring the error bars, which in some cases were likely so large as to encompass the entire possible range (Fig. 8). In contrast, our larger sample size allows for fairly precise estimates (Fig. 8). These simulations support Vanhove’s (2013) contention that thousands of subjects are required to provide reliable conclusions about ultimate attainment. Note that we cannot conclude that differences in stimuli or population do not matter for ultimate attainment, only that studying such effects requires very large datasets. We return to this issue in the General Discussion.

#### 4. General discussion

Taken together, the analyses above all point to a grammar-learning ability that is preserved throughout childhood and declines rapidly in late adolescence. This model provided a better fit to the data than did a wide range of alternatives, including models with declines that were earlier or later, faster or slower, sharper or smoother.

In addition to providing the first empirical estimate of how language-learning ability changes with age, we addressed two related issues. First, we found that native and non-native learners both require around 30 years to reach asymptotic performance, at least in immersion

settings. While this question has not been previously addressed, these findings are compatible with what is known about the initial period of learning.

Second, we found that ultimate attainment—that is, the level of asymptotic performance—is fairly consistent for learners who begin prior to 10–12 years of age. We found no evidence that the ultimate attainment curve reaches a floor at around puberty, as has been previously proposed (Johnson & Newport, 1989). While these results differed from the conclusions of some prior studies, our simulations showed that the prior findings were in fact too noisy to provide precise estimates.<sup>3</sup> To provide reliable results about ultimate attainment, a study should have in excess of 10,000 subjects (see also Vanhove, 2013). This suggests that the results of those prior studies, all but one of which has fewer than 250 subjects, largely reflect statistical noise. The remaining study had many subjects but uncertain validity (see discussion above).

This set of results is internally consistent, adding credibility to the whole. However, our conclusions—like any conclusions—are only as good as the data supporting them. Below, we address a number of possible concerns. These include both methodological concerns about the data and how they were collected but also more theoretical

<sup>3</sup> We also noted a number of limitations and confounds in prior studies, such as how ultimate attainment was defined, which would have biased results. However, detailed investigation shows that the resulting biases and imprecisions were likely swamped by the effect of low power (see Supplementary Materials, “Effect of Analysis Decisions”).

concerns, like the possibility that results differ across subsets of subjects or items. We then conclude by discussing the implications of our results, should they prove valid and robust.

#### 4.1. Potential concerns and complications

##### 4.1.1. Familiarity with the testing procedure

One possible concern is that differences across subjects were due to age-related differences in familiarity with the Internet. Prior comparisons of Internet-based and offline datasets have found little support for this concern (Hartshorne & Germine, 2015). Similarly, some of the differences between children and adults could conceivably be due to general test-taking ability. In order to better understand interactions between subject age and test method, if any, it would be ideal to gather data from a variety of tests in a variety of modalities.

Crucially, however, most of our analyses did not depend on the *current* age of the subject but on their age at first exposure, which should weaken any effects of current age. Moreover, we can compare the learning trajectories of learners who started at different ages (see Figs. 4 and 7 but especially Figs. S21–S22 in the Supplementary Materials). If older subjects are substantially better at taking our test, this should appear as more rapid early learning. As inspection of the figures indicates, any such effect is inconsistent and small.

##### 4.1.2. Test modality

Our use of a written comprehension test was dictated by our methodology. Comprehension studies can be scored automatically (which is crucial when there are over half a million subjects), and written tests do not require high-quality audio equipment or sound booths. Nonetheless, one might ask how these choices affected our results.

Certainly, differences between production and comprehension and between written and oral modalities can affect comparisons between native and non-native speakers (Bialystok & Miller, 1999). Listening places high demands on speed and memory (one can re-read but not re-hear), and the speech must be analyzed by non-native acoustic phonetics and phonology, which we do not test here. Written tests require literacy. Production allows one to strategically avoid difficult and imperfectly learned words and constructions.

Whether any of these factors affect estimates of a critical period depends on whether they interact with the variables that define critical period effects, namely age at first exposure, current age, and years of experience. While the necessary studies are not currently feasible, this is likely to change as technology improves. (For instance, we are exploring the use of machine learning to characterize the nativeness of a written text.)

Importantly, none of these considerations would make the study of critical periods in written comprehension uninteresting or uninformative, merely complex. Results from any modality must reflect underlying grammatical ability at least to some degree, and reading comprehension is important in its own right, given the importance of reading in many modern societies. (In fact, for many non-native speakers, this may be their primary use for the non-native language.)

