

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

Acoustic correlates and listener ratings of function word reduction

in child versus adult speech

Melissa A. Redford

The University of Oregon

and

Phil J. Howson

Leibniz-Zentrum Allgemeine Sprachwissenschaft

Running title:

Acoustic correlates and listener ratings of reduction

Please address correspondence to:

Melissa A. Redford, Professor; Department of Linguistics | 1290 University of Oregon | Eugene, OR
97403; redford@uoregon.edu.

1 **Abstract**

2 The present study investigated “the” reduction in phrase-medial Verb-*the*-Noun sequences elicited from
3 five-year-old children and college-aged adults. Several measures of reduction were calculated based on
4 acoustic measurement of these sequences. Analyses on the measures indicated that the determiner vowel
5 was reduced in both child and adult speech relative to content word vowels, but it was reduced less in
6 children’s speech compared to adults’ speech. Listener ratings on the sequences indicated a preference for
7 adult speech over children's speech. Acoustic measures of reduction also predicted goodness ratings.
8 Listeners preferred sequences with shorter and lower amplitude determiner vowels relative to content
9 word vowels. They also preferred a more neutral schwa over more coarticulated versions. In sequences
10 where ratings differed by age group, the effect of coarticulation was limited to adult speech and the effect
11 of relative schwa duration was limited to children’s speech. The results are discussed with reference to
12 communicative pressures on function word production in speech, including the rhythmic and semantic
13 pressure towards reduction versus the pressure to convey adequate information in the acoustic signal. It is
14 argued that competing pressures on production may delay the acquisition of adult-like function word
15 reduction.

16
17 *Keywords:* speech production; speech acquisition; speech rhythm; speech motor skills

1 I. INTRODUCTION

2 Function words, like determiners, refine the message and help to define the grammatical structure
3 of a sentence. As high frequency items with minimal semantic weight, function words are typically
4 unstressed and phonetically reduced relative to lower frequency content words with maximal semantic
5 weight (Bell et al., 2009). The 10 most frequent function words in English are monosyllabic (Jurafsky et
6 al., 1998). These alternate with lexically-stressed content words, many of which are also monosyllabic
7 (Cutler and Carter, 1987). In spoken language, the recurrent alternation of unstressed function words with
8 stressed content words, though not itself periodic, contributes substantially to the rhythm of English
9 (Allen and Hawkins, 1978; Dauer, 1983; Deterding, 2001). Given the importance of lexical stress and
10 rhythm to speech processing (e.g., Cutler and Butterfield, 1992; Mattys et al., 2005; Dilley and McAuley,
11 2008), it is reasonable to think that listeners come to expect function word reduction along the temporal
12 and amplitude dimensions that define speech rhythm (Grabe and Low, 2002; He, 2012; Tilsen and
13 Arvaniti, 2013). But reduction along these dimensions is also typically associated with greater
14 coarticulation (i.e., greater gestural “overlap” or hypo-articulated speech; Agwuele et al., 2008; Moon and
15 Lindblom, 1994). When coarticulation is extreme, vowel quality may be impacted to the point of
16 distorting the phonetic shape of the determiner vowel, rendering it more difficult to process. And, of
17 course, reduction-related decreases along the temporal and amplitude dimensions can also render function
18 words inaudible to the listener (see, e.g., Dilley and Pitt, 2010). These observations suggest that the
19 rhythmic requirements of speech production may compete with the functional pressure to produce
20 intelligible speech. This competition may complicate the acquisition of function word reduction during
21 spoken language development. It is an interest in children’s acquisition of function word reduction that
22 motivates the present study, which investigates the relationship between the acoustic correlates of “the”
23 reduction and adult listener ratings of speech rhythmicity.

24 A. The developmental context

25 In adult speech, reduced unstressed syllables are shorter, quieter, and more coarticulated with
26 adjacent speech sounds than unreduced stressed syllables (Dauer, 1983; Fourakis, 1991; Fowler, 1981;

1 Plag et al., 2011). This is especially true for monosyllabic function words, which are typically more
2 reduced than unstressed syllables in content words (Fuchs, 2016; van Bergem, 1993). Although children
3 acquire the overall temporal and amplitude patterns associated with lexical stress by age 2 or 3 years
4 (Ballard et al., 2012; Kehoe et al., 1995; Schwartz et al., 1996), they elide unstressed function words in
5 extra-metrical positions in early speech (see, e.g., Gerken, 1996) and do not reduce these in running
6 speech to the same extent as adults until sometime in middle childhood (Allen and Hawkins, 1978;
7 Nittrouer, 1993; Goffman, 2004; Redford, 2018). For example, Nittrouer (1993) found that the vowel of
8 the indefinite determiner “a” was longer in speech produced by 3-, 5-, and 7-year-old children than in
9 adult speech, but that their stressed vowel durations were similar to adults’ stressed vowel durations.
10 Redford (2018) found that, in comparison to adults, 5- and 8-year-old children produced longer vowels in
11 the definite determiner “the” relative to adjacent content word vowels. In addition, she found that the
12 determiner vowel was louder relative to the content word vowel in children’s speech compared to adults’
13 speech, but that determiner vowel formant frequencies varied as a function of the following content word
14 onset in both child and adult speech. One aim of the present study is to confirm these limited findings by
15 investigating “the” reduction in speech elicited from different English-speaking 5-year-old children and
16 college-aged adults.

17 Function word reduction, or lack thereof, likely impacts speech rhythm. Interval-based studies of
18 rhythm indicate that English-speaking children’s speech is both more vocalic overall than adults’ speech
19 until late middle childhood and also that the vocalic intervals in children’s speech vary less in duration
20 than in adults’ speech (Payne et al., 2011; Polyanskaya and Ordin, 2015). Whereas the greater vocalicness
21 of children’s speech might be explained by age-related differences in articulation rate (see below),
22 reduced variability in vowel durations suggests a pattern that speech-language pathologists might
23 characterize as “excessive equal stress.” Such a pattern impedes speech intelligibility (e.g., Shriberg et al.,
24 2003). Since the relative duration and amplitude patterns typical of lexical stress is mastered early
25 (Ballard et al., 2012), the reduced variability of older children’s speech requires another explanation.
26 Hawkins and Allen (1978) long ago suggested that the explanation might be incomplete function word

1 reduction. Sirsa and Redford (2011) provided some support for this explanation when they found that
2 interval-based measures of school-aged children’s speech rhythm were better predicted by a ratio of
3 determiner vowel duration to a subsequent noun vowel duration (an acoustic measure of function word
4 reduction) than by a ratio of unstressed vowel duration to stressed vowel duration within a disyllabic
5 word (an acoustic measure of lexical stress patterning) or by a ratio of final vowel duration to the mean of
6 non-final vowel durations (an acoustic measure of final lengthening).

7 Of course, school-aged children’s speech differs from adults’ speech along a variety of
8 dimensions that intersect with function word reduction and, by extension, with speech rhythm. The most
9 obvious of these is articulation rate, which is slower in children’s speech than in adults’ speech until at
10 least age 12 years (Lee et al., 1999). Whereas rate changes in adults’ speech are associated with targeted
11 changes in stressed vowel production (Gay, 1981), children’s distinctively slower articulation rate is due
12 to slower articulatory movements into and out of all segmental targets, including unstressed vowels
13 (Redford, 2014). These slower movements are part of an overall pattern of articulation that is associated
14 with immature motor skills, including larger amplitude articulatory movements relative to oral-facial size
15 (Riely and Smith, 2003) and greater spatial-temporal variability (Smith and Zelaznik, 2004). Since these
16 patterns would seem to conspire against the production of very short, quiet, coarticulated unstressed
17 vowels, it is likely that age-related differences in function word reduction are due to immature speech
18 motor skills, not to immature prosodic representations. This conclusion is consistent with results from
19 kinematic and acoustic studies on the production of different metrical structures in child versus adult
20 speech (e.g., Goffman, 2004; Redford, 2018): both children and adults produce trochaic and iambic
21 patterns; the patterns produced by children are simply less contrastive than those produced by adults. It is
22 also consistent with the evidence from cross-linguistic studies on the acquisition of prosody, which
23 suggests the acquisition of language-specific rhythmic structures and intonational patterns by age 3 years
24 (see Fikkert, Liu, and Ota, 2021) – long before the adult-like production of these patterns.

25 But even if the acquisition of adult-like function word reduction is “merely” a product of speech
26 motor development, it must still be learned. This learning begins with patterns extracted from the ambient

1 language. And findings from studies of early child language strongly suggest that this happens early. In
2 particular, developmental studies of speech processing indicate that English-speaking toddlers make use
3 of function words to identify noun-picture correspondences (Kedar et al., 2006); soon thereafter, they
4 regularly produce these words in syntactically correct sentences (~ age 3 years; Abu-Akel et al., 2004).
5 The question is: Once children use function words correctly, how do they know that their speech is still
6 not phonetically accurate? Relatedly, what motivates children to adjust their production of function words
7 across developmental time until they achieve adult-like patterns of reduction? If the answer to the second
8 question is that children strive to communicate with others (adults included), then the answer to the first is
9 that they are not always successful in doing so. Specifically, both the elision of function words and their
10 too-fulsome production may impede communication because it disrupts speech rhythm. Communication
11 failure motivates the child to try again, leading to the adjustments that characterize speech motor learning
12 and the acquisition of adult-like speech patterns.