##### 4.1.3. Item selection and quiz difficulty

Another potential worry is that our results may depend on smallish differences among subjects who are already near the ceiling (for relevant discussion, see: Abrahamsson & Hyltenstam, 2009; Birdsong, 2006). Mitigating this concern is that, as we argued in the Introduction, the ceiling is where all the action is. What is remarkable about language is that we are (nearly) all extremely good at it, including adult learners. For reference, we noted that over-regularizations of irregular verbs, which are among the most salient errors in the speech of preschoolers, occur in only 0.75% of their utterances. On a continuum of linguistic ability that includes apes and machines at one end, preschoolers and reasonably diligent late learners are clustered at the other end, near

native-speaking adults. Indeed, the question in the critical period literature has never been why adults are incapable of learning a new language—obviously they are—but why adult learners so rarely (if ever) achieve native-like mastery. Likewise, asking whether adult learners can master basic syntax may be theoretically interesting but distracts from the original motivation for this literature: adult learners rarely, if ever, achieve the same level of mastery as those who started in childhood. In order to study that phenomenon, the relevant yardstick is the asymptotic performance of native speakers.

Still, we can ask whether our results hold for both items mastered early in typical development and for items mastered only in adolescence or adulthood. We found no evidence of such a difference: In the best-fitting models of learning, the learning rate began to slow at approximately the same time for the 47 items that are mastered by the youngest monolingual English-speakers in the sample (ages 7–8) as for the 48 items that are mastered only by the older ones: 17.3 years old and 18.2 years old, respectively. Moreover, if there were substantial interactions between item and age of first exposure, we would expect to see substantial differences in terms of which items were more or less difficult for early and late learners. However, item difficulty was strongly correlated across learners regardless of age of first exposure (for details of these analyses, see Supplementary Materials, “Item Effects”).

We might similarly ask whether results vary based on the type of syntactic construction tested. Prior analyses of ultimate attainment have provided conflicting results, likely due to the power issues discussed above (Coppieters, 1987; Flege et al., 1999; Johnson & Newport, 1989, 1991; McDonald, 2000; Weber-Fox & Neville, 1996) and the theoretical issues raised below. Our just-discussed analyses of item difficulty provide some initial evidence against substantial differences across syntactic phenomena. More precise analyses would involve the direct comparison of different types of constructions. Unfortunately, our quiz was designed to cover a wide range of phenomena, and thus we have few items of any given type, making it difficult to distinguish differences between *items* and differences between *item types*. In any case, such analyses raise thorny theoretical questions: different theories of syntactic processing categorize phenomena differently, and any given sentence involves many different phenomena. Thus, classifying items by syntactic phenomena is far from trivial and may not even be the right approach. Progress on this question will require a significant amount of further research.<sup>4</sup> If it turns out that different aspects of syntax do indeed have different critical periods, the conclusions presented here would need to be revised. Design of follow-up studies may be informed by comparing items in our dataset, which is available at <http://osf.io/pyb8s>.

##### 4.1.4. The effect of the first language

Our results are unlikely to be specific to any one language or language family: Participants listed more than 6000 native languages or combinations of them. The best-represented language families among immersion and non-immersion learners were Uralic (N = 54,664), Slavic (N = 41,640), West Germanic (N = 38,385), Romance (N = 40,476), Turkic (N = 29,816), and Chinese (N = 15,161). The remaining 29% of participants either had multiple native languages or had native languages belonging to a different family. Thus, no language contributed more than a small fraction of the immersion or non-

<sup>4</sup> We note a further difficulty. All research in this domain has treated items as fixed effects, averaging across them. This simplifies calculation, but at a cost: such statistical analyses do not directly assess the question of whether the results generalize beyond the items used (Baayen, Davidson, & Bates, 2008; Clark, 1973). This problem is mitigated somewhat when using a large and representative set of items—as we do—but is particularly problematic when looking at smaller samples of items. The standard solution currently is to use mixed effects modeling (Baayen et al., 2008). However, mixed effects modeling requires significant computational power. We have so far been unable to identify a tractable method of applying mixed effects modeling to a dataset the size of the present one.

immersion learners (Fig. 3C). However, this leaves the possibility that our results reflect an epiphenomenal average of very different trajectories for very different types of learners (Bialystok & Miller, 1999; McDonald, 2000).

It is uncontroversial that speakers of different native languages make characteristic mistakes when speaking English (Schachter, 1990, among others); indeed, the algorithm we used as part of our recruitment strategy depended on this fact (see Section 2.2). However, that is logically distinct from the question as to whether critical periods differ across native languages. Ideally, we would compare the results of our model for speakers of different native languages. However, our samples of individual languages are too small. Specifically, because our data are unevenly distributed across ages and learner conditions, we risk overfitting certain conditions (such as monolinguals) at the expense of others. As described in the Method, we circumvented this issue by averaging across subjects in each bin prior to running the model. This is not applied easily to subsets of the data: too many bins have few or no subjects. In any case, we lack a computationally tractable method for comparing model fits for different datasets. Thus, we must leave this for future research.

We can, however, address a related question. It could be that speakers of different native languages learn English more or less quickly and to a greater or lesser degree. At best, this would add noise to our analyses. At worst, to the extent that native language is confounded with other variables of interest in our sample (e.g., age of first exposure), it could have distorted our results. Anecdotally, many people perceive that speakers of certain languages are better or worse at English, though it is hard to know how much this is confounded with accent (which likely has a critical period distinct from that of syntax), cultural variation in age at first exposure, and differences in the types of exposure (e.g., songs, movies, tourism, coursework) and instructional methods. For instance, in our dataset, speakers of Chinese and Western Germanic languages tended to start learning English in immersion settings earlier than speakers of Turkic or Uralic languages (5.2 and 5.9 years old vs. 13.4 and 14.8 years old, respectively). More systematically, some studies have suggested different patterns of ultimate attainment for speakers of different native languages (Bialystok & Miller, 1999), though caution is warranted given the extremely low power for such studies (see Fig. 8 and surrounding discussion).

We considered the effect of native language on three different metrics of learning success: the level of ultimate attainment (how well the most advanced learners do), the age at the end of the optimal period (the last age to start learning in order to reach native-like performance), and the shape of the learning curve (performance as a function of years of experience). In keeping with our earlier analyses, ultimate attainment was defined as the average performance for subjects no older than 70 years old and with at least 30 years of experience with English. To increase power, we grouped subjects into Uralic, Slavic, West Germanic, Romance, and Chinese language groups (no other language group had nearly as many speakers at similarly wide ranges of years of experience and ages of first exposure). For each measurement, we assessed the level of evidence that speakers of one language group differed from the others using Bayes Factor model comparison with the BIC approximation (Wagenmakers, 2007). Details for all analyses are provided in the Supplementary Materials, under “Item Effects.”

By looking at ultimate attainment, we can assess whether speakers of different languages have greater or lesser success in learning English, equating for years of experience. In fact, the differences across language groups were small (see Fig. S14) and generally not reliable. In most cases, analyses favored the null hypothesis (no difference between the target language and the other languages), and differences across language groups were inconsistent: among learners who began at age 0, the best-performing language group was Romance, for learners beginning at 1–5 years old, it was West Germanic, and for learners who began at 6–10 years old, it was Chinese. Likewise, analysis indicated that the length of the optimal period does not vary across language groups. We

found slightly more evidence for differences in learning curves. In particular, simultaneous English-Chinese speakers could be distinguished from the rest, whereas simultaneous bilinguals who spoke Romance or West Germanic languages both matched the overall pattern. However, the actual differences are subtle and seem to reflect slightly faster initial learning by the Chinese speakers (Fig. S18). Most other comparisons were not possible due to insufficiently many subjects (see Supplementary Materials).

Thus, although speakers of different languages make different mistakes, we find only limited evidence of differences in learning once learning context (immersion vs. non-immersion), years of experience, and age at first exposure are taken into account. That said, power analyses suggest that we only had sufficient subjects to detect relatively large effects, meaning that we cannot rule out more subtle differences (see Supplementary Materials, under “Item Effects”). These power analyses should, however, provide guidance on sample sizes for future research along these lines.

Whatever these analyses say about language-learning in general, they do not provide any evidence that our findings were heavily confounded by differences across the native languages in our sample.

## 4.2. Implications

The analyses above suggest that our findings are reasonably robust, particularly in comparison to those of previous studies. While this inspires confidence, it should also suggest caution: future work that successfully addresses the limitations of the present study may similarly prompt significant revisions in what we believe to be true. Science is the process of becoming less wrong, and while hopefully the revisions are smaller and smaller after each step, there is no way of knowing that this is the case in advance. Thus, confirmation and extension of the present results is crucial, particularly given the novelty of our questions, methods, models, and results.