13 **B. The present study**

14 The hypothesis that competing pressures on function word production delays the acquisition of
15 adult-like reduction predicts that children do not reduce function words to the same extent as adults. The
16 hypothesis that communicative pressures help shape children's speech motor learning, including learning
17 that underpins the acquisition of function word reduction, predicts that adults prefer adult-like speech
18 rhythm patterns over children's speech rhythm patterns. This prediction is in line with the evidence that
19 adult listeners prefer more intelligible children's speech than less intelligible children's speech (e.g.,
20 Redford et al., 2018). In the context of the current study, the more specific prediction is that adult listeners
21 will prefer adult-like reduction of function words over incomplete reduction of these words. These
22 predictions motivate the current study. Here, we sought to identify the acoustic correlates that best
23 distinguished child from adult productions of a determiner, and then asked whether these same correlates
24 account for adult listener ratings of speech rhythmicity on sequences that contained the determiner.

1 II. EXPERIMENT 1

2 The aims of Experiment 1 were to confirm previous study findings on age-dependent differences
3 in function word reduction and to identify those acoustic correlates of reduction that are most salient to
4 adult listeners and so are most likely to influence the predicted preference for adult speech rhythm over
5 children’s speech rhythm. Simple Subject-Verb-*the*-Noun sentences were elicited from a group of 5-year-
6 old children and college-aged adults. Both the verb and the object noun were monosyllabic. The sentence
7 frame included a final “today” in order to avoid phrase-final lengthening effects on the object noun. The
8 consonantal context on either side of “the” was controlled. The stressed vowels in the verb and noun were
9 manipulated to investigate vowel-to-vowel coarticulatory effects on the unstressed determiner vowel.
10 Several measures of reduction were calculated from the acoustic data and analyzed for an effect of age.
11 These were the duration of schwa divided by the duration of adjacent content word vowels (i.e., relative
12 duration), the amplitude of schwa divided by the amplitude of adjacent content word vowels (i.e., relative
13 amplitude), and the effect of vowel context on schwa formant frequencies (i.e., coarticulation). Based on
14 previous findings, the predictions were that children would produce relatively longer and higher
15 amplitude determiner vowels compared to adults, but that the unstressed vowels would be similarly
16 coarticulated with adjacent vowels in children’s and adults’ speech.

17 A. Methods

18 1. *Participants*

19 Participants were 12 school-aged children (7 female) and 12 college-aged adults (6 female) drawn
20 from a larger project on speech rhythm acquisition. Children ranged in age from 5;3 to 6;2 with a mean
21 age of 5;8 (SD = 3 months). They were recruited via a database built and maintained by a group of
22 developmental labs at the University of Oregon and from summer camps run by the YMCA in Eugene,
23 Oregon. Typical development in children was determined based on parental report and on an in-
24 laboratory assessment of speech-language skills using the Diagnostic Evaluation of Articulation and
25 Phonology (DEAP; Dodd et al., 2002) and the Clinical Evaluation of Language Fundamentals (CELF-5;
26 Wiig et al., 2013). Inclusion criteria were standardized scores within 1 standard deviation of the mean on

1 both the DEAP and the CELF-5. Further selection criteria were based on age (the larger sample included
2 8-year-olds), the order in which children participated in the study (recruitment was on-going), and the
3 quality of the video recording (not relevant to the present study). The college-aged adults were recruited
4 by word-of-mouth from the University of Oregon student body. None of the adult participants reported a
5 history of speech-language therapy. All participants, including adults, completed and passed a pure-tone
6 hearing screen (1000, 2000, and 4000 Hz at 25 dB). All participants were financially compensated for
7 their time. Children also earned a small prize at the end of the study session.

8 **2. *Speech Elicitation***

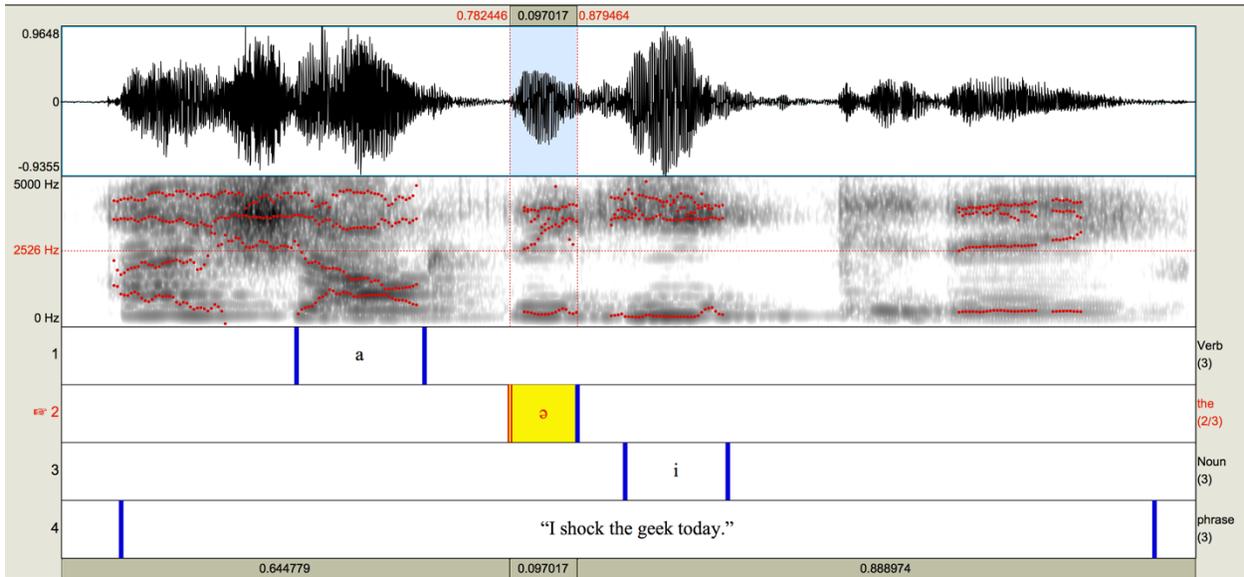
9 Speech materials were designed to investigate carryover and anticipatory V-to-V effects on the
10 production of “the” in simple Subject-Verb-Object sentences, where the verb and object noun were varied
11 to create different stressed vowel environments for the determiner. A total of 16 sentences were created.
12 Here, we focus on a subset of 4 sentences that were designed to investigate the effect of vowel height and
13 backness on schwa production using the vowels [a] and [i]. The verbs were “shock” and “tweak” and the
14 nouns were “god” and “geek.” The target sentences, which had a first person singular subject and the
15 object noun in penultimate position, were as follows: *I shock the god today; I shock the geek today; I*
16 *tweak the god today; I tweak the geek today.* An adult female speaker of west coast American English
17 recorded the target sentences with a second person singular subject (e.g., *You shock the geek today*),
18 followed by the question: *What do you do today?* Care was taken to produce the model sentence under a
19 single intonational contour, followed by a clear prosodic break before the question.

20 Participants were introduced to the object nouns with different cartoon pictures; namely, a
21 studious looking boy for *geek* and an unfamiliar mythological deity for *god*. Verbs were associated with
22 different hand gestures (a flicking gesture for *tweak* and a whole-hand expansion gesture for *shock*),
23 which the experimenter deployed next to the cartoon picture during elicitation when each stimulus
24 sentence was played. The participant’s task was to repeat back the model sentence in first person after the
25 question prompt. The experimenter controlled the pace of sentence elicitation and provided feedback on
26 production during practice. If a repetition was deemed errorful or disfluent, the experimenter elicited a

1 new repetition by replaying the stimuli at the end of a block. The sentences were elicited in random order
 2 four times, once per repetition block. Speech was audio-visually (AV) recorded with a Panasonic AJ-
 3 PX270. Audio for the video was recorded with Shure SM81 Condenser at 44100 Hz.

4 **3. Acoustic Segmentation and Measurement**

5 The first three good repetitions of each sentence from every speaker were selected for
 6 measurement, resulting in a total of 288 sentences for analysis. A good repetition was defined as a fluent
 7 utterance that the speaker produced while looking face-on at the video camera. The AV recording of each
 8 sentence was isolated and saved. Audio was stripped from the files and displayed for vowel segmentation
 9 as an oscillogram and a spectrogram in Praat (Boersma and Weenink, 2019). Utterance boundaries were
 10 identified by the onset and offset of acoustic energy in the sentence and its duration measured. The verb,
 11 determiner, and noun vowel were then segmented and durations measured based on repeated listening,
 12 visible periodicity and abrupt changes in the oscillogram and/or the presence of formant structure.



13
 14 *Figure 1. Example of vowel segmentations for a Verb-the-Noun sequence embedded in a simple S-V-O*
 15 *sentence elicited from a 5-year-old child.*

16 Figure 1 illustrates the segmentation criteria for the sentence, *I shock the geek today*, produced by
 17 a 5-year-old boy. As in Figure 1, stressed vowels in the monosyllabic verb and noun were typically

1 produced in such a way that the expected visible cues were robustly present. Determiner vowels were also
2 easily identified based on some subset of the cues. The reliability of the authors' segmentation was
3 assessed on 25% of the data. A research assistant, blind to the purpose of the experiment, was asked to
4 segment 3 sentences chosen at random from each speaker according to the above criteria. Interval
5 durations from the new segmentations were correlated with those based on the original segmentations on
6 the same data. The bivariate correlation indicated very high inter-rater reliability, $r(216) = 0.922$.