Nonetheless, we believe it is useful to consider the implications of the present findings, on the presumption that they prove to be (reasonably) robust:

### 4.2.1. The nature of the critical period for second language acquisition

On the assumption that the present results apply broadly to syntax acquisition by diverse learners, they have profound theoretical implications. Most importantly, they clarify the shape of the well-attested critical period for second-language acquisition: a plateau followed by a continuous decline. The end of the plateau period must be due to changes in late adolescence rather than childhood, whether they are biological, social, or environmental. Thus the critical period cannot be attributed to neuronal death or syntactic pruning in the first few years of life, nor to hormonal changes surrounding adrenarche or puberty (Johnson & Newport, 1989; Lenneberg, 1967; Pinker, 1994). Also casting doubt on the effect of hormones is our finding that girls do not show a decline in learning ability before boys do, despite their earlier age of puberty (see Supplementary Materials). Likewise, the critical period cannot be explained by documented developmental changes in working memory, episodic memory, reasoning ability, processing speed, or social cognition (Hakuta et al., 2003; Hartshorne & Germine, 2015; Klindt, Devaine, & Daunizeau, 2017; Morgan-Short & Ullman, 2012; Newport, 1988), to the diminished likelihood that adolescent and adult immigrants will be immersed in an environment of native speakers and identify with the new culture,<sup>5</sup> or to gradually accumulating interference from a first language (Hernandez et al., 2005; Jia et al., 2002; Sebastián-Gallés et al., 2005).

In short, these data are inconsistent with any hypothesis that places

<sup>5</sup> Note that while critical period researchers widely assume that there are age-related effects on cultural identification among immigrant groups, this may not in fact be the case (Chudek, Cheung, & Heine, 2015).

the decline in childhood—which is to say, every prior *specific* hypothesis that we know of. What, then, *could* explain the critical period? There are a number of possibilities. For instance, it remains possible that the critical period is an epiphenomenon of culture: the age we identified (17–18 years old) coincides with a number of social changes, any of which could diminish one’s ability, opportunity, or willingness to learn a new language. In many cultures, this age marks the transition to the workforce or to professional education, which may diminish opportunities to learn. Note that causality (if any) could run the other direction: cultures may have chosen this age for certain transitions because of age-dependent changes in neural plasticity. Further traction on these issues could come from cross-cultural comparison, or comparison of individuals within a culture who are on different educational tracks.

Alternatively, the critical period could reflect interference from the first language, so long as this interference is non-linear rather than gradually accumulating. While it has generally been assumed that interference from the first language would be proportional to the amount of first language learned—something inconsistent with our data—we cannot rule out the possibility of non-linear interference. Neural network models, which are capable of showing interference from a first language (Hernandez et al., 2005), can exhibit surprising nonlinearities (Haykin, 1999; Hernandez et al., 2005). It remains to be seen whether they can successfully model the nonlinearities we actually observed.

Finally, the end of the critical period might reflect late-emerging neural maturation processes that compromise the circuitry responsible for successful language acquisition (whether specific to language or not). While language acquisition researchers often focus on neural development in the childhood years, the brain undergoes significant changes through adolescence and early adulthood (Blakemore & Mills, 2014; Mills, Lalonde, Clasen, Giedd, & Blakemore, 2014; Pinto, Hornby, Jones, & Murphy, 2010; Shafee, Buckner, & Fischl, 2015; Tamnes et al., 2010). While continued development of the prefrontal cortex is perhaps the most familiar, changes occur throughout the brain and along multiple dimensions. Drawing on these and other findings, some researchers have suggested that adolescence may involve a number of different biologically-driven critical periods (Crews, He, & Hodge, 2007; Fuhrmann, Knoll, & Blakemore, 2015; see also Ghitza & Gelman, 2014).

Little is certain about the relationship between neural maturation and behavioral maturation, other than the likelihood it is complex. Current evidence suggests that critical periods in perception involve a complex interplay of neurochemical and epigenetic promoters and brakes for both synaptic pruning and outgrowth (Werker & Hensch, 2015). Given this complexity, and the relative sparseness of the data on neural maturation, it is hard to say whether any of the identified neural maturation processes might correspond to the changes in syntax acquisition that we observed.

Nor can we do much more than speculate as to whether these maturational process (if any) are specific to structures subserving language acquisition. It is notable that language-learning ability is, out of every cognitive ability whose developmental trajectory has been characterized behaviorally, the only one that is stable through childhood and declines sharply in late adolescence (Hartshorne & Germine, 2015). This observation is consistent with the possibility of language-specific maturation. However, the developmental trajectories of some cognitive abilities, such as procedural memory, have not been well characterized (Fuhrmann et al., 2015; Hartshorne & Germine, 2015). Moreover, cognitive testing has largely focused on simple abilities that can be measured in a single, short session (e.g., working memory). In contrast, syntax acquisition takes place over much longer intervals and involves learning a complex, interlocking system. Thus, progress on this question will require characterization of a broader range of cognitive abilities, as well as acquisition of other complex systems (e.g., music or chess).

In attempting to gain traction on these issues, there are additional complexities, which future studies should seek to clarify. The duration

of the critical period may differ for other aspects of language, like phonology and vocabulary. Moreover, we cannot be certain that syntax learning ability is a unitary construct rather than the combination of multiple factors potentially operating on distinct timelines and affecting different aspects of syntax differently. Second, the exact timing of the critical period may be obfuscated by older learners deploying conscious learning strategies, absorbing explicit instruction, or transferring knowledge from the first language. Some purchase on these issues may come from additional studies, potentially using different methods (e.g., online processing, production, ERP, or longitudinal studies), should obtaining sufficiently many subjects become feasible. Finally, because our dataset consists of people’s performance in a second language, it does not directly address the question of how age affects the learning of a first language. It is possible that exposure to linguistic input delays the atrophy of language learning circuitry, in which case the decline in learning ability we have documented would represent the prolongation of a critical period that terminates sooner in people who have been deprived of all language input (Curtiss, 1994; de Villiers, 2007; Mayberry, 1993; Newport, 1990). Because delayed first-language acquisition is fortunately rare, it would be impossible to achieve a sample size similar to the one here, but our results could be used to guide smaller, targeted studies.

Crucially, the investigation of these issues—all of which have long been of interest but difficult to address—can now be guided by the finding that the ability to learn the grammar of a new language, though indeed compromised in adults compared to children, is largely or entirely preserved up to the cusp of adulthood.

#### 4.2.2. Additional implications

The dataset bears on many issues beyond those discussed in detail above. For instance, the data contain a rich source of information about dialect variation and L1 transfer effects. We briefly mention a few other issues. First, prior work has indicated that simultaneous bilinguals do not reach the same level of proficiency in phonology as individuals with a single first language (Sebastián-Gallés et al., 2005). We extend this finding to syntax, where it is apparent throughout the lifespan (Fig. 5B). This finding is consistent with some earlier work suggesting that a sufficiently sensitive test can distinguish even highly proficient bilinguals from monolinguals (Abrahamsson & Hyltenstam, 2008, 2009).<sup>6</sup> Our model captures this difference as one of exposure, estimating that simultaneous bilinguals receive only 63% as much English input as monolinguals (see Fig. S6). Though parsimonious, this is not the only possible explanation; alternatives include the effects of suppression of the non-target language and influences of each language on the other (Birdsong & Gertken, 2013).

Similarly, there are a number of interesting demographic effects. We confirm prior findings of a main effect of education on ultimate attainment, with post-secondary education resulting in higher accuracy (see Supplementary Materials, “Education Differences”) (Birdsong, 2014; Hakuta et al., 2003). We likewise find a main effect for gender, with higher accuracy by females (see Supplementary Materials, “Gender Differences”). In neither case do these main effects appear to interact with age at first exposure, and so they are unlikely to be relevant for critical periods. However, they likely have implications for other aspects of language learning.

We have made the data available (<http://osf.io/pyb8s>) in the hopes they will prove informative for investigation of these and other questions.

<sup>6</sup> This finding also has practical consequences for research. Many researchers have argued that if later learners can reach monolingual levels of performance, that would be evidence against critical periods (and conversely, the failure of later learners to match monolinguals would be evidence for critical periods) (e.g., Abrahamsson & Hyltenstam, 2009). This standard, in conjunction with our results, leads to the unlikely conclusion that the critical period for syntax closes prior to birth. For additional discussion, see Birdsong and Gertken (2013).

## Acknowledgements

We are indebted to David Barner, David Birdsong, Kenji Hakuta, Elissa Newport, Laura-Ann Petitto, and Michael Ullman for comments, to Tanya Ivonchik and Brandon Benson for help with developing the quiz, and to the hundreds of thousands of volunteers who participated in the study. This research was supported by an NIH NRSA award to JKH (5F32HD072748) and the Center for Minds, Brains, & Machines (NSF STC CCF-1231216).

## Contributions

JKH designed the study, collected the data, and performed the analyses. All three authors contributed to designing the analyses and to writing the paper.

## Appendix A. Supplementary material

Supplementary data associated with this article can be found, in the online version, at <http://dx.doi.org/10.1016/j.cognition.2018.04.007>.