7 Amplitude (dB) and formant frequency were measured at 10 evenly-spaced intervals across the
8 entire schwa duration, and at 5 evenly-spaced intervals in the latter half and first half of the verb and noun
9 vowel, respectively. The spectrogram settings were as follows: the view range was set from 0 to 7000 Hz;
10 the window length was 0.005 s; the dynamic range was 50 dB. The standard pre-emphasis view setting of
11 6 dB per octave was used. Amplitude used the standard Praat settings including a minimum pitch setting
12 of 75 Hz, a time step that was one quarter the spectrogram window length divided by the minimum pitch,
13 and a cubic interpolation method. If the automatic track was visibly inaccurate relative to the
14 spectrogram, tracking was adjusted by changing the number of formants tracked or the range in Hz used
15 for picking peaks. The solution for improving tracking typically differed by vowel type: for example, the
16 number of formants increased for low vowels if F1 and F2 were not separately tracked; the range in Hz
17 increased for high vowels if F2 was poorly tracked.

18 **4. Analyses**

19 Articulation rate was calculated for each target sentence (syll/s). Sentences with pauses were
20 excluded from this calculation ($N=36$). Relative schwa duration and amplitude was calculated by dividing
21 the determiner vowel duration/amplitude by the duration/amplitude of the verb (Det:V) or by the
22 duration/amplitude of the noun (Det:N). Lower ratios signaled greater reduction than higher ratio.
23 Although all participants generally produced fluent sentences, a few sentences were produced with a
24 pause between the verb and object noun phrase ($N=24$) or between the object noun phrase and the
25 adverbial ($N=2$). Since a prosodic break is associated with pre-final lengthening, Det:V and Det:N vowel
26 durations were calculated only when the content word was not also at a prosodic boundary. All vowels

1 from one hyperarticulated sentence, where all elements were uttered under narrow prosodic focus, was
2 also excluded from the analyses.

3 A linear mixed effects model tested for the fixed effects of age group (Group, 2 levels: child,
4 adult) and vowel context (Context, 4 levels: [a]_[a], [a]_[i], [i]_[a], [i]_[i]) on the temporal and
5 amplitude measures of reduction. The model was built using the lme4 package (Bates, Mächler, Bolker,
6 and Walker, 2015). Speaker and repetition were included as random effects. Repetition was removed
7 when shown to have no significant effect on the results. The lmerTest package (Kuznetsova et al., 2017)
8 was used to estimate the degrees of freedom with Satterthwaite's method (Satterthwaite, 1946). Box and
9 whisker plots present all data used in the analyses. The whiskers represent 1.5 times the interquartile
10 range. The potential outliers (circles) and extreme values (stars) that are shown in the plots were not
11 excluded from the analyses.

12 Schwa coarticulation was assessed based on formant frequencies in two vowel contexts – the
13 [a]_[i] and [i]_[a], contexts. These contexts were chosen to test both the effect the stressed vowels on
14 determiner production and the direction of this effect. Formant frequencies were analysed using a
15 smoothing spline (SS)ANOVA (Gu, 2014) and an R script (Mielke, 2015; R Core Team, 2019). For each
16 of the first three formants, best-fit frequency trajectories and 95% confidence intervals were calculated.
17 Though not presented here, the first three formants of the stressed vowels were analysed in the same way
18 and found to conform to expectations consistent with the low vowel versus high vowel target.

19 **B. Results**

20 As expected, children spoke more slowly than adults: mean articulation rate was 3.22 syllables
21 per second in child speech ($SD = 0.49$ syll) and 4.37 syllables per second in adult speech ($SD = 0.63$ syll).
22 The absolute duration of the determiner vowel was therefore longer in child speech [$M(144) = 115$ ms,
23 $SD = 26$ ms] compared to adult speech ($M(144) = 74$ ms, $SD = 18$ ms). More importantly, the mixed
24 effects models indicated that both children and adults produced a longer, higher amplitude schwa relative
25 to adjacent stressed vowels, though the effect of age group on temporal reduction interacted with the
26 particular Verb-*the*-Noun sequence produced. Specifically, the data in Figure 2 show that schwa duration

1 was generally longer in child speech after [a] but not after [i]. This was true whether relative duration was
2 measured in relation to the verb or noun. Accordingly, there was a significant interaction between age
3 group and vowel context on both Det:V duration [$F(3,257) = 2.73, p = 0.044$] and on Det:N duration
4 [$F(3,257) = 2.88, p = 0.037$]. The simple effect of context was also significant on both Det:V duration
5 [$F(3,257) = 12.64, p < 0.001$] and Det:N duration [$F(3,257) = 72.26, p < .001$]. The simple effect of age
6 group was not significant on either measure.

7 The main effect of vowel context was also significant on relative schwa amplitude [Det:V
8 amplitude: $F(3,257) = 4.51, p = 0.004$; Det:N amplitude: $F(3,257) = 26.28, p < 0.001$]. But the data in
9 Figure 3 show that children produced a higher amplitude schwa compared to adults when this measure
10 was calculated in relation to the verb [$F(1,22) = 12.28, p = 0.002$]. The interaction between age group and
11 vowel context was not significant no matter how relative amplitude was calculated. The effect of age
12 group on Det:V and not on Det:N amplitude could reflect an effect of phrase position on child speech or
13 the more consistent modulation of amplitude across the phrase in adult speech.

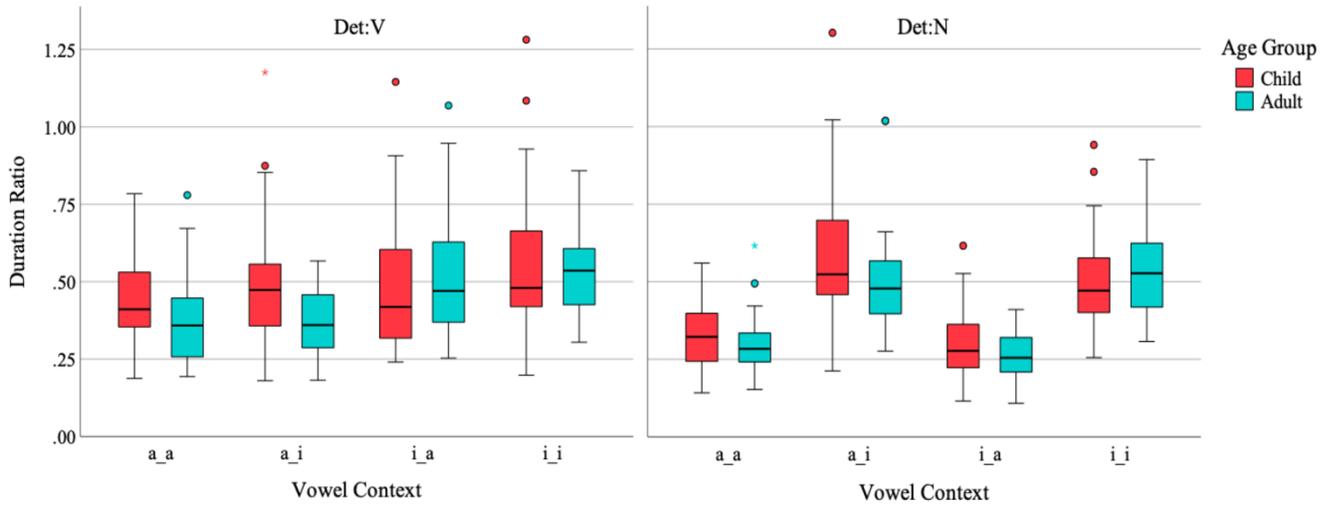
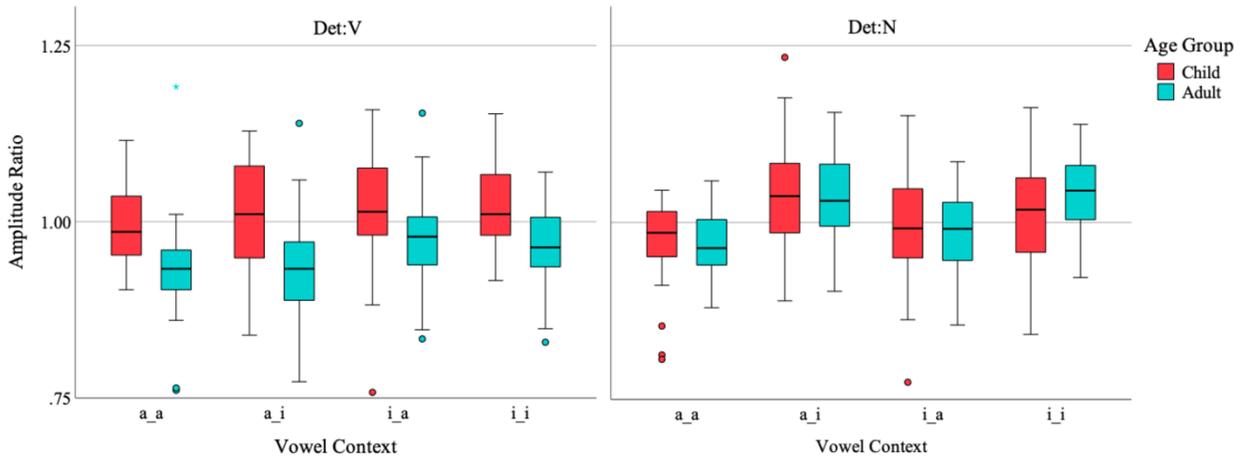