## References

- Abrahamsson, N. (2012). Age of onset and nativelike L2 ultimate attainment of morphosyntactic and phonetic intuition. *Studies in Second Language Acquisition*, 34(02), 187–214.
- Abrahamsson, N., & Hyltenstam, K. (2008). The robustness of aptitude effects in near-native second language acquisition. *Studies in Second Language Acquisition*, 30(4), 481–509.
- Abrahamsson, N., & Hyltenstam, K. (2009). Age of onset and nativelikeness in a second language: Listener perception versus linguistic scrutiny. *Language Learning*, 59(2), 249–306.
- Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language*, 59, 390–412.
- Berman, R. A. (Ed.). (2004). *Language development across childhood and adolescence*. Philadelphia, PA: John Benjamins Publishing Company.
- Berman, R. A. (2007). Developing linguistic knowledge and language use across adolescence. In E. Hoff, & M. Shatz (Eds.), *Blackwell handbook of language development* (pp. 347–367). Malden, MA: Blackwell Publishing.
- Bialystok, E., & Miller, B. (1999). The problem of age in second-language acquisition: Influences from language, structure, and task. *Bilingualism: Language and Cognition*, 2(02), 127–145.
- Birdsong, D. (2006). Age and second language acquisition and processing: A selective overview. *Language Learning*, 56, 9–49.
- Birdsong, D. (2014). The critical period hypothesis for second language acquisition: Tailoring the coat of many colors. In M. Pawlak & L. Aronin (Eds.), *Essential topics in applied linguistics and multilingualism. Studies in honor of David Singleton* (pp. 43–50). Berlin and New York: Springer.
- Birdsong, D. (2017). Critical periods. In M. Aronoff (Ed.), *Oxford bibliographies in linguistics*. New York: Oxford University Press.
- Birdsong, D., & Gertken, L. M. (2013). In faint praise of folly: A critical review of native/non-native speaker comparisons, with examples from native and bilingual processing of French complex syntax. *Language, Interaction and Acquisition*, 4(2), 107–133.
- Birdsong, D., & Molis, M. (2001). On the evidence for maturational constraints in second-language acquisition. *Journal of Memory and Language*, 44(2), 235–249.
- Blakemore, S.-J., & Mills, K. L. (2014). Is adolescence a sensitive period for sociocultural processing? *Annual Review of Psychology*, 65, 9.1–9.21.
- Bruer, J. T. (1999). *The myth of the first three years*. New York: Free Press.
- Chudek, M., Cheung, B. Y., & Heine, S. J. (2015). US immigrants' patterns of acculturation are sensitive to their age, language, and cultural contact but show no evidence of a sensitive window for acculturation. *Journal of Cognition and Culture*, 15, 174–190.
- Clark, H. H. (1973). The language-as-fixed-effect fallacy: A critique of language statistics in psychological research. *Journal of Verbal Learning and Verbal Behavior*, 12(4), 335–359.
- Coppieters, R. (1987). Competence differences between native and near-native speakers. *Language*, 544–573.
- Crain, S., & Thornton, R. (2011). Syntax acquisition. *WIREs Cognitive Science*, 3(2), 185–203.
- Crews, F., He, J., & Hodge, C. (2007). Adolescent cortical development: A critical period of vulnerability for addiction. *Pharmacology, Biochemistry, and Behavior*, 86, 189–199.
- Curtiss, S. (1994). Learning as a cognitive system: Its independence and selective vulnerability. In C. P. Otero (Vol. Ed.), *Noam Chomsky: Critical assessments: Vol. 1*, (pp. 227–228). New York, NY: Routledge.
- de Villiers, J. G. (2007). The interface of language and theory of mind. *Lingua*, 117, 1858–1878. <http://dx.doi.org/10.1016/j.lingua.2006.11.006>.
- DeKeyser, R. M. (2000). The robustness of critical period effects in second language acquisition. *Studies in Second Language Acquisition*, 22(04), 499–533.
- DeKeyser, R. M., Alfi-Shabtay, I., & Ravid, D. (2010). Cross-linguistic evidence for the nature of age effects in second language acquisition. *Applied Psycholinguistics*, 31(03), 413–438.
- Efron, B., & Tibshirani, R. (1993). *An introduction to the bootstrap*. Boca Raton, FL: Chapman & Hall/CRC.
- Flege, J. E., Yeni-Komshian, G. H., & Liu, S. (1999). Age constraints on second-language acquisition. *Journal of Memory and Language*, 41(1), 78–104.
- Friedman, J. H. (1991). Multivariate adaptive regression splines. *The Annals of Statistics*, 19(1), 1–141.
- Fuhrmann, D., Knoll, L. J., & Blakemore, S.-J. (2015). Adolescence as a sensitive period of brain development. *Trends in Cognitive Sciences*, 19(10), 558–566.
- Germine, L. T., Duchaine, B., & Nakayama, K. (2011). Where cognitive development and aging meet: face learning ability peaks after age 30. *Cognition*, 118(2), 201–210. <http://dx.doi.org/10.1016/j.cognition.2010.11.002>.
- Ghitza, Y., & Gelman, A. (2014). *The Great Society, Reagan's Revolution, and generations of presidential voting*.
- Granena, G., & Long, M. H. (2013). Age of onset, length of residence, language aptitude, and ultimate L2 attainment in three linguistic domains. *Second Language Research*, 29(3), 311–343.