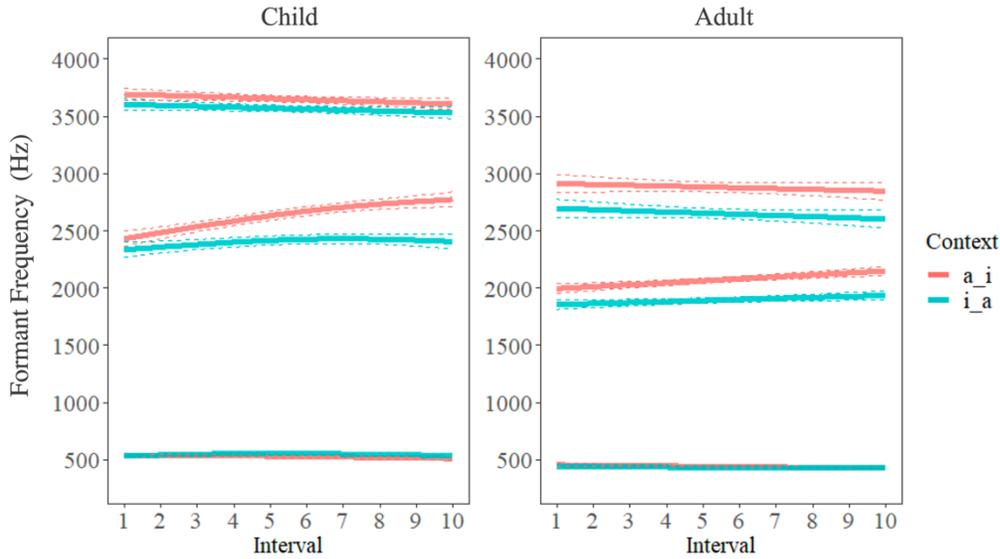
Figure 2. Vowel duration ratios for *Det:V* (left) and *Det:N* (right) shown by age group and the vowel context within which the determiner appeared.



1
2 Figure 3. Vowel amplitude ratios for *Det:V* (left) and *Det:N* (right) are shown by age group and the
3 vowel context within which the determiner appeared.

4 The effect of age group on function word reduction extended to coarticulation. Figure 4 shows
5 that both children and adults produced schwa differently as a function of the upcoming stressed vowel,
6 but not the preceding one. In child speech, the largest effect of vowel context is seen on F2, which was
7 higher before [i] than before [a]. In adult speech, the data indicate an additional effect of context on F3,
8 which was also higher before [i] compared to [a]. These results suggest that children move the

1 constriction location forward during schwa ahead of [i]; adults may also front the tongue body more
 2 before [i] than before [a], but changes in F3 also suggest greater lip spreading and a higher overall
 3 constriction degree (= larger subapical aperture) in this context (Lindblom et al., 2011).



4
 5 *Figure 4. Best-fit curves for schwa F1, F2, and F3 are shown as a function of the stressed vowel context*
 6 *within which the determiner was produced in child (left) and adult (right) speech. Dotted lines represent*
 7 *the 95% confidence interval.*

8 C. Discussion

9 The results from Experiment 1 indicate that when group differences are observed, relative schwa
 10 duration and amplitude is higher in child speech compared to adult speech. This result replicates previous
 11 findings, including from studies where the elicited speech was less well controlled and so also less
 12 contrived (e.g., Sirsa and Redford, 2011). The results also indicate that both children and adults
 13 coarticulated the determiner vowel with the following noun vowel, but how this was done differed by age.
 14 Whereas the results from child speech suggest an adjustment along the front-back dimension (i.e., F2)
 15 during schwa articulation when it was produced in advance of [i], the results from adult speech suggest
 16 that a single posture was adopted across the entire schwa duration that adjusted for both constriction
 17 location as well as the shape of the front cavity (i.e., F3) in advance of [i]. These age-related differences

1 in vowel-to-vowel coarticulation may be related to children’s overall slower articulation rate and longer
2 determiner vowels (see, e.g., Agwuele et al., 2008; Moon and Lindblom, 1994) or to differences in the
3 speech plan representation (see General Discussion). Either way, slower movements into and out of a
4 vowel target cannot explain why the *relative* duration and amplitude of schwa varied with age. Instead,
5 the Det:V and Det:N results are consistent with the interpretation that 5-year-old children do not reduce
6 function words to the same extent as adults, at least along the temporal and amplitude domains most
7 closely associated with speech rhythm.

8 **III. EXPERIMENT 2**

9 In Experiment 2, we investigate whether age-related differences in function word reduction
10 influence adult listeners’ ratings of speech rhythmicity. Verb-*the*-Noun sequences were excised from the
11 subset of sentences with different vowel contexts. The sequences were blocked by speaker and age group
12 and presented to listeners, who were instructed to rate the rhythmic quality of the sequence on a goodness
13 scale. The blocked design was used to encourage listeners to attend to within speaker variability in
14 production rather than to the many acoustic characteristics that distinguish individual speakers from one
15 another and children from adults. Despite this encouragement, we expected listeners to rate children’s
16 speech as less good overall than adults’ speech – either because children do not reduce function words to
17 the same extent as adults or because listeners make global intelligibility judgments even when instructed
18 to attend to speech rhythm. We also expected that the acoustic correlates of function word reduction
19 would account for listener ratings of rhythmicity independently from a preference for adult speech,
20 assuming that listeners are indeed as sensitive to speech rhythm as the psycholinguistic literature would
21 suggest. We were particularly interested in the character of this sensitivity: Do certain correlates of
22 function word reduction matter more to listeners than others? Do the same correlates that influence
23 ratings of child speech also account for ratings of adult speech? Under the hypothesis that communicative
24 pressures help shape children’s speech motor learning, including the learning that underpins the
25 acquisition of function word reduction, we expected that the answer to both questions would be “yes” and
26 that the specifics of this answer would follow from the results of Experiment 1. In particular, we expected

1 that measures of reduction that systematically varied with age group (especially, Det:V duration and
2 Det:V amplitude; see Experiment 1) would best predict listener ratings of goodness independently from a
3 preference for adult speech over child speech.

4 **A. Methods**

5 **1. Participants**

6 A total of 100 adult listeners participated in Experiment 2. Their mean age was 35.5 years ($SD =$
7 10.7 years); 81 self-identified as female. Listeners were recruited using Amazon's Mechanical Turk
8 crowdsourcing platform (MTurk; Buhrmester et al., 2011). Recruitment was limited to self-reported
9 native speakers of English residing in either the United States or Canada. It was further restricted to just
10 those MTurkers who had previously completed a minimum of 5000 HITs ("Human Intelligence Tasks")
11 with an acceptance rate of at least 95%. Each listener was compensated for their time upon completing the
12 task.

13 **2. Stimuli**

14 The sentences that were elicited from children and adults in Experiment 1 were used to
15 investigate the effect of age group and determiner vowel reduction on listener preference. Stimuli were
16 those sentences with maximally contrastive vowel contexts (*I shock the geek today* and *I tweak the god*
17 *today*). The Verb-*the*-Noun sequence from these sentences was excised and amplitude normalized to
18 70dB. This resulted in 72 sequences for the [a]_[i] and [i]_[a] vowel context (= 1 sentence \times 3 elicitations
19 \times 24 speakers).

20 **3. Procedure**

21 One group of 50 adult listeners rated the goodness of all [a]_[i] stimuli, another group rated the
22 goodness of all [i]_[a] stimuli. Each stimulus was presented 3 times to the listener for a total of 216
23 stimuli (= 24 speakers \times 3 elicitations \times 3 repetitions). Stimuli were blocked by speaker, and speaker was
24 blocked by age group. Stimuli within speaker and speaker within age block were randomized for each
25 listener, as was the order of the age blocks. Listeners were instructed to wear headphones set to a
26 comfortable listening volume based on a preliminary task, which was to listen to and then type 3 different

1 words. They were then instructed on the main task. The instructions, given in full below, sought to draw
2 listeners' attention to the rhythmic aspects of the sequence and thus away from other features that are
3 specific to child versus adult speech (e.g., pitch):

4 *This experiment is broken into 2 parts. You will be prompted at the start of Part 1 to begin, and*
5 *then again at the start of Part 2. Feel free to take a break between parts if you need one. In both*
6 *parts of the experiment, you will hear the same 3-word stretch of speech produced by different*
7 *talkers. The sequence of words has been taken out of a single sentence context. Your task is to*
8 *rate the rhythmic quality of the sequence. On each trial, a cross will appear on the screen and*
9 *then disappear. Then you will hear audio play. After the audio plays, a scale will appear with*
10 *numbers from 1 to 7. When the scale appears, rate the rhythmic quality of the 3-word sequence*
11 *from 1 to 7: 1 = Sounds Weird. 7 = Sounds Great. Please press the keyboard button*
12 *corresponding to your rating. Before the study begins, there will be practice trials to familiarize*
13 *you with the task.*

14 The instructions were presented on the screen with breaks between thoughts and between steps
15 (e.g., “Then you will hear audio play.” ¶ “After the audio plays, a scale will appear...” ¶ “When the scale
16 appears...”). Listeners were given a few practice trials with the task-specific manipulation using a sham
17 sequence. These trials were not meant to teach listeners about rhythm; they were meant to familiarize
18 them with the task. After the practice trials, listeners were presented with the experimental stimuli. The
19 task took an average of 18 minutes to complete (+/- 5.7 minutes).

20 **4. Analyses**

21 Ratings that were provided too quickly (< 300 ms) or too slowly (> 2400 ms) were excluded
22 from the analyses (= 14% of the data), following best practices (see Ratcliff, 1993). Ratings were then
23 averaged within listener across repetitions to generate a single score for each unique stimulus that the
24 listener heard. This procedure resulted in a total of 3,600 ratings per vowel context (50 listeners x 24
25 speakers x 3 unique elicitations per speaker). In a first set of analyses, a linear mixed-effects model was

1 used to test the fixed effects of age group (Group: 2 levels) and vowel context (Context: 2 levels) and
2 their interaction on ratings. The model was built using the lme4 package (Bates et al., 2015) in R (R
3 Development Core Team, 2019), and included a random intercept for every combination of the levels of
4 listener and speaker. The lmerTest package (Kuznetsova et al., 2017) was used to estimate the degrees of
5 freedom with Satterthwaite's method (Satterthwaite, 1946).

6 In a second set of analyses, multiple linear regression was used to evaluate the specific influence
7 of determiner vowel reduction on goodness ratings. The models were implemented in SPSS (IBM SPSS
8 Statistics Version 27). The ratings for the [a]_[i] and [i]_[a] contexts were fit separately. To maximize the
9 independence of residuals in the model, ratings were standardized across listeners within a context based
10 on individual listener means and standard deviations (i.e., z-scored). The Durbin-Watson statistic
11 indicated a minor case of positive autocorrelation in each model ($d = 1.50$ in the [a]_[i] model; $d = 1.43$ in
12 the [i]_[a] model). This was corrected by adding a lag-1 of the dependent variable, which was entered
13 first in the model. The assumption of homoscedasticity was also met in the data: scatterplots of predicted
14 values versus residuals showed no relationship in either the [a]_[i] or [i]_[a] model. Q-Q plots of the
15 residuals indicated that the assumption of normality was also met in both models.

16 The predictor variables in the multiple regression analyses were the acoustic correlates of
17 temporal and amplitude reduction from Experiment 1 (Det:V duration, Det:V amplitude, Det:N duration,
18 Det:N amplitude) and a single measure of schwa coarticulation. This measure was based on the mean
19 formant frequency measures from Experiment 1. It was the Euclidean distance of each schwa produced by
20 the speaker from the mean schwa for that speaker. All predictor values were log transformed to minimize
21 the influence of extreme values on the results. Tests for correlations between the acoustic predictor
22 variables entered into the models indicated expected significant pairwise correlations between relative
23 duration and amplitude for Det:V and Det:N as well as an unsurprising relationship between relative
24 duration and amplitude. The very strongest relationship, which was between Det:N and Det:V duration,
25 had a coefficient of nearly 0.8 (Pearson's $r = 0.797$) in the [a]_[i] data and of nearly 0.7 in the [i]_[a] data

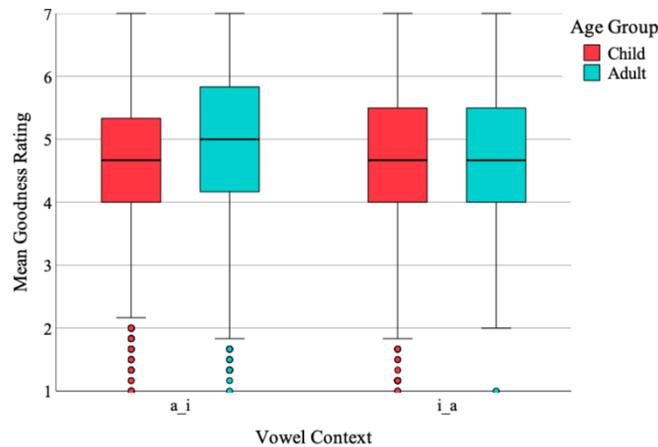
1 (Pearson's $r = 0.660$). Despite this, collinearity in the models was low; the largest VIF value was 2.51 in
2 the [a]_[i] model and it was 2.30 in the [i]_[a] model).

3 After the predictor variables of interest were included in the model, age group was added to
4 control for age-related effects that were not of interest in the experiment (e.g., the effect of F0). Speaker
5 was initially included for the same reason, but then eliminated because its effect was not significant.

6 B. Results

7 The data in Figure 5 show that goodness ratings were lower overall for child speech than for adult
8 speech [$F(1, 22) = 6.10, p = 0.022$], but this effect interacted with vowel context [$F(1, 7076) = 35.87, p <$
9 0.001].

10

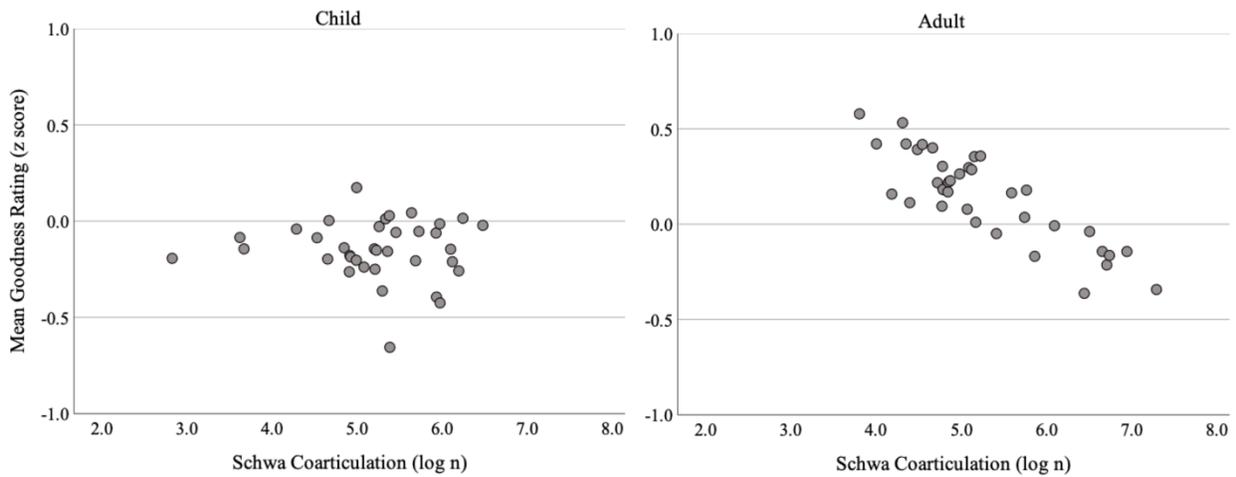


11

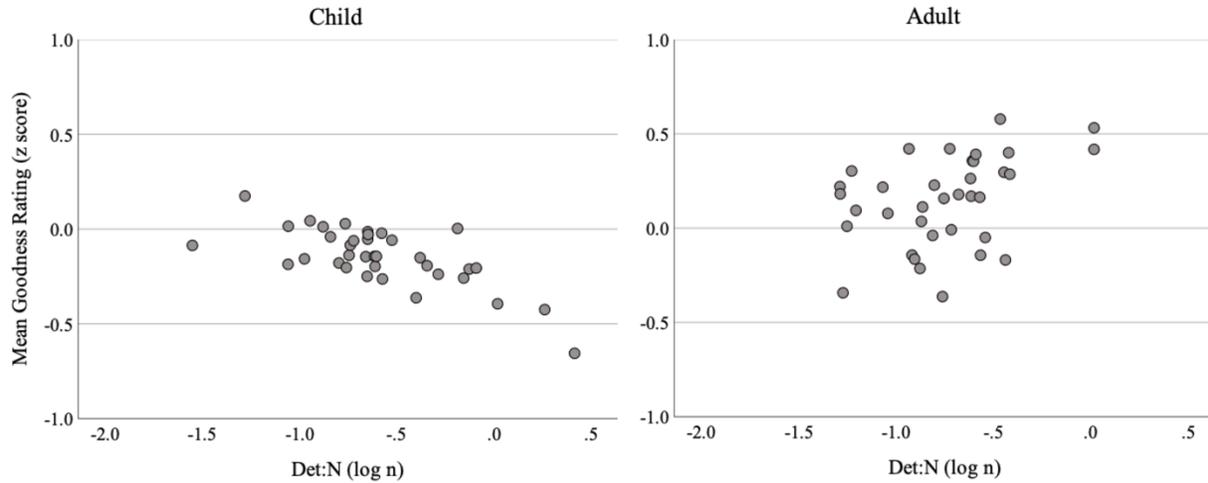
12 *Figure 5. Goodness ratings on V-the-N sequences in child and adult speech as a function of vowel*
13 *context.*

14 The second set of analyses tested for the independent influence of reduction on listener ratings.
15 As expected based on the mixed effects model results, age group accounted for a significant proportion of
16 the variance in the full [a]_[i] model [$b = -.143, t(3591) = -7.207, p < 0.001$], but not in the [i]_[a] model
17 [$b = -0.018, t(3591) = -1.03, p = 0.302$]. The strongest predictor variable of interest in both models was
18 schwa coarticulation [a_i model: $b = -.125, t(3591) = -7.70, p < 0.001$; i_a model: $b = -0.055, t(3591) =$
19 $-3.36, p < 0.001$]. Det:N duration was also significant in the [a]_[i] model [$b = -0.076, t(3591) = -3.28, p$

1 < 0.001] and nearly so in the [i]_[a] model [$b = -0.043$, $t(3591) = -1.96$, $p = 0.05$]. Det:N amplitude was
2 significant in the [i]_[a] model [$b = 0.042$, $t(3591) = 2.15$, $p = 0.032$]. The overall model of ratings on
3 stimuli from [a]_[i] elicitation accounted for 10% of the variance ($R = 0.317$; Adjusted $R^2 = 0.099$),
4 which represents a significant improvement over the null model [$F(7,3591) = 57.33$, $p < 0.001$]; the
5 overall model of ratings on stimuli from [i]_[a] elicitation accounted for 9% of the variance ($R = 0.302$;
6 Adjusted $R^2 = 0.089$), which was also significant [$F(7,3591) = 51.43$, $p < 0.001$].