- Guion, S. G., Flege, J. E., Liu, S. H., & Yeni-Komshian, G. H. (2000). Age of learning effects on the duration of sentences produced in a second language. *Applied Psycholinguistics*, 21(02), 205–228.
- Hakuta, K., Bialystok, E., & Wiley, E. (2003). Critical evidence: A test of the critical-period hypothesis for second-language acquisition. *Psychological Science*, 14(1), 31–38.
- Halberda, J., Ly, R., Wilmer, J. B., Naiman, D. Q., & Germine, L. (2012). Number sense across the lifespan as revealed by a massive Internet-based sample. *Proceedings of the National Academy of Sciences*, 109(28), 11116–11120. <http://dx.doi.org/10.1073/pnas.1200196109>.
- Hartshorne, J. K., & Germine, L. T. (2015). When does cognitive functioning peak? The asynchronous rise and fall of different cognitive abilities across the life span. *Psychological Science*, 26(4), 433–443.
- Haykin, S. (1999). *Neural networks: A comprehensive guide* (2nd ed.). Upper Saddle River, NJ: Prentice Hall.
- Hernandez, A. E., Li, P., & MacWhinney, B. (2005). The emergence of competing modules in bilingualism. *Trends in Cognitive Sciences*, 9(5), 220–225.
- Huang, B. H. (2015). A synthesis of empirical research on the linguistic outcomes of early foreign language instruction. *International Journal of Multilingualism*, 13(3), 257–273. <http://dx.doi.org/10.1080/14790718.2015.1066792>.
- Jaeger, T. F. (2008). Categorical data analysis: Away from ANOVAs (transformation or not) and towards logit mixed models. *Journal of Memory and Language*, 59(4), 434–446. <http://dx.doi.org/10.1016/j.jml.2007.11.007>.
- Jia, G., Aaronson, D., & Wu, Y. (2002). Long-term language attainment of bilingual immigrants: Predictive variables and language group differences. *Applied Psycholinguistics*, 23(04), 599–621.
- Johnson, J. S., & Newport, E. L. (1989). Critical period effects in second language learning: The influence of maturational state on the acquisition of English as a second language. *Cognitive Psychology*, 21(1), 60–99.
- Johnson, J. S., & Newport, E. L. (1991). Critical period effects on universal properties of language: The status of subadjacency in the acquisition of a second language. *Cognition*, 39(3), 215–258.
- Kidd, E., & Bavin, E. L. (2002). English-speaking children's comprehension of relative clauses: Evidence for general-cognitive and language-specific constraints on development. *Journal of Psycholinguistic Research*, 31(6), 599–617.
- Kidd, E., & Lum, J. A. G. (2008). Sex differences in past tense overregularization. *Developmental Science*, 11(6), 882–889.
- Klindt, D., Devaine, M., & Daunizeau, J. (2017). Does the way we read others' mind change over the lifespan? Insights from a massive Web poll of cognitive skills from childhood to late adulthood. *Cortex*, 86, 205–215.
- Krashen, S. D., Long, M. A., & Scarcella, R. C. (1979). Age, rate, and eventual attainment in second language acquisition. *TESOL Quarterly*, 573–582.
- Lenneberg, E. (1967). *Biological foundations of language*. New York: Wiley.
- Marcus, G. F., Pinker, S., Ullman, M. T., Hollander, M., Rosen, T. J., & Xu, F. (1992). Overregularization in language acquisition. *Monographs of the Society for Research in Child Development*, 57(4), 1–182.
- Mayberry, R. I. (1993). First-Language acquisition after childhood differs from second-language acquisition: The case of American sign language. *Journal of Speech, Language, and Hearing Research*, 36(6), 1258–1270.
- Mayberry, R. I., & Lock, E. (2003). Age constraints on first versus second language acquisition: Evidence for linguistic plasticity and epigenesis. *Brain and Language*, 87(3), 369–384.
- Mayberry, R. I., Lock, E., & Kazmi, H. (2002). Development: Linguistic ability and early language exposure. *Nature*, 417(6884), 38.
- McDonald, J. L. (2000). Grammaticality judgments in a second language: Influences of age of acquisition and native language. *Applied Psycholinguistics*, 21(03), 395–423.
- Messenger, K., Branigan, H. P., McLean, J. F., & Sorace, A. (2012). Is young children's passive syntax semantically constrained? Evidence from syntactic priming. *Journal of Memory and Language*, 66, 568–587.
- Milborrow, S. (2014). *Earth: Multivariate adaptive regression spline models*. R package version 3.2-7. < <http://cran.r-project.org/web/packages/earth/index.html> > .
- Mills, K. L., Lalonde, F., Clasen, L. S., Giedd, J. N., & Blakemore, S.-J. (2014). Developmental changes in the structure of the social brain in late childhood and adolescence. *Social Cognitive Affective Neuroscience*, 9(1), 123–131.
- Morgan-Short, K., & Ullman, M. T. (2012). The neurocognition of second language. In A. Mackey, & S. Gass (Eds.), *Handbook of second language acquisition* (pp. 1–18). Routledge.
- Mullen, K., Aridia, D., Gil, D., Windover, D., & Cline, J. (2011). DEoptim: An R package for global optimization by differential evolution. *Journal of Statistical Software*, 40(6),