7
8 *Figure 6. The effect of schwa coarticulation on model predicted goodness ratings for [a]_[i] sequences*
9 *as a function of age group.*



1
 2 *Figure 7. The effect of relative schwa duration (Det:N) on model predicted goodness ratings for [a]_[i]*
 3 *sequences as a function of age group.*

4 When the analyses on ratings of [a]_[i] sequences were conducted separately by age group, the
 5 relationship between the acoustic predictors and goodness ratings was found to differ for child and adult
 6 speech. For example, Figure 6 shows the relationship between determiner vowel coarticulation (log
 7 transformed) and goodness ratings (z-scored); Figure 7 shows the relationship for Det:N duration and
 8 goodness ratings. In adult speech (right panels), the data show that goodness rating varied inversely with
 9 degree of coarticulation in adult speech, but not with relative duration. Specifically, when the determiner
 10 vowel was further away from an adult speaker’s average schwa in $F1 \times F2 \times F3$ space, the overall
 11 sequence was rated as less good than when it was closer to their average schwa; this relationship
 12 accounted for a significant proportion of the variance in the adult [a]_[i] model [$b = -0.170$, $t(1793) =$
 13 -7.36 , $p < 0.001$]. In child speech (left panels), goodness rating was uncorrelated with schwa
 14 coarticulation, but it was higher when schwa duration was shorter than when it was longer; this
 15 relationship also accounted for a significant proportion of the variance in the child [a]_[i] model [Det:N
 16 duration, $b = -0.211$, $t(1792) = -6.01$, $p < 0.001$]. The significant inverse relationship in child speech but
 17 not adult speech is likely dependent on the greater range of relative durations, at both extremes, in the
 18 child speech data.

1 An unexpected positive relationship between Det:V and rating goodness was also found in child
2 speech for [a]_[i] sequences [$b = 0.105, t(1792) = 3.00, p = 0.003$]. Yet, a straightforward bivariate
3 correlation between the variables shows the expected inverse relationship [$r(1800) = -0.05, p = 0.03$].
4 Also, Det:N and Det:V duration are positively correlated (see Methods). This results is therefore
5 interpreted to indicate that Det:N duration captured all of the shared variance due to vowel reduction in
6 the child [a]_[i] model. Some additional variance in ratings was then accounted for by an increase in the
7 determiner vowel relative to [a] in the child model – a finding that could also indicate an influence on
8 ratings from the stressed vowel itself. Overall, the adult [a]_[i] model accounted for 13% of the variance
9 in goodness ratings ($R = 0.362$; Adjusted $R^2 = 0.131$) and the child [a]_[i] model accounted for 7% of the
10 variance ($R = 0.265$, Adjusted $R^2 = 0.067$). Both models were significantly different from the null models
11 [Child [a]_[i] model: $F(6,1793) = 22.50, p < 0.001$; Adult [a]_[i] model: $F(6,1793) = 44.94, p < 0.001$].

12 Consistent with the child [a]_[i] model results, both Det:N duration and Det:N amplitude were
13 significant predictors of goodness rating in the [i]_[a] model when the non-significant effect of age group
14 was removed [Det:N duration, $b = -0.05, t(3591) = -2.36, p = 0.019$; Det:N amplitude, $b = 0.046, t(3591)$
15 $= 2.38, p = 0.017$]. The signs of the coefficients indicated that, as expected, goodness ratings were higher
16 when the determiner vowel was shorter and of lower amplitude relative to the noun vowel. The effect of
17 schwa coarticulation remained strong in the partial model as well [$b = -0.056, t(3591) = -3.37, p <$
18 0.001]. As in the full model, when schwa was further away from its mean value in $F1 \times F2 \times F3$ space
19 listeners rated the sequence as less good than when it is closer to its mean value. The model R^2 is as
20 before: $R = 0.301$; Adjusted $R^2 = 0.089$. Overall, the change in coefficients from the full [i]_[a] model
21 that controlled for age group and the partial [i]_[a] model where this nonsignificant variable was removed
22 suggests that listeners were sensitive to the overlap between the temporal and amplitude correlates of
23 determiner vowel reduction and age group.

24 C. Discussion

25 The results from Experiment 2 upheld the expectation that adults would rate child speech as less
26 good than adult speech. Also upheld was the expectation that function word reduction would predict

1 goodness ratings independently of the preference for adult speech. Although duration and amplitude
2 measures of reduction were expected to predict goodness ratings better than schwa coarticulation, schwa
3 coarticulation was the stronger overall predictor. Relative schwa duration and amplitude did not combine
4 to explain additional variance in ratings within each model. Instead, relative schwa duration combined
5 with schwa coarticulation to predict goodness ratings on [a][i] sequences on top of the effect of age, and
6 relative schwa amplitude combined with schwa coarticulation to predict ratings on [i][a] sequences. On
7 the other hand, the analyses by age group on ratings of [a][i] sequences suggested that listeners attended
8 to different correlates of reduction depending on the speaker's age: the global measure of coarticulation
9 predicted goodness ratings of adult speech only; the relative duration of schwa in [a][i] sequences was
10 the only predictor of rating variance in child speech.

11 **V. GENERAL DISCUSSION**

12 The strong-weak pattern of English rhythm leads to the reduction of function words, which are
13 typically unstressed even when stranded in prosodic positions where they are considered extra-metrical
14 (Selkirk, 1996). The high frequency with which function words occur in spoken language coupled with
15 their minimal semantic weight contributes to their reduction (Bell et al., 2009; Jurafsky et al., 1998). But
16 function words also provide critical grammatical information – to the point where they are treated as
17 heads of phrases in some current theories of syntax (Müller, 2016). For this reason, if function words are
18 reduced to the point of not being heard, the speaker's message may either be ungrammatical or the
19 information they wish to communicate may be compromised (see, e.g., Baese-Berk et al., 2016). We
20 argue that the slow development of adult-like function word reduction results from this tension between a
21 rhythmic and semantic pressure towards maximal reduction and a listener-oriented pressure towards
22 maximal intelligibility. The present results are consistent with findings from prior acoustic and kinematic
23 studies of immature function word production in child speech: results from Experiment 1 were that 5-
24 year-old speakers do not reduce determiners to the same extent as adult speakers in all contexts. More
25 specifically, when child and adult determiners differs, they differs in the direction of relatively longer and
26 higher amplitude schwa in child speech compared to adult speech.

1 Although schwa was relatively longer in child speech compared to adult speech, it was still much
2 shorter than the vowels in the adjacent content words within the same sequence. In fact, it was often just
3 half the duration of the adjacent content word vowel. This finding is consistent with children’s relatively
4 early acquisition of lexical stress patterns (Ballard et al., 2012). A possible implication of the finding is
5 that children reduce monosyllabic function words to the same degree that they reduce weak syllables in
6 content words, whereas adults reduce function words to a greater degree than they reduce weak syllables
7 in content words. This possibility is consistent with empirical findings (Goffman, 2004; Fuchs, 2016; van
8 Bergem, 1993). For example, Goffman (2004) compared child (4 to 7 years old) and adult production of
9 weak syllables in two-syllable sequences. The syllables were embedded in a discourse context that led
10 speakers to produce them either as a determiner-noun sequence or as an iambically stressed content word.
11 Vertical movements of the lip-jaw complex were measured. In adult speech, the results were that
12 movement duration and amplitude was smaller when the syllables were treated as a determiner-noun
13 sequence than when it was treated as a disyllabic content word. In child speech, the results were that the
14 weak syllable was produced with shorter and lower amplitude than the strong syllable, but there was no
15 effect of morphosyntactic environment. This again suggests that children do not need to learn to reduce
16 function words; they need to learn to reduce them to the same extent as adults.

17 The results from Experiment 1 also showed that both children and adults coarticulated the
18 determiner with a following noun. If vowel-to-vowel coarticulation can be used to index chunking in the
19 speech plan, this result is consistent with chunking along morphosyntactic lines rather than along metrical
20 ones. Recall that the elicited sequences had a strong-weak-strong stress pattern; that is, strong
21 monosyllabic content words separated by a weak function word. In this context, the weak function word
22 should adhere to the preceding strong syllable to form a trochaic foot (Selkirk, 1996). Instead, the
23 coarticulatory pattern found here suggests an adherence to the following strong syllable. This adherence
24 pattern promotes a coherent determiner noun phrase. Given that listeners are known to use coarticulatory
25 cues to aid in speech segmentation (Mattys, 2004), chunking along morphosyntactic lines has functional
26 value. Of course, the distributional patterns of lexical stress in English are such that listeners also use a

1 trochaic pattern to aid in speech segmentation (Cutler and Butterfield, 1992; Mattys et al., 2005; Dilley
2 and McAuley, 2008). This again suggests competing pressures on production: the chunking of
3 determiner-noun sequences will (usually) result in the production of an iambic pattern; the most common
4 disyllabic nouns and adjectives in English are produced with a trochaic pattern (see Cutler and Carter,
5 1987). Conflicting pressures such as these may also influence the production of determiner noun phrases.

6 Although the noun vowel influenced the production of schwa in the determiner in both child and
7 adult speech, the transition towards the noun was only evident in children's speech. In adult speech,
8 schwa was simply produced with a different quality depending on the identity of the noun vowel. It could
9 be that these differences are merely an epiphenomenon of age-related differences in articulation rate. In
10 particular, children's slower rate of articulation could provide them with more time to adjust articulatory
11 movements into and out of a vowel target compared to adults. Such an explanation assumes that schwa
12 coarticulation in adult speech is due to target undershoot (i.e., hypoarticulation; Agwuele et al., 2008;
13 Moon and Lindblom, 1994). Target undershoot increases as the time allotted for target attainment is
14 decreased so long as articulatory effort is held constant. If we explain the effect of age on coarticulation in
15 these terms, it would suggest that children's slower articulation rates reflect a production strategy
16 designed to maximize sequential target attainment. There is some evidence to support this possibility.
17 Recall that children's speech movements are in fact larger relative to their oral-facial size (Riely and
18 Smith, 2003), suggesting that they expend more effort in speaking compared to adults – presumably, in
19 order to achieve acoustic targets.

20 Yet, if we explain the effect of age on coarticulation as an epiphenomenon of speech rate, we
21 must still explain why the relative duration and amplitude of schwa is (typically) less than that of an
22 adjacent full vowel. Saying that it is “reduced” only says that, at some level of representation, its timing is
23 due to language factors, which is to say that its timing is planned before execution. And, if timing is part
24 of the plan, then coarticulation may not be due to target undershoot, but instead to target overlap (Fowler,
25 1980). This view of coarticulation suggests an alternative explanation for the present results – one that is
26 more consistent with the idea of competing rhythmic and morphosyntactic pressures on production: if

1 children are less able than adults to resolve these pressures, their determiner-noun sequences are less
2 tightly bound than in adult speech. If the determiner-noun sequence is less tightly bound in child speech
3 than in adult speech, then the vowels in children's sequences will also be less overlapped and schwa less
4 subject to "truncation" allowing for its fuller realization (Harrington et al., 1995).

5 Regardless of exactly why children do not reduce function words to the same extent as adults,
6 they eventually learn to do so. Our suggestion is that learning requires a honing process that depends on
7 feedback in the form of communicative success or failure. This hypothesis is consistent with an adult
8 listener preference for adult speech over child speech. It also predicts that listener ratings of speech
9 rhythmicity will be most influenced by those measures of reduction that best distinguish between child
10 and adult speech. The results from Experiment 2 confirmed the predicted listener preference for adult
11 speech over child speech, but the influence of specific measures of reduction on listener preference was
12 more complicated than we had expected. Based on the results from Experiment 1, we had expected that
13 measures of relative duration and amplitude of schwa (especially, Det:V duration and Det:V amplitude)
14 would predict listener goodness ratings on *V-the-N* sequences. Instead, the strongest predictor of
15 goodness ratings was a global measure of schwa coarticulation.

16 The global measure of schwa coarticulation that was used as a predictor variable in analyses of
17 goodness ratings in Experiment 2 was calculated as the Euclidean distance of schwa from the speaker
18 mean schwa in $F1 \times F2 \times F3$ space. Goodness ratings on *V-the-N* sequences decreased when schwa was
19 more distant from the speaker mean. The assumption is that the larger the distance from the mean, the
20 more schwa varied as a function of context. But production variability may be due to other sources as
21 well, including immature motor skills (Smith and Zelaznik, 2004). The result could therefore indicate that
22 listeners were attending to segmental target attainment, that schwa has a context-sensitive target
23 (Browman and Goldstein, 1992), and that reduction per se was not relevant to the judgments that were
24 made. Then again, the possibility that listeners rated sequences based on how similar the realization of
25 schwa was to an expected schwa target could also be due to the experimental design.

1 The speech stimuli were blocked both by speaker and by age group. The goal of blocking was to
2 attune listeners to the rhythmicity of the sequence and away from global differences in production due to
3 individual or group differences (e.g., differences in F0, differences in segmental articulation). But the
4 blocked design also likely attuned listeners to that which was most variable within a speaker and age
5 group. It could be that temporal and amplitude patterns are simply more stable than the realization of
6 schwa within a speaker and age group. But this explanation does not account for why relative schwa
7 duration was the only significant predictor of ratings on [a]_[i] sequences elicited from children or why
8 the global measure of coarticulation was the only significant predictor of ratings on [a]_[i] sequences in
9 adult speech. Also, child speech is notoriously variable and this variability is especially significant in the
10 spatial-temporal domain (Smith and Zelaznik, 2004). It therefore seems likely that listener expectations
11 for reduction provide a better explanation for why their ratings on adult sequences were influenced by
12 different correlates than their ratings on child sequences. For example, it could be that listeners expect
13 some average amount of reduction in the temporal and amplitude domains and variability around that
14 amount is less important than variability around an average that is already higher than the amount
15 expected.

16 Although the results conform in several ways to the predictions made, the gap between the
17 present results and the working hypothesis that communication success shapes speech motor learning is
18 still admittedly large. Much more work is needed to firmly establish a relationship between child speech
19 production and adult speech processing. Future research could start by confirming the relationship
20 between function word reduction and children's speech rhythm. This relationship is assumed based on
21 perceptual analyses of children's speech (e.g., Allen and Hawkins, 1978) and on a correlation between
22 interval-based rhythm measures and the relative duration measures reported here (e.g., Sirsa and Redford,
23 2011). Combined acoustic and perceptual studies on large samples of spontaneous speech would go some
24 way towards confirming the relationship. If confirmed, it would also be worth directly testing the link
25 between violations in rhythm due to more or less reduced function words in adult speech processing.
26 More generally, the relationship between child speech production and adult speech perception requires

1 further study. Although the link between immature motor skills and speech intelligibility is well-
2 established, the relationship is rarely investigated in detail. For example, it is not clear how adult listeners
3 weigh the relative contribution of variable segmental articulation and immature speech rhythm when
4 processing child speech. Finally, ecologically-valid descriptions of adult listener response to child speech
5 and child responses to adult behaviors is needed to better characterize the details of how communicative
6 interactions may drive speech motor learning.

7

8 **ACKNOWLEDGMENTS**

9 This research was wholly supported by the Eunice Kennedy Shriver National Institute of Child Health
10 and Human Development (NICHD) under grant R01HD087452 (PI: Redford). The content is solely the
11 authors' responsibility and does not necessarily reflect the views of NICHD.

1 **REFERENCES**

- 2 Abu-Akel, A., Bailey, A. L., and Thum, Y. M. (2004). “Describing the acquisition of determiners in
3 English: A growth modeling approach,” *J. Psycholinguist. Res.* **33**(5), 407-424.
- 4 Agwuele, A., Sussman, H. M., and Lindblom, B. (2008). “The effect of speaking rate on consonant vowel
5 coarticulation,” *Phonetica.* **65**(4), 194-209.
- 6 Allen, G., and Hawkins. S. (1978). “The development of phonological rhythm,” in *Syllables and*
7 *Segments*, edited by A. Bell and J. Bybee Hooper (North-Holland Publishing: New York), pp.
8 173–185.
- 9 Baese-Berk, M. M., Dilley, L. C., Schmidt, S., Morrill, T. H., and Pitt, M. A. (2016). “Revisiting Neil
10 Armstrong’s moon-landing quote: implications for speech perception, function word reduction,
11 and acoustic ambiguity,” *PloS one.* **11**(9), e0155975.
- 12 Ballard, K. J., Djaja, D., Arciuli, J., James, D. G., and van Doorn, J. (2012). “Developmental trajectory
13 for production of prosody: lexical stress contrastivity in children ages 3 to 7 years and in adults,”
14 *J. Speech Lang. Hear. Res.* **55**, 1822-1835.
- 15 Bates, D. and Mächler, M., Bolker, B., and Walker, S. (2015). “Fitting linear mixed-effects models using
16 lme4,” *J. Stat. Softw.* **67**(1), 1-48.
- 17 Bell, A., Brenier, J. M., Gregory, M., Girand, C., and Jurafsky, D. (2009). “Predictability effects on
18 durations of content and function words in conversational English,” *J. Mem. Lang.* **60**(1), 92-111.
- 19 Boersma, P. and Weenink, D. (2019). Praat: doing phonetics by computer [Computer program]. Version
20 6.0.49. <http://www.praat.org/>
- 21 Browman, C. P., and Goldstein, L. (1992). “Targetless schwa: An articulatory analysis.” In *Papers in*
22 *laboratory phonology II: Gesture, segment, prosody*, edited by G.J. Docherty and D.R. Ladd
23 (Cambridge University Press, Cambridge), pp. 26–56.

- 1 Buhrmester, M., Kwang, T., and Gosling, S. D. (2011). "Amazon's mechanical turk: a new source of
2 inexpensive, yet high-quality, data?" *Perspect. Psychol. Sci.* **6**(1), 3-5.
- 3 Cutler, A. and Butterfield, S. (1992). "Rhythmic cues to speech segmentation: Evidence from juncture
4 misperception," *J. Mem. Lang.* **31**(2), 218-236.
- 5 Cutler, A. and Carter, D. (1987). "The predominance of strong initial syllables in the English
6 vocabulary," *Comput. Speech Lang.* **2**, 133-142.
- 7 Dauer, R. M. (1983). "Stress-timing and syllable-timing reanalyzed," *J. Phon.* **11**, 51-62.
- 8 Deterding, D. (2001). "The measurement of rhythm: A comparison of Singapore and British English," *J.*
9 *Phon.* **29**(2), 217-230.
- 10 Dilley, L. C. and McAuley, J. D. (2008). "Distal prosodic context affects word segmentation and lexical
11 processing," *J. Mem. Lang.* **59**, 294-311.
- 12 Dilley, L. C. and Pitt, M. A. (2010). "Altering context speech rate can cause words to appear or
13 disappear," *Psychol. Sci.* **21**(11), 1664-1670.
- 14 Dodd, B., Zhu, H., Crosbie, S., Holm, A., and Ozanne, A. (2002). *Diagnostic evaluation of articulation*
15 *and phonology (DEAP)*. (Psychological Corporation, London).
- 16 Fikkert, P., Liu, L., and Mitsuhiro, O. (2021) "The acquisition of word prosody," in *The Oxford*
17 *Handbook of Language Prosody*, edited by C. Gussenhoven and A. Chen (Oxford University
18 Press, London), pp. 541-552.
- 19 Fourakis, M. (1991). "Tempo, stress, and vowel reduction in American English," *J. Acoust. Soc. Am.*
20 **90**(4), 1816-1827.
- 21 Fowler, C. A. (1980). "Coarticulation and theories of extrinsic timing," *J. Phon.* **8**(1), 113-133.