- 1–26.
- Newport, E. L. (1988). Constraints on learning and their role in language acquisition: Studies of the acquisition of American Sign Language. *Language Sciences*, 10(1), 147–172.
- Newport, E. L. (1990). Maturation constraints on language learning. *Cognitive Science*, 14(1), 11–28.
- Nippold, M. A. (2007). *Later language development: School-age children, adolescents, and young adults* (3rd ed.). Austin, TX: Pro-Ed.
- Patkowski, M. S. (1980). The sensitive period for the acquisition of syntax in a secondary language. *Language Learning*, 30(2), 449–468.
- Pinker, S. (1994). *The language instinct*. New York: William Morrow.
- Pinker, S. (1999). *Words and rules: The ingredients of language*. New York, NY: HarperCollins.
- Pinto, J. G. A., Hornby, K. R., Jones, D. G., & Murphy, K. M. (2010). Developmental changes in GABAergic mechanisms in human visual cortex across the lifespan. *Frontiers in Cellular Neuroscience*, 4(16), <http://dx.doi.org/10.3389/fncel.2010.00016>.
- Pulvermüller, F., & Schumann, J. H. (1994). Neurobiological mechanisms of language acquisition. *Language Learning*, 44, 681–734. <http://dx.doi.org/10.1111/j.1467-1770.1994.tb00635.x>.
- Qureshi, M. A. (2016). A meta-analysis: Age and second language grammar acquisition. *System*, 60, 147–160. <http://dx.doi.org/10.1016/j.system.2016.06.001>.
- Rowland, C. F., & Pine, J. M. (2000). Subject-auxiliary inversion errors and wh-question acquisition: 'what children do know?'. *Journal of Child Language*, 27(1), 157–181.
- Schachter, J. (1990). On the issue of completeness in second language acquisition. *Second Language Research*, 6(2), 93–124. <http://dx.doi.org/10.1177/026765839000600201>.
- Sebastián-Gallés, N., Echeverría, S., & Bosch, L. (2005). The influence of initial exposure on lexical representation: Comparing early and simultaneous bilinguals. *Journal of Memory and Language*, 52(2), 240–255.
- Shafee, R., Buckner, R. L., & Fischl, B. (2015). Gray matter myelination of 1555 human brains using partial volume corrected MRI images. *NeuroImage*, 105, 473–485.
- Snow, C. E., & Hoefnagel-Höhle, M. (1978). The critical period for language acquisition: Evidence from second language learning. *Child Development*, 1114–1128.
- Stevens, G. (1999). Age at immigration and second language proficiency among foreign-born adults. *Language in Society*, 28(04), 555–578.
- Stewart, N., Ungemach, C., Harris, A. J., Bartels, D. M., Newell, B. R., Paolacci, G., & Chandler, J. (2015). The average laboratory samples a population of 7,300 Amazon Mechanical Turk workers. *Judgment and Decision Making*, 10(5), 479–491.
- Tamnes, C. K., Ostby, Y., Fjell, A. M., Westlye, L. T., Due-Tønnessen, P., & Walhovd, K. B. (2010). Brain maturation in adolescence and young adulthood: Regional age-related changes in cortical thickness and white matter volume and microstructure. *Cerebral Cortex*, 20, 534–548.
- Vanhove, J. (2013). The critical period hypothesis in second language acquisition: A statistical critique and a reanalysis. *PLoS ONE*, 8(7), e69172. <http://dx.doi.org/10.1371/journal.pone.0069172.s003>.
- Wagenmakers, E.-J. (2007). A practical solution to the pervasive problems of p values. *Psychonomic Bulletin & Review*, 14(5), 779–804.
- Weber-Fox, C., & Neville, H. (1996). Maturation constraints on functional specializations for language processing: ERP and behavioral evidence in bilingual speakers. *Journal of Cognitive Neuroscience*, 8(3), 231–256.
- Werker, J. F., & Hensch, T. (2015). Critical periods in speech perception: New directions. *Annual Review of Psychology*, 66, 173–196.