- 1 Fowler, C. A. (1981). "Production and perception of coarticulation among stressed and unstressed
2 vowels," J. Speech Hear. Res. **24**(1), 127-139.
- 3 Fuchs, R. (2016). "The acoustic correlates of stress and accent in English content and function words,"
4 Proc. Speech Prosody, **8**, pp. 290-294.
- 5 Gay, T. (1981). "Mechanisms in the control of speech rate," *Phonetica*. **38**(1-3), 148-158.
- 6 Gerken, L. (1996). "Prosodic structure in young children's language production," *Language*. **72**(4), 683-
7 712.
- 8 Goffman, L. (2004). "Kinematic differentiation of prosodic categories in normal and disordered language
9 development," J. Speech Lang. Hear. Res., **47**(5), 1088-1102.
- 10 Grabe, E., and Low, E. L. (2002). "Durational variability in speech and the rhythm class hypothesis," in
11 C. Gussenhoven and N. Warner (eds.), *Laboratory Phonology 7* (Mouton de Gruyter, Berlin), pp.
12 515-546.
- 13 Gu, C. (2014). "Smoothing spline ANOVA models: R package gss," J. Stat. Softw. **58**(5), 1-25.
- 14 Harrington, J., Fletcher, J., and Roberts, C. (1995). "Coarticulation and the accented/unaccented
15 distinction: evidence from jaw movement data," J. Phon. **23**(3), 305-322.
- 16 He, L. (2012). "Syllabic intensity variations as quantification of speech rhythm: evidence from both L1
17 and L2," Proc. Speech Prosody, pp. 466-469.
- 18 Jurafsky, D., Bell, A., Fosler-Lussier, E., Girand, C., and Raymond, W. (1998). "Reduction of English
19 function words in switchboard," in ICSLP-98, 3111-3114.
- 20 Kedar, Y., Casasola, M., and Lust, B. (2006). "Getting there faster: 18- and 24-month-old infants' use of
21 function words to determine reference," *Child Dev.* **77**(2), 325-338.

- 1 Kehoe, M., Stoel-Gammon, C. and Buder, E.H. (1995). "Acoustic correlates of stress in young children's
2 speech," J. Speech Hear. Res. **38**, 338-350.
- 3 Kuznetsova, A., Brockhoff, P. B., Christensen, R. H. B. (2017). "lmerTest package: test in linear mixed
4 effects models," J. Stat. Softw. **82**(13), 1-26.
- 5 Lee, S., Potamianos, A., and Narayanan, S. (1999). "Acoustics of children's speech: Developmental
6 changes of temporal and spectral parameters," J. Acoust. Soc. Am. **105**(3), 1455-1468.
- 7 Lindblom, B., Sundberg, J., Branderud, P., and Djamshidpey, H. (2011). "Articulatory modeling and front
8 cavity acoustics," Proc. Fonetik TMH-QPSR. **51**(1), 17-20.
- 9 Mattys, S. L. (2004). "Stress versus coarticulation: toward an integrated approach to explicit speech
10 segmentation," J. Exp. Psychol. Hum. Percept. Perform., **30**(2), 397.
- 11 Mattys, S.L., White, L., and Melhorn, J.F. (2005). "Integration of multiple speech segmentation cues: a
12 hierarchical framework," J. Exp. Psychol. Gen. **134**(4), 477.
- 13 Mielke, J. (2015). "An ultrasound study of Canadian French rhotic vowels with polar smoothing spline
14 comparisons," J. Acoust. Soc. Am. **137**(5), 2858-2869.
- 15 Moon, S. J., and Lindblom, B. (1994). "Interaction between duration, context, and speaking style in
16 English stressed vowels," J. Acoust. Soc. Am. **96**(1), 40-55.
- 17 Müller, Stefan (2016). *Grammatical Theory: From Transformational Grammar to Constraint-based*
18 *Approaches*. (Language Science Press, Berlin).
- 19 Nittrouer, S. (1993). "The emergence of mature gestural patterns is not uniform: Evidence from an
20 acoustic study," J. Speech Lang. Hear. Res. **36**(5), 959-972.
- 21 Payne, E., Post, B., Astruc, L., Prieto, and Vanrell, M. (2011). "Measuring child rhythm," Lang. Speech
22 **55**(2), 203-229.

- 1 Plag, I., Kunter, G., and Schramm, M. (2011). "Acoustic correlates of primary and secondary stress in
2 North American English," J. Phon. **39**(3), 362-374.
- 3 Polyanskaya, L. and Ordin, M. (2015). "Acquisition of speech rhythm in first language," J. Acoust. Soc.
4 Am. **138**(3), EL199-EL204.
- 5 R Core Team. (2019). R: A Language and Environment for Statistical Computing. R Foundation for
6 Statistical Computing. Vienna, Austria.
- 7 Ratcliff, R. (1993). "Methods for dealing with reaction time outliers," Psychol. Bull. **114**(3), 510-532.
- 8 Redford, M. A. (2018). "Grammatical word production across metrical contexts in school-aged children's
9 and adults' speech," J. Speech Lang. Hear. Res. **61**, 1339-1354.
- 10 Redford, M.A. (2014). "The perceived clarity of children's speech varies as a function of their default
11 speech rate," J. Acoust. Soc. Am. **135**, 2952-2963.
- 12 Redford, M. A., Kapatsinski, V., and Cornell-Fabiano, J. (2018). "Lay listener classification and
13 evaluation of typical and atypical children's speech," Lang. Speech, **61**, 277-302.
- 14 Riely, R.R., Smith, A. (2003). "Speech movements do not scale by orofacial structure size," J. Appl.
15 Physiol. **94**, 2119-2126.
- 16 Satterthwaite, F. E. (1946). "An approximate distribution of estimates of variance components,"
17 Biometrics **2**, 110-114.
- 18 Schwartz, R. G., Petinou, K., Goffman, L., Lazowski, G. and Cartusciello, C. (1996). "Young children's
19 production of syllable stress: An acoustic analysis," J. Acoust. Soc. Am. **99**(5), 3192-3200.
- 20 Selkirk, E. (1996). "The prosodic structure of function words," in *Signal to Syntax: Bootstrapping from*
21 *Speech to Grammar in Early Acquisition*, edited by J.L. Morgan and K. Demuth (eds.), (Erlbaum,
22 Mahwah), pp. 187-214.

- 1 Shriberg, L. D., Campbell, T. F., Karlsson, H. B., Brown, R. L., McSweeny, J. L., and Nadler, C. J.
2 (2003). "A diagnostic marker for childhood apraxia of speech: The lexical stress ratio," Clin.
3 Linguist. Phon. **17**(7), 549-574.
- 4 Smith, A., and Zelaznik, H. N. (2004). "Development of functional synergies for speech motor
5 coordination in childhood and adolescence," Dev. Psychobio. **45**(1), 22-33.
- 6 Sirsa, H. and Redford, M. A. (2011). "Towards understanding the protracted acquisition of English
7 rhythm," Proc. Int. Congress Phon. Sci. 1862-1865.
- 8 Tilsen, S., and Arvaniti, A. (2013). "Speech rhythm analysis with decomposition of the amplitude
9 envelope: characterizing rhythmic patterns within and across languages," J. Acoust. Soc.
10 Am. **134**(1), 628-639.
- 11 Van Bergem, D. R. (1993). "Acoustic vowel reduction as a function of sentence accent, word stress, and
12 word class," Speech Commun. **12**(1), 1-23.
- 13 Wiig, E. H., Secord, W., and Semel, E. (2013). *Clinical Evaluation of Language Fundamentals (CELF-*
14 *5)*, (Harcourt Brace Jovanovich, San Antonio, TX).

1 **FIGURE LEGENDS**

2 **Figure 1.** Example of vowel segmentations for a Verb-the-Noun sequence embedded in a simple S-V-O
3 sentence elicited from a 5-year-old child.

4 **Figure 2.** Vowel duration ratios for Det:V (left) and Det:N (right) shown by age group and the vowel
5 context within which the determiner appeared.

6 **Figure 3.** Vowel amplitude ratios for Det:V (left) and Det:N (right) are shown by age group and the
7 vowel context within which the determiner appeared.

8 **Figure 4.** Best-fit curves for schwa F1, F2, and F3 are shown as a function of the stressed vowel context
9 within which the determiner was produced in child (left) and adult (right) speech. Dotted lines
10 represent the 95% confidence interval.

11 **Figure 5.** Goodness ratings on V-the-N sequences in child and adult speech as a function of vowel
12 context.

13 **Figure 6.** The effect of schwa coarticulation on model predicted goodness ratings for [a]_[i] sequences as
14 a function of age group.

15 **Figure 7.** The effect of relative schwa duration (Det:N) on model predicted goodness ratings for [a]_[i]
16 sequences as a function of age group